



(86) **Date de dépôt PCT/PCT Filing Date:** 2007/10/19
(87) **Date publication PCT/PCT Publication Date:** 2008/05/02
(45) **Date de délivrance/Issue Date:** 2016/02/23
(85) **Entrée phase nationale/National Entry:** 2009/04/09
(86) **N° demande PCT/PCT Application No.:** US 2007/022414
(87) **N° publication PCT/PCT Publication No.:** 2008/051511
(30) **Priorité/Priority:** 2006/10/20 (US60/853,284)

(51) **Cl.Int./Int.Cl. C12Q 1/68** (2006.01),
C07H 21/00 (2006.01), **C07H 21/04** (2006.01),
G01N 33/53 (2006.01), **A61K 39/395** (2006.01)

(72) **Inventeurs/Inventors:**
BARE, LANCE, US;
DEVLIN, JAMES J., US;
ROSENDAAL, FRITS R., US;
REITSMA, PIETER H., US;
BEZEMER, IRENE D., US

(73) **Propriétaires/Owners:**
CELERA CORPORATION, US;

(54) **Titre : POLYMORPHISMES GENETIQUES ASSOCIES A LA THROMBOSE VEINEUSE, PROCEDES POUR LES DETECTER ET UTILISATIONS**
(54) **Title: GENETIC POLYMORPHISMS ASSOCIATED WITH VENOUS THROMBOSIS, METHODS OF DETECTION AND USES THEREOF**

(57) **Abrégé/Abstract:**

The present invention is based on the discovery of genetic polymorphisms that are associated with coronary heart disease and in particular VT and response to drug treatment. In particular, the present invention relates to nucleic acid molecules containing the polymorphisms, variant proteins encoded by such nucleic acid molecules, reagents for detecting the polymorphic nucleic acid molecules and proteins, and methods of using the nucleic acid and proteins as well as methods of using reagents for their detection.



(73) **Propriétaires(suite)/Owners(continued):**

LEIDEN UNIVERSITY MEDICAL CENTER (LUMC) ACTING ON BEHALF OF ACADEMIC HOSPITAL LEIDEN (AZL), NL

(74) **Agent:** FETHERSTONHAUGH & CO.

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
2 May 2008 (02.05.2008)

PCT

(10) International Publication Number
WO 2008/051511 A3

(51) International Patent Classification:

C12Q 1/68 (2006.01) *A61K 48/00* (2006.01)
C07H 21/00 (2006.01)

(21) International Application Number:

PCT/US2007/022414

(22) International Filing Date: 19 October 2007 (19.10.2007)

(25) Filing Language:

English

(26) Publication Language:

English

(30) Priority Data:

60/853,284 20 October 2006 (20.10.2006) US

(71) Applicant (for all designated States except US): **APPLERA CORPORATION** [US/US]; c/o Celera, An Applera Business Unit, 1401 Harbor Bay Parkway, Alameda, CA 94502 (US).

(72) Inventor; and

(75) Inventor/Applicant (for US only): **BARE, Lance** [US/US]; c/o Celera, An Applera Business Unit, 1401 Harbor Bay Parkway, Alameda, CA 94502 (US).

(74) Agent: **LEE, Victor**; Celera, An Applera Business Unit, 1401 Harbor Bay Parkway, Alameda, CA 94502 (US).

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RS, RU, SC, SD, SE, SG, SK, SL, SM, SV, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, MT, NL, PL, PT, RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

- with international search report
- before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments

(88) Date of publication of the international search report:

30 October 2008

(54) Title: GENETIC POLYMORPHISMS ASSOCIATED WITH VENOUS THROMBOSIS, METHODS OF DETECTION AND USES THEREOF

(57) Abstract: The present invention is based on the discovery of genetic polymorphisms that are associated with coronary heart disease and in particular VT and response to drug treatment. In particular, the present invention relates to nucleic acid molecules containing the polymorphisms, variant proteins encoded by such nucleic acid molecules, reagents for detecting the polymorphic nucleic acid molecules and proteins, and methods of using the nucleic acid and proteins as well as methods of using reagents for their detection.



WO 2008/051511 A3

GENETIC POLYMORPHISMS ASSOCIATED WITH VENOUS THROMBOSIS, METHODS OF DETECTION AND USES THEREOF

FIELD OF THE INVENTION

The present invention is in the field of thrombosis diagnosis and therapy. In particular, the present invention relates to specific single nucleotide polymorphisms (SNPs) in the human genome, and their association with venous thrombosis and related pathologies. Based on differences in allele frequencies in the patient population relative to normal individuals, the naturally-occurring SNPs disclosed herein can be used as targets for the design of diagnostic reagents and the development of therapeutic agents, as well as for disease association and linkage analyses. In particular, the SNPs of the present invention are useful for identifying an individual who is at an increased or decreased risk of developing venous thrombosis and for early detection of the disease, for providing clinically important information for the prevention and/or treatment of venous thrombosis, for screening and selecting therapeutic agents, and for predicting a patient's response to therapeutic agents. The SNPs disclosed herein are also useful for human identification applications. Methods, assays, kits, and reagents for detecting the presence of these polymorphisms and their encoded products are provided.

BACKGROUND OF THE INVENTION

Venous Thrombosis

The development of a blood clot is known as thrombosis. Venous thrombosis (VT) is the formation of a blood clot in the veins. Several conditions can lead to an increased tendency to develop blood clots in the veins or arteries (National Hemophilia Foundation, *HemAware* newsletter, Vol. 6 (5), 2001), and such conditions may be inherited or acquired. Examples of acquired conditions are surgery and trauma, prolonged immobilization, cancer, myeloproliferative disorders, and even pregnancy, all of which may result in thrombosis (U. Seligsohn and A. Lubetsky, *New Eng J Med* 344(16):1222-1231, 2001). Inherited causes include mutations in any of several different clotting, anticoagulant, or thrombolytic factors, such as the

factor V gene (the factor V Leiden mutation), prothrombin gene (factor II), and methylenetetrahydrofolate reductase gene (MTHFR). Other likely inherited causes are an increase in the expression levels of the factors VIII, IX or XI, or fibrinogen genes (U. Seligsohn and A. Lubetsky, *New Eng J Med* 344(16):1222-1231, 2001). VT may result from a genetic mutation alone or in concert with environmental factors, such as smoking.

VT is considered a chronic and polygenic disease (AI Schafer, *New Engl J Med* 340:955-956, 1999). There is evidence to suggest that patients with a first episode of VT be treated with anticoagulant agents (C. Kearon, JA Julian *et al.*, *New Engl J Med* 340:901-907, 1999).

Venous thrombotic events occur for the first time in about 100 per 100,000 people each year in the United States. About one-third of patients with symptomatic VT manifest pulmonary embolism (PE), whereas two-thirds manifest deep vein thrombosis (DVT) (RH White, *Circulation* 107(23 Suppl 1):I4-8 Review, 2003). DVT is an acute VT in a deep vein, usually in the thigh, legs, or pelvic, and it is a serious and potentially fatal disorder that can arise as a complication for hospital patients, but may also affect otherwise healthy people (AWA Lensing, HR Buller *et al.*, *Lancet* 353:479-485, 1999). Large blood clots in VT may interfere with blood circulation and impede normal blood flow. In some instances, blood clots may break off and travel to distant major organs such as the brain, heart or lungs as in PE and result in fatality.

VT is a multifactorial disease resulting from both acquired and genetic factors (FR Rosendaal, *Lancet*, 1999, 353:1167-1173; MD Silverstein, JA Heit *et al.*, *Arch Intern Med*, 1998; 158:585-593). Over 200,000 new cases of VT occur annually. Of these, 30 percent of patients die within three days; one in five suffer sudden death due to PE (*Seminars in Thrombosis and Hemostasis*, 2002, Vol. 28, Suppl. 2). Caucasians and African-Americans have a significantly higher incidence than Hispanics, Asians or Pacific Islanders (RH White, *Circulation* 107(23 Suppl 1):I4-8 Review, 2003).

Thus, there is an urgent need for novel genetic markers that are predictive of predisposition to VT, particularly for individuals who are unrecognized as having a predisposition to developing the disease. Such genetic markers may enable screening of VT in much larger populations compared with the populations that can currently be evaluated by using existing risk factors and biomarkers. The availability of a genetic test may allow, for example, appropriate preventive treatments for acute venous thrombotic events to be provided for high risk individuals (such preventive treatments may include, for example, anticoagulant agents). Moreover, the discovery of genetic markers associated with VT may provide novel targets for therapeutic intervention or preventive treatments.

SNPs

The genomes of all organisms undergo spontaneous mutation in the course of their continuing evolution, generating variant forms of progenitor genetic sequences (Gusella, *Ann. Rev. Biochem.* 55, 831-854 (1986)). A variant form may confer an evolutionary advantage or disadvantage relative to a progenitor form or may be neutral. In some instances, a variant form confers an evolutionary advantage to the species and is eventually incorporated into the DNA of many or most members of the species and effectively becomes the progenitor form. Additionally, the effects of a variant form may be both beneficial and detrimental, depending on the circumstances. For example, a heterozygous sickle cell mutation confers resistance to malaria, but a homozygous sickle cell mutation is usually lethal. In many cases, both progenitor and variant forms survive and co-exist in a species population. The coexistence of multiple forms of a genetic sequence gives rise to genetic polymorphisms, including SNPs.

Approximately 90% of all genetic polymorphisms in the human genome are SNPs. SNPs are single base positions in DNA at which different alleles, or alternative nucleotides, exist in a population. The SNP position (interchangeably referred to herein as SNP, SNP site, SNP locus, SNP marker, or marker) is usually preceded by and followed by highly conserved sequences of the allele (*e.g.*, sequences that vary in less than 1/100 or 1/1000 members of the populations). An individual may be homozygous or heterozygous for an allele at each SNP position. A SNP can, in some instances, be referred to as a "cSNP" to denote that the nucleotide sequence containing the SNP is an amino acid coding sequence.

A SNP may arise from a substitution of one nucleotide for another at the polymorphic site. Substitutions can be transitions or transversions. A transition is the replacement of one purine nucleotide by another purine nucleotide, or one pyrimidine by another pyrimidine. A transversion is the replacement of a purine by a pyrimidine, or vice versa. A SNP may also be a single base insertion or deletion variant referred to as an "indel" (Weber *et al.*, "Human diallelic insertion/deletion polymorphisms," *Am J Hum Genet* 2002 Oct.; 71(4):854-62).

A synonymous codon change, or silent mutation/SNP (terms such as "SNP," "polymorphism," "mutation," "mutant," "variation," and "variant" are used herein interchangeably), is one that does not result in a change of amino acid due to the degeneracy of the genetic code. A substitution that changes a codon coding for one amino acid to a codon coding for a different amino acid (*i.e.*, a non-synonymous codon change) is referred to as a missense mutation. A nonsense mutation results in a type of non-synonymous codon change in which a stop codon is formed, thereby leading to premature termination of a polypeptide chain

and a truncated protein. A read-through mutation is another type of non-synonymous codon change that causes the destruction of a stop codon, thereby resulting in an extended polypeptide product. While SNPs can be bi-, tri-, or tetra- allelic, the vast majority of the SNPs are bi-allelic, and are thus often referred to as “bi-allelic markers,” or “di-allelic markers.”

5 As used herein, references to SNPs and SNP genotypes include individual SNPs and/or haplotypes, which are groups of SNPs that are generally inherited together. Haplotypes can have stronger correlations with diseases or other phenotypic effects compared with individual SNPs, and therefore may provide increased diagnostic accuracy in some cases (Stephens *et al. Science* 293, 489-493, 20 July 2001).

10 Causative SNPs are those SNPs that produce alterations in gene expression or in the expression, structure, and/or function of a gene product, and therefore are most predictive of a possible clinical phenotype. One such class includes SNPs falling within regions of genes encoding a polypeptide product, *i.e.* cSNPs. These SNPs may result in an alteration of the amino acid sequence of the polypeptide product (*i.e.*, non-synonymous codon changes) and give rise to
15 the expression of a defective or other variant protein. Furthermore, in the case of nonsense mutations, a SNP may lead to premature termination of a polypeptide product. Such variant products can result in a pathological condition, *e.g.*, genetic disease. Examples of genes in which a SNP within a coding sequence causes a genetic disease include sickle cell anemia and cystic fibrosis.

20 Causative SNPs do not necessarily have to occur in coding regions; causative SNPs can occur in, for example, any genetic region that can ultimately affect the expression, structure, and/or activity of the protein encoded by a nucleic acid. Such genetic regions include, for example, those involved in transcription, such as SNPs in transcription factor binding domains, SNPs in promoter regions, in areas involved in transcript processing, such as SNPs at intron-exon
25 boundaries that may cause defective splicing, or SNPs in mRNA processing signal sequences such as polyadenylation signal regions. Some SNPs that are not causative SNPs nevertheless are in close association with, and therefore segregate with, a disease-causing sequence. In this situation, the presence of a SNP correlates with the presence of, or predisposition to, or an increased risk in developing the disease. These SNPs, although not causative, are nonetheless
30 also useful for diagnostics, disease predisposition screening, and other uses.

An association study of a SNP and a specific disorder involves determining the presence or frequency of the SNP allele in biological samples from individuals with the disorder of interest, such as VT, and comparing the information to that of controls (*i.e.*, individuals who do not have the disorder; controls may be also referred to as “healthy” or “normal” individuals) who

are preferably of similar age and race. The appropriate selection of patients and controls is important to the success of SNP association studies. Therefore, a pool of individuals with well-characterized phenotypes is extremely desirable.

5 A SNP may be screened in diseased tissue samples or any biological sample obtained from a diseased individual, and compared to control samples, and selected for its increased (or decreased) occurrence in a specific pathological condition, such as pathologies related to VT. Once a statistically significant association is established between one or more SNP(s) and a pathological condition (or other phenotype) of interest, then the region around the SNP can optionally be thoroughly screened to identify the causative genetic locus/sequence(s) (*e.g.*,
10 causative SNP/mutation, gene, regulatory region, etc.) that influences the pathological condition or phenotype. Association studies may be conducted within the general population and are not limited to studies performed on related individuals in affected families (linkage studies).

Clinical trials have shown that patient response to treatment with pharmaceuticals is often heterogeneous. There is a continuing need to improve pharmaceutical agent design and therapy.
15 In that regard, SNPs can be used to identify patients most suited to therapy with particular pharmaceutical agents (this is often termed "pharmacogenomics"). Similarly, SNPs can be used to exclude patients from certain treatment due to the patient's increased likelihood of developing toxic side effects or their likelihood of not responding to the treatment. Pharmacogenomics can also be used in pharmaceutical research to assist the drug development and selection process.
20 (Linder *et al.* (1997), *Clinical Chemistry*, 43, 254; Marshall (1997), *Nature Biotechnology*, 15, 1249; International Patent Application WO 97/40462, Spectra Biomedical; and Schafer *et al.* (1998), *Nature Biotechnology*, 16: 3).

SUMMARY OF THE INVENTION

25 The present invention relates to the identification of novel SNPs, unique combinations of such SNPs, and haplotypes of SNPs that are associated with VT. The polymorphisms disclosed herein are directly useful as targets for the design of diagnostic reagents and the development of therapeutic agents for use in the diagnosis and treatment of VT.

30 Based on the identification of SNPs associated with VT, the present invention also provides methods of detecting these variants as well as the design and preparation of detection reagents needed to accomplish this task. The invention specifically provides, for example, novel SNPs in genetic sequences involved in VT, isolated nucleic acid molecules (including, for example, DNA and RNA molecules) containing these SNPs, variant proteins encoded by nucleic acid molecules containing such SNPs, antibodies to the encoded variant proteins, computer-based

and data storage systems containing the novel SNP information, methods of detecting these SNPs in a test sample, methods of identifying individuals who have an altered (*i.e.*, increased or decreased) risk of developing VT based on the presence or absence of one or more particular nucleotides (alleles) at one or more SNP sites disclosed herein or the detection of one or more
5 encoded variant products (*e.g.*, variant mRNA transcripts or variant proteins), methods of identifying individuals who are more or less likely to respond to a treatment (or more or less likely to experience undesirable side effects from a treatment, etc.), methods of screening for compounds useful in the treatment of a disorder associated with a variant gene/protein, compounds identified by these methods, methods of treating disorders mediated by a variant
10 gene/protein, methods of using the novel SNPs of the present invention for human identification, etc.

In Tables 1-2, the present invention provides gene information, references to the identification of transcript sequences (SEQ ID NOS: 1-20), encoded amino acid sequences (SEQ ID NOS: 21-40), genomic sequences (SEQ ID NOS: 69-81), transcript-based context sequences
15 (SEQ ID NOS: 41-68) and genomic-based context sequences (SEQ ID NOS: 82-199) that contain the SNPs of the present invention, and extensive SNP information that includes observed alleles, allele frequencies, populations/ethnic groups in which alleles have been observed, information about the type of SNP and corresponding functional effect, and, for cSNPs, information about the encoded polypeptide product. The actual transcript sequences (SEQ ID
20 NOS: 1-20), amino acid sequences (SEQ ID NOS: 21-40), genomic sequences (SEQ ID NOS: 69-81), transcript-based SNP context sequences (SEQ ID NOS: 41-68), and genomic-based SNP context sequences (SEQ ID NOS: 82-199), together with primer sequences (SEQ ID NOS: 200-223) are provided in the Sequence Listing.

In one embodiment of the invention, applicants teach a method for identifying an
25 individual who has an altered risk for developing VT, comprising detecting a single nucleotide polymorphism (SNP) in any one of the nucleotide sequences of SEQ ID NOS: 1-20, SEQ ID NOS: 41-68, SEQ ID NOS: 69-81, and SEQ ID NOS: 82-199 in said individual's nucleic acids, wherein the SNP is as specified in Table 1 and Table 2, respectively, and the presence of the SNP is correlated with an altered risk for VT in said individual. In a specific embodiment of the
30 present invention, SNPs that occur naturally in the human genome are provided as isolated nucleic acid molecules. These SNPs are associated with VT, such that they can have a variety of uses in the diagnosis and/or treatment of VT and related pathologies. In an alternative embodiment, a nucleic acid of the invention is an amplified polynucleotide, which is produced by amplification of a SNP-containing nucleic acid template. In another embodiment, the invention

provides for a variant protein that is encoded by a nucleic acid molecule containing a SNP disclosed herein.

In yet another embodiment of the invention, a reagent for detecting a SNP in the context of its naturally-occurring flanking nucleotide sequences (which can be, *e.g.*, either DNA or mRNA) is provided. In particular, such a reagent may be in the form of, for example, a hybridization probe or an amplification primer that is useful in the specific detection of a SNP of interest. In an alternative embodiment, a protein detection reagent is used to detect a variant protein that is encoded by a nucleic acid molecule containing a SNP disclosed herein. A preferred embodiment of a protein detection reagent is an antibody or an antigen-reactive antibody fragment.

Various embodiments of the invention also provide kits comprising SNP detection reagents, and methods for detecting the SNPs disclosed herein by employing detection reagents. In a specific embodiment, the present invention provides for a method of identifying an individual having an increased or decreased risk of developing VT by detecting the presence or absence of one or more SNP alleles disclosed herein. In another embodiment, a method for diagnosis of VT by detecting the presence or absence of one or more SNP alleles disclosed herein is provided.

The nucleic acid molecules of the invention can be inserted in an expression vector, such as to produce a variant protein in a host cell. Thus, the present invention also provides for a vector comprising a SNP-containing nucleic acid molecule, genetically-engineered host cells containing the vector, and methods for expressing a recombinant variant protein using such host cells. In another specific embodiment, the host cells, SNP-containing nucleic acid molecules, and/or variant proteins can be used as targets in a method for screening and identifying therapeutic agents or pharmaceutical compounds useful in the treatment of VT.

An aspect of this invention is a method for treating VT in a human subject wherein said human subject harbors a SNP, gene, transcript, and/or encoded protein identified in Tables 1-2, which method comprises administering to said human subject a therapeutically or prophylactically effective amount of one or more agents counteracting the effects of the disease, such as by inhibiting (or stimulating) the activity of the gene, transcript, and/or encoded protein identified in Tables 1-2.

Another aspect of this invention is a method for identifying an agent useful in therapeutically or prophylactically treating VT in a human subject wherein said human subject harbors a SNP, gene, transcript, and/or encoded protein identified in Tables 1-2, which method comprises contacting the gene, transcript, or encoded protein with a candidate agent under

conditions suitable to allow formation of a binding complex between the gene, transcript, or encoded protein and the candidate agent and detecting the formation of the binding complex, wherein the presence of the complex identifies said agent.

Another aspect of this invention is a method for treating VT in a human subject, which method comprises:

(i) determining that said human subject harbors a SNP, gene, transcript, and/or encoded protein identified in Tables 1-2, and

(ii) administering to said subject a therapeutically or prophylactically effective amount of one or more agents counteracting the effects of the disease such as anticoagulant agents.

Many other uses and advantages of the present invention will be apparent to those skilled in the art upon review of the detailed description of the preferred embodiments herein. Solely for clarity of discussion, the invention is described in the sections below by way of non-limiting examples.

Various embodiments of the invention provide a method of determining whether a human has an increased risk for venous thrombosis (VT), the method comprising testing nucleic acid from said human for the presence or absence of a polymorphism in gene *F9* at position 101 of SEQ ID NO:82 or its complement, wherein A at position 101 of SEQ ID NO:82 or T at position 101 of its complement indicates said human has said increased risk for VT.

Description of the Files Contained On the CD-R Named CD000012ORD

Each of the CD-Rs contains the following file:

File SEQLIST_CD000012PCT.txt provides the Sequence Listing. The Sequence Listing provides the transcript sequences (SEQ ID NOS: 1-20) and protein sequences (SEQ ID NOS: 21-40) as referred to in Table 1, and genomic sequences (SEQ ID NOS: 69-81) as referred to in Table 2, for each VT-associated gene or genomic region (for intergenic SNPs) that contains one or more SNPs of the present invention. Also provided in the Sequence Listing are context sequences flanking each SNP, including both transcript-based context sequences as referred to in Table 1 (SEQ ID NOS: 41-68) and genomic-based context sequences as referred to in Table 2 (SEQ ID NOS: 82-199). In addition, the Sequence Listing provides the primer sequences from Table 3 (SEQ ID NOS: 200-223), which are oligonucleotides that have been synthesized and used in the laboratory to assay the SNPs disclosed in Tables 4-5 during the course of association studies to verify the association of these SNPs with VT. The context sequences generally provide 100bp upstream (5') and 100bp downstream (3') of each SNP, with the SNP in the middle of the context sequence, for a total of 200bp of context sequence surrounding each SNP.

File SEQLIST_CD000012PCT.txt is 1,794KB in size, and was created on October 18, 2007. A computer readable format of the sequence listing is also submitted herein on a separate CDR labeled CRF. The

information recorded in the CRF CDR is identical to the sequence listing as provided on the CDR Duplicate Copy 1 and Duplicate Copy 2.

DESCRIPTION OF TABLE 1 AND TABLE 2

Table 1 and Table 2 (both provided on the CD-R) disclose the SNP and associated gene/transcript/protein information of the present invention. For each gene, Table 1 provides a header containing gene, transcript and protein information, followed by a transcript and protein sequence identifier (SEQ ID), and then SNP information regarding each SNP found in that gene/transcript including the transcript context sequence. For each gene in Table 2, a header is provided that contains gene and genomic information, followed by a genomic sequence identifier (SEQ ID) and then SNP information regarding each SNP found in that gene, including the genomic context sequence.

Note that SNP markers may be included in both Table 1 and Table 2; Table 1 presents the SNPs relative to their transcript sequences and encoded protein sequences, whereas Table 2 presents the SNPs relative to their genomic sequences. In some instances Table 2 may also include, after the last gene sequence, genomic sequences of one or more intergenic regions, as well as SNP context sequences and other SNP information for any SNPs that lie within these intergenic regions. Additionally, in either Table 1 or 2 a "Related Interrogated SNP" may be listed following a SNP which is determined to be in LD with that interrogated SNP according to the given Power value. SNPs can readily be cross-referenced between all Tables based on their Celera hCV (or, in some instances, hDV) identification numbers, and to the Sequence Listing based on their corresponding SEQ ID NOS.

The gene/transcript/protein information includes:

- a gene number (1 through n, where n = the total number of genes in the Table)
- a Celera hCG and UID internal identification numbers for the gene
- a Celera hCT and UID internal identification numbers for the transcript (Table 1 only)
- a public Genbank accession number (*e.g.*, RefSeq NM number) for the transcript (Table 1 only)
- a Celera hCP and UID internal identification numbers for the protein encoded by the hCT transcript (Table 1 only)
- a public Genbank accession number (*e.g.*, RefSeq NP number) for the protein (Table 1 only)
- an art-known gene symbol
- an art-known gene/protein name

- Celera genomic axis position (indicating start nucleotide position-stop nucleotide position)

- the chromosome number of the chromosome on which the gene is located

5 - an OVTM (Online Mendelian Inheritance in Man; Johns Hopkins University/NCBI) public reference number for obtaining further information regarding the medical significance of each gene

- alternative gene/protein name(s) and/or symbol(s) in the OVTEM entry

Note that, due to the presence of alternative splice forms, multiple transcript/protein entries may be provided for a single gene entry in Table 1; *i.e.*, for a single Gene Number, 10 multiple entries may be provided in series that differ in their transcript/protein information and sequences.

Following the gene/transcript/protein information is a transcript context sequence and (Table 1), or a genomic context sequence (Table 2), for each SNP within that gene.

15 After the last gene sequence, Table 2 may include additional genomic sequences of intergenic regions (in such instances, these sequences are identified as "Intergenic region:" followed by a numerical identification number), as well as SNP context sequences and other SNP information for any SNPs that lie within each intergenic region (such SNPs are identified as "INTERGENIC" for SNP type).

Note that the transcript, protein, and transcript-based SNP context sequences are all 20 provided in the Sequence Listing. The transcript-based SNP context sequences are provided in both Table 1 and also in the Sequence Listing. The genomic and genomic-based SNP context sequences are provided in the Sequence Listing. The genomic-based SNP context sequences are provided in both Table 2 and in the Sequence Listing. SEQ ID NOS are indicated in Table 1 for the transcript-based context sequences (SEQ ID NOS: 41-68); SEQ ID NOS are indicated in 25 Table 2 for the genomic-based context sequences (SEQ ID NOS: 82-199).

The SNP information includes:

- context sequence (taken from the transcript sequence in Table 1, the genomic sequence in Table 2) with the SNP represented by its IUB code, including 100 bp upstream (5') of the SNP position plus 100 bp downstream (3') of the SNP position (the transcript-based SNP context 30 sequences in Table 1 are provided in the Sequence Listing as SEQ ID NOS: 41-68; the genomic-based SNP context sequences in Table 2 are provided in the Sequence Listing as SEQ ID NOS: 82-199).

- Celera hCV internal identification number for the SNP (in some instances, an "hDV" number is given instead of an "hCV" number).

- The corresponding public identification number for the SNP, the RS number.
- SNP position (position of the SNP within the given transcript sequence (Table 1) or within the given genomic sequence (Table 2)).
- “Related Interrogated SNP” is as the interrogated SNP with which the listed SNP is in LD at the given value of Power.
- SNP source (may include any combination of one or more of the following five codes, depending on which internal sequencing projects and/or public databases the SNP has been observed in: “Applera” = SNP observed during the re-sequencing of genes and regulatory regions of 39 individuals, “Celera” = SNP observed during shotgun sequencing and assembly of the Celera human genome sequence, “Celera Diagnostics” = SNP observed during re-sequencing of nucleic acid samples from individuals who have a disease, “dbSNP” = SNP observed in the dbSNP public database, “HGBASE” = SNP observed in the HGBASE public database, “HGMD” = SNP observed in the Human Gene Mutation Database (HGMD) public database, “HapMap” = SNP observed in the International HapMap Project public database, “CSNP” = SNP observed in an internal Applied Biosystems (Foster City, CA) database of coding SNPS (cSNPs)). Note that multiple “Applera” source entries for a single SNP indicate that the same SNP was covered by multiple overlapping amplification products and the re-sequencing results (*e.g.*, observed allele counts) from each of these amplification products is being provided.
- Population/allele/allele count information in the format of [population1(first_allele,count|second_allele,count)population2(first_allele,count|second_allele,count) total (first_allele,total count|second_allele,total count)]. The information in this field includes populations/ethnic groups in which particular SNP alleles have been observed (“cau” = Caucasian, “his” = Hispanic, “chn” = Chinese, and “afr” = African-American, “jpn” = Japanese, “ind” = Indian, “mex” = Mexican, “ain” = “American Indian, “cra” = Celera donor, “no_pop” = no population information available), identified SNP alleles, and observed allele counts (within each population group and total allele counts), where available [“-“ in the allele field represents a deletion allele of an insertion/deletion (“indel”) polymorphism (in which case the corresponding insertion allele, which may be comprised of one or more nucleotides, is indicated in the allele field on the opposite side of the “|”); “-“ in the count field indicates that allele count information is not available]. For certain SNPs from the public dbSNP database, population/ethnic information is indicated as follows (this population information is publicly available in dbSNP): “HISP1” = human individual DNA (anonymized samples) from 23 individuals of self-described HISPANIC heritage; “PAC1” = human individual DNA (anonymized samples) from 24 individuals of self-described PACIFIC RIM heritage; “CAUC1” = human individual DNA

(anonymized samples) from 31 individuals of self-described CAUCASIAN heritage; “AFR1” = human individual DNA (anonymized samples) from 24 individuals of self-described AFRICAN/AFRICAN AMERICAN heritage; “P1” = human individual DNA (anonymized samples) from 102 individuals of self-described heritage; “PA130299515”; “SC_12_A” = SANGER 12 DNAs of Asian origin from Coriell cell repositories, 6 of which are male and 6 female; “SC_12_C” = SANGER 12 DNAs of Caucasian origin from Coriell cell repositories from the CEPH/UTAH library. Six male and 6 female; “SC_12_AA” = SANGER 12 DNAs of African-American origin from Coriell cell repositories 6 of which are male and 6 female; “SC_95_C” = SANGER 95 DNAs of Caucasian origin from Coriell cell repositories from the CEPH/UTAH library; and “SC_12_CA” = Caucasians - 12 DNAs from Coriell cell repositories that are from the CEPH/UTAH library.

Note that for SNPs of “Applera” SNP source, genes/regulatory regions of 39 individuals (20 Caucasians and 19 African Americans) were re-sequenced and, since each SNP position is represented by two chromosomes in each individual (with the exception of SNPs on X and Y chromosomes in males, for which each SNP position is represented by a single chromosome), up to 78 chromosomes were genotyped for each SNP position. Thus, the sum of the African-American (“afr”) allele counts is up to 38, the sum of the Caucasian allele counts (“cau”) is up to 40, and the total sum of all allele counts is up to 78.

Note that semicolons separate population/allele/count information corresponding to each indicated SNP source; *i.e.*, if four SNP sources are indicated, such as “Celera,” “dbSNP,” “HGBASE,” and “HGMD,” then population/allele/count information is provided in four groups which are separated by semicolons and listed in the same order as the listing of SNP sources, with each population/allele/count information group corresponding to the respective SNP source based on order; thus, in this example, the first population/allele/count information group would correspond to the first listed SNP source (Celera) and the third population/allele/count information group separated by semicolons would correspond to the third listed SNP source (HGBASE); if population/allele/count information is not available for any particular SNP source, then a pair of semicolons is still inserted as a place-holder in order to maintain correspondence between the list of SNP sources and the corresponding listing of population/allele/count information.

- SNP type (*e.g.*, location within gene/transcript and/or predicted functional effect)
 [“VTES-SENSE MUTATION” = SNP causes a change in the encoded amino acid (*i.e.*, a non-synonymous coding SNP); “SILENT MUTATION” = SNP does not cause a change in the encoded amino acid (*i.e.*, a synonymous coding SNP); “STOP CODON MUTATION” = SNP is

located in a stop codon; “NONSENSE MUTATION” = SNP creates or destroys a stop codon; “UTR 5” = SNP is located in a 5’ UTR of a transcript; “UTR 3” = SNP is located in a 3’ UTR of a transcript; “PUTATIVE UTR 5” = SNP is located in a putative 5’ UTR; “PUTATIVE UTR 3” = SNP is located in a putative 3’ UTR; “DONOR SPLICE SITE” = SNP is located in a donor splice site (5’ intron boundary); “ACCEPTOR SPLICE SITE” = SNP is located in an acceptor splice site (3’ intron boundary); “CODING REGION” = SNP is located in a protein-coding region of the transcript; “EXON” = SNP is located in an exon; “INTRON” = SNP is located in an intron; “hmCS” = SNP is located in a human-mouse conserved segment; “TFBS” = SNP is located in a transcription factor binding site; “UNKNOWN” = SNP type is not defined; “INTERGENIC” = SNP is intergenic, *i.e.*, outside of any gene boundary].

- Protein coding information (Table 1 only), where relevant, in the format of [protein SEQ ID NO: #, amino acid position, (amino acid-1, codon1) (amino acid-2, codon2)]. The information in this field includes SEQ ID NO of the encoded protein sequence, position of the amino acid residue within the protein identified by the SEQ ID NO that is encoded by the codon containing the SNP, amino acids (represented by one-letter amino acid codes) that are encoded by the alternative SNP alleles (in the case of stop codons, “X” is used for the one-letter amino acid code), and alternative codons containing the alternative SNP nucleotides which encode the amino acid residues (thus, for example, for missense mutation-type SNPs, at least two different amino acids and at least two different codons are generally indicated; for silent mutation-type SNPs, one amino acid and at least two different codons are generally indicated, etc.). In instances where the SNP is located outside of a protein-coding region (*e.g.*, in a UTR region), “None” is indicated following the protein SEQ ID NO.

DESCRIPTION OF TABLE 3

Table 3 provides sequences (SEQ ID NOS: 200-223) of oligonucleotides that have been synthesized and used in the laboratory to assay the SNPs disclosed in Tables 4-5 during the course of association studies to verify the association of these SNPs with VT. The experiments that were conducted using these primers are explained in detail in Example 1, below.

Table 3 provides the following:

- the column labeled “Marker” lists the Celera identifier hCV number for each SNP marker.

- the column labeled “Alleles” designates the two alternative alleles at the SNP site identified by the hCV identification number that are targeted by the allele-specific oligonucleotides.

- allele-specific oligonucleotides with their respective SEQ ID numbers are shown in the next two columns, "Sequence A (allele-specific primer)" and "Sequence B (allele-specific primer)." These two primers were used in conjunction with a common primer in each PCR assay to genotype DNA samples for each SNP marker. Note that alleles may be presented in Table 3 based on a different orientation (i.e., the reverse complement) relative to how the same alleles are presented in Tables 1 and 2.

- common oligonucleotides with their respective SEQ ID numbers are shown in the column, "Sequence C (common primer)." Each common primer was used in conjunction with the two allele-specific primers to genotype DNA samples for each SNP marker.

All sequences are given in the 5' to 3' direction.

DESCRIPTION OF TABLE 4

Table 4 provides results of statistical analyses for certain SNPs disclosed in Tables 1 and 2 (SNPs can be cross-referenced between tables based on their hCV identification numbers), and the association of these SNPs with VT based on a genotyping analysis, unstratified by patient phenotype. The experiment that provided this data is explained in detail in Example 1, below.

SNP association with VT was found when cases with VT were compared to controls who did not have VT. From patients at the Leiden Thrombophilia Study (LETS), 866 samples were obtained; from the Multiple Environmental and Genetic Assessment (MEGA) study, 4,383 samples were obtained. The numbers of cases and controls genotyped for each assay are provided under the respective columns "Case" and "Cont."

The statistical results provided in Table 4 show that the association of these SNPs with VT is supported by P values < 0.05 in an allelic association test, based on either a dominant/recessive mode of inheritance, or homozygous/heterozygous (Mode column). "Rec" in the Mode column indicates the trait is recessive; "Hom" indicates the risk association is observed in individuals homozygous for the risk allele.

In Table 4, the column labeled "Marker" presents each SNP as identified by its unique Celera hCV identification number. The column labeled "RS" presents each SNP as identified by its reference sequence (rs) number as assigned by NCBI. The column "Gene Symbol" presents the standard symbol for the gene containing the SNP; i.e., the symbol approved by the Human Genome Organization (HUGO) Gene Nomenclature Committee. The column labeled "Risk Allele" presents the variant nucleotide for each of the identified SNPs that is associated with risk. The column labeled "Non Risk Allele" presents the wild type nucleotide for each of the

identified SNPs that is not associated with risk. The allele may be presented in Table 4 as the reverse complement relative to how the same allele is presented in Tables 1 and/or 2.

The column labeled "Mode" indicates the genetic mode under which the P value for association was calculated. Under a genotypic analysis (described in examples below), when two
5 copies of the SNP are required to see the observed effect, the mode is recessive, or "Rec." When one or two copies of the SNP are required to see the association, the mode is dominant, or "Dom." When the association is found by simply comparing the frequency of the allele in the case population to the control population, the mode is "Allelic." The allelic mode closely approximates an additive mode. The column labeled "P val" indicates the results of either the
10 asymptotic chi-square test for genotypic association (Rec or Dom), or the Fisher Exact test (Allelic) to determine if the qualitative phenotype is a function of the SNP genotype. The column labeled "OR" (odds ratio) indicates an approximation of the relative risk for an individual for the defined endpoint associated with the SNP. An OR of less than one indicates that the allele is protective for VT, and an OR greater than one indicates the allele increases the risk of
15 VT.

Note that SNPs can be cross-referenced between the tables herein based on their hCV identification numbers. Some of the SNPs that are included in the tables may possess two different hCV identification numbers. For instance, hCV916107 represents the same SNP as hCV26887450.

20

DESCRIPTION OF TABLE 5 (stratified by sex, female)

The column labeled "Stratum" lists the subgroups of individuals from cases and controls in which VT association was observed. Reference is made to the Stratum Key below the table for an explanation of symbols used.

25

Table 5 provides the results of statistical analyses for certain SNPs disclosed in Tables 1 and 2; namely, the association of SNP alleles with a risk for VT based on case-control studies. The experiment that provided this data is explained in detail in Example 1, below. Note that SNPs can be cross-referenced between tables based on their hCV identification numbers. The statistical results provided in Table 5 show that the association of these SNPs with VT is
30 supported by P values <0.05 in allelic association tests. The data presented were obtained from individually genotyped samples. Case samples were limited to patients that had a history of VT, while controls had no history of VT.

In Table 5, the column labeled "Marker" presents each SNP as identified by its unique Celera hCV identification number. The column labeled "RS" presents each SNP as identified by

its reference sequence (rs) number as assigned by NCBI. The column "Gene Symbol" presents the standard symbol for the gene containing the SNP; i.e., the symbol approved by the Human Genome Organization (HUGO) Gene Nomenclature Committee. The column labeled "Risk Allele" presents the variant nucleotide for each of the identified SNPs that is associated with risk.

5 The column labeled "Non Risk Allele" presents the wild type nucleotide for each of the identified SNPs that is not associated with risk. Each allele may be presented in Tables 1 and/or 2 as the complement of the allele presented in Table 5; e.g., "G" may be presented as its complement, "C." The column labeled "P val" indicates the results of the Fisher Exact test, to determine the association of one allele with risk for VT. The column labeled "OR" (odds ratio) shows an approximation of the relative VT risk for individuals with the risk allele, based on the observed frequencies of alleles in cases vs. controls. An OR less than one would indicate an allele is protective for VT, and an OR greater than one indicates the allele is associated with an increased risk of VT.

15 DESCRIPTION OF TABLE 6

Table 6 provides a list of the sample LD SNPs that are related to and derived from an interrogated SNP. These LD SNPs are provided as an example of the groups of SNPs which can also serve as markers for disease association based on their being in LD with the interrogated SNP. The criteria and process of selecting such LD SNPs, including the calculation of the r^2 value and the r^2 threshold value, are described in Example, below.

In Table 6, the column labeled "Interrogated SNP" presents each marker as identified by its unique identifier, the hCV number. The column labeled "Interrogated rs" presents the publicly known identifier rs number for the corresponding hCV number. The column labeled "LD SNP" presents the hCV numbers of the LD SNPs that are derived from their corresponding interrogated SNPs. The column labeled "LD SNP rs" presents the publicly known rs number for the corresponding hCV number. The column labeled "Power (T)" presents the level of power where the r^2 threshold is set. For example, when power is set at 51%, the threshold r^2 value calculated therefrom is the minimum r^2 that an LD SNP must have in reference to an interrogated SNP, in order for the LD SNP to be classified as a marker capable of being associated with a disease phenotype at greater than 51 % probability. The column labeled "Threshold r^2 " presents the minimum value of r^2 that an LD SNP must meet in reference to an interrogated SNP in order to qualify as an LD SNP. The column labeled " r^2 " presents the actual r^2 value of the LD SNP in reference to the interrogated SNP to which it is related.

DETAILED DESCRIPTION OF THE INVENTION

The present invention provides SNPs associated with VT, nucleic acid molecules containing SNPs, methods and reagents for the detection of the SNPs disclosed herein, uses of these SNPs for the development of detection reagents, and assays or kits that utilize such reagents. The VT-associated SNPs disclosed herein are useful for diagnosing, screening for, and evaluating predisposition to VT and related pathologies in humans. Furthermore, such SNPs and their encoded products are useful targets for the development of therapeutic agents.

A large number of SNPs have been identified from re-sequencing DNA from 39 individuals, and they are indicated as "Applera" SNP source in Tables 1-2. Their allele frequencies observed in each of the Caucasian and African-American ethnic groups are provided. Additional SNPs included herein were previously identified during shotgun sequencing and assembly of the human genome, and they are indicated as "Celera" SNP source in Tables 1-2. Furthermore, the information provided in Table 1-2, particularly the allele frequency information obtained from 39 individuals and the identification of the precise position of each SNP within each gene/transcript, allows haplotypes (*i.e.*, groups of SNPs that are co-inherited) to be readily inferred. The present invention encompasses SNP haplotypes, as well as individual SNPs.

Thus, the present invention provides individual SNPs associated with VT, as well as combinations of SNPs and haplotypes in genetic regions associated with VT, polymorphic/variant transcript sequences (SEQ ID NOS: 1-20) and genomic sequences (SEQ ID NOS: 69-81) containing SNPs, encoded amino acid sequences (SEQ ID NOS: 21-40), and both transcript-based SNP context sequences (SEQ ID NOS: 41-68) and genomic-based SNP context sequences (SEQ ID NOS: 82-199) (transcript sequences, protein sequences, and transcript-based SNP context sequences are provided in Table 1 and the Sequence Listing; genomic sequences and genomic-based SNP context sequences are provided in Table 2 and the Sequence Listing), methods of detecting these polymorphisms in a test sample, methods of determining the risk of an individual of having or developing VT, methods of screening for compounds useful for treating disorders associated with a variant gene/protein such as VT, compounds identified by these screening methods, methods of using the disclosed SNPs to select a treatment strategy, methods of treating a disorder associated with a variant gene/protein (*i.e.*, therapeutic methods), methods of determining if an individual is likely to respond to a specific treatment and methods of using the SNPs of the present invention for human identification.

The present invention provides novel SNPs associated with VT, as well as SNPs that were previously known in the art, but were not previously known to be associated with VT. Accordingly, the present invention provides novel compositions and methods based on the novel

SNPs disclosed herein, and also provides novel methods of using the known, but previously unassociated, SNPs in methods relating to VT (*e.g.*, for diagnosing VT, etc.). In Tables 1-2, known SNPs are identified based on the public database in which they have been observed, which is indicated as one or more of the following SNP types: “dbSNP” = SNP observed in dbSNP, “HGBASE” = SNP observed in HGBASE, and “HGMD” = SNP observed in the Human Gene Mutation Database (HGMD). Novel SNPs for which the SNP source is only “Applera” and none other, *i.e.*, those that have not been observed in any public databases and which were also not observed during shotgun sequencing and assembly of the Celera human genome sequence (*i.e.*, “Celera” SNP source), are also noted in the tables.

Particular SNP alleles of the present invention can be associated with either an increased risk of having or developing VT, or a decreased risk of having or developing VT. SNP alleles that are associated with a decreased risk of having or developing VT may be referred to as “protective” alleles, and SNP alleles that are associated with an increased risk of having or developing VT may be referred to as “susceptibility” alleles, “risk” alleles, or “risk factors”.

Thus, whereas certain SNPs (or their encoded products) can be assayed to determine whether an individual possesses a SNP allele that is indicative of an increased risk of having or developing VT (*i.e.*, a susceptibility allele), other SNPs (or their encoded products) can be assayed to determine whether an individual possesses a SNP allele that is indicative of a decreased risk of having or developing VT (*i.e.*, a protective allele). Similarly, particular SNP alleles of the present invention can be associated with either an increased or decreased likelihood of responding to a particular treatment or therapeutic compound, or an increased or decreased likelihood of experiencing toxic effects from a particular treatment or therapeutic compound. The term “altered” may be used herein to encompass either of these two possibilities (*e.g.*, an increased or a decreased risk/likelihood).

Those skilled in the art will readily recognize that nucleic acid molecules may be double-stranded molecules and that reference to a particular site on one strand refers, as well, to the corresponding site on a complementary strand. In defining a SNP position, SNP allele, or nucleotide sequence, reference to an adenine, a thymine (uridine), a cytosine, or a guanine at a particular site on one strand of a nucleic acid molecule also defines the thymine (uridine), adenine, guanine, or cytosine (respectively) at the corresponding site on a complementary strand of the nucleic acid molecule. Thus, reference may be made to either strand in order to refer to a particular SNP position, SNP allele, or nucleotide sequence. Probes and primers, may be designed to hybridize to either strand and SNP genotyping methods disclosed herein may

generally target either strand. Throughout the specification, in identifying a SNP position, reference is generally made to the protein-encoding strand, only for the purpose of convenience.

References to variant peptides, polypeptides, or proteins of the present invention include peptides, polypeptides, proteins, or fragments thereof, that contain at least one amino acid residue that differs from the corresponding amino acid sequence of the art-known peptide/polypeptide/protein (the art-known protein may be interchangeably referred to as the “wild-type,” “reference,” or “normal” protein). Such variant peptides/polypeptides/proteins can result from a codon change caused by a nonsynonymous nucleotide substitution at a protein-coding SNP position (*i.e.*, a missense mutation) disclosed by the present invention. Variant peptides/polypeptides/proteins of the present invention can also result from a nonsense mutation, *i.e.*, a SNP that creates a premature stop codon, a SNP that generates a read-through mutation by abolishing a stop codon, or due to any SNP disclosed by the present invention that otherwise alters the structure, function/activity, or expression of a protein, such as a SNP in a regulatory region (*e.g.* a promoter or enhancer) or a SNP that leads to alternative or defective splicing, such as a SNP in an intron or a SNP at an exon/intron boundary. As used herein, the terms “polypeptide,” “peptide,” and “protein” are used interchangeably.

ISOLATED NUCLEIC ACID MOLECULES AND SNP DETECTION REAGENTS & KITS

Tables 1 and 2 provide a variety of information about each SNP of the present invention that is associated with VT, including the transcript sequences (SEQ ID NOS: 1-20), genomic sequences (SEQ ID NOS: 69-81), and protein sequences (SEQ ID NOS: 21-40) of the encoded gene products (with the SNPs indicated by IUB codes in the nucleic acid sequences). In addition, Tables 1 and 2 include SNP context sequences, which generally include 100 nucleotide upstream (5') plus 100 nucleotides downstream (3') of each SNP position (SEQ ID NOS: 41-68 correspond to transcript-based SNP context sequences disclosed in Table 1, and SEQ ID NOS: 82-199 correspond to genomic-based context sequences disclosed in Table 2), the alternative nucleotides (alleles) at each SNP position, and additional information about the variant where relevant, such as SNP type (coding, missense, splice site, UTR, etc.), human populations in which the SNP was observed, observed allele frequencies, information about the encoded protein, etc.

Isolated Nucleic Acid Molecules

The present invention provides isolated nucleic acid molecules that contain one or more SNPs disclosed Table 1 and/or Table 2. Isolated nucleic acid molecules containing one or more SNPs disclosed in at least one of Tables 1-2 may be interchangeably referred to throughout the present text as “SNP-containing nucleic acid molecules”. Isolated nucleic acid molecules may optionally encode a full-length variant protein or fragment thereof. The isolated nucleic acid molecules of the present invention also include probes and primers (which are described in greater detail below in the section entitled “SNP Detection Reagents”), which may be used for assaying the disclosed SNPs, and isolated full-length genes, transcripts, cDNA molecules, and fragments thereof, which may be used for such purposes as expressing an encoded protein.

As used herein, an “isolated nucleic acid molecule” generally is one that contains a SNP of the present invention or one that hybridizes to such molecule such as a nucleic acid with a complementary sequence, and is separated from most other nucleic acids present in the natural source of the nucleic acid molecule. Moreover, an “isolated” nucleic acid molecule, such as a cDNA molecule containing a SNP of the present invention, can be substantially free of other cellular material, or culture medium when produced by recombinant techniques, or chemical precursors or other chemicals when chemically synthesized. A nucleic acid molecule can be fused to other coding or regulatory sequences and still be considered “isolated”. Nucleic acid molecules present in non-human transgenic animals, which do not naturally occur in the animal, are also considered “isolated”. For example, recombinant DNA molecules contained in a vector are considered “isolated”. Further examples of “isolated” DNA molecules include recombinant DNA molecules maintained in heterologous host cells, and purified (partially or substantially) DNA molecules in solution. Isolated RNA molecules include in vivo or in vitro RNA transcripts of the isolated SNP-containing DNA molecules of the present invention. Isolated nucleic acid molecules according to the present invention further include such molecules produced synthetically.

Generally, an isolated SNP-containing nucleic acid molecule comprises one or more SNP positions disclosed by the present invention with flanking nucleotide sequences on either side of the SNP positions. A flanking sequence can include nucleotide residues that are naturally associated with the SNP site and/or heterologous nucleotide sequences. Preferably the flanking sequence is up to about 500, 300, 100, 60, 50, 30, 25, 20, 15, 10, 8, or 4 nucleotides (or any other length in-between) on either side of a SNP position, or as long as the full-length gene or entire protein-coding sequence (or any portion thereof such as an exon), especially if the SNP-containing nucleic acid molecule is to be used to produce a protein or protein fragment.

For full-length genes and entire protein-coding sequences, a SNP flanking sequence can be, for example, up to about 5KB, 4KB, 3KB, 2KB, 1KB on either side of the SNP. Furthermore, in

such instances, the isolated nucleic acid molecule comprises exonic sequences (including protein-coding and/or non-coding exonic sequences), but may also include intronic sequences. Thus, any protein coding sequence may be either contiguous or separated by introns. The important point is that the nucleic acid is isolated from remote and unimportant flanking sequences and is of
5 appropriate length such that it can be subjected to the specific manipulations or uses described herein such as recombinant protein expression, preparation of probes and primers for assaying the SNP position, and other uses specific to the SNP-containing nucleic acid sequences.

An isolated SNP-containing nucleic acid molecule can comprise, for example, a full-length gene or transcript, such as a gene isolated from genomic DNA (*e.g.*, by cloning or PCR
10 amplification), a cDNA molecule, or an mRNA transcript molecule. Polymorphic transcript sequences are referred to in Table 1 and provided in the Sequence Listing (SEQ ID NOS: 1-20), and polymorphic genomic sequences are referred to in Table 2 and provided in the Sequence Listing (SEQ ID NOS: 69-81). Furthermore, fragments of such full-length genes and transcripts that contain one or more SNPs disclosed herein are also encompassed by the present invention, and such
15 fragments may be used, for example, to express any part of a protein, such as a particular functional domain or an antigenic epitope.

Thus, the present invention also encompasses fragments of the nucleic acid sequences as disclosed in Tables 1-2 (transcript sequences are referred to in Table 1 as SEQ ID NOS: 1-20, genomic sequences are referred to in Table 2 as SEQ ID NOS: 69-81, transcript-based SNP context
20 sequences are referred to in Table 1 as SEQ ID NO: 41-68, and genomic-based SNP context sequences are referred to in Table 2 as SEQ ID NO: 82-199) and their complements. The actual sequences referred to in the tables are provided in the Sequence Listing. A fragment typically comprises a contiguous nucleotide sequence at least about 8 or more nucleotides, more preferably at least about 12 or more nucleotides, and even more preferably at least about 16 or more nucleotides.
25 Further, a fragment could comprise at least about 18, 20, 22, 25, 30, 40, 50, 60, 80, 100, 150, 200, 250 or 500 (or any other number in-between) nucleotides in length. The length of the fragment will be based on its intended use. For example, the fragment can encode epitope-bearing regions of a variant peptide or regions of a variant peptide that differ from the normal/wild-type protein, or can be useful as a polynucleotide probe or primer. Such fragments can be isolated using the nucleotide
30 sequences provided in Table 1 and/or Table 2 for the synthesis of a polynucleotide probe. A labeled probe can then be used, for example, to screen a cDNA library, genomic DNA library, or mRNA to isolate nucleic acid corresponding to the coding region. Further, primers can be used in amplification reactions, such as for purposes of assaying one or more SNPs sites or for cloning specific regions of a gene.

An isolated nucleic acid molecule of the present invention further encompasses a SNP-containing polynucleotide that is the product of any one of a variety of nucleic acid amplification methods, which are used to increase the copy numbers of a polynucleotide of interest in a nucleic acid sample. Such amplification methods are well known in the art, and they include but are not limited to, polymerase chain reaction (PCR) (U.S. Patent Nos. 4,683,195; and 4,683,202; *PCR Technology: Principles and Applications for DNA Amplification*, ed. H.A. Erlich, Freeman Press, NY, NY, 1992), ligase chain reaction (LCR) (Wu and Wallace, *Genomics* 4:560, 1989; Landegren *et al.*, *Science* 241:1077, 1988), strand displacement amplification (SDA) (U.S. Patent Nos. 5,270,184; and 5,422,252), transcription-mediated amplification (TMA) (U.S. Patent No. 5,399,491), linked linear amplification (LLA) (U.S. Patent No. 6,027,923), and the like, and isothermal amplification methods such as nucleic acid sequence based amplification (NASBA), and self-sustained sequence replication (Guatelli *et al.*, *Proc. Natl. Acad. Sci. USA* 87: 1874, 1990). Based on such methodologies, a person skilled in the art can readily design primers in any suitable regions 5' and 3' to a SNP disclosed herein. Such primers may be used to amplify DNA of any length so long that it contains the SNP of interest in its sequence.

As used herein, an "amplified polynucleotide" of the invention is a SNP-containing nucleic acid molecule whose amount has been increased at least two fold by any nucleic acid amplification method performed *in vitro* as compared to its starting amount in a test sample. In other preferred embodiments, an amplified polynucleotide is the result of at least ten fold, fifty fold, one hundred fold, one thousand fold, or even ten thousand fold increase as compared to its starting amount in a test sample. In a typical PCR amplification, a polynucleotide of interest is often amplified at least fifty thousand fold in amount over the unamplified genomic DNA, but the precise amount of amplification needed for an assay depends on the sensitivity of the subsequent detection method used.

Generally, an amplified polynucleotide is at least about 16 nucleotides in length. More typically, an amplified polynucleotide is at least about 20 nucleotides in length. In a preferred embodiment of the invention, an amplified polynucleotide is at least about 30 nucleotides in length. In a more preferred embodiment of the invention, an amplified polynucleotide is at least about 32, 40, 45, 50, or 60 nucleotides in length. In yet another preferred embodiment of the invention, an amplified polynucleotide is at least about 100, 200, 300, 400, or 500 nucleotides in length. While the total length of an amplified polynucleotide of the invention can be as long as an exon, an intron or the entire gene where the SNP of interest resides, an amplified product is typically up to about 1,000 nucleotides in length (although certain amplification methods may generate amplified products greater than 1000 nucleotides in length). More preferably, an

amplified polynucleotide is not greater than about 600-700 nucleotides in length. It is understood that irrespective of the length of an amplified polynucleotide, a SNP of interest may be located anywhere along its sequence.

5 In a specific embodiment of the invention, the amplified product is at least about 201 nucleotides in length, comprises one of the transcript-based context sequences or the genomic-based context sequences shown in Tables 1-2. Such a product may have additional sequences on its 5' end or 3' end or both. In another embodiment, the amplified product is about 101
10 nucleotides in length, and it contains a SNP disclosed herein. Preferably, the SNP is located at the middle of the amplified product (*e.g.*, at position 101 in an amplified product that is 201 nucleotides in length, or at position 51 in an amplified product that is 101 nucleotides in length), or within 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12, 15, or 20 nucleotides from the middle of the amplified product (however, as indicated above, the SNP of interest may be located anywhere along the length of the amplified product).

15 The present invention provides isolated nucleic acid molecules that comprise, consist of, or consist essentially of one or more polynucleotide sequences that contain one or more SNPs disclosed herein, complements thereof, and SNP-containing fragments thereof.

Accordingly, the present invention provides nucleic acid molecules that consist of any of the nucleotide sequences shown in Table 1 and/or Table 2 (transcript sequences are referred to in Table 1 as SEQ ID NOS: 1-20, genomic sequences are referred to in Table 2 as SEQ ID NOS: 69-81,
20 transcript-based SNP context sequences are referred to in Table 1 as SEQ ID NO: 41-68, and genomic-based SNP context sequences are referred to in Table 2 as SEQ ID NO: 82-199), or any nucleic acid molecule that encodes any of the variant proteins referred to in Table 1 (SEQ ID NOS: 21-40). The actual sequences referred to in the tables are provided in the Sequence Listing. A nucleic acid molecule consists of a nucleotide sequence when the nucleotide sequence is the
25 complete nucleotide sequence of the nucleic acid molecule.

The present invention further provides nucleic acid molecules that consist essentially of any of the nucleotide sequences referred to in Table 1 and/or Table 2 (transcript sequences are referred to in Table 1 as SEQ ID NOS:1-20, genomic sequences are referred to in Table 2 as SEQ ID NOS: 69-81, transcript-based SNP context sequences are referred to in Table 1 as SEQ ID NO: 41-68, and
30 genomic-based SNP context sequences are referred to in Table 2 as SEQ ID NO: 82-199), or any nucleic acid molecule that encodes any of the variant proteins referred to in Table 1 (SEQ ID NOS: 21-40). The actual sequences referred to in the tables are provided in the Sequence Listing. A nucleic acid molecule consists essentially of a nucleotide sequence when such a nucleotide sequence is present with only a few additional nucleotide residues in the final nucleic acid molecule.

The present invention further provides nucleic acid molecules that comprise any of the nucleotide sequences shown in Table 1 and/or Table 2 or a SNP-containing fragment thereof (transcript sequences are referred to in Table 1 as SEQ ID NOS: 1-20, genomic sequences are referred to in Table 2 as SEQ ID NOS: 69-81, transcript-based SNP context sequences are referred to in Table 1 as SEQ ID NO: 41-68, and genomic-based SNP context sequences are referred to in Table 2 as SEQ ID NO: 82-199), or any nucleic acid molecule that encodes any of the variant proteins provided in Table 1 (SEQ ID NOS: 21-40). The actual sequences referred to in the tables are provided in the Sequence Listing. A nucleic acid molecule comprises a nucleotide sequence when the nucleotide sequence is at least part of the final nucleotide sequence of the nucleic acid molecule. In such a fashion, the nucleic acid molecule can be only the nucleotide sequence or have additional nucleotide residues, such as residues that are naturally associated with it or heterologous nucleotide sequences. Such a nucleic acid molecule can have one to a few additional nucleotides or can comprise many more additional nucleotides. A brief description of how various types of these nucleic acid molecules can be readily made and isolated is provided below, and such techniques are well known to those of ordinary skill in the art (Sambrook and Russell, 2000, *Molecular Cloning: A Laboratory Manual*, Cold Spring Harbor Press, NY).

The isolated nucleic acid molecules can encode mature proteins plus additional amino or carboxyl-terminal amino acids or both, or amino acids interior to the mature peptide (when the mature form has more than one peptide chain, for instance). Such sequences may play a role in processing of a protein from precursor to a mature form, facilitate protein trafficking, prolong or shorten protein half-life, or facilitate manipulation of a protein for assay or production. As generally is the case *in situ*, the additional amino acids may be processed away from the mature protein by cellular enzymes.

Thus, the isolated nucleic acid molecules include, but are not limited to, nucleic acid molecules having a sequence encoding a peptide alone, a sequence encoding a mature peptide and additional coding sequences such as a leader or secretory sequence (*e.g.*, a pre-pro or pro-protein sequence), a sequence encoding a mature peptide with or without additional coding sequences, plus additional non-coding sequences, for example introns and non-coding 5' and 3' sequences such as transcribed but untranslated sequences that play a role in, for example, transcription, mRNA processing (including splicing and polyadenylation signals), ribosome binding, and/or stability of mRNA. In addition, the nucleic acid molecules may be fused to heterologous marker sequences encoding, for example, a peptide that facilitates purification.

Isolated nucleic acid molecules can be in the form of RNA, such as mRNA, or in the form DNA, including cDNA and genomic DNA, which may be obtained, for example, by molecular

cloning or produced by chemical synthetic techniques or by a combination thereof (Sambrook and Russell, 2000, *Molecular Cloning: A Laboratory Manual*, Cold Spring Harbor Press, NY).

Furthermore, isolated nucleic acid molecules, particularly SNP detection reagents such as probes and primers, can also be partially or completely in the form of one or more types of nucleic acid analogs, such as peptide nucleic acid (PNA) (U.S. Patent Nos. 5,539,082; 5,527,675; 5,623,049; 5,714,331). The nucleic acid, especially DNA, can be double-stranded or single-stranded.

Single-stranded nucleic acid can be the coding strand (sense strand) or the complementary non-coding strand (anti-sense strand). DNA, RNA, or PNA segments can be assembled, for example, from fragments of the human genome (in the case of DNA or RNA) or single nucleotides, short oligonucleotide linkers, or from a series of oligonucleotides, to provide a synthetic nucleic acid molecule. Nucleic acid molecules can be readily synthesized using the sequences provided herein as a reference; oligonucleotide and PNA oligomer synthesis techniques are well known in the art (see, *e.g.*, Corey, "Peptide nucleic acids: expanding the scope of nucleic acid recognition", *Trends Biotechnol.* 1997 Jun;15(6):224-9, and Hyrup et al., "Peptide nucleic acids (PNA): synthesis, properties and potential applications", *Bioorg Med Chem.* 1996 Jan;4(1):5-23).

Furthermore, large-scale automated oligonucleotide/PNA synthesis (including synthesis on an array or bead surface or other solid support) can readily be accomplished using commercially available nucleic acid synthesizers, such as the Applied Biosystems (Foster City, CA) 3900 High-Throughput DNA Synthesizer or Expedite 8909 Nucleic Acid Synthesis System, and the sequence information provided herein.

The present invention encompasses nucleic acid analogs that contain modified, synthetic, or non-naturally occurring nucleotides or structural elements or other alternative/modified nucleic acid chemistries known in the art. Such nucleic acid analogs are useful, for example, as detection reagents (*e.g.*, primers/probes) for detecting one or more SNPs identified in Table 1 and/or Table 2. Furthermore, kits/systems (such as beads, arrays, etc.) that include these analogs are also encompassed by the present invention. For example, PNA oligomers that are based on the polymorphic sequences of the present invention are specifically contemplated. PNA oligomers are analogs of DNA in which the phosphate backbone is replaced with a peptide-like backbone (Lagriffoul *et al.*, *Bioorganic & Medicinal Chemistry Letters*, 4: 1081-1082 (1994), Petersen *et al.*, *Bioorganic & Medicinal Chemistry Letters*, 6: 793-796 (1996), Kumar *et al.*, *Organic Letters* 3(9): 1269-1272 (2001), WO96/04000). PNA hybridizes to complementary RNA or DNA with higher affinity and specificity than conventional oligonucleotides and oligonucleotide analogs. The properties of PNA enable novel molecular biology and biochemistry applications unachievable with traditional oligonucleotides and peptides.

Additional examples of nucleic acid modifications that improve the binding properties and/or stability of a nucleic acid include the use of base analogs such as inosine, intercalators (U.S. Patent No. 4,835,263) and the minor groove binders (U.S. Patent No. 5,801,115). Thus, references herein to nucleic acid molecules, SNP-containing nucleic acid molecules, SNP
5 detection reagents (*e.g.*, probes and primers), oligonucleotides/polynucleotides include PNA oligomers and other nucleic acid analogs. Other examples of nucleic acid analogs and alternative/modified nucleic acid chemistries known in the art are described in *Current Protocols in Nucleic Acid Chemistry*, John Wiley & Sons, N.Y. (2002).

The present invention further provides nucleic acid molecules that encode fragments of
10 the variant polypeptides disclosed herein as well as nucleic acid molecules that encode obvious variants of such variant polypeptides. Such nucleic acid molecules may be naturally occurring, such as paralogs (different locus) and orthologs (different organism), or may be constructed by recombinant DNA methods or by chemical synthesis. Non-naturally occurring variants may be
15 made by mutagenesis techniques, including those applied to nucleic acid molecules, cells, or organisms. Accordingly, the variants can contain nucleotide substitutions, deletions, inversions and insertions (in addition to the SNPs disclosed in Tables 1-2). Variation can occur in either or both the coding and non-coding regions. The variations can produce conservative and/or non-conservative amino acid substitutions.

Further variants of the nucleic acid molecules disclosed in Tables 1-2, such as naturally
20 occurring allelic variants (as well as orthologs and paralogs) and synthetic variants produced by mutagenesis techniques, can be identified and/or produced using methods well known in the art. Such further variants can comprise a nucleotide sequence that shares at least 70-80%, 80-85%, 85-90%, 91%, 92%, 93%, 94%, 95%, 96%, 97%, 98%, or 99% sequence identity with a nucleic acid sequence disclosed in Table 1 and/or Table 2 (or a fragment thereof) and that includes a
25 novel SNP allele disclosed in Table 1 and/or Table 2. Further, variants can comprise a nucleotide sequence that encodes a polypeptide that shares at least 70-80%, 80-85%, 85-90%, 91%, 92%, 93%, 94%, 95%, 96%, 97%, 98%, or 99% sequence identity with a polypeptide sequence disclosed in Table 1 (or a fragment thereof) and that includes a novel SNP allele disclosed in
30 Table 1 and/or Table 2. Thus, an aspect of the present invention that is specifically contemplated are isolated nucleic acid molecules that have a certain degree of sequence variation compared with the sequences shown in Tables 1-2, but that contain a novel SNP allele disclosed herein. In other words, as long as an isolated nucleic acid molecule contains a novel SNP allele disclosed herein, other portions of the nucleic acid molecule that flank the novel SNP allele can vary to some degree from the specific transcript, genomic, and context sequences referred to and shown

in Tables 1-2, and can encode a polypeptide that varies to some degree from the specific polypeptide sequences referred to in Table 1.

To determine the percent identity of two amino acid sequences or two nucleotide sequences of two molecules that share sequence homology, the sequences are aligned for optimal comparison purposes (*e.g.*, gaps can be introduced in one or both of a first and a second amino acid or nucleic acid sequence for optimal alignment and non-homologous sequences can be disregarded for comparison purposes). In a preferred embodiment, at least 30%, 40%, 50%, 60%, 70%, 80%, or 90% or more of the length of a reference sequence is aligned for comparison purposes. The amino acid residues or nucleotides at corresponding amino acid positions or nucleotide positions are then compared. When a position in the first sequence is occupied by the same amino acid residue or nucleotide as the corresponding position in the second sequence, then the molecules are identical at that position (as used herein, amino acid or nucleic acid "identity" is equivalent to amino acid or nucleic acid "homology"). The percent identity between the two sequences is a function of the number of identical positions shared by the sequences, taking into account the number of gaps, and the length of each gap, which need to be introduced for optimal alignment of the two sequences.

The comparison of sequences and determination of percent identity between two sequences can be accomplished using a mathematical algorithm. (*Computational Molecular Biology*, Lesk, A.M., ed., Oxford University Press, New York, 1988; *Biocomputing: Informatics and Genome Projects*, Smith, D.W., ed., Academic Press, New York, 1993; *Computer Analysis of Sequence Data, Part 1*, Griffin, A.M., and Griffin, H.G., eds., Humana Press, New Jersey, 1994; *Sequence Analysis in Molecular Biology*, von Heinje, G., Academic Press, 1987; and *Sequence Analysis Primer*, Gribskov, M. and Devereux, J., eds., M Stockton Press, New York, 1991). In a preferred embodiment, the percent identity between two amino acid sequences is determined using the Needleman and Wunsch algorithm (*J. Mol. Biol.* (48):444-453 (1970)) which has been incorporated into the GAP program in the GCG software package, using either a Blossom 62 matrix or a PAM250 matrix, and a gap weight of 16, 14, 12, 10, 8, 6, or 4 and a length weight of 1, 2, 3, 4, 5, or 6.

In yet another preferred embodiment, the percent identity between two nucleotide sequences is determined using the GAP program in the GCG software package (Devereux, J., *et al.*, *Nucleic Acids Res.* 12(1):387 (1984)), using a NWSgapdna.CMP matrix and a gap weight of 40, 50, 60, 70, or 80 and a length weight of 1, 2, 3, 4, 5, or 6. In another embodiment, the percent identity between two amino acid or nucleotide sequences is determined using the algorithm of E. Myers and W. Miller (CABIOS, 4:11-17 (1989)) which has been incorporated into the ALIGN

program (version 2.0), using a PAM120 weight residue table, a gap length penalty of 12, and a gap penalty of 4.

The nucleotide and amino acid sequences of the present invention can further be used as a "query sequence" to perform a search against sequence databases to, for example, identify other family members or related sequences. Such searches can be performed using the NBLAST and XBLAST programs (version 2.0) of Altschul, *et al.* (*J. Mol. Biol.* 215:403-10 (1990)). BLAST nucleotide searches can be performed with the NBLAST program, score = 100, wordlength = 12 to obtain nucleotide sequences homologous to the nucleic acid molecules of the invention. BLAST protein searches can be performed with the XBLAST program, score = 50, wordlength = 3 to obtain amino acid sequences homologous to the proteins of the invention. To obtain gapped alignments for comparison purposes, Gapped BLAST can be utilized as described in Altschul *et al.* (*Nucleic Acids Res.* 25(17):3389-3402 (1997)). When utilizing BLAST and gapped BLAST programs, the default parameters of the respective programs (*e.g.*, XBLAST and NBLAST) can be used. In addition to BLAST, examples of other search and sequence comparison programs used in the art include, but are not limited to, FASTA (Pearson, *Methods Mol. Biol.* 25, 365-389 (1994)) and KERR (Dufresne *et al.*, *Nat Biotechnol* 2002 Dec;20(12):1269-71). For further information regarding bioinformatics techniques, see *Current Protocols in Bioinformatics*, John Wiley & Sons, Inc., N.Y.

The present invention further provides non-coding fragments of the nucleic acid molecules disclosed in Table 1 and/or Table 2. Preferred non-coding fragments include, but are not limited to, promoter sequences, enhancer sequences, intronic sequences, 5' untranslated regions (UTRs), 3' untranslated regions, gene modulating sequences and gene termination sequences. Such fragments are useful, for example, in controlling heterologous gene expression and in developing screens to identify gene-modulating agents.

SNP Detection Reagents

In a specific aspect of the present invention, the SNPs disclosed in Table 1 and/or Table 2, and their associated transcript sequences (referred to in Table 1 as SEQ ID NOS: 1-20), genomic sequences (referred to in Table 2 as SEQ ID NOS: 69-81), and context sequences (transcript-based context sequences are referred to in Table 1 as SEQ ID NOS: 41-68; genomic-based context sequences are provided in Table 2 as SEQ ID NOS: 82-199), can be used for the design of SNP detection reagents. The actual sequences referred to in the tables are provided in the Sequence Listing. As used herein, a "SNP detection reagent" is a reagent that specifically detects a specific target SNP position disclosed herein, and that is preferably specific for a particular nucleotide

(allele) of the target SNP position (*i.e.*, the detection reagent preferably can differentiate between different alternative nucleotides at a target SNP position, thereby allowing the identity of the nucleotide present at the target SNP position to be determined). Typically, such detection reagent hybridizes to a target SNP-containing nucleic acid molecule by complementary base-pairing in a sequence specific manner, and discriminates the target variant sequence from other nucleic acid sequences such as an art-known form in a test sample. An example of a detection reagent is a probe that hybridizes to a target nucleic acid containing one or more of the SNPs referred to in Table 1 and/or Table 2. In a preferred embodiment, such a probe can differentiate between nucleic acids having a particular nucleotide (allele) at a target SNP position from other nucleic acids that have a different nucleotide at the same target SNP position. In addition, a detection reagent may hybridize to a specific region 5' and/or 3' to a SNP position, particularly a region corresponding to the context sequences referred to in Table 1 and/or Table 2 (transcript-based context sequences are referred to in Table 1 as SEQ ID NOS: 41-68; genomic-based context sequences are referred to in Table 2 as SEQ ID NOS: 82-199). Another example of a detection reagent is a primer that acts as an initiation point of nucleotide extension along a complementary strand of a target polynucleotide. The SNP sequence information provided herein is also useful for designing primers, *e.g.* allele-specific primers, to amplify (*e.g.*, using PCR) any SNP of the present invention.

In one preferred embodiment of the invention, a SNP detection reagent is an isolated or synthetic DNA or RNA polynucleotide probe or primer or PNA oligomer, or a combination of DNA, RNA and/or PNA, that hybridizes to a segment of a target nucleic acid molecule containing a SNP identified in Table 1 and/or Table 2. A detection reagent in the form of a polynucleotide may optionally contain modified base analogs, intercalators or minor groove binders. Multiple detection reagents such as probes may be, for example, affixed to a solid support (*e.g.*, arrays or beads) or supplied in solution (*e.g.*, probe/primer sets for enzymatic reactions such as PCR, RT-PCR, TaqMan assays, or primer-extension reactions) to form a SNP detection kit.

A probe or primer typically is a substantially purified oligonucleotide or PNA oligomer. Such oligonucleotide typically comprises a region of complementary nucleotide sequence that hybridizes under stringent conditions to at least about 8, 10, 12, 16, 18, 20, 22, 25, 30, 40, 50, 55, 60, 65, 70, 80, 90, 100, 120 (or any other number in-between) or more consecutive nucleotides in a target nucleic acid molecule. Depending on the particular assay, the consecutive nucleotides can either include the target SNP position, or be a specific region in close enough proximity 5' and/or 3' to the SNP position to carry out the desired assay.

Other preferred primer and probe sequences can readily be determined using the transcript sequences (SEQ ID NOS: 1-20), genomic sequences (SEQ ID NOS: 69-81), and SNP context sequences (transcript-based context sequences are referred to in Table 1 as SEQ ID NOS: 41-68; genomic-based context sequences are referred to in Table 2 as SEQ ID NOS: 82-199) disclosed in the Sequence Listing and in Tables 1-2. The actual sequences referred to in the tables are provided in the Sequence Listing. It will be apparent to one of skill in the art that such primers and probes are directly useful as reagents for genotyping the SNPs of the present invention, and can be incorporated into any kit/system format.

In order to produce a probe or primer specific for a target SNP-containing sequence, the gene/transcript and/or context sequence surrounding the SNP of interest is typically examined using a computer algorithm that starts at the 5' or at the 3' end of the nucleotide sequence. Typical algorithms will then identify oligomers of defined length that are unique to the gene/SNP context sequence, have a GC content within a range suitable for hybridization, lack predicted secondary structure that may interfere with hybridization, and/or possess other desired characteristics or that lack other undesired characteristics.

A primer or probe of the present invention is typically at least about 8 nucleotides in length. In one embodiment of the invention, a primer or a probe is at least about 10 nucleotides in length. In a preferred embodiment, a primer or a probe is at least about 12 nucleotides in length. In a more preferred embodiment, a primer or probe is at least about 16, 17, 18, 19, 20, 21, 22, 23, 24 or 25 nucleotides in length. While the maximal length of a probe can be as long as the target sequence to be detected, depending on the type of assay in which it is employed, it is typically less than about 50, 60, 65, or 70 nucleotides in length. In the case of a primer, it is typically less than about 30 nucleotides in length. In a specific preferred embodiment of the invention, a primer or a probe is within the length of about 18 and about 28 nucleotides.

However, in other embodiments, such as nucleic acid arrays and other embodiments in which probes are affixed to a substrate, the probes can be longer, such as on the order of 30-70, 75, 80, 90, 100, or more nucleotides in length (see the section below entitled "SNP Detection Kits and Systems").

For analyzing SNPs, it may be appropriate to use oligonucleotides specific for alternative SNP alleles. Such oligonucleotides that detect single nucleotide variations in target sequences may be referred to by such terms as "allele-specific oligonucleotides", "allele-specific probes", or "allele-specific primers". The design and use of allele-specific probes for analyzing polymorphisms is described in, *e.g.*, *Mutation Detection A Practical Approach*, ed. Cotton *et al.* Oxford University

Press, 1998; Saiki *et al.*, *Nature* 324, 163-166 (1986); Dattagupta, EP235,726; and Saiki, WO 89/11548.

While the design of each allele-specific primer or probe depends on variables such as the precise composition of the nucleotide sequences flanking a SNP position in a target nucleic acid molecule, and the length of the primer or probe, another factor in the use of primers and probes is the stringency of the condition under which the hybridization between the probe or primer and the target sequence is performed. Higher stringency conditions utilize buffers with lower ionic strength and/or a higher reaction temperature, and tend to require a more perfect match between probe/primer and a target sequence in order to form a stable duplex. If the stringency is too high, however, hybridization may not occur at all. In contrast, lower stringency conditions utilize buffers with higher ionic strength and/or a lower reaction temperature, and permit the formation of stable duplexes with more mismatched bases between a probe/primer and a target sequence. By way of example and not limitation, exemplary conditions for high stringency hybridization conditions using an allele-specific probe are as follows: prehybridization with a solution containing 5X standard saline phosphate EDTA (SSPE), 0.5% NaDodSO₄ (SDS) at 55°C, and incubating probe with target nucleic acid molecules in the same solution at the same temperature, followed by washing with a solution containing 2X SSPE, and 0.1% SDS at 55°C or room temperature.

Moderate stringency hybridization conditions may be used for allele-specific primer extension reactions with a solution containing, *e.g.*, about 50mM KCl at about 46°C. Alternatively, the reaction may be carried out at an elevated temperature such as 60°C. In another embodiment, a moderately stringent hybridization condition suitable for oligonucleotide ligation assay (OLA) reactions wherein two probes are ligated if they are completely complementary to the target sequence may utilize a solution of about 100mM KCl at a temperature of 46°C.

In a hybridization-based assay, allele-specific probes can be designed that hybridize to a segment of target DNA from one individual but do not hybridize to the corresponding segment from another individual due to the presence of different polymorphic forms (*e.g.*, alternative SNP alleles/nucleotides) in the respective DNA segments from the two individuals. Hybridization conditions should be sufficiently stringent that there is a significant detectable difference in hybridization intensity between alleles, and preferably an essentially binary response, whereby a probe hybridizes to only one of the alleles or significantly more strongly to one allele. While a probe may be designed to hybridize to a target sequence that contains a SNP site such that the SNP site aligns anywhere along the sequence of the probe, the probe is preferably designed to

hybridize to a segment of the target sequence such that the SNP site aligns with a central position of the probe (*e.g.*, a position within the probe that is at least three nucleotides from either end of the probe). This design of probe generally achieves good discrimination in hybridization between different allelic forms.

5 In another embodiment, a probe or primer may be designed to hybridize to a segment of target DNA such that the SNP aligns with either the 5' most end or the 3' most end of the probe or primer. In a specific preferred embodiment that is particularly suitable for use in a oligonucleotide ligation assay (U.S. Patent No. 4,988,617), the 3' most nucleotide of the probe aligns with the SNP position in the target sequence.

10 Oligonucleotide probes and primers may be prepared by methods well known in the art. Chemical synthetic methods include, but are limited to, the phosphotriester method described by Narang *et al.*, 1979, *Methods in Enzymology* 68:90; the phosphodiester method described by Brown *et al.*, 1979, *Methods in Enzymology* 68:109, the diethylphosphoamidate method described by Beaucage *et al.*, 1981, *Tetrahedron Letters* 22:1859; and the solid support method
15 described in U.S. Patent No. 4,458,066.

Allele-specific probes are often used in pairs (or, less commonly, in sets of 3 or 4, such as if a SNP position is known to have 3 or 4 alleles, respectively, or to assay both strands of a nucleic acid molecule for a target SNP allele), and such pairs may be identical except for a one nucleotide mismatch that represents the allelic variants at the SNP position. Commonly, one
20 member of a pair perfectly matches a reference form of a target sequence that has a more common SNP allele (*i.e.*, the allele that is more frequent in the target population) and the other member of the pair perfectly matches a form of the target sequence that has a less common SNP allele (*i.e.*, the allele that is rarer in the target population). In the case of an array, multiple pairs of probes can be immobilized on the same support for simultaneous analysis of multiple different
25 polymorphisms.

In one type of PCR-based assay, an allele-specific primer hybridizes to a region on a target nucleic acid molecule that overlaps a SNP position and only primes amplification of an allelic form to which the primer exhibits perfect complementarity (Gibbs, 1989, *Nucleic Acid Res.* 17 2427-2448). Typically, the primer's 3'-most nucleotide is aligned with and
30 complementary to the SNP position of the target nucleic acid molecule. This primer is used in conjunction with a second primer that hybridizes at a distal site. Amplification proceeds from the two primers, producing a detectable product that indicates which allelic form is present in the test sample. A control is usually performed with a second pair of primers, one of which shows a single base mismatch at the polymorphic site and the other of which exhibits perfect

complementarity to a distal site. The single-base mismatch prevents amplification or substantially reduces amplification efficiency, so that either no detectable product is formed or it is formed in lower amounts or at a slower pace. The method generally works most effectively when the mismatch is at the 3'-most position of the oligonucleotide (*i.e.*, the 3'-most position of the oligonucleotide aligns with the target SNP position) because this position is most destabilizing to elongation from the primer (see, *e.g.*, WO 93/22456). This PCR-based assay can be utilized as part of the TaqMan™ assay, described below.

In a specific embodiment of the invention, a primer of the invention contains a sequence substantially complementary to a segment of a target SNP-containing nucleic acid molecule except that the primer has a mismatched nucleotide in one of the three nucleotide positions at the 3'-most end of the primer, such that the mismatched nucleotide does not base pair with a particular allele at the SNP site. In a preferred embodiment, the mismatched nucleotide in the primer is the second from the last nucleotide at the 3'-most position of the primer. In a more preferred embodiment, the mismatched nucleotide in the primer is the last nucleotide at the 3'-most position of the primer.

In another embodiment of the invention, a SNP detection reagent of the invention is labeled with a fluorogenic reporter dye that emits a detectable signal. While the preferred reporter dye is a fluorescent dye, any reporter dye that can be attached to a detection reagent such as an oligonucleotide probe or primer is suitable for use in the invention. Such dyes include, but are not limited to, Acridine, AMCA, BODIPY, Cascade Blue, Cy2, Cy3, Cy5, Cy7, Dabcyl, Edans, Eosin, Erythrosin, Fluorescein, 6-Fam, Tet, Joe, Hex, Oregon Green, Rhodamine, Rhodol Green, Tamra, Rox, and Texas Red.

In yet another embodiment of the invention, the detection reagent may be further labeled with a quencher dye such as Tamra, especially when the reagent is used as a self-quenching probe such as a TaqMan™ (U.S. Patent Nos. 5,210,015 and 5,538,848) or Molecular Beacon probe (U.S. Patent Nos. 5,118,801 and 5,312,728), or other stemless or linear beacon probe (Livak *et al.*, 1995, PCR Method Appl. 4:357-362; Tyagi *et al.*, 1996, Nature Biotechnology 14: 303-308; Nazarenko *et al.*, 1997, Nucl. Acids Res. 25:2516-2521; U.S. Patent Nos. 5,866,336 and 6,117,635).

The detection reagents of the invention may also contain other labels, including but not limited to, biotin for streptavidin binding, hapten for antibody binding, and oligonucleotide for binding to another complementary oligonucleotide such as pairs of zipcodes.

The present invention also contemplates reagents that do not contain (or that are complementary to) a SNP nucleotide identified herein but that are used to assay one or more SNPs disclosed herein. For example, primers that flank, but do not hybridize directly to a target SNP position provided herein are useful in primer extension reactions in which the primers

hybridize to a region adjacent to the target SNP position (*i.e.*, within one or more nucleotides from the target SNP site). During the primer extension reaction, a primer is typically not able to extend past a target SNP site if a particular nucleotide (allele) is present at that target SNP site, and the primer extension product can be detected in order to determine which SNP allele is present at the target SNP site. For example, particular ddNTPs are typically used in the primer extension reaction to terminate primer extension once a ddNTP is incorporated into the extension product (a primer extension product which includes a ddNTP at the 3'-most end of the primer extension product, and in which the ddNTP is a nucleotide of a SNP disclosed herein, is a composition that is specifically contemplated by the present invention). Thus, reagents that bind to a nucleic acid molecule in a region adjacent to a SNP site and that are used for assaying the SNP site, even though the bound sequences do not necessarily include the SNP site itself, are also contemplated by the present invention.

SNP Detection Kits and Systems

A person skilled in the art will recognize that, based on the SNP and associated sequence information disclosed herein, detection reagents can be developed and used to assay any SNP of the present invention individually or in combination, and such detection reagents can be readily incorporated into one of the established kit or system formats which are well known in the art. The terms "kits" and "systems", as used herein in the context of SNP detection reagents, are intended to refer to such things as combinations of multiple SNP detection reagents, or one or more SNP detection reagents in combination with one or more other types of elements or components (*e.g.*, other types of biochemical reagents, containers, packages such as packaging intended for commercial sale, substrates to which SNP detection reagents are attached, electronic hardware components, etc.). Accordingly, the present invention further provides SNP detection kits and systems, including but not limited to, packaged probe and primer sets (*e.g.*, TaqMan probe/primer sets), arrays/microarrays of nucleic acid molecules, and beads that contain one or more probes, primers, or other detection reagents for detecting one or more SNPs of the present invention. The kits/systems can optionally include various electronic hardware components; for example, arrays ("DNA chips") and microfluidic systems ("lab-on-a-chip" systems) provided by various manufacturers typically comprise hardware components. Other kits/systems (*e.g.*, probe/primer sets) may not include electronic hardware components, but may be comprised of, for example, one or more SNP detection reagents (along with, optionally, other biochemical reagents) packaged in one or more containers.

In some embodiments, a SNP detection kit typically contains one or more detection reagents and other components (*e.g.*, a buffer, enzymes such as DNA polymerases or ligases, chain extension nucleotides such as deoxynucleotide triphosphates, and in the case of Sanger-type DNA sequencing reactions, chain terminating nucleotides, positive control sequences, negative control sequences, and the like) necessary to carry out an assay or reaction, such as amplification and/or detection of a SNP-containing nucleic acid molecule. A kit may further contain means for determining the amount of a target nucleic acid, and means for comparing the amount with a standard, and can comprise instructions for using the kit to detect the SNP-containing nucleic acid molecule of interest. In one embodiment of the present invention, kits are provided which contain the necessary reagents to carry out one or more assays to detect one or more SNPs disclosed herein. In a preferred embodiment of the present invention, SNP detection kits/systems are in the form of nucleic acid arrays, or compartmentalized kits, including microfluidic/lab-on-a-chip systems.

SNP detection kits/systems may contain, for example, one or more probes, or pairs of probes, that hybridize to a nucleic acid molecule at or near each target SNP position. Multiple pairs of allele-specific probes may be included in the kit/system to simultaneously assay large numbers of SNPs, at least one of which is a SNP of the present invention. In some kits/systems, the allele-specific probes are immobilized to a substrate such as an array or bead. For example, the same substrate can comprise allele-specific probes for detecting at least 1; 10; 100; 1000; 10,000; 100,000 (or any other number in-between) or substantially all of the SNPs shown in Table 1 and/or Table 2.

The terms "arrays," "microarrays," and "DNA chips" are used herein interchangeably to refer to an array of distinct polynucleotides affixed to a substrate, such as glass, plastic, paper, nylon or other type of membrane, filter, chip, or any other suitable solid support. The polynucleotides can be synthesized directly on the substrate, or synthesized separate from the substrate and then affixed to the substrate. In one embodiment, the microarray is prepared and used according to the methods described in U.S. Patent No. 5,837,832, Chee *et al.*, PCT application W095/11995 (Chee *et al.*), Lockhart, D. J. *et al.* (1996; *Nat. Biotech.* 14: 1675-1680) and Schena, M. *et al.* (1996; *Proc. Natl. Acad. Sci.* 93: 10614-10619), all of which are incorporated herein in their entirety by reference. In other embodiments, such arrays are produced by the methods described by Brown *et al.*, U.S. Patent No. 5,807,522.

Nucleic acid arrays are reviewed in the following references: Zammattéo *et al.*, "New chips for molecular biology and diagnostics", *Biotechnol Annu Rev.* 2002;8:85-101; Sosnowski *et al.*, "Active microelectronic array system for DNA hybridization, genotyping and

pharmacogenomic applications”, *Psychiatr Genet.* 2002 Dec;12(4):181-92; Heller, “DNA microarray technology: devices, systems, and applications”, *Annu Rev Biomed Eng.* 2002;4:129-53. Epub 2002 Mar 22; Kolchinsky *et al.*, “Analysis of SNPs and other genomic variations using gel-based chips”, *Hum Mutat.* 2002 Apr;19(4):343-60; and McGall *et al.*, “High-density
5 genechip oligonucleotide probe arrays”, *Adv Biochem Eng Biotechnol.* 2002;77:21-42.

Any number of probes, such as allele-specific probes, may be implemented in an array, and each probe or pair of probes can hybridize to a different SNP position. In the case of polynucleotide probes, they can be synthesized at designated areas (or synthesized separately and then affixed to designated areas) on a substrate using a light-directed chemical process. Each DNA chip can
10 contain, for example, thousands to millions of individual synthetic polynucleotide probes arranged in a grid-like pattern and miniaturized (*e.g.*, to the size of a dime). Preferably, probes are attached to a solid support in an ordered, addressable array.

A microarray can be composed of a large number of unique, single-stranded polynucleotides, usually either synthetic antisense polynucleotides or fragments of cDNAs, fixed
15 to a solid support. Typical polynucleotides are preferably about 6-60 nucleotides in length, more preferably about 15-30 nucleotides in length, and most preferably about 18-25 nucleotides in length. For certain types of microarrays or other detection kits/systems, it may be preferable to use oligonucleotides that are only about 7-20 nucleotides in length. In other types of arrays, such as arrays used in conjunction with chemiluminescent detection technology, preferred probe
20 lengths can be, for example, about 15-80 nucleotides in length, preferably about 50-70 nucleotides in length, more preferably about 55-65 nucleotides in length, and most preferably about 60 nucleotides in length. The microarray or detection kit can contain polynucleotides that cover the known 5' or 3' sequence of a gene/transcript or target SNP site, sequential polynucleotides that cover the full-length sequence of a gene/transcript; or unique
25 polynucleotides selected from particular areas along the length of a target gene/transcript sequence, particularly areas corresponding to one or more SNPs disclosed in Table 1 and/or Table 2. Polynucleotides used in the microarray or detection kit can be specific to a SNP or SNPs of interest (*e.g.*, specific to a particular SNP allele at a target SNP site, or specific to particular SNP alleles at multiple different SNP sites), or specific to a polymorphic
30 gene/transcript or genes/transcripts of interest.

Hybridization assays based on polynucleotide arrays rely on the differences in hybridization stability of the probes to perfectly matched and mismatched target sequence variants. For SNP genotyping, it is generally preferable that stringency conditions used in hybridization assays are high enough such that nucleic acid molecules that differ from one another at

as little as a single SNP position can be differentiated (*e.g.*, typical SNP hybridization assays are designed so that hybridization will occur only if one particular nucleotide is present at a SNP position, but will not occur if an alternative nucleotide is present at that SNP position). Such high stringency conditions may be preferable when using, for example, nucleic acid arrays of allele-specific probes for SNP detection. Such high stringency conditions are described in the preceding section, and are well known to those skilled in the art and can be found in, for example, *Current Protocols in Molecular Biology*, John Wiley & Sons, N.Y. (1989), 6.3.1-6.3.6.

In other embodiments, the arrays are used in conjunction with chemiluminescent detection technology. The following patents and patent applications, provide additional information pertaining to chemiluminescent detection: U.S. patent publication nos. US 2005/049778 and US 2005/0026151 describe chemiluminescent approaches for microarray detection; U.S. Patent Nos. 6124478, 6107024, 5994073, 5981768, 5871938, 5843681, 5800999, and 5773628 describe methods and compositions of dioxetane for performing chemiluminescent detection; and U.S. published application US2002/0110828 discloses methods and compositions for microarray controls.

In one embodiment of the invention, a nucleic acid array can comprise an array of probes of about 15-25 nucleotides in length. In further embodiments, a nucleic acid array can comprise any number of probes, in which at least one probe is capable of detecting one or more SNPs disclosed in Table 1 and/or Table 2, and/or at least one probe comprises a fragment of one of the sequences selected from the group consisting of those disclosed in Table 1, Table 2, the Sequence Listing, and sequences complementary thereto, said fragment comprising at least about 8 consecutive nucleotides, preferably 10, 12, 15, 16, 18, 20, more preferably 22, 25, 30, 40, 47, 50, 55, 60, 65, 70, 80, 90, 100, or more consecutive nucleotides (or any other number in-between) and containing (or being complementary to) a novel SNP allele disclosed in Table 1 and/or Table 2. In some embodiments, the nucleotide complementary to the SNP site is within 5, 4, 3, 2, or 1 nucleotide from the center of the probe, more preferably at the center of said probe.

A polynucleotide probe can be synthesized on the surface of the substrate by using a chemical coupling procedure and an ink jet application apparatus, as described in PCT application W095/251116 (Baldeschweiler *et al.*). In another aspect, a "gridded" array analogous to a dot (or slot) blot may be used to arrange and link cDNA fragments or oligonucleotides to the surface of a substrate using a vacuum system, thermal, UV, mechanical or chemical bonding procedures. An array, such as those described above, may be produced by hand or by using available devices (slot blot or dot blot apparatus), materials (any suitable solid support), and machines (including robotic instruments), and may contain 8, 24, 96, 384, 1536, 6144 or more polynucleotides, or any other number which lends itself to the efficient use of commercially available instrumentation.

Using such arrays or other kits/systems, the present invention provides methods of identifying the SNPs disclosed herein in a test sample. Such methods typically involve incubating a test sample of nucleic acids with an array comprising one or more probes corresponding to at least one SNP position of the present invention, and assaying for binding of a nucleic acid from the test sample with one or more of the probes. Conditions for incubating a SNP detection reagent (or a kit/system that employs one or more such SNP detection reagents) with a test sample vary. Incubation conditions depend on such factors as the format employed in the assay, the detection methods employed,

and the type and nature of the detection reagents used in the assay. One skilled in the art will recognize that any one of the commonly available hybridization, amplification and array assay formats can readily be adapted to detect the SNPs disclosed herein.

A SNP detection kit/system of the present invention may include components that are used to prepare nucleic acids from a test sample for the subsequent amplification and/or detection of a SNP-containing nucleic acid molecule. Such sample preparation components can be used to produce nucleic acid extracts (including DNA and/or RNA), proteins or membrane extracts from any bodily fluids (such as blood, serum, plasma, urine, saliva, phlegm, gastric juices, semen, tears, sweat, etc.), skin, hair, cells (especially nucleated cells), biopsies, buccal swabs or tissue specimens. The test samples used in the above-described methods will vary based on such factors as the assay format, nature of the detection method, and the specific tissues, cells or extracts used as the test sample to be assayed. Methods of preparing nucleic acids, proteins, and cell extracts are well known in the art and can be readily adapted to obtain a sample that is compatible with the system utilized. Automated sample preparation systems for extracting nucleic acids from a test sample are commercially available, and examples are Qiagen's BioRobot 9600, Applied Biosystems' PRISM™ 6700 sample preparation system, and Roche Molecular Systems' COBAS™ AmpliPrep System.

Another form of kit contemplated by the present invention is a compartmentalized kit. A compartmentalized kit includes any kit in which reagents are contained in separate containers. Such containers include, for example, small glass containers, plastic containers, strips of plastic, glass or paper, or arraying material such as silica. Such containers allow one to efficiently transfer reagents from one compartment to another compartment such that the test samples and reagents are not cross-contaminated, or from one container to another vessel not included in the kit, and the agents or solutions of each container can be added in a quantitative fashion from one compartment to another or to another vessel. Such containers may include, for example, one or

more containers which will accept the test sample, one or more containers which contain at least one probe or other SNP detection reagent for detecting one or more SNPs of the present invention, one or more containers which contain wash reagents (such as phosphate buffered saline, Tris-buffers, etc.), and one or more containers which contain the reagents used to reveal
5 the presence of the bound probe or other SNP detection reagents. The kit can optionally further comprise compartments and/or reagents for, for example, nucleic acid amplification or other enzymatic reactions such as primer extension reactions, hybridization, ligation, electrophoresis (preferably capillary electrophoresis), mass spectrometry, and/or laser-induced fluorescent detection. The kit may also include instructions for using the kit. Exemplary compartmentalized kits include
10 microfluidic devices known in the art (see, *e.g.*, Weigl *et al.*, “Lab-on-a-chip for drug development”, *Adv Drug Deliv Rev.* 2003 Feb 24;55(3):349-77). In such microfluidic devices, the containers may be referred to as, for example, microfluidic “compartments”, “chambers”, or “channels”.

Microfluidic devices, which may also be referred to as “lab-on-a-chip” systems, biomedical micro-electro-mechanical systems (bioMEMs), or multicomponent integrated
15 systems, are exemplary kits/systems of the present invention for analyzing SNPs. Such systems miniaturize and compartmentalize processes such as probe/target hybridization, nucleic acid amplification, and capillary electrophoresis reactions in a single functional device. Such microfluidic devices typically utilize detection reagents in at least one aspect of the system, and such detection reagents may be used to detect one or more SNPs of the present invention. One
20 example of a microfluidic system is disclosed in U.S. Patent No. 5,589,136, which describes the integration of PCR amplification and capillary electrophoresis in chips. Exemplary microfluidic systems comprise a pattern of microchannels designed onto a glass, silicon, quartz, or plastic wafer included on a microchip. The movements of the samples may be controlled by electric, electroosmotic or hydrostatic forces applied across different areas of the microchip to create
25 functional microscopic valves and pumps with no moving parts. Varying the voltage can be used as a means to control the liquid flow at intersections between the micro-machined channels and to change the liquid flow rate for pumping across different sections of the microchip. See, for example, U.S. Patent Nos. 6,153,073, Dubrow *et al.*, and 6,156,181, Parce *et al.*

For genotyping SNPs, an exemplary microfluidic system may integrate, for example,
30 nucleic acid amplification, primer extension, capillary electrophoresis, and a detection method such as laser induced fluorescence detection. In a first step of an exemplary process for using such an exemplary system, nucleic acid samples are amplified, preferably by PCR. Then, the amplification products are subjected to automated primer extension reactions using ddNTPs (specific fluorescence for each ddNTP) and the appropriate oligonucleotide primers to carry out

primer extension reactions which hybridize just upstream of the targeted SNP. Once the extension at the 3' end is completed, the primers are separated from the unincorporated fluorescent ddNTPs by capillary electrophoresis. The separation medium used in capillary electrophoresis can be, for example, polyacrylamide, polyethyleneglycol or dextran. The incorporated ddNTPs in the single nucleotide primer extension products are identified by laser-induced fluorescence detection. Such an exemplary microchip can be used to process, for example, at least 96 to 384 samples, or more, in parallel.

USES OF NUCLEIC ACID MOLECULES

10 The nucleic acid molecules of the present invention have a variety of uses, especially in the diagnosis and treatment of VT. For example, the nucleic acid molecules are useful as hybridization probes, such as for genotyping SNPs in messenger RNA, transcript, cDNA, genomic DNA, amplified DNA or other nucleic acid molecules, and for isolating full-length cDNA and genomic clones encoding the variant peptides disclosed in Table 1 as well as their orthologs.

15 A probe can hybridize to any nucleotide sequence along the entire length of a nucleic acid molecule referred to in Table 1 and/or Table 2. Preferably, a probe of the present invention hybridizes to a region of a target sequence that encompasses a SNP position indicated in Table 1 and/or Table 2. More preferably, a probe hybridizes to a SNP-containing target sequence in a sequence-specific manner such that it distinguishes the target sequence from other nucleotide sequences which vary from the target sequence only by which nucleotide is present at the SNP site. Such a probe is particularly useful for detecting the presence of a SNP-containing nucleic acid in a test sample, or for determining which nucleotide (allele) is present at a particular SNP site (*i.e.*, genotyping the SNP site).

20 A nucleic acid hybridization probe may be used for determining the presence, level, form, and/or distribution of nucleic acid expression. The nucleic acid whose level is determined can be DNA or RNA. Accordingly, probes specific for the SNPs described herein can be used to assess the presence, expression and/or gene copy number in a given cell, tissue, or organism. These uses are relevant for diagnosis of disorders involving an increase or decrease in gene expression relative to normal levels. *In vitro* techniques for detection of mRNA include, for example, Northern blot hybridizations and *in situ* hybridizations. *In vitro* techniques for detecting DNA include Southern blot hybridizations and *in situ* hybridizations (Sambrook and Russell, 2000, Molecular Cloning: A Laboratory Manual, Cold Spring Harbor Press, Cold Spring Harbor, NY).

30 Probes can be used as part of a diagnostic test kit for identifying cells or tissues in which a variant protein is expressed, such as by measuring the level of a variant protein-encoding nucleic

acid (*e.g.*, mRNA) in a sample of cells from a subject or determining if a polynucleotide contains a SNP of interest.

Thus, the nucleic acid molecules of the invention can be used as hybridization probes to detect the SNPs disclosed herein, thereby determining whether an individual with the polymorphisms is at risk for VT or has developed early stage VT. Detection of a SNP associated with a disease phenotype provides a diagnostic tool for an active disease and/or genetic predisposition to the disease.

Furthermore, the nucleic acid molecules of the invention are therefore useful for detecting a gene (gene information is disclosed in Table 2, for example) which contains a SNP disclosed herein and/or products of such genes, such as expressed mRNA transcript molecules (transcript information is disclosed in Table 1, for example), and are thus useful for detecting gene expression. The nucleic acid molecules can optionally be implemented in, for example, an array or kit format for use in detecting gene expression.

The nucleic acid molecules of the invention are also useful as primers to amplify any given region of a nucleic acid molecule, particularly a region containing a SNP identified in Table 1 and/or Table 2.

The nucleic acid molecules of the invention are also useful for constructing recombinant vectors (described in greater detail below). Such vectors include expression vectors that express a portion of, or all of, any of the variant peptide sequences referred to in Table 1. Vectors also include insertion vectors, used to integrate into another nucleic acid molecule sequence, such as into the cellular genome, to alter *in situ* expression of a gene and/or gene product. For example, an endogenous coding sequence can be replaced via homologous recombination with all or part of the coding region containing one or more specifically introduced SNPs.

The nucleic acid molecules of the invention are also useful for expressing antigenic portions of the variant proteins, particularly antigenic portions that contain a variant amino acid sequence (*e.g.*, an amino acid substitution) caused by a SNP disclosed in Table 1 and/or Table 2.

The nucleic acid molecules of the invention are also useful for constructing vectors containing a gene regulatory region of the nucleic acid molecules of the present invention.

The nucleic acid molecules of the invention are also useful for designing ribozymes corresponding to all, or a part, of an mRNA molecule expressed from a SNP-containing nucleic acid molecule described herein.

The nucleic acid molecules of the invention are also useful for constructing host cells expressing a part, or all, of the nucleic acid molecules and variant peptides.

The nucleic acid molecules of the invention are also useful for constructing transgenic animals expressing all, or a part, of the nucleic acid molecules and variant peptides. The production of recombinant cells and transgenic animals having nucleic acid molecules which contain the SNPs disclosed in Table 1 and/or Table 2 allow, for example, effective clinical design of treatment compounds and dosage regimens.

The nucleic acid molecules of the invention are also useful in assays for drug screening to identify compounds that, for example, modulate nucleic acid expression.

The nucleic acid molecules of the invention are also useful in gene therapy in patients whose cells have aberrant gene expression. Thus, recombinant cells, which include a patient's cells that have been engineered *ex vivo* and returned to the patient, can be introduced into an individual where the recombinant cells produce the desired protein to treat the individual.

SNP Genotyping Methods

The process of determining which specific nucleotide (*i.e.*, allele) is present at each of one or more SNP positions, such as a SNP position in a nucleic acid molecule disclosed in Table 1 and/or Table 2, is referred to as SNP genotyping. The present invention provides methods of SNP genotyping, such as for use in screening for VT or related pathologies, or determining predisposition thereto, or determining responsiveness to a form of treatment, or in genome mapping or SNP association analysis, etc.

Nucleic acid samples can be genotyped to determine which allele(s) is/are present at any given genetic region (*e.g.*, SNP position) of interest by methods well known in the art. The neighboring sequence can be used to design SNP detection reagents such as oligonucleotide probes, which may optionally be implemented in a kit format. Exemplary SNP genotyping methods are described in Chen *et al.*, "Single nucleotide polymorphism genotyping: biochemistry, protocol, cost and throughput", *Pharmacogenomics J.* 2003;3(2):77-96; Kwok *et al.*, "Detection of single nucleotide polymorphisms", *Curr Issues Mol Biol.* 2003 Apr;5(2):43-60; Shi, "Technologies for individual genotyping: detection of genetic polymorphisms in drug targets and disease genes", *Am J Pharmacogenomics.* 2002;2(3):197-205; and Kwok, "Methods for genotyping single nucleotide polymorphisms", *Annu Rev Genomics Hum Genet* 2001;2:235-58. Exemplary techniques for high-throughput SNP genotyping are described in Marnellos, "High-throughput SNP analysis for genetic association studies", *Curr Opin Drug Discov Devel.* 2003 May;6(3):317-21. Common SNP genotyping methods include, but are not limited to, TaqMan assays, molecular beacon assays, nucleic acid arrays, allele-specific primer extension, allele-specific PCR, arrayed primer extension, homogeneous primer extension assays, primer extension with detection by mass

spectrometry, pyrosequencing, multiplex primer extension sorted on genetic arrays, ligation with rolling circle amplification, homogeneous ligation, OLA (U.S. Patent No. 4,988,167), multiplex ligation reaction sorted on genetic arrays, restriction-fragment length polymorphism, single base extension-tag assays, and the Invader assay. Such methods may be used in combination with
5 detection mechanisms such as, for example, luminescence or chemiluminescence detection, fluorescence detection, time-resolved fluorescence detection, fluorescence resonance energy transfer, fluorescence polarization, mass spectrometry, and electrical detection.

Various methods for detecting polymorphisms include, but are not limited to, methods in which protection from cleavage agents is used to detect mismatched bases in RNA/RNA or
10 RNA/DNA duplexes (Myers *et al.*, *Science* 230:1242 (1985); Cotton *et al.*, *PNAS* 85:4397 (1988); and Saleeba *et al.*, *Meth. Enzymol.* 217:286-295 (1992)), comparison of the electrophoretic mobility of variant and wild type nucleic acid molecules (Orita *et al.*, *PNAS* 86:2766 (1989); Cotton *et al.*, *Mutat. Res.* 285:125-144 (1993); and Hayashi *et al.*, *Genet. Anal. Tech. Appl.* 9:73-79 (1992)), and assaying the movement of polymorphic or wild-type fragments in polyacrylamide gels containing a
15 gradient of denaturant using denaturing gradient gel electrophoresis (DGGE) (Myers *et al.*, *Nature* 313:495 (1985)). Sequence variations at specific locations can also be assessed by nuclease protection assays such as RNase and S1 protection or chemical cleavage methods.

In a preferred embodiment, SNP genotyping is performed using the TaqMan assay, which is also known as the 5' nuclease assay (U.S. Patent Nos. 5,210,015 and 5,538,848). The TaqMan
20 assay detects the accumulation of a specific amplified product during PCR. The TaqMan assay utilizes an oligonucleotide probe labeled with a fluorescent reporter dye and a quencher dye. The reporter dye is excited by irradiation at an appropriate wavelength, it transfers energy to the quencher dye in the same probe via a process called fluorescence resonance energy transfer (FRET). When attached to the probe, the excited reporter dye does not emit a signal. The
25 proximity of the quencher dye to the reporter dye in the intact probe maintains a reduced fluorescence for the reporter. The reporter dye and quencher dye may be at the 5' most and the 3' most ends, respectively, or vice versa. Alternatively, the reporter dye may be at the 5' or 3' most end while the quencher dye is attached to an internal nucleotide, or vice versa. In yet another embodiment, both the reporter and the quencher may be attached to internal nucleotides
30 at a distance from each other such that fluorescence of the reporter is reduced.

During PCR, the 5' nuclease activity of DNA polymerase cleaves the probe, thereby separating the reporter dye and the quencher dye and resulting in increased fluorescence of the reporter. Accumulation of PCR product is detected directly by monitoring the increase in fluorescence of the reporter dye. The DNA polymerase cleaves the probe between the reporter

dye and the quencher dye only if the probe hybridizes to the target SNP-containing template which is amplified during PCR, and the probe is designed to hybridize to the target SNP site only if a particular SNP allele is present.

Preferred TaqMan primer and probe sequences can readily be determined using the SNP and associated nucleic acid sequence information provided herein. A number of computer programs, such as Primer Express (Applied Biosystems, Foster City, CA), can be used to rapidly obtain optimal primer/probe sets. It will be apparent to one of skill in the art that such primers and probes for detecting the SNPs of the present invention are useful in diagnostic assays for VT and related pathologies, and can be readily incorporated into a kit format. The present invention also includes modifications of the Taqman assay well known in the art such as the use of Molecular Beacon probes (U.S. Patent Nos. 5,118,801 and 5,312,728) and other variant formats (U.S. Patent Nos. 5,866,336 and 6,117,635).

Another preferred method for genotyping the SNPs of the present invention is the use of two oligonucleotide probes in an OLA (see, *e.g.*, U.S. Patent No. 4,988,617). In this method, one probe hybridizes to a segment of a target nucleic acid with its 3' most end aligned with the SNP site. A second probe hybridizes to an adjacent segment of the target nucleic acid molecule directly 3' to the first probe. The two juxtaposed probes hybridize to the target nucleic acid molecule, and are ligated in the presence of a linking agent such as a ligase if there is perfect complementarity between the 3' most nucleotide of the first probe with the SNP site. If there is a mismatch, ligation would not occur. After the reaction, the ligated probes are separated from the target nucleic acid molecule, and detected as indicators of the presence of a SNP.

The following patents, patent applications, and published international patent applications provide additional information pertaining to techniques for carrying out various types of OLA: U.S. Patent Nos. 6027889, 6268148, 5494810, 5830711, and 6054564 describe OLA strategies for performing SNP detection; WO 97/31256 and WO 00/56927 describe OLA strategies for performing SNP detection using universal arrays, wherein a zipcode sequence can be introduced into one of the hybridization probes, and the resulting product, or amplified product, hybridized to a universal zip code array; wherein zipcodes are incorporated into OLA probes, and amplified PCR products are determined by electrophoretic or universal zipcode array readout; U.S. applications 60/427818, 60/445636, and 60/445494 describe SNPlax methods and software for multiplexed SNP detection using OLA followed by PCR, wherein zipcodes are incorporated into OLA probes, and amplified PCR products are hybridized with a zipchute reagent, and the identity of

the SNP determined from electrophoretic readout of the zipchute. In some embodiments, OLA is carried out prior to PCR (or another method of nucleic acid amplification). In other embodiments, PCR (or another method of nucleic acid amplification) is carried out prior to OLA.

Another method for SNP genotyping is based on mass spectrometry. Mass spectrometry takes advantage of the unique mass of each of the four nucleotides of DNA. SNPs can be unambiguously genotyped by mass spectrometry by measuring the differences in the mass of nucleic acids having alternative SNP alleles. MALDI-TOF (Matrix Assisted Laser Desorption Ionization – Time of Flight) mass spectrometry technology is preferred for extremely precise determinations of molecular mass, such as SNPs. Numerous approaches to SNP analysis have been developed based on mass spectrometry. Preferred mass spectrometry-based methods of SNP genotyping include primer extension assays, which can also be utilized in combination with other approaches, such as traditional gel-based formats and microarrays.

Typically, the primer extension assay involves designing and annealing a primer to a template PCR amplicon upstream (5') from a target SNP position. A mix of dideoxynucleotide triphosphates (ddNTPs) and/or deoxynucleotide triphosphates (dNTPs) are added to a reaction mixture containing template (*e.g.*, a SNP-containing nucleic acid molecule which has typically been amplified, such as by PCR), primer, and DNA polymerase. Extension of the primer terminates at the first position in the template where a nucleotide complementary to one of the ddNTPs in the mix occurs. The primer can be either immediately adjacent (*i.e.*, the nucleotide at the 3' end of the primer hybridizes to the nucleotide next to the target SNP site) or two or more nucleotides removed from the SNP position. If the primer is several nucleotides removed from the target SNP position, the only limitation is that the template sequence between the 3' end of the primer and the SNP position cannot contain a nucleotide of the same type as the one to be detected, or this will cause premature termination of the extension primer. Alternatively, if all four ddNTPs alone, with no dNTPs, are added to the reaction mixture, the primer will always be extended by only one nucleotide, corresponding to the target SNP position. In this instance, primers are designed to bind one nucleotide upstream from the SNP position (*i.e.*, the nucleotide at the 3' end of the primer hybridizes to the nucleotide that is immediately adjacent to the target SNP site on the 5' side of the target SNP site). Extension by only one nucleotide is preferable, as it minimizes the overall mass of the extended primer, thereby increasing the resolution of mass differences between alternative SNP nucleotides. Furthermore, mass-tagged ddNTPs can be employed in the primer extension reactions in place of unmodified ddNTPs. This increases the mass difference between primers extended with these ddNTPs, thereby providing increased sensitivity and accuracy, and is particularly useful for typing heterozygous base positions. Mass-

tagging also alleviates the need for intensive sample-preparation procedures and decreases the necessary resolving power of the mass spectrometer.

The extended primers can then be purified and analyzed by MALDI-TOF mass spectrometry to determine the identity of the nucleotide present at the target SNP position. In one method of analysis, the products from the primer extension reaction are combined with light absorbing crystals that form a matrix. The matrix is then hit with an energy source such as a laser to ionize and desorb the nucleic acid molecules into the gas-phase. The ionized molecules are then ejected into a flight tube and accelerated down the tube towards a detector. The time between the ionization event, such as a laser pulse, and collision of the molecule with the detector is the time of flight of that molecule. The time of flight is precisely correlated with the mass-to-charge ratio (m/z) of the ionized molecule. Ions with smaller m/z travel down the tube faster than ions with larger m/z and therefore the lighter ions reach the detector before the heavier ions. The time-of-flight is then converted into a corresponding, and highly precise, m/z . In this manner, SNPs can be identified based on the slight differences in mass, and the corresponding time of flight differences, inherent in nucleic acid molecules having different nucleotides at a single base position. For further information regarding the use of primer extension assays in conjunction with MALDI-TOF mass spectrometry for SNP genotyping, see, *e.g.*, Wise *et al.*, "A standard protocol for single nucleotide primer extension in the human genome using matrix-assisted laser desorption/ionization time-of-flight mass spectrometry", *Rapid Commun Mass Spectrom.* 2003;17(11):1195-202.

The following references provide further information describing mass spectrometry-based methods for SNP genotyping: Bocker, "SNP and mutation discovery using base-specific cleavage and MALDI-TOF mass spectrometry", *Bioinformatics.* 2003 Jul;19 Suppl 1:I44-I53; Storm *et al.*, "MALDI-TOF mass spectrometry-based SNP genotyping", *Methods Mol Biol.* 2003;212:241-62; Jurinke *et al.*, "The use of Mass ARRAY technology for high throughput genotyping", *Adv Biochem Eng Biotechnol.* 2002;77:57-74; and Jurinke *et al.*, "Automated genotyping using the DNA MassArray technology", *Methods Mol Biol.* 2002;187:179-92.

SNPs can also be scored by direct DNA sequencing. A variety of automated sequencing procedures can be utilized ((1995) *Biotechniques* 19:448), including sequencing by mass spectrometry (see, *e.g.*, PCT International Publication No. WO94/16101; Cohen *et al.*, *Adv. Chromatogr.* 36:127-162 (1996); and Griffin *et al.*, *Appl. Biochem. Biotechnol.* 38:147-159 (1993)). The nucleic acid sequences of the present invention enable one of ordinary skill in the art to readily design sequencing primers for such automated sequencing procedures. Commercial

instrumentation, such as the Applied Biosystems 377, 3100, 3700, 3730, and 3730xl DNA Analyzers (Foster City, CA), is commonly used in the art for automated sequencing.

Other methods that can be used to genotype the SNPs of the present invention include single-strand conformational polymorphism (SSCP), and denaturing gradient gel electrophoresis (DGGE) (Myers *et al.*, *Nature* 313:495 (1985)). SSCP identifies base differences by alteration in electrophoretic migration of single stranded PCR products, as described in Orita *et al.*, *Proc. Nat. Acad.* Single-stranded PCR products can be generated by heating or otherwise denaturing double stranded PCR products. Single-stranded nucleic acids may refold or form secondary structures that are partially dependent on the base sequence. The different electrophoretic mobilities of single-stranded amplification products are related to base-sequence differences at SNP positions. DGGE differentiates SNP alleles based on the different sequence-dependent stabilities and melting properties inherent in polymorphic DNA and the corresponding differences in electrophoretic migration patterns in a denaturing gradient gel (Erlich, ed., *PCR Technology, Principles and Applications for DNA Amplification*, W.H. Freeman and Co, New York, 1992, Chapter 7).

Sequence-specific ribozymes (U.S. Patent No. 5,498,531) can also be used to score SNPs based on the development or loss of a ribozyme cleavage site. Perfectly matched sequences can be distinguished from mismatched sequences by nuclease cleavage digestion assays or by differences in melting temperature. If the SNP affects a restriction enzyme cleavage site, the SNP can be identified by alterations in restriction enzyme digestion patterns, and the corresponding changes in nucleic acid fragment lengths determined by gel electrophoresis

SNP genotyping can include the steps of, for example, collecting a biological sample from a human subject (*e.g.*, sample of tissues, cells, fluids, secretions, etc.), isolating nucleic acids (*e.g.*, genomic DNA, mRNA or both) from the cells of the sample, contacting the nucleic acids with one or more primers which specifically hybridize to a region of the isolated nucleic acid containing a target SNP under conditions such that hybridization and amplification of the target nucleic acid region occurs, and determining the nucleotide present at the SNP position of interest, or, in some assays, detecting the presence or absence of an amplification product (assays can be designed so that hybridization and/or amplification will only occur if a particular SNP allele is present or absent). In some assays, the size of the amplification product is detected and compared to the length of a control sample; for example, deletions and insertions can be detected by a change in size of the amplified product compared to a normal genotype.

SNP genotyping is useful for numerous practical applications, as described below. Examples of such applications include, but are not limited to, SNP-disease association analysis,

disease predisposition screening, disease diagnosis, disease prognosis, disease progression monitoring, determining therapeutic strategies based on an individual's genotype ("pharmacogenomics"), developing therapeutic agents based on SNP genotypes associated with a disease or likelihood of responding to a drug, stratifying a patient population for clinical trial for a treatment regimen, predicting the likelihood that an individual will experience toxic side effects from a therapeutic agent, and human identification applications such as forensics.

Analysis of Genetic Association Between SNPs and Phenotypic Traits

SNP genotyping for disease diagnosis, disease predisposition screening, disease prognosis, determining drug responsiveness (pharmacogenomics), drug toxicity screening, and other uses described herein, typically relies on initially establishing a genetic association between one or more specific SNPs and the particular phenotypic traits of interest.

Different study designs may be used for genetic association studies (*Modern Epidemiology*, Lippincott Williams & Wilkins (1998), 609-622). Observational studies are most frequently carried out in which the response of the patients is not interfered with. The first type of observational study identifies a sample of persons in whom the suspected cause of the disease is present and another sample of persons in whom the suspected cause is absent, and then the frequency of development of disease in the two samples is compared. These sampled populations are called cohorts, and the study is a prospective study. The other type of observational study is case-control or a retrospective study. In typical case-control studies, samples are collected from individuals with the phenotype of interest (cases) such as certain manifestations of a disease, and from individuals without the phenotype (controls) in a population (target population) that conclusions are to be drawn from. Then the possible causes of the disease are investigated retrospectively. As the time and costs of collecting samples in case-control studies are considerably less than those for prospective studies, case-control studies are the more commonly used study design in genetic association studies, at least during the exploration and discovery stage.

In both types of observational studies, there may be potential confounding factors that should be taken into consideration. Confounding factors are those that are associated with both the real cause(s) of the disease and the disease itself, and they include demographic information such as age, gender, ethnicity as well as environmental factors. When confounding factors are not matched in cases and controls in a study, and are not controlled properly, spurious association results can arise. If potential confounding factors are identified, they should be controlled for by analysis methods explained below.

In a genetic association study, the cause of interest to be tested is a certain allele or a SNP or a combination of alleles or a haplotype from several SNPs. Thus, tissue specimens (*e.g.*, whole blood) from the sampled individuals may be collected and genomic DNA genotyped for the SNP(s) of interest. In addition to the phenotypic trait of interest, other information such as demographic (*e.g.*, age, gender, ethnicity, etc.), clinical, and environmental information that may influence the outcome of the trait can be collected to further characterize and define the sample set. In many cases, these factors are known to be associated with diseases and/or SNP allele frequencies. There are likely gene-environment and/or gene-gene interactions as well. Analysis methods to address gene-environment and gene-gene interactions (for example, the effects of the presence of both susceptibility alleles at two different genes can be greater than the effects of the individual alleles at two genes combined) are discussed below.

After all the relevant phenotypic and genotypic information has been obtained, statistical analyses are carried out to determine if there is any significant correlation between the presence of an allele or a genotype with the phenotypic characteristics of an individual. Preferably, data inspection and cleaning are first performed before carrying out statistical tests for genetic association. Epidemiological and clinical data of the samples can be summarized by descriptive statistics with tables and graphs. Data validation is preferably performed to check for data completion, inconsistent entries, and outliers. Chi-squared tests and t-tests (Wilcoxon rank-sum tests if distributions are not normal) may then be used to check for significant differences between cases and controls for discrete and continuous variables, respectively. To ensure genotyping quality, Hardy-Weinberg disequilibrium tests can be performed on cases and controls separately. Significant deviation from Hardy-Weinberg equilibrium (HWE) in both cases and controls for individual markers can be indicative of genotyping errors. If HWE is violated in a majority of markers, it is indicative of population substructure that should be further investigated. Moreover, Hardy-Weinberg disequilibrium in cases only can indicate genetic association of the markers with the disease (*Genetic Data Analysis*, Weir B., Sinauer (1990)).

To test whether an allele of a single SNP is associated with the case or control status of a phenotypic trait, one skilled in the art can compare allele frequencies in cases and controls. Standard chi-squared tests and Fisher exact tests can be carried out on a 2x2 table (2 SNP alleles x 2 outcomes in the categorical trait of interest). To test whether genotypes of a SNP are associated, chi-squared tests can be carried out on a 3x2 table (3 genotypes x 2 outcomes). Score tests are also carried out for genotypic association to contrast the three genotypic frequencies (major homozygotes, heterozygotes and minor homozygotes) in cases and controls, and to look for trends using 3 different modes of inheritance, namely dominant (with contrast coefficients 2,

-
-1, -1), additive (with contrast coefficients 1, 0, -1) and recessive (with contrast coefficients 1, 1, -2). Odds ratios for minor versus major alleles, and odds ratios for heterozygote and homozygote variants versus the wild type genotypes are calculated with the desired confidence limits, usually 95%.

5 In order to control for confounders and to test for interaction and effect modifiers, stratified analyses may be performed using stratified factors that are likely to be confounding, including demographic information such as age, ethnicity, and gender, or an interacting element or effect modifier, such as a known major gene (*e.g.*, APOE for Alzheimer's disease or HLA genes for autoimmune diseases), or environmental factors such as smoking in lung cancer.

10 Stratified association tests may be carried out using Cochran-Mantel-Haenszel tests that take into account the ordinal nature of genotypes with 0, 1, and 2 variant alleles. Exact tests by StatXact may also be performed when computationally possible. Another way to adjust for confounding effects and test for interactions is to perform stepwise multiple logistic regression analysis using statistical packages such as SAS or R. Logistic regression is a model-building technique in

15 which the best fitting and most parsimonious model is built to describe the relation between the dichotomous outcome (for instance, getting a certain disease or not) and a set of independent variables (for instance, genotypes of different associated genes, and the associated demographic and environmental factors). The most common model is one in which the logit transformation of the odds ratios is expressed as a linear combination of the variables (main effects) and their

20 cross-product terms (interactions) (*Applied Logistic Regression*, Hosmer and Lemeshow, Wiley (2000)). To test whether a certain variable or interaction is significantly associated with the outcome, coefficients in the model are first estimated and then tested for statistical significance of their departure from zero.

 In addition to performing association tests one marker at a time, haplotype association

25 analysis may also be performed to study a number of markers that are closely linked together. Haplotype association tests can have better power than genotypic or allelic association tests when the tested markers are not the disease-causing mutations themselves but are in linkage disequilibrium with such mutations. The test will even be more powerful if the disease is indeed caused by a combination of alleles on a haplotype (*e.g.*, APOE is a haplotype formed by 2 SNPs

30 that are very close to each other). In order to perform haplotype association effectively, marker-marker linkage disequilibrium measures, both D' and r^2 , are typically calculated for the markers within a gene to elucidate the haplotype structure. Recent studies (Daly *et al*, *Nature Genetics*, 29, 232-235, 2001) in linkage disequilibrium indicate that SNPs within a gene are organized in block pattern, and a high degree of linkage disequilibrium exists within blocks and very little

linkage disequilibrium exists between blocks. Haplotype association with the disease status can be performed using such blocks once they have been elucidated.

Haplotype association tests can be carried out in a similar fashion as the allelic and genotypic association tests. Each haplotype in a gene is analogous to an allele in a multi-allelic marker. One skilled in the art can either compare the haplotype frequencies in cases and controls or test genetic association with different pairs of haplotypes. It has been proposed (Schaid *et al*, *Am. J. Hum. Genet.*, 70, 425-434, 2002) that score tests can be done on haplotypes using the program "haplo.score." In that method, haplotypes are first inferred by EM algorithm and score tests are carried out with a generalized linear model (GLM) framework that allows the adjustment of other factors.

An important decision in the performance of genetic association tests is the determination of the significance level at which significant association can be declared when the P value of the tests reaches that level. In an exploratory analysis where positive hits will be followed up in subsequent confirmatory testing, an unadjusted P value < 0.2 (a significance level on the lenient side), for example, may be used for generating hypotheses for significant association of a SNP with certain phenotypic characteristics of a disease. It is preferred that a p-value < 0.05 (a significance level traditionally used in the art) is achieved in order for a SNP to be considered to have an association with a disease. It is more preferred that a p-value < 0.01 (a significance level on the stringent side) is achieved for an association to be declared. When hits are followed up in confirmatory analyses in more samples of the same source or in different samples from different sources, adjustment for multiple testing will be performed as to avoid excess number of hits while maintaining the experiment-wide error rates at 0.05. While there are different methods to adjust for multiple testing to control for different kinds of error rates, a commonly used but rather conservative method is Bonferroni correction to control the experiment-wise or family-wise error rate (*Multiple comparisons and multiple tests*, Westfall *et al*, SAS Institute (1999)). Permutation tests to control for the false discovery rates, FDR, can be more powerful (Benjamini and Hochberg, *Journal of the Royal Statistical Society, Series B* 57, 1289-1300, 1995, *Resampling-based Multiple Testing*, Westfall and Young, Wiley (1993)). Such methods to control for multiplicity would be preferred when the tests are dependent and controlling for false discovery rates is sufficient as opposed to controlling for the experiment-wise error rates.

In replication studies using samples from different populations after statistically significant markers have been identified in the exploratory stage, meta-analyses can then be performed by combining evidence of different studies (*Modern Epidemiology*, Lippincott

Williams & Wilkins, 1998, 643-673). If available, association results known in the art for the same SNPs can be included in the meta-analyses.

5 Since both genotyping and disease status classification can involve errors, sensitivity analyses may be performed to see how odds ratios and p-values would change upon various estimates on genotyping and disease classification error rates.

10 It has been well known that subpopulation-based sampling bias between cases and controls can lead to spurious results in case-control association studies (Ewens and Spielman, *Am. J. Hum. Genet.* 62, 450-458, 1995) when prevalence of the disease is associated with different subpopulation groups. Such bias can also lead to a loss of statistical power in genetic association studies. To detect population stratification, Pritchard and Rosenberg (Pritchard *et al.* *Am. J. Hum. Gen.* 1999, 65:220-228) suggested typing markers that are unlinked to the disease and using results of association tests on those markers to determine whether there is any population stratification. When stratification is detected, the genomic control (GC) method as proposed by Devlin and Roeder (Devlin *et al. Biometrics* 1999, 55:997-1004) can be used to
15 adjust for the inflation of test statistics due to population stratification. GC method is robust to changes in population structure levels as well as being applicable to DNA pooling designs (Devlin *et al. Genet. Epidem.* 20001, 21:273-284).

20 While Pritchard's method recommended using 15-20 unlinked microsatellite markers, it suggested using more than 30 biallelic markers to get enough power to detect population stratification. For the GC method, it has been shown (Bacanu *et al. Am. J. Hum. Genet.* 2000, 66:1933-1944) that about 60-70 biallelic markers are sufficient to estimate the inflation factor for the test statistics due to population stratification. Hence, 70 intergenic SNPs can be chosen in unlinked regions as indicated in a genome scan (Kehoe *et al. Hum. Mol. Genet.* 1999, 8:237-245).

25 Once individual risk factors, genetic or non-genetic, have been found for the predisposition to disease, the next step is to set up a classification/prediction scheme to predict the category (for instance, disease or no-disease) that an individual will be in depending on his genotypes of associated SNPs and other non-genetic risk factors. Logistic regression for discrete trait and linear regression for continuous trait are standard techniques for such tasks (*Applied Regression Analysis*, Draper and Smith, Wiley (1998)).
30 Moreover, other techniques can also be used for setting up classification. Such techniques include, but are not limited to, MART, CART, neural network, and discriminant analyses that are suitable for use in comparing the performance of different methods (*The Elements of Statistical Learning*, Hastie, Tibshirani & Friedman, Springer (2002)).

Disease Diagnosis and Predisposition Screening

Information on association/correlation between genotypes and disease-related phenotypes can be exploited in several ways. For example, in the case of a highly statistically significant association between one or more SNPs with predisposition to a disease for which treatment is available, detection of such a genotype pattern in an individual may justify immediate administration of treatment, or at least the institution of regular monitoring of the individual. Detection of the susceptibility alleles associated with serious disease in a couple contemplating having children may also be valuable to the couple in their reproductive decisions. In the case of a weaker but still statistically significant association between a SNP and a human disease, immediate therapeutic intervention or monitoring may not be justified after detecting the susceptibility allele or SNP. Nevertheless, the subject can be motivated to begin simple life-style changes (*e.g.*, diet, exercise) that can be accomplished at little or no cost to the individual but would confer potential benefits in reducing the risk of developing conditions for which that individual may have an increased risk by virtue of having the risk allele(s).

The SNPs of the invention may contribute to the development of VT in an individual in different ways. Some polymorphisms occur within a protein coding sequence and contribute to disease phenotype by affecting protein structure. Other polymorphisms occur in noncoding regions but may exert phenotypic effects indirectly via influence on, for example, replication, transcription, and/or translation. A single SNP may affect more than one phenotypic trait. Likewise, a single phenotypic trait may be affected by multiple SNPs in different genes.

As used herein, the terms “diagnose,” “diagnosis,” and “diagnostics” include, but are not limited to any of the following: detection of VT that an individual may presently have, predisposition/susceptibility screening (*i.e.*, determining the increased risk of an individual in developing VT in the future, or determining whether an individual has a decreased risk of developing VT in the future), determining a particular type or subclass of VT in an individual known to have VT, confirming or reinforcing a previously made diagnosis of VT, pharmacogenomic evaluation of an individual to determine which therapeutic strategy that individual is most likely to positively respond to or to predict whether a patient is likely to respond to a particular treatment such as statins, predicting whether a patient is likely to experience toxic effects from a particular treatment or therapeutic compound, and evaluating the future prognosis of an individual having VT. Such diagnostic uses are based on the SNPs individually or in a unique combination or SNP haplotypes of the present invention.

Haplotypes are particularly useful in that, for example, fewer SNPs can be genotyped to determine if a particular genomic region harbors a locus that influences a particular phenotype, such as in linkage disequilibrium-based SNP association analysis.

5 Linkage disequilibrium (LD) refers to the co-inheritance of alleles (*e.g.*, alternative nucleotides) at two or more different SNP sites at frequencies greater than would be expected from the separate frequencies of occurrence of each allele in a given population. The expected frequency of co-occurrence of two alleles that are inherited independently is the frequency of the first allele multiplied by the frequency of the second allele. Alleles that co-occur at expected frequencies are said to be in "linkage equilibrium." In contrast, LD refers to any non-random genetic association between allele(s) at two or more different SNP sites, which is generally due to the physical proximity of the two loci along a chromosome. LD can occur when two or more SNPs sites are in close physical proximity to each other on a given chromosome and therefore alleles at these SNP sites will tend to remain unseparated for multiple generations with the consequence that a particular nucleotide (allele) at one SNP site will show a non-random association with a particular nucleotide (allele) at a different SNP site located nearby. Hence, genotyping one of the SNP sites will give almost the same information as genotyping the other SNP site that is in LD.

15 Various degrees of LD can be encountered between two or more SNPs with the result being that some SNPs are more closely associated (*i.e.*, in stronger LD) than others. Furthermore, the physical distance over which LD extends along a chromosome differs between different regions of the genome, and therefore the degree of physical separation between two or more SNP sites necessary for LD to occur can differ between different regions of the genome.

20 For diagnostic purposes and similar uses, if a particular SNP site is found to be useful for diagnosing VT (*e.g.*, has a significant statistical association with the condition and/or is recognized as a causative polymorphism for the condition), then the skilled artisan would recognize that other SNP sites which are in LD with this SNP site would also be useful for diagnosing the condition. Thus, polymorphisms (*e.g.*, SNPs and/or haplotypes) that are not the actual disease-causing (causative) polymorphisms, but are in LD with such causative polymorphisms, are also useful. In such instances, the genotype of the polymorphism(s) that is/are in LD with the causative polymorphism is predictive of the genotype of the causative polymorphism and, consequently, predictive of the phenotype (*e.g.*, VT) that is influenced by the causative SNP(s). Therefore, polymorphic markers that are in LD with causative polymorphisms are useful as diagnostic markers, and are particularly useful when the actual causative polymorphism(s) is/are unknown.

Examples of polymorphisms that can be in LD with one or more causative polymorphisms (and/or in LD with one or more polymorphisms that have a significant statistical association with a condition) and therefore useful for diagnosing the same condition that the causative/associated SNP(s) is used to diagnose, include other SNPs in the same gene, protein-coding, or mRNA transcript-coding region as the causative/associated SNP, other SNPs in the same exon or same intron as the causative/associated SNP, other SNPs in the same haplotype block as the causative/associated SNP, other SNPs in the same intergenic region as the causative/associated SNP, SNPs that are outside but near a gene (*e.g.*, within 6kb on either side, 5' or 3', of a gene boundary) that harbors a causative/associated SNP, etc. Such useful LD SNPs can be selected from among the SNPs disclosed in Tables 1-2, for example.

Linkage disequilibrium in the human genome is reviewed in: Wall *et al.*, "Haplotype blocks and linkage disequilibrium in the human genome", *Nat Rev Genet.* 2003 Aug;4(8):587-97; Garner *et al.*, "On selecting markers for association studies: patterns of linkage disequilibrium between two and three diallelic loci", *Genet Epidemiol.* 2003 Jan;24(1):57-67; Ardlie *et al.*, "Patterns of linkage disequilibrium in the human genome", *Nat Rev Genet.* 2002 Apr;3(4):299-309 (erratum in *Nat Rev Genet* 2002 Jul;3(7):566); and Remm *et al.*, "High-density genotyping and linkage disequilibrium in the human genome using chromosome 22 as a model"; *Curr Opin Chem Biol.* 2002 Feb;6(1):24-30; Haldane JBS (1919) The combination of linkage values, and the calculation of distances between the loci of linked factors. *J Genet* 8:299-309; Mendel, G. (1866) Versuche über Pflanzen-Hybriden. Verhandlungen des naturforschenden Vereines in Brünn [Proceedings of the Natural History Society of Brünn]; Lewin B (1990) *Genes IV.* Oxford University Press, New York, USA; Hartl DL and Clark AG (1989) *Principles of Population Genetics 2nd ed.* Sinauer Associates, Inc. Sunderland, Mass., USA; Gillespie JH (2004) *Population Genetics: A Concise Guide. 2nd ed.* Johns Hopkins University Press. USA; Lewontin RC (1964) The interaction of selection and linkage. I. General considerations; heterotic models. *Genetics* 49:49-67; Hoel PG (1954) *Introduction to Mathematical Statistics 2nd ed.* John Wiley & Sons, Inc. New York, USA; Hudson RR (2001) Two-locus sampling distributions and their application. *Genetics* 159:1805-1817; Dempster AP, Laird NM, Rubin DB (1977) Maximum likelihood from incomplete data via the EM algorithm. *J R Stat Soc* 39:1-38; Excoffier L, Slatkin M (1995) Maximum-likelihood estimation of molecular haplotype frequencies in a diploid population. *Mol Biol Evol* 12(5):921-927; Tregouet DA, Escolano S, Tiret L, Mallet A, Golmard JL (2004) A new algorithm for haplotype-based association analysis: the Stochastic-EM algorithm. *Ann Hum Genet* 68(Pt 2):165-177; Long AD and Langley CH (1999) The power of association studies to detect the contribution of candidate genetic loci to variation in complex

traits. *Genome Research* 9:720-731; Agresti A (1990) *Categorical Data Analysis*. John Wiley & Sons, Inc. New York, USA; Lange K (1997) *Mathematical and Statistical Methods for Genetic Analysis*. Springer-Verlag New York, Inc. New York, USA; The International HapMap Consortium (2003) The International HapMap Project. *Nature* 426:789-796; The International HapMap Consortium (2005) A haplotype map of the human genome. *Nature* 437:1299-1320; Thorisson GA, Smith AV, Krishnan L, Stein LD (2005), The International HapMap Project Web Site. *Genome Research* 15:1591-1593; McVean G, Spencer CCA, Chaix R (2005) Perspectives on human genetic variation from the HapMap project. *PLoS Genetics* 1(4):413-418; Hirschhorn JN, Daly MJ (2005) Genome-wide association studies for common diseases and complex traits. *Nat Genet* 6:95-108; Schrodin SJ (2005) A probabilistic approach to large-scale association scans: a semi-Bayesian method to detect disease-predisposing alleles. *SAGMB* 4(1):31; Wang WYS, Barratt BJ, Clayton DG, Todd JA (2005) Genome-wide association studies: theoretical and practical concerns. *Nat Rev Genet* 6:109-118. Pritchard JK, Przeworski M (2001) Linkage disequilibrium in humans: models and data. *Am J Hum Genet* 69:1-14.

As discussed above, one aspect of the present invention is the discovery that SNPs which are in certain LD distance with the interrogated SNP can also be used as valid markers for identifying an increased or decreased risks of having or developing VT. As used herein, the term "interrogated SNP" refers to SNPs that have been found to be associated with an increased or decreased risk of disease using genotyping results and analysis, or other appropriate experimental method as exemplified in the working examples described in this application. As used herein, the term "LD SNP" refers to a SNP that has been characterized as a SNP associating with an increased or decreased risk of diseases due to their being in LD with the "interrogated SNP" under the methods of calculation described in the application. Below, applicants describe the methods of calculation with which one of ordinary skilled in the art may determine if a particular SNP is in LD with an interrogated SNP. The parameter r^2 is commonly used in the genetics art to characterize the extent of linkage disequilibrium between markers (Hudson, 2001). As used herein, the term "in LD with" refers to a particular SNP that is measured at above the threshold of a parameter such as r^2 with an interrogated SNP.

It is now common place to directly observe genetic variants in a sample of chromosomes obtained from a population. Suppose one has genotype data at two genetic markers located on the same chromosome, for the markers A and B . Further suppose that two alleles segregate at each of these two markers such that alleles A_1 and A_2 can be found at marker A and alleles B_1 and B_2 at marker B . Also assume that these two markers are on a human autosome. If one is to examine a specific individual and find that they are heterozygous at both markers, such that their

two-marker genotype is $A_1A_2B_1B_2$, then there are two possible configurations: the individual in question could have the alleles A_1B_1 on one chromosome and A_2B_2 on the remaining chromosome; alternatively, the individual could have alleles A_1B_2 on one chromosome and A_2B_1 on the other. The arrangement of alleles on a chromosome is called a haplotype. In this illustration, the individual could have haplotypes A_1B_1/A_2B_2 or A_1B_2/A_2B_1 (see Hartl and Clark (1989) for a more complete description). The concept of linkage equilibrium relates the frequency of haplotypes to the allele frequencies.

Assume that a sample of individuals is selected from a larger population. Considering the two markers described above, each having two alleles, there are four possible haplotypes: A_1B_1 , A_1B_2 , A_2B_1 and A_2B_2 . Denote the frequencies of these four haplotypes with the following notation.

$$P_{11} = \text{freq}(A_1B_1) \quad (1)$$

$$P_{12} = \text{freq}(A_1B_2) \quad (2)$$

$$P_{21} = \text{freq}(A_2B_1) \quad (3)$$

$$P_{22} = \text{freq}(A_2B_2) \quad (4)$$

The allele frequencies at the two markers are then the sum of different haplotype frequencies, it is straightforward to write down a similar set of equations relating single-marker allele frequencies to two-marker haplotype frequencies:

$$p_1 = \text{freq}(A_1) = P_{11} + P_{12} \quad (5)$$

$$p_2 = \text{freq}(A_2) = P_{21} + P_{22} \quad (6)$$

$$q_1 = \text{freq}(B_1) = P_{11} + P_{21} \quad (7)$$

$$q_2 = \text{freq}(B_2) = P_{12} + P_{22} \quad (8)$$

Note that the four haplotype frequencies and the allele frequencies at each marker must sum to a frequency of 1.

$$P_{11} + P_{12} + P_{21} + P_{22} = 1 \quad (9)$$

$$p_1 + p_2 = 1 \quad (10)$$

$$q_1 + q_2 = 1 \quad (11)$$

If there is no correlation between the alleles at the two markers, one would expect that the frequency of the haplotypes would be approximately the product of the composite alleles. Therefore,

$$P_{11} \approx p_1q_1 \quad (12)$$

$$P_{12} \approx p_1q_2 \quad (13)$$

$$P_{21} \approx p_2q_1 \quad (14)$$

$$P_{22} \approx p_2q_2 \quad (15)$$

5 These approximating equations (12)-(15) represent the concept of linkage equilibrium where there is independent assortment between the two markers – the alleles at the two markers occur together at random. These are represented as approximations because linkage equilibrium and linkage disequilibrium are concepts typically thought of as properties of a sample of chromosomes; and as such they are susceptible to stochastic fluctuations due to the sampling
10 process. Empirically, many pairs of genetic markers will be in linkage equilibrium, but certainly not all pairs.

Having established the concept of linkage equilibrium above, applicants can now describe the concept of linkage disequilibrium (LD), which is the deviation from linkage equilibrium. Since the frequency of the A_1B_1 haplotype is approximately the product of the allele frequencies
15 for A_1 and B_1 under the assumption of linkage equilibrium as stated mathematically in (12), a simple measure for the amount of departure from linkage equilibrium is the difference in these two quantities, D ,

$$D = P_{11} - p_1q_1 \quad (16)$$

20 $D = 0$ indicates perfect linkage equilibrium. Substantial departures from $D = 0$ indicates LD in the sample of chromosomes examined. Many properties of D are discussed in Lewontin (1964) including the maximum and minimum values that D can take. Mathematically, using basic algebra, it can be shown that D can also be written solely in terms of haplotypes:

$$D = P_{11}P_{22} - P_{12}P_{21} \quad (17)$$

25 If one transforms D by squaring it and subsequently dividing by the product of the allele frequencies of A_1 , A_2 , B_1 and B_2 , the resulting quantity, called r^2 , is equivalent to the square of
30 the Pearson's correlation coefficient commonly used in statistics (e.g. Hoel, 1954).

$$r^2 = \frac{D^2}{p_1p_2q_1q_2} \quad (18)$$

35 As with D , values of r^2 close to 0 indicate linkage equilibrium between the two markers examined in the sample set. As values of r^2 increase, the two markers are said to be in linkage

disequilibrium. The range of values that r^2 can take are from 0 to 1. $r^2 = 1$ when there is a perfect correlation between the alleles at the two markers.

In addition, the quantities discussed above are sample-specific. And as such, it is necessary to formulate notation specific to the samples studied. In the approach discussed here, three types of samples are of primary interest: (i) a sample of chromosomes from individuals affected by a disease-related phenotype (cases), (ii) a sample of chromosomes obtained from individuals not affected by the disease-related phenotype (controls), and (iii) a standard sample set used for the construction of haplotypes and calculation pairwise linkage disequilibrium. For the allele frequencies used in the development of the method described below, an additional subscript will be added to denote either the case or control sample sets.

$$p_{1,cs} = \text{freq}(A_1 \text{ in cases}) \quad (19)$$

$$p_{2,cs} = \text{freq}(A_2 \text{ in cases}) \quad (20)$$

$$q_{1,cs} = \text{freq}(B_1 \text{ in cases}) \quad (21)$$

$$q_{2,cs} = \text{freq}(B_2 \text{ in cases}) \quad (22)$$

Similarly,

$$p_{1,ct} = \text{freq}(A_1 \text{ in controls}) \quad (23)$$

$$p_{2,ct} = \text{freq}(A_2 \text{ in controls}) \quad (24)$$

$$q_{1,ct} = \text{freq}(B_1 \text{ in controls}) \quad (25)$$

$$q_{2,ct} = \text{freq}(B_2 \text{ in controls}) \quad (26)$$

As a well-accepted sample set is necessary for robust linkage disequilibrium calculations, data obtained from the International HapMap project (The International HapMap Consortium 2003, 2005; Thorisson et al, 2005; McVean et al, 2005) can be used for the calculation of pairwise r^2 values. Indeed, the samples genotyped for the International HapMap Project were selected to be representative examples from various human sub-populations with sufficient numbers of chromosomes examined to draw meaningful and robust conclusions from the patterns of genetic variation observed. The International HapMap project website contains a description of the project, methods utilized and samples examined. It is useful to examine empirical data to get a sense of the patterns present in such data.

Haplotype frequencies were explicit arguments in equation (18) above. However, knowing the 2-marker haplotype frequencies requires that phase to be determined for doubly heterozygous samples. When phase is unknown in the data examined, various algorithms can be used to infer phase from the genotype data. This issue was discussed earlier where the doubly heterozygous individual with a 2-SNP genotype of $A_1A_2B_1B_2$ could have one of two different

sets of chromosomes: A_1B_1/A_2B_2 or A_1B_2/A_2B_1 . One such algorithm to estimate haplotype frequencies is the expectation-maximization (EM) algorithm first formalized by Dempster et al (1977). This algorithm is often used in genetics to infer haplotype frequencies from genotype data (e.g. Excoffier and Slatkin, 1995; Tregouet et al, 2004). It should be noted that for the two-
 5 SNP case explored here, EM algorithms have very little error provided that the allele frequencies and sample sizes are not too small. The impact on r^2 values is typically negligible.

As correlated genetic markers share information, interrogation of SNP markers in LD with a disease-associated SNP marker can also have sufficient power to detect disease association (Long and Langley, 1999). The relationship between the power to directly find
 10 disease-associated alleles and the power to indirectly detect disease-association was investigated by Pritchard and Przeworski (2001). In a straight-forward derivation, it can be shown that the power to detect disease association indirectly at a marker locus in linkage disequilibrium with a disease-association locus is approximately the same as the power to detect disease-association directly at the disease- association locus if the sample size is increased by a factor of $\frac{1}{r^2}$ (the
 15 reciprocal of equation 18) at the marker in comparison with the disease- association locus.

Therefore, if one calculated the power to detect disease-association indirectly with an experiment having N samples, then equivalent power to directly detect disease-association (at the actual disease-susceptibility locus) would necessitate an experiment using approximately
 r^2N samples. This elementary relationship between power, sample size and linkage
 20 disequilibrium can be used to derive an r^2 threshold value useful in determining whether or not genotyping markers in linkage disequilibrium with a SNP marker directly associated with disease status has enough power to indirectly detect disease-association.

To commence a derivation of the power to detect disease-associated markers through an indirect process, define the effective chromosomal sample size as
 25

$$n = \frac{4N_{cs}N_{ct}}{N_{cs} + N_{ct}}; \quad (27)$$

where N_{cs} and N_{ct} are the numbers of diploid cases and controls, respectively. This is necessary to handle situations where the numbers of cases and controls are not equivalent. For
 30 equal case and control sample sizes, $N_{cs} = N_{ct} = N$, the value of the effective number of chromosomes is simply $n = 2N$ – as expected. Let power be calculated for a significance level α (such that traditional P-values below α will be deemed statistically significant). Define the standard Gaussian distribution function as $\Phi(\bullet)$. Mathematically,

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{\theta^2}{2}} d\theta \quad (28)$$

Alternatively, the following error function notation (Erf) may also be used,

5

$$\Phi(x) = \frac{1}{2} \left[1 + \operatorname{Erf} \left(\frac{x}{\sqrt{2}} \right) \right] \quad (29)$$

For example, $\Phi(1.644854) = 0.95$. The value of r^2 may be derived to yield a pre-specified minimum amount of power to detect disease association through indirect interrogation.

10 Noting that the LD SNP marker could be the one that is carrying the disease-association allele, therefore that this approach constitutes a lower-bound model where all indirect power results are expected to be at least as large as those interrogated.

Denote by β the error rate for not detecting truly disease-associated markers. Therefore, $1 - \beta$ is the classical definition of statistical power. Substituting the Pritchard-Pzreworski result into the sample size, the power to detect disease association at a significance level of α is given by the approximation

$$1 - \beta \cong \Phi \left[\frac{|q_{1,cs} - q_{1,ct}|}{\sqrt{\frac{q_{1,cs}(1 - q_{1,cs}) + q_{1,ct}(1 - q_{1,ct})}{r^2 n}}} - Z_{1-\alpha/2} \right]; \quad (30)$$

where Z_u is the inverse of the standard normal cumulative distribution evaluated at u ($u \in (0,1)$).

20 $Z_u = \Phi^{-1}(u)$, where $\Phi(\Phi^{-1}(u)) = \Phi^{-1}(\Phi(u)) = u$. For example, setting $\alpha = 0.05$, and therefore $1 - \alpha/2 = 0.975$, we obtain $Z_{0.975} = 1.95996$. Next, setting power equal to a threshold of a minimum power of T ,

$$T = \Phi \left[\frac{|q_{1,cs} - q_{1,ct}|}{\sqrt{\frac{q_{1,cs}(1 - q_{1,cs}) + q_{1,ct}(1 - q_{1,ct})}{r^2 n}}} - Z_{1-\alpha/2} \right] \quad (31)$$

25

and solving for r^2 , the following threshold r^2 is obtained:

$$r_T^2 = \frac{[q_{1,cs}(1-q_{1,cs}) + q_{1,ct}(1-q_{1,ct})]}{n(q_{1,cs} - q_{1,ct})^2} [\Phi^{-1}(T) + Z_{1-\alpha/2}] \quad (32)$$

5 Or,

$$r_T^2 = \left(\frac{Z_T + Z_{1-\alpha/2}}{n} \right) \left[\frac{q_{1,cs} - (q_{1,cs})^2 + q_{1,ct} - (q_{1,ct})^2}{(q_{1,cs} - q_{1,ct})^2} \right] \quad (33)$$

Suppose that r^2 is calculated between an interrogated SNP and a number of other SNPs
 10 with varying levels of LD with the interrogated SNP. The threshold value r_T^2 is the minimum
 value of linkage disequilibrium between the interrogated SNP and the potential LD SNPs such
 that the LD SNP still retains a power greater or equal to T for detecting disease-association. For
 example, suppose that SNP rs200 is genotyped in a case-control disease-association study and it
 is found to be associated with a disease phenotype. Further suppose that the minor allele
 15 frequency in 1,000 case chromosomes was found to be 16% in contrast with a minor allele
 frequency of 10% in 1,000 control chromosomes. Given those measurements one could have
 predicted, prior to the experiment, that the power to detect disease association at a significance
 level of 0.05 was quite high – approximately 98% using a test of allelic association. Applying
 equation (32) one can calculate a minimum value of r^2 to indirectly assess disease association
 20 assuming that the minor allele at SNP rs200 is truly disease-predisposing for a threshold level of
 power. If one sets the threshold level of power to be 80%, then $r_T^2 = 0.489$ given the same
 significance level and chromosome numbers as above. Hence, any SNP with a pairwise r^2 value
 with rs200 greater than 0.489 is expected to have greater than 80% power to detect the disease
 association. Further, this is assuming the conservative model where the LD SNP is disease-
 25 associated only through linkage disequilibrium with the interrogated SNP rs200.

The contribution or association of particular SNPs and/or SNP haplotypes with disease
 phenotypes, such as VT, enables the SNPs of the present invention to be used to develop superior
 diagnostic tests capable of identifying individuals who express a detectable trait, such as VT, as
 the result of a specific genotype, or individuals whose genotype places them at an increased or
 30 decreased risk of developing a detectable trait at a subsequent time as compared to individuals
 who do not have that genotype. As described herein, diagnostics may be based on a single SNP
 or a group of SNPs. Combined detection of a plurality of SNPs (for example, 2, 3, 4, 5, 6, 7, 8, 9,
 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 24, 25, 30, 32, 48, 50, 64, 96, 100, or any other number
 in-between, or more, of the SNPs provided in Table 1 and/or Table 2) typically increases the

probability of an accurate diagnosis. For example, the presence of a single SNP known to correlate with VT might indicate a probability of 20% that an individual has or is at risk of developing VT, whereas detection of five SNPs, each of which correlates with VT, might indicate a probability of 80% that an individual has or is at risk of developing VT. To further
5 increase the accuracy of diagnosis or predisposition screening, analysis of the SNPs of the present invention can be combined with that of other polymorphisms or other risk factors of VT, such as disease symptoms, pathological characteristics, family history, diet, environmental factors or lifestyle factors.

It will, of course, be understood by practitioners skilled in the treatment or diagnosis of
10 VT that the present invention generally does not intend to provide an absolute identification of individuals who are at risk (or less at risk) of developing VT, and/or pathologies related to VT, but rather to indicate a certain increased (or decreased) degree or likelihood of developing the disease based on statistically significant association results. However, this information is extremely valuable as it can be used to, for example, initiate preventive treatments or to allow an
15 individual carrying one or more significant SNPs or SNP haplotypes to foresee warning signs such as minor clinical symptoms, or to have regularly scheduled physical exams to monitor for appearance of a condition in order to identify and begin treatment of the condition at an early stage. Particularly with diseases that are extremely debilitating or fatal if not treated on time, the knowledge of a potential predisposition, even if this predisposition is not absolute, would likely
20 contribute in a very significant manner to treatment efficacy.

The diagnostic techniques of the present invention may employ a variety of methodologies to determine whether a test subject has a SNP or a SNP pattern associated with an increased or decreased risk of developing a detectable trait or whether the individual suffers from a detectable trait as a result of a particular polymorphism/mutation, including, for example,
25 methods which enable the analysis of individual chromosomes for haplotyping, family studies, single sperm DNA analysis, or somatic hybrids. The trait analyzed using the diagnostics of the invention may be any detectable trait that is commonly observed in pathologies and disorders related to VT.

Another aspect of the present invention relates to a method of determining whether an
30 individual is at risk (or less at risk) of developing one or more traits or whether an individual expresses one or more traits as a consequence of possessing a particular trait-causing or trait-influencing allele. These methods generally involve obtaining a nucleic acid sample from an individual and assaying the nucleic acid sample to determine which nucleotide(s) is/are present at one or more SNP positions, wherein the assayed nucleotide(s) is/are indicative of an increased or

decreased risk of developing the trait or indicative that the individual expresses the trait as a result of possessing a particular trait-causing or trait-influencing allele.

In another embodiment, the SNP detection reagents of the present invention are used to determine whether an individual has one or more SNP allele(s) affecting the level (*e.g.*, the concentration of mRNA or protein in a sample, etc.) or pattern (*e.g.*, the kinetics of expression, rate of decomposition, stability profile, K_m , V_{max} , etc.) of gene expression (collectively, the "gene response" of a cell or bodily fluid). Such a determination can be accomplished by screening for mRNA or protein expression (*e.g.*, by using nucleic acid arrays, RT-PCR, TaqMan assays, or mass spectrometry), identifying genes having altered expression in an individual, genotyping SNPs disclosed in Table 1 and/or Table 2 that could affect the expression of the genes having altered expression (*e.g.*, SNPs that are in and/or around the gene(s) having altered expression, SNPs in regulatory/control regions, SNPs in and/or around other genes that are involved in pathways that could affect the expression of the gene(s) having altered expression, or all SNPs could be genotyped), and correlating SNP genotypes with altered gene expression. In this manner, specific SNP alleles at particular SNP sites can be identified that affect gene expression.

Pharmacogenomics and Therapeutics/Drug Development

The present invention provides methods for assessing the pharmacogenomics of a subject harboring particular SNP alleles or haplotypes to a particular therapeutic agent or pharmaceutical compound, or to a class of such compounds. Pharmacogenomics deals with the roles which clinically significant hereditary variations (*e.g.*, SNPs) play in the response to drugs due to altered drug disposition and/or abnormal action in affected persons. See, *e.g.*, Roses, *Nature* 405, 857-865 (2000); Gould Rothberg, *Nature Biotechnology* 19, 209-211 (2001); Eichelbaum, *Clin. Exp. Pharmacol. Physiol.* 23(10-11):983-985 (1996); and Linder, *Clin. Chem.* 43(2):254-266 (1997). The clinical outcomes of these variations can result in severe toxicity of therapeutic drugs in certain individuals or therapeutic failure of drugs in certain individuals as a result of individual variation in metabolism. Thus, the SNP genotype of an individual can determine the way a therapeutic compound acts on the body or the way the body metabolizes the compound. For example, SNPs in drug metabolizing enzymes can affect the activity of these enzymes, which in turn can affect both the intensity and duration of drug action, as well as drug metabolism and clearance.

The discovery of SNPs in drug metabolizing enzymes, drug transporters, proteins for pharmaceutical agents, and other drug targets has explained why some patients do not obtain the expected drug effects, show an exaggerated drug effect, or experience serious toxicity from standard

drug dosages. SNPs can be expressed in the phenotype of the extensive metabolizer and in the phenotype of the poor metabolizer. Accordingly, SNPs may lead to allelic variants of a protein in which one or more of the protein functions in one population are different from those in another population. SNPs and the encoded variant peptides thus provide targets to ascertain a genetic predisposition that can affect treatment modality. For example, in a ligand-based treatment, SNPs may give rise to amino terminal extracellular domains and/or other ligand-binding regions of a receptor that are more or less active in ligand binding, thereby affecting subsequent protein activation. Accordingly, ligand dosage would necessarily be modified to maximize the therapeutic effect within a given population containing particular SNP alleles or haplotypes.

As an alternative to genotyping, specific variant proteins containing variant amino acid sequences encoded by alternative SNP alleles could be identified. Thus, pharmacogenomic characterization of an individual permits the selection of effective compounds and effective dosages of such compounds for prophylactic or therapeutic uses based on the individual's SNP genotype, thereby enhancing and optimizing the effectiveness of the therapy. Furthermore, the production of recombinant cells and transgenic animals containing particular SNPs/haplotypes allow effective clinical design and testing of treatment compounds and dosage regimens. For example, transgenic animals can be produced that differ only in specific SNP alleles in a gene that is orthologous to a human disease susceptibility gene.

Pharmacogenomic uses of the SNPs of the present invention provide several significant advantages for patient care, particularly in treating VT. Pharmacogenomic characterization of an individual, based on an individual's SNP genotype, can identify those individuals unlikely to respond to treatment with a particular medication and thereby allows physicians to avoid prescribing the ineffective medication to those individuals. On the other hand, SNP genotyping of an individual may enable physicians to select the appropriate medication and dosage regimen that will be most effective based on an individual's SNP genotype. This information increases a physician's confidence in prescribing medications and motivates patients to comply with their drug regimens. Furthermore, pharmacogenomics may identify patients predisposed to toxicity and adverse reactions to particular drugs or drug dosages. Adverse drug reactions lead to more than 100,000 avoidable deaths per year in the United States alone and therefore represent a significant cause of hospitalization and death, as well as a significant economic burden on the healthcare system (Pfof *et. al.*, *Trends in Biotechnology*, Aug. 2000.). Thus, pharmacogenomics based on the SNPs disclosed herein has the potential to both save lives and reduce healthcare costs substantially.

Pharmacogenomics in general is discussed further in Rose *et al.*, "Pharmacogenetic analysis of clinically relevant genetic polymorphisms", *Methods Mol Med.* 2003;85:225-37.

Pharmacogenomics as it relates to Alzheimer's disease and other neurodegenerative disorders is discussed in Cacabelos, "Pharmacogenomics for the treatment of dementia", *Ann Med.* 2002;34(5):357-79, Maimone *et al.*, "Pharmacogenomics of neurodegenerative diseases", *Eur J Pharmacol.* 2001 Feb 9;413(1):11-29, and Poirier, "Apolipoprotein E: a pharmacogenetic target for the treatment of Alzheimer's disease", *Mol Diagn.* 1999 Dec;4(4):335-41. Pharmacogenomics as it relates to cardiovascular disorders is discussed in Siest *et al.*, "Pharmacogenomics of drugs affecting the cardiovascular system", *Clin Chem Lab Med.* 2003 Apr;41(4):590-9, Mukherjee *et al.*, "Pharmacogenomics in cardiovascular diseases", *Prog Cardiovasc Dis.* 2002 May-Jun;44(6):479-98, and Mooser *et al.*, "Cardiovascular pharmacogenetics in the SNP era", *J Thromb Haemost.* 2003 Jul;1(7):1398-402. Pharmacogenomics as it relates to cancer is discussed in McLeod *et al.*, "Cancer pharmacogenomics: SNPs, chips, and the individual patient", *Cancer Invest.* 2003;21(4):630-40 and Watters *et al.*, "Cancer pharmacogenomics: current and future applications", *Biochim Biophys Acta.* 2003 Mar 17;1603(2):99-111.

The SNPs of the present invention also can be used to identify novel therapeutic targets for VT. For example, genes containing the disease-associated variants ("variant genes") or their products, as well as genes or their products that are directly or indirectly regulated by or interacting with these variant genes or their products, can be targeted for the development of therapeutics that, for example, treat the disease or prevent or delay disease onset. The therapeutics may be composed of, for example, small molecules, proteins, protein fragments or peptides, antibodies, nucleic acids, or their derivatives or mimetics which modulate the functions or levels of the target genes or gene products.

The SNP-containing nucleic acid molecules disclosed herein, and their complementary nucleic acid molecules, may be used as antisense constructs to control gene expression in cells, tissues, and organisms. Antisense technology is well established in the art and extensively reviewed in *Antisense Drug Technology: Principles, Strategies, and Applications*, Crooke (ed.), Marcel Dekker, Inc.: New York (2001). An antisense nucleic acid molecule is generally designed to be complementary to a region of mRNA expressed by a gene so that the antisense molecule hybridizes to the mRNA and thereby blocks translation of mRNA into protein. Various classes of antisense oligonucleotides are used in the art, two of which are cleavers and blockers. Cleavers, by binding to target RNAs, activate intracellular nucleases (*e.g.*, RNaseH or RNase L) that cleave the target RNA. Blockers, which also bind to target RNAs, inhibit protein translation through steric hindrance of ribosomes. Exemplary blockers include peptide nucleic acids, morpholinos, locked nucleic acids, and methylphosphonates (see, *e.g.*, Thompson, *Drug Discovery Today*, 7 (17): 912-917 (2002)). Antisense oligonucleotides are directly useful as

therapeutic agents, and are also useful for determining and validating gene function (*e.g.*, in gene knock-out or knock-down experiments).

Antisense technology is further reviewed in: Lavery *et al.*, "Antisense and RNAi: powerful tools in drug target discovery and validation", *Curr Opin Drug Discov Devel.* 2003 Jul;6(4):561-9; Stephens *et al.*, "Antisense oligonucleotide therapy in cancer", *Curr Opin Mol Ther.* 2003 Apr;5(2):118-22; Kurreck, "Antisense technologies. Improvement through novel chemical modifications", *Eur J Biochem.* 2003 Apr;270(8):1628-44; Dias *et al.*, "Antisense oligonucleotides: basic concepts and mechanisms", *Mol Cancer Ther.* 2002 Mar;1(5):347-55; Chen, "Clinical development of antisense oligonucleotides as anti-cancer therapeutics", *Methods Mol Med.* 2003;75:621-36; Wang *et al.*, "Antisense anticancer oligonucleotide therapeutics", *Curr Cancer Drug Targets.* 2001 Nov;1(3):177-96; and Bennett, "Efficiency of antisense oligonucleotide drug discovery", *Antisense Nucleic Acid Drug Dev.* 2002 Jun;12(3):215-24.

The SNPs of the present invention are particularly useful for designing antisense reagents that are specific for particular nucleic acid variants. Based on the SNP information disclosed herein, antisense oligonucleotides can be produced that specifically target mRNA molecules that contain one or more particular SNP nucleotides. In this manner, expression of mRNA molecules that contain one or more undesired polymorphisms (*e.g.*, SNP nucleotides that lead to a defective protein such as an amino acid substitution in a catalytic domain) can be inhibited or completely blocked. Thus, antisense oligonucleotides can be used to specifically bind a particular polymorphic form (*e.g.*, a SNP allele that encodes a defective protein), thereby inhibiting translation of this form, but which do not bind an alternative polymorphic form (*e.g.*, an alternative SNP nucleotide that encodes a protein having normal function).

Antisense molecules can be used to inactivate mRNA in order to inhibit gene expression and production of defective proteins. Accordingly, these molecules can be used to treat a disorder, such as VT, characterized by abnormal or undesired gene expression or expression of certain defective proteins. This technique can involve cleavage by means of ribozymes containing nucleotide sequences complementary to one or more regions in the mRNA that attenuate the ability of the mRNA to be translated. Possible mRNA regions include, for example, protein-coding regions and particularly protein-coding regions corresponding to catalytic activities, substrate/ligand binding, or other functional activities of a protein.

The SNPs of the present invention are also useful for designing RNA interference reagents that specifically target nucleic acid molecules having particular SNP variants. RNA interference (RNAi), also referred to as gene silencing, is based on using double-stranded RNA (dsRNA) molecules to turn genes off. When introduced into a cell, dsRNAs are processed by the

cell into short fragments (generally about 21, 22, or 23 nucleotides in length) known as small interfering RNAs (siRNAs) which the cell uses in a sequence-specific manner to recognize and destroy complementary RNAs (Thompson, *Drug Discovery Today*, 7 (17): 912-917 (2002)). Accordingly, an aspect of the present invention specifically contemplates isolated nucleic acid molecules that are about 18-26 nucleotides in length, preferably 19-25 nucleotides in length, and more preferably 20, 21, 22, or 23 nucleotides in length, and the use of these nucleic acid molecules for RNAi. Because RNAi molecules, including siRNAs, act in a sequence-specific manner, the SNPs of the present invention can be used to design RNAi reagents that recognize and destroy nucleic acid molecules having specific SNP alleles/nucleotides (such as deleterious alleles that lead to the production of defective proteins), while not affecting nucleic acid molecules having alternative SNP alleles (such as alleles that encode proteins having normal function). As with antisense reagents, RNAi reagents may be directly useful as therapeutic agents (*e.g.*, for turning off defective, disease-causing genes), and are also useful for characterizing and validating gene function (*e.g.*, in gene knock-out or knock-down experiments).

The following references provide a further review of RNAi: Reynolds *et al.*, "Rational siRNA design for RNA interference", *Nat Biotechnol.* 2004 Mar;22(3):326-30. Epub 2004 Feb 01; Chi *et al.*, "Genomewide view of gene silencing by small interfering RNAs", *PNAS* 100(11):6343-6346, 2003; Vickers *et al.*, "Efficient Reduction of Target RNAs by Small Interfering RNA and RNase H-dependent Antisense Agents", *J. Biol. Chem.* 278: 7108-7118, 2003; Agami, "RNAi and related mechanisms and their potential use for therapy", *Curr Opin Chem Biol.* 2002 Dec;6(6):829-34; Lavery *et al.*, "Antisense and RNAi: powerful tools in drug target discovery and validation", *Curr Opin Drug Discov Devel.* 2003 Jul;6(4):561-9; Shi, "Mammalian RNAi for the masses", *Trends Genet* 2003 Jan;19(1):9-12), Shuey *et al.*, "RNAi: gene-silencing in therapeutic intervention", *Drug Discovery Today* 2002 Oct;7(20):1040-1046; McManus *et al.*, *Nat Rev Genet* 2002 Oct;3(10):737-47; Xia *et al.*, *Nat Biotechnol* 2002 Oct;20(10):1006-10; Plasterk *et al.*, *Curr Opin Genet Dev* 2000 Oct;10(5):562-7; Boshier *et al.*, *Nat Cell Biol* 2000 Feb;2(2):E31-6; and Hunter, *Curr Biol* 1999 Jun 17;9(12):R440-2).

A subject suffering from a pathological condition, such as VT, ascribed to a SNP may be treated so as to correct the genetic defect (see Kren *et al.*, *Proc. Natl. Acad. Sci. USA* 96:10349-10354 (1999)). Such a subject can be identified by any method that can detect the polymorphism in a biological sample drawn from the subject. Such a genetic defect may be permanently corrected by administering to such a subject a nucleic acid fragment incorporating a repair sequence that supplies the normal/wild-type nucleotide at the position of the SNP. This site-

specific repair sequence can encompass an RNA/DNA oligonucleotide that operates to promote endogenous repair of a subject's genomic DNA. The site-specific repair sequence is administered in an appropriate vehicle, such as a complex with polyethylenimine, encapsulated in anionic liposomes, a viral vector such as an adenovirus, or other pharmaceutical composition that
5 promotes intracellular uptake of the administered nucleic acid. A genetic defect leading to an inborn pathology may then be overcome, as the chimeric oligonucleotides induce incorporation of the normal sequence into the subject's genome. Upon incorporation, the normal gene product is expressed, and the replacement is propagated, thereby engendering a permanent repair and therapeutic enhancement of the clinical condition of the subject.

10 In cases in which a cSNP results in a variant protein that is ascribed to be the cause of, or a contributing factor to, a pathological condition, a method of treating such a condition can include administering to a subject experiencing the pathology the wild-type/normal cognate of the variant protein. Once administered in an effective dosing regimen, the wild-type cognate provides complementation or remediation of the pathological condition.

15 The invention further provides a method for identifying a compound or agent that can be used to treat VT. The SNPs disclosed herein are useful as targets for the identification and/or development of therapeutic agents. A method for identifying a therapeutic agent or compound typically includes assaying the ability of the agent or compound to modulate the activity and/or expression of a SNP-containing nucleic acid or the encoded product and thus identifying an agent or
20 a compound that can be used to treat a disorder characterized by undesired activity or expression of the SNP-containing nucleic acid or the encoded product. The assays can be performed in cell-based and cell-free systems. Cell-based assays can include cells naturally expressing the nucleic acid molecules of interest or recombinant cells genetically engineered to express certain nucleic acid molecules.

25 Variant gene expression in a VT patient can include, for example, either expression of a SNP-containing nucleic acid sequence (for instance, a gene that contains a SNP can be transcribed into an mRNA transcript molecule containing the SNP, which can in turn be translated into a variant protein) or altered expression of a normal/wild-type nucleic acid sequence due to one or more SNPs (for instance, a regulatory/control region can contain a SNP that affects the level or pattern of
30 expression of a normal transcript).

Assays for variant gene expression can involve direct assays of nucleic acid levels (*e.g.*, mRNA levels), expressed protein levels, or of collateral compounds involved in a signal pathway. Further, the expression of genes that are up- or down-regulated in response to the signal pathway

can also be assayed. In this embodiment, the regulatory regions of these genes can be operably linked to a reporter gene such as luciferase.

Modulators of variant gene expression can be identified in a method wherein, for example, a cell is contacted with a candidate compound/agent and the expression of mRNA determined. The level of expression of mRNA in the presence of the candidate compound is compared to the level of expression of mRNA in the absence of the candidate compound. The candidate compound can then be identified as a modulator of variant gene expression based on this comparison and be used to treat a disorder such as VT that is characterized by variant gene expression (*e.g.*, either expression of a SNP-containing nucleic acid or altered expression of a normal/wild-type nucleic acid molecule due to one or more SNPs that affect expression of the nucleic acid molecule) due to one or more SNPs of the present invention. When expression of mRNA is statistically significantly greater in the presence of the candidate compound than in its absence, the candidate compound is identified as a stimulator of nucleic acid expression. When nucleic acid expression is statistically significantly less in the presence of the candidate compound than in its absence, the candidate compound is identified as an inhibitor of nucleic acid expression.

The invention further provides methods of treatment, with the SNP or associated nucleic acid domain (*e.g.*, catalytic domain, ligand/substrate-binding domain, regulatory/control region, etc.) or gene, or the encoded mRNA transcript, as a target, using a compound identified through drug screening as a gene modulator to modulate variant nucleic acid expression. Modulation can include either up-regulation (*i.e.*, activation or agonization) or down-regulation (*i.e.*, suppression or antagonization) of nucleic acid expression.

Expression of mRNA transcripts and encoded proteins, either wild type or variant, may be altered in individuals with a particular SNP allele in a regulatory/control element, such as a promoter or transcription factor binding domain, that regulates expression. In this situation, methods of treatment and compounds can be identified, as discussed herein, that regulate or overcome the variant regulatory/control element, thereby generating normal, or healthy, expression levels of either the wild type or variant protein.

The SNP-containing nucleic acid molecules of the present invention are also useful for monitoring the effectiveness of modulating compounds on the expression or activity of a variant gene, or encoded product, in clinical trials or in a treatment regimen. Thus, the gene expression pattern can serve as an indicator for the continuing effectiveness of treatment with the compound, particularly with compounds to which a patient can develop resistance, as well as an indicator for toxicities. The gene expression pattern can also serve as a marker indicative of a physiological response of the affected cells to the compound. Accordingly, such monitoring would allow either

increased administration of the compound or the administration of alternative compounds to which the patient has not become resistant. Similarly, if the level of nucleic acid expression falls below a desirable level, administration of the compound could be commensurately decreased.

In another aspect of the present invention, there is provided a pharmaceutical pack
5 comprising a therapeutic agent (*e.g.*, a small molecule drug, antibody, peptide, antisense or RNAi nucleic acid molecule, etc.) and a set of instructions for administration of the therapeutic agent to humans diagnostically tested for one or more SNPs or SNP haplotypes provided by the present invention.

The SNPs/haplotypes of the present invention are also useful for improving many
10 different aspects of the drug development process. For instance, an aspect of the present invention includes selecting individuals for clinical trials based on their SNP genotype. For example, individuals with SNP genotypes that indicate that they are likely to positively respond to a drug can be included in the trials, whereas those individuals whose SNP genotypes indicate that they are less likely to or would not respond to the drug, or who are at risk for suffering toxic
15 effects or other adverse reactions, can be excluded from the clinical trials. This not only can improve the safety of clinical trials, but also can enhance the chances that the trial will demonstrate statistically significant efficacy. Furthermore, the SNPs of the present invention may explain why certain previously developed drugs performed poorly in clinical trials and may help identify a subset of the population that would benefit from a drug that had previously
20 performed poorly in clinical trials, thereby “rescuing” previously developed drugs, and enabling the drug to be made available to a particular VT patient population that can benefit from it.

SNPs have many important uses in drug discovery, screening, and development. A high probability exists that, for any gene/protein selected as a potential drug target, variants of that gene/protein will exist in a patient population. Thus, determining the impact of gene/protein
25 variants on the selection and delivery of a therapeutic agent should be an integral aspect of the drug discovery and development process. (*Jazwinska, A Trends Guide to Genetic Variation and Genomic Medicine*, 2002 Mar; S30-S36).

Knowledge of variants (*e.g.*, SNPs and any corresponding amino acid polymorphisms) of a particular therapeutic target (*e.g.*, a gene, mRNA transcript, or protein) enables parallel
30 screening of the variants in order to identify therapeutic candidates (*e.g.*, small molecule compounds, antibodies, antisense or RNAi nucleic acid compounds, etc.) that demonstrate efficacy across variants (*Rothberg, Nat Biotechnol* 2001 Mar;19(3):209-11). Such therapeutic candidates would be expected to show equal efficacy across a larger segment of the patient population, thereby leading to a larger potential market for the therapeutic candidate.

Furthermore, identifying variants of a potential therapeutic target enables the most common form of the target to be used for selection of therapeutic candidates, thereby helping to ensure that the experimental activity that is observed for the selected candidates reflects the real activity expected in the largest proportion of a patient population (Jazwinska, *A Trends Guide to Genetic Variation and Genomic Medicine*, 2002 Mar; S30-S36).

Additionally, screening therapeutic candidates against all known variants of a target can enable the early identification of potential toxicities and adverse reactions relating to particular variants. For example, variability in drug absorption, distribution, metabolism and excretion (ADME) caused by, for example, SNPs in therapeutic targets or drug metabolizing genes, can be identified, and this information can be utilized during the drug development process to minimize variability in drug disposition and develop therapeutic agents that are safer across a wider range of a patient population. The SNPs of the present invention, including the variant proteins and encoding polymorphic nucleic acid molecules provided in Tables 1-2, are useful in conjunction with a variety of toxicology methods established in the art, such as those set forth in *Current Protocols in Toxicology*, John Wiley & Sons, Inc., N.Y.

Furthermore, therapeutic agents that target any art-known proteins (or nucleic acid molecules, either RNA or DNA) may cross-react with the variant proteins (or polymorphic nucleic acid molecules) disclosed in Table 1, thereby significantly affecting the pharmacokinetic properties of the drug. Consequently, the protein variants and the SNP-containing nucleic acid molecules disclosed in Tables 1-2 are useful in developing, screening, and evaluating therapeutic agents that target corresponding art-known protein forms (or nucleic acid molecules). Additionally, as discussed above, knowledge of all polymorphic forms of a particular drug target enables the design of therapeutic agents that are effective against most or all such polymorphic forms of the drug target.

Pharmaceutical Compositions and Administration Thereof

Any of the VT-associated proteins, and encoding nucleic acid molecules, disclosed herein can be used as therapeutic targets (or directly used themselves as therapeutic compounds) for treating VT and related pathologies, and the present disclosure enables therapeutic compounds (*e.g.*, small molecules, antibodies, therapeutic proteins, RNAi and antisense molecules, etc.) to be developed that target (or are comprised of) any of these therapeutic targets.

In general, a therapeutic compound will be administered in a therapeutically effective amount by any of the accepted modes of administration for agents that serve similar utilities. The actual amount of the therapeutic compound of this invention, *i.e.*, the active ingredient, will

depend upon numerous factors such as the severity of the disease to be treated, the age and relative health of the subject, the potency of the compound used, the route and form of administration, and other factors.

5 Therapeutically effective amounts of therapeutic compounds may range from, for example, approximately 0.01-50 mg per kilogram body weight of the recipient per day; preferably about 0.1-20 mg/kg/day. Thus, as an example, for administration to a 70 kg person, the dosage range would most preferably be about 7 mg to 1.4 g per day.

10 In general, therapeutic compounds will be administered as pharmaceutical compositions by any one of the following routes: oral, systemic (*e.g.*, transdermal, intranasal, or by suppository), or parenteral (*e.g.*, intramuscular, intravenous, or subcutaneous) administration. The preferred manner of administration is oral or parenteral using a convenient daily dosage regimen, which can be adjusted according to the degree of affliction. Oral compositions can take the form of tablets, pills, capsules, semisolids, powders, sustained release formulations, solutions, suspensions, elixirs, aerosols, or any other appropriate compositions.

15 The choice of formulation depends on various factors such as the mode of drug administration (*e.g.*, for oral administration, formulations in the form of tablets, pills, or capsules are preferred) and the bioavailability of the drug substance. Recently, pharmaceutical formulations have been developed especially for drugs that show poor bioavailability based upon the principle that bioavailability can be increased by increasing the surface area, *i.e.*, decreasing
20 particle size. For example, U.S. Patent No. 4,107,288 describes a pharmaceutical formulation having particles in the size range from 10 to 1,000 nm in which the active material is supported on a cross-linked matrix of macromolecules. U.S. Patent No. 5,145,684 describes the production of a pharmaceutical formulation in which the drug substance is pulverized to nanoparticles (average particle size of 400 nm) in the presence of a surface modifier and then dispersed in a
25 liquid medium to give a pharmaceutical formulation that exhibits remarkably high bioavailability.

Pharmaceutical compositions are comprised of, in general, a therapeutic compound in combination with at least one pharmaceutically acceptable excipient. Acceptable excipients are non-toxic, aid administration, and do not adversely affect the therapeutic benefit of the
30 therapeutic compound. Such excipients may be any solid, liquid, semi-solid or, in the case of an aerosol composition, gaseous excipient that is generally available to one skilled in the art.

Solid pharmaceutical excipients include starch, cellulose, talc, glucose, lactose, sucrose, gelatin, malt, rice, flour, chalk, silica gel, magnesium stearate, sodium stearate, glycerol monostearate, sodium chloride, dried skim milk and the like. Liquid and semisolid excipients

may be selected from glycerol, propylene glycol, water, ethanol and various oils, including those of petroleum, animal, vegetable or synthetic origin, *e.g.*, peanut oil, soybean oil, mineral oil, sesame oil, etc. Preferred liquid carriers, particularly for injectable solutions, include water, saline, aqueous dextrose, and glycols.

5 Compressed gases may be used to disperse a compound of this invention in aerosol form. Inert gases suitable for this purpose are nitrogen, carbon dioxide, etc.

Other suitable pharmaceutical excipients and their formulations are described in Remington's Pharmaceutical Sciences, edited by E. W. Martin (Mack Publishing Company, 18th ed., 1990).

10 The amount of the therapeutic compound in a formulation can vary within the full range employed by those skilled in the art. Typically, the formulation will contain, on a weight percent (wt %) basis, from about 0.01-99.99 wt % of the therapeutic compound based on the total formulation, with the balance being one or more suitable pharmaceutical excipients. Preferably, the compound is present at a level of about 1-80 wt %.

15 Therapeutic compounds can be administered alone or in combination with other therapeutic compounds or in combination with one or more other active ingredient(s). For example, an inhibitor or stimulator of an VT-associated protein can be administered in combination with another agent that inhibits or stimulates the activity of the same or a different VT-associated protein to thereby counteract the affects of VT.

20 For further information regarding pharmacology, see *Current Protocols in Pharmacology*, John Wiley & Sons, Inc., N.Y.

Human Identification Applications

In addition to their diagnostic and therapeutic uses in VT and related pathologies, the
25 SNPs provided by the present invention are also useful as human identification markers for such applications as forensics, paternity testing, and biometrics (see, *e.g.*, Gill, "An assessment of the utility of single nucleotide polymorphisms (SNPs) for forensic purposes", *Int J Legal Med.* 2001;114(4-5):204-10). Genetic variations in the nucleic acid sequences between individuals can be used as genetic markers to identify individuals and to associate a biological sample with an
30 individual. Determination of which nucleotides occupy a set of SNP positions in an individual identifies a set of SNP markers that distinguishes the individual. The more SNP positions that are analyzed, the lower the probability that the set of SNPs in one individual is the same as that in an unrelated individual. Preferably, if multiple sites are analyzed, the sites are unlinked (*i.e.*, inherited independently). Thus, preferred sets of SNPs can be selected from among the SNPs

disclosed herein, which may include SNPs on different chromosomes, SNPs on different chromosome arms, and/or SNPs that are dispersed over substantial distances along the same chromosome arm.

Furthermore, among the SNPs disclosed herein, preferred SNPs for use in certain forensic/human identification applications include SNPs located at degenerate codon positions (i.e., the third position in certain codons which can be one of two or more alternative nucleotides and still encode the same amino acid), since these SNPs do not affect the encoded protein. SNPs that do not affect the encoded protein are expected to be under less selective pressure and are therefore expected to be more polymorphic in a population, which is typically an advantage for forensic/human identification applications. However, for certain forensics/human identification applications, such as predicting phenotypic characteristics (e.g., inferring ancestry or inferring one or more physical characteristics of an individual) from a DNA sample, it may be desirable to utilize SNPs that affect the encoded protein.

For many of the SNPs disclosed in Tables 1-2 (which are identified as “Applera” SNP source), Tables 1-2 provide SNP allele frequencies obtained by re-sequencing the DNA of chromosomes from 39 individuals (Tables 1-2 also provide allele frequency information for “Celera” source SNPs and, where available, public SNPs from dbEST, HGBASE, and/or HGMD). The allele frequencies provided in Tables 1-2 enable these SNPs to be readily used for human identification applications. Although any SNP disclosed in Table 1 and/or Table 2 could be used for human identification, the closer that the frequency of the minor allele at a particular SNP site is to 50%, the greater the ability of that SNP to discriminate between different individuals in a population since it becomes increasingly likely that two randomly selected individuals would have different alleles at that SNP site. Using the SNP allele frequencies provided in Tables 1-2, one of ordinary skill in the art could readily select a subset of SNPs for which the frequency of the minor allele is, for example, at least 1%, 2%, 5%, 10%, 20%, 25%, 30%, 40%, 45%, or 50%, or any other frequency in-between. Thus, since Tables 1-2 provide allele frequencies based on the re-sequencing of the chromosomes from 39 individuals, a subset of SNPs could readily be selected for human identification in which the total allele count of the minor allele at a particular SNP site is, for example, at least 1, 2, 4, 8, 10, 16, 20, 24, 30, 32, 36, 38, 39, 40, or any other number in-between.

Furthermore, Tables 1-2 also provide population group (interchangeably referred to herein as ethnic or racial groups) information coupled with the extensive allele frequency information. For example, the group of 39 individuals whose DNA was re-sequenced was made-up of 20 Caucasians and 19 African-Americans. This population group information enables

further refinement of SNP selection for human identification. For example, preferred SNPs for human identification can be selected from Tables 1-2 that have similar allele frequencies in both the Caucasian and African-American populations; thus, for example, SNPs can be selected that have equally high discriminatory power in both populations. Alternatively, SNPs can be selected for which there is a statistically significant difference in allele frequencies between the Caucasian and African-American populations (as an extreme example, a particular allele may be observed only in either the Caucasian or the African-American population group but not observed in the other population group); such SNPs are useful, for example, for predicting the race/ethnicity of an unknown perpetrator from a biological sample such as a hair or blood stain recovered at a crime scene. For a discussion of using SNPs to predict ancestry from a DNA sample, including statistical methods, see Frudakis *et al.*, "A Classifier for the SNP-Based Inference of Ancestry," *Journal of Forensic Sciences* 2003; 48(4):771-782.

SNPs have numerous advantages over other types of polymorphic markers, such as short tandem repeats (STRs). For example, SNPs can be easily scored and are amenable to automation, making SNPs the markers of choice for large-scale forensic databases. SNPs are found in much greater abundance throughout the genome than repeat polymorphisms. Population frequencies of two polymorphic forms can usually be determined with greater accuracy than those of multiple polymorphic forms at multi-allelic loci. SNPs are mutationally more stable than repeat polymorphisms. SNPs are not susceptible to artefacts such as stutter bands that can hinder analysis. Stutter bands are frequently encountered when analyzing repeat polymorphisms, and are particularly troublesome when analyzing samples such as crime scene samples that may contain mixtures of DNA from multiple sources. Another significant advantage of SNP markers over STR markers is the much shorter length of nucleic acid needed to score a SNP. For example, STR markers are generally several hundred base pairs in length. A SNP, on the other hand, comprises a single nucleotide, and generally a short conserved region on either side of the SNP position for primer and/or probe binding. This makes SNPs more amenable to typing in highly degraded or aged biological samples that are frequently encountered in forensic casework in which DNA may be fragmented into short pieces.

SNPs also are not subject to microvariant and "off-ladder" alleles frequently encountered when analyzing STR loci. Microvariants are deletions or insertions within a repeat unit that change the size of the amplified DNA product so that the amplified product does not migrate at the same rate as reference alleles with normal sized repeat units. When separated by size, such as by electrophoresis on a polyacrylamide gel, microvariants do not align with a reference allelic ladder of standard sized repeat units, but rather migrate between the reference alleles. The

reference allelic ladder is used for precise sizing of alleles for allele classification; therefore alleles that do not align with the reference allelic ladder lead to substantial analysis problems. Furthermore, when analyzing multi-allelic repeat polymorphisms, occasionally an allele is found that consists of more or less repeat units than has been previously seen in the population, or more
5 or less repeat alleles than are included in a reference allelic ladder. These alleles will migrate outside the size range of known alleles in a reference allelic ladder, and therefore are referred to as “off-ladder” alleles. In extreme cases, the allele may contain so few or so many repeats that it migrates well out of the range of the reference allelic ladder. In this situation, the allele may not even be observed, or, with multiplex analysis, it may migrate within or close to the size range for
10 another locus, further confounding analysis.

SNP analysis avoids the problems of microvariants and off-ladder alleles encountered in STR analysis. Importantly, microvariants and off-ladder alleles may provide significant problems, and may be completely missed, when using analysis methods such as oligonucleotide hybridization arrays, which utilize oligonucleotide probes specific for certain known alleles.
15 Furthermore, off-ladder alleles and microvariants encountered with STR analysis, even when correctly typed, may lead to improper statistical analysis, since their frequencies in the population are generally unknown or poorly characterized, and therefore the statistical significance of a matching genotype may be questionable. All these advantages of SNP analysis are considerable in light of the consequences of most DNA identification cases, which may lead
20 to life imprisonment for an individual, or re-association of remains to the family of a deceased individual.

DNA can be isolated from biological samples such as blood, bone, hair, saliva, or semen, and compared with the DNA from a reference source at particular SNP positions. Multiple SNP markers can be assayed simultaneously in order to increase the power of discrimination and the
25 statistical significance of a matching genotype. For example, oligonucleotide arrays can be used to genotype a large number of SNPs simultaneously. The SNPs provided by the present invention can be assayed in combination with other polymorphic genetic markers, such as other SNPs known in the art or STRs, in order to identify an individual or to associate an individual with a particular biological sample.

30 Furthermore, the SNPs provided by the present invention can be genotyped for inclusion in a database of DNA genotypes, for example, a criminal DNA databank such as the FBI’s Combined DNA Index System (CODIS) database. A genotype obtained from a biological sample of unknown source can then be queried against the database to find a matching genotype, with the SNPs of the present invention providing nucleotide positions at which to compare the

known and unknown DNA sequences for identity. Accordingly, the present invention provides a database comprising novel SNPs or SNP alleles of the present invention (*e.g.*, the database can comprise information indicating which alleles are possessed by individual members of a population at one or more novel SNP sites of the present invention), such as for use in forensics, biometrics, or other human identification applications. Such a database typically comprises a computer-based system in which the SNPs or SNP alleles of the present invention are recorded on a computer readable medium (see the section of the present specification entitled “Computer-Related Embodiments”).

The SNPs of the present invention can also be assayed for use in paternity testing. The object of paternity testing is usually to determine whether a male is the father of a child. In most cases, the mother of the child is known and thus, the mother's contribution to the child's genotype can be traced. Paternity testing investigates whether the part of the child's genotype not attributable to the mother is consistent with that of the putative father. Paternity testing can be performed by analyzing sets of polymorphisms in the putative father and the child, with the SNPs of the present invention providing nucleotide positions at which to compare the putative father's and child's DNA sequences for identity. If the set of polymorphisms in the child attributable to the father does not match the set of polymorphisms of the putative father, it can be concluded, barring experimental error, that the putative father is not the father of the child. If the set of polymorphisms in the child attributable to the father match the set of polymorphisms of the putative father, a statistical calculation can be performed to determine the probability of coincidental match, and a conclusion drawn as to the likelihood that the putative father is the true biological father of the child.

In addition to paternity testing, SNPs are also useful for other types of kinship testing, such as for verifying familial relationships for immigration purposes, or for cases in which an individual alleges to be related to a deceased individual in order to claim an inheritance from the deceased individual, etc. For further information regarding the utility of SNPs for paternity testing and other types of kinship testing, including methods for statistical analysis, see Krawczak, “Informativity assessment for biallelic single nucleotide polymorphisms”, *Electrophoresis* 1999 Jun;20(8):1676-81.

The use of the SNPs of the present invention for human identification further extends to various authentication systems, commonly referred to as biometric systems, which typically convert physical characteristics of humans (or other organisms) into digital data. Biometric systems include various technological devices that measure such unique anatomical or physiological characteristics as finger, thumb, or palm prints; hand geometry; vein patterning on the back of the hand; blood

vessel patterning of the retina and color and texture of the iris; facial characteristics; voice patterns; signature and typing dynamics; and DNA. Such physiological measurements can be used to verify identity and, for example, restrict or allow access based on the identification. Examples of applications for biometrics include physical area security, computer and network security, aircraft passenger check-in and boarding, financial transactions, medical records access, government benefit distribution, voting, law enforcement, passports, visas and immigration, prisons, various military applications, and for restricting access to expensive or dangerous items, such as automobiles or guns (see, for example, O'Connor, *Stanford Technology Law Review* and U.S. Patent No. 6,119,096).

Groups of SNPs, particularly the SNPs provided by the present invention, can be typed to uniquely identify an individual for biometric applications such as those described above. Such SNP typing can readily be accomplished using, for example, DNA chips/arrays. Preferably, a minimally invasive means for obtaining a DNA sample is utilized. For example, PCR amplification enables sufficient quantities of DNA for analysis to be obtained from buccal swabs or fingerprints, which contain DNA-containing skin cells and oils that are naturally transferred during contact.

Further information regarding techniques for using SNPs in forensic/human identification applications can be found in, for example, *Current Protocols in Human Genetics*, John Wiley & Sons, N.Y. (2002), 14.1-14.7.

VARIANT PROTEINS, ANTIBODIES, VECTORS, HOST CELLS, & USES THEREOF

Variant Proteins Encoded by SNP-Containing Nucleic Acid Molecules

The present invention provides SNP-containing nucleic acid molecules, many of which encode proteins having variant amino acid sequences as compared to the art-known (*i.e.*, wild-type) proteins. Amino acid sequences encoded by the polymorphic nucleic acid molecules of the present invention are referred to as SEQ ID NOS: 21-40 in Table 1 and provided in the Sequence Listing. These variants will generally be referred to herein as variant proteins/peptides/polypeptides, or polymorphic proteins/peptides/polypeptides of the present invention. The terms "protein," "peptide," and "polypeptide" are used herein interchangeably.

A variant protein of the present invention may be encoded by, for example, a nonsynonymous nucleotide substitution at any one of the cSNP positions disclosed herein. In addition, variant proteins may also include proteins whose expression, structure, and/or function is altered by a SNP disclosed herein, such as a SNP that creates or destroys a stop codon, a SNP that affects splicing, and a SNP in control/regulatory elements, *e.g.* promoters, enhancers, or transcription factor binding domains.

As used herein, a protein or peptide is said to be "isolated" or "purified" when it is substantially free of cellular material or chemical precursors or other chemicals. The variant proteins of the present invention can be purified to homogeneity or other lower degrees of purity. The level of purification will be based on the intended use. The key feature is that the preparation
5 allows for the desired function of the variant protein, even if in the presence of considerable amounts of other components.

As used herein, "substantially free of cellular material" includes preparations of the variant protein having less than about 30% (by dry weight) other proteins (*i.e.*, contaminating protein), less than about 20% other proteins, less than about 10% other proteins, or less than about 5% other
10 proteins. When the variant protein is recombinantly produced, it can also be substantially free of culture medium, *i.e.*, culture medium represents less than about 20% of the volume of the protein preparation.

The language "substantially free of chemical precursors or other chemicals" includes preparations of the variant protein in which it is separated from chemical precursors or other
15 chemicals that are involved in its synthesis. In one embodiment, the language "substantially free of chemical precursors or other chemicals" includes preparations of the variant protein having less than about 30% (by dry weight) chemical precursors or other chemicals, less than about 20% chemical precursors or other chemicals, less than about 10% chemical precursors or other chemicals, or less than about 5% chemical precursors or other chemicals.

20 An isolated variant protein may be purified from cells that naturally express it, purified from cells that have been altered to express it (recombinant host cells), or synthesized using known protein synthesis methods. For example, a nucleic acid molecule containing SNP(s) encoding the variant protein can be cloned into an expression vector, the expression vector introduced into a host cell, and the variant protein expressed in the host cell. The variant protein can then be isolated from
25 the cells by any appropriate purification scheme using standard protein purification techniques. Examples of these techniques are described in detail below (Sambrook and Russell, 2000, Molecular Cloning: A Laboratory Manual, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY).

The present invention provides isolated variant proteins that comprise, consist of or consist essentially of amino acid sequences that contain one or more variant amino acids encoded
30 by one or more codons that contain a SNP of the present invention.

Accordingly, the present invention provides variant proteins that consist of amino acid sequences that contain one or more amino acid polymorphisms (or truncations or extensions due to creation or destruction of a stop codon, respectively) encoded by the SNPs provided in Table 1

and/or Table 2. A protein consists of an amino acid sequence when the amino acid sequence is the entire amino acid sequence of the protein.

The present invention further provides variant proteins that consist essentially of amino acid sequences that contain one or more amino acid polymorphisms (or truncations or extensions due to
5 creation or destruction of a stop codon, respectively) encoded by the SNPs provided in Table 1 and/or Table 2. A protein consists essentially of an amino acid sequence when such an amino acid sequence is present with only a few additional amino acid residues in the final protein.

The present invention further provides variant proteins that comprise amino acid sequences that contain one or more amino acid polymorphisms (or truncations or extensions due to creation or
10 destruction of a stop codon, respectively) encoded by the SNPs provided in Table 1 and/or Table 2. A protein comprises an amino acid sequence when the amino acid sequence is at least part of the final amino acid sequence of the protein. In such a fashion, the protein may contain only the variant amino acid sequence or have additional amino acid residues, such as a contiguous encoded sequence that is naturally associated with it or heterologous amino acid residues. Such a protein can have a
15 few additional amino acid residues or can comprise many more additional amino acids. A brief description of how various types of these proteins can be made and isolated is provided below.

The variant proteins of the present invention can be attached to heterologous sequences to form chimeric or fusion proteins. Such chimeric and fusion proteins comprise a variant protein operatively linked to a heterologous protein having an amino acid sequence not substantially
20 homologous to the variant protein. "Operatively linked" indicates that the coding sequences for the variant protein and the heterologous protein are ligated in-frame. The heterologous protein can be fused to the N-terminus or C-terminus of the variant protein. In another embodiment, the fusion protein is encoded by a fusion polynucleotide that is synthesized by conventional techniques including automated DNA synthesizers. Alternatively, PCR amplification of gene
25 fragments can be carried out using anchor primers which give rise to complementary overhangs between two consecutive gene fragments which can subsequently be annealed and re-amplified to generate a chimeric gene sequence (see Ausubel *et al.*, *Current Protocols in Molecular Biology*, 1992). Moreover, many expression vectors are commercially available that already encode a fusion moiety (*e.g.*, a GST protein). A variant protein-encoding nucleic acid can be
30 cloned into such an expression vector such that the fusion moiety is linked in-frame to the variant protein.

In many uses, the fusion protein does not affect the activity of the variant protein. The fusion protein can include, but is not limited to, enzymatic fusion proteins, for example, beta-galactosidase fusions, yeast two-hybrid GAL fusions, poly-His fusions, MYC-tagged, HI-tagged

and Ig fusions. Such fusion proteins, particularly poly-His fusions, can facilitate their purification following recombinant expression. In certain host cells (*e.g.*, mammalian host cells), expression and/or secretion of a protein can be increased by using a heterologous signal sequence. Fusion proteins are further described in, for example, Terpe, "Overview of tag protein fusions: from
5 molecular and biochemical fundamentals to commercial systems", *Appl Microbiol Biotechnol.* 2003 Jan;60(5):523-33. Epub 2002 Nov 07; Graddis *et al.*, "Designing proteins that work using recombinant technologies", *Curr Pharm Biotechnol.* 2002 Dec;3(4):285-97; and Nilsson *et al.*, "Affinity fusion strategies for detection, purification, and immobilization of recombinant proteins", *Protein Expr Purif.* 1997 Oct;11(1):1-16.

10 The present invention also relates to further obvious variants of the variant polypeptides of the present invention, such as naturally-occurring mature forms (*e.g.*, allelic variants), non-naturally occurring recombinantly-derived variants, and orthologs and paralogs of such proteins that share sequence homology. Such variants can readily be generated using art-known techniques in the fields of recombinant nucleic acid technology and protein biochemistry. It is understood, however,
15 that variants exclude those known in the prior art before the present invention.

Further variants of the variant polypeptides disclosed in Table 1 can comprise an amino acid sequence that shares at least 70-80%, 80-85%, 85-90%, 91%, 92%, 93%, 94%, 95%, 96%, 97%, 98%, or 99% sequence identity with an amino acid sequence disclosed in Table 1 (or a fragment thereof) and that includes a novel amino acid residue (allele) disclosed in Table 1
20 (which is encoded by a novel SNP allele). Thus, an aspect of the present invention that is specifically contemplated are polypeptides that have a certain degree of sequence variation compared with the polypeptide sequences shown in Table 1, but that contain a novel amino acid residue (allele) encoded by a novel SNP allele disclosed herein. In other words, as long as a polypeptide contains a novel amino acid residue disclosed herein, other portions of the
25 polypeptide that flank the novel amino acid residue can vary to some degree from the polypeptide sequences shown in Table 1.

Full-length pre-processed forms, as well as mature processed forms, of proteins that comprise one of the amino acid sequences disclosed herein can readily be identified as having complete sequence identity to one of the variant proteins of the present invention as well as being
30 encoded by the same genetic locus as the variant proteins provided herein.

Orthologs of a variant peptide can readily be identified as having some degree of significant sequence homology/identity to at least a portion of a variant peptide as well as being encoded by a gene from another organism. Preferred orthologs will be isolated from non-human mammals, preferably primates, for the development of human therapeutic targets and agents. Such orthologs

can be encoded by a nucleic acid sequence that hybridizes to a variant peptide-encoding nucleic acid molecule under moderate to stringent conditions depending on the degree of relatedness of the two organisms yielding the homologous proteins.

Variant proteins include, but are not limited to, proteins containing deletions, additions and
5 substitutions in the amino acid sequence caused by the SNPs of the present invention. One class of substitutions is conserved amino acid substitutions in which a given amino acid in a polypeptide is substituted for another amino acid of like characteristics. Typical conservative substitutions are replacements, one for another, among the aliphatic amino acids Ala, Val, Leu, and Ile; interchange
10 of the hydroxyl residues Ser and Thr; exchange of the acidic residues Asp and Glu; substitution between the amide residues Asn and Gln; exchange of the basic residues Lys and Arg; and replacements among the aromatic residues Phe and Tyr. Guidance concerning which amino acid changes are likely to be phenotypically silent are found in, for example, Bowie *et al.*, *Science* 247:1306-1310 (1990).

Variant proteins can be fully functional or can lack function in one or more activities, *e.g.*
15 ability to bind another molecule, ability to catalyze a substrate, ability to mediate signaling, etc. Fully functional variants typically contain only conservative variations or variations in non-critical residues or in non-critical regions. Functional variants can also contain substitution of similar amino acids that result in no change or an insignificant change in function. Alternatively, such substitutions may positively or negatively affect function to some degree. Non-functional
20 variants typically contain one or more non-conservative amino acid substitutions, deletions, insertions, inversions, truncations or extensions, or a substitution, insertion, inversion, or deletion of a critical residue or in a critical region.

Amino acids that are essential for function of a protein can be identified by methods known in the art, such as site-directed mutagenesis or alanine-scanning mutagenesis (Cunningham *et al.*,
25 *Science* 244:1081-1085 (1989)), particularly using the amino acid sequence and polymorphism information provided in Table 1. The latter procedure introduces single alanine mutations at every residue in the molecule. The resulting mutant molecules are then tested for biological activity such as enzyme activity or in assays such as an *in vitro* proliferative activity. Sites that are critical for binding partner/substrate binding can also be determined by structural analysis such as
30 crystallization, nuclear magnetic resonance or photoaffinity labeling (Smith *et al.*, *J. Mol. Biol.* 224:899-904 (1992); de Vos *et al.* *Science* 255:306-312 (1992)).

Polypeptides can contain amino acids other than the 20 amino acids commonly referred to as the 20 naturally occurring amino acids. Further, many amino acids, including the terminal amino acids, may be modified by natural processes, such as processing and other post-

translational modifications, or by chemical modification techniques well known in the art.

Accordingly, the variant proteins of the present invention also encompass derivatives or analogs in which a substituted amino acid residue is not one encoded by the genetic code, in which a substituent group is included, in which the mature polypeptide is fused with another compound, such as a compound to increase the half-life of the polypeptide (*e.g.*, polyethylene glycol), or in which additional amino acids are fused to the mature polypeptide, such as a leader or secretory sequence or a sequence for purification of the mature polypeptide or a pro-protein sequence.

Known protein modifications include, but are not limited to, acetylation, acylation, ADP-ribosylation, amidation, covalent attachment of flavin, covalent attachment of a heme moiety, covalent attachment of a nucleotide or nucleotide derivative, covalent attachment of a lipid or lipid derivative, covalent attachment of phosphatidylinositol, cross-linking, cyclization, disulfide bond formation, demethylation, formation of covalent crosslinks, formation of cystine, formation of pyroglutamate, formylation, gamma carboxylation, glycosylation, GPI anchor formation, hydroxylation, iodination, methylation, myristoylation, oxidation, proteolytic processing, phosphorylation, prenylation, racemization, selenoylation, sulfation, transfer-RNA mediated addition of amino acids to proteins such as arginylation, and ubiquitination.

Such protein modifications are well known to those of skill in the art and have been described in great detail in the scientific literature. Several particularly common modifications, glycosylation, lipid attachment, sulfation, gamma-carboxylation of glutamic acid residues, hydroxylation and ADP-ribosylation, for instance, are described in most basic texts, such as *Proteins - Structure and Molecular Properties*, 2nd Ed., T.E. Creighton, W. H. Freeman and Company, New York (1993); Wold, F., *Posttranslational Covalent Modification of Proteins*, B.C. Johnson, Ed., Academic Press, New York 1-12 (1983); Seifter *et al.*, *Meth. Enzymol.* 182: 626-646 (1990); and Rattan *et al.*, *Ann. N.Y. Acad. Sci.* 663:48-62 (1992).

The present invention further provides fragments of the variant proteins in which the fragments contain one or more amino acid sequence variations (*e.g.*, substitutions, or truncations or extensions due to creation or destruction of a stop codon) encoded by one or more SNPs disclosed herein. The fragments to which the invention pertains, however, are not to be construed as encompassing fragments that have been disclosed in the prior art before the present invention.

As used herein, a fragment may comprise at least about 4, 8, 10, 12, 14, 16, 18, 20, 25, 30, 50, 100 (or any other number in-between) or more contiguous amino acid residues from a variant protein, wherein at least one amino acid residue is affected by a SNP of the present invention, *e.g.*, a variant amino acid residue encoded by a nonsynonymous nucleotide substitution at a cSNP position provided by the present invention. The variant amino acid encoded by a cSNP may occupy any

residue position along the sequence of the fragment. Such fragments can be chosen based on the ability to retain one or more of the biological activities of the variant protein or the ability to perform a function, *e.g.*, act as an immunogen. Particularly important fragments are biologically active fragments. Such fragments will typically comprise a domain or motif of a variant protein of the present invention, *e.g.*, active site, transmembrane domain, or ligand/substrate binding domain. Other fragments include, but are not limited to, domain or motif-containing fragments, soluble peptide fragments, and fragments containing immunogenic structures. Predicted domains and functional sites are readily identifiable by computer programs well known to those of skill in the art (*e.g.*, PROSITE analysis) (*Current Protocols in Protein Science*, John Wiley & Sons, N.Y. (2002)).

10

Uses of Variant Proteins

The variant proteins of the present invention can be used in a variety of ways, including but not limited to, in assays to determine the biological activity of a variant protein, such as in a panel of multiple proteins for high-throughput screening; to raise antibodies or to elicit another type of immune response; as a reagent (including the labeled reagent) in assays designed to quantitatively determine levels of the variant protein (or its binding partner) in biological fluids; as a marker for cells or tissues in which it is preferentially expressed (either constitutively or at a particular stage of tissue differentiation or development or in a disease state); as a target for screening for a therapeutic agent; and as a direct therapeutic agent to be administered into a human subject. Any of the variant proteins disclosed herein may be developed into reagent grade or kit format for commercialization as research products. Methods for performing the uses listed above are well known to those skilled in the art (see, *e.g.*, *Molecular Cloning: A Laboratory Manual*, Cold Spring Harbor Laboratory Press, Sambrook and Russell, 2000, and *Methods in Enzymology: Guide to Molecular Cloning Techniques*, Academic Press, Berger, S. L. and A. R. Kimmel eds., 1987).

In a specific embodiment of the invention, the methods of the present invention include detection of one or more variant proteins disclosed herein. Variant proteins are disclosed in Table 1 and in the Sequence Listing as SEQ ID NOS: 21-40. Detection of such proteins can be accomplished using, for example, antibodies, small molecule compounds, aptamers, ligands/substrates, other proteins or protein fragments, or other protein-binding agents. Preferably, protein detection agents are specific for a variant protein of the present invention and can therefore discriminate between a variant protein of the present invention and the wild-type protein or another variant form. This can generally be accomplished by, for example, selecting or designing detection agents that bind to the region of a protein that differs between the variant

30

and wild-type protein, such as a region of a protein that contains one or more amino acid substitutions that is/are encoded by a non-synonymous cSNP of the present invention, or a region of a protein that follows a nonsense mutation-type SNP that creates a stop codon thereby leading to a shorter polypeptide, or a region of a protein that follows a read-through mutation-type SNP that destroys a stop codon thereby leading to a longer polypeptide in which a portion of the polypeptide is present in one version of the polypeptide but not the other.

In another specific aspect of the invention, the variant proteins of the present invention are used as targets for diagnosing VT or for determining predisposition to VT in a human. Accordingly, the invention provides methods for detecting the presence of, or levels of, one or more variant proteins of the present invention in a cell, tissue, or organism. Such methods typically involve contacting a test sample with an agent (*e.g.*, an antibody, small molecule compound, or peptide) capable of interacting with the variant protein such that specific binding of the agent to the variant protein can be detected. Such an assay can be provided in a single detection format or a multi-detection format such as an array, for example, an antibody or aptamer array (arrays for protein detection may also be referred to as “protein chips”). The variant protein of interest can be isolated from a test sample and assayed for the presence of a variant amino acid sequence encoded by one or more SNPs disclosed by the present invention. The SNPs may cause changes to the protein and the corresponding protein function/activity, such as through non-synonymous substitutions in protein coding regions that can lead to amino acid substitutions, deletions, insertions, and/or rearrangements; formation or destruction of stop codons; or alteration of control elements such as promoters. SNPs may also cause inappropriate post-translational modifications.

One preferred agent for detecting a variant protein in a sample is an antibody capable of selectively binding to a variant form of the protein (antibodies are described in greater detail in the next section). Such samples include, for example, tissues, cells, and biological fluids isolated from a subject, as well as tissues, cells and fluids present within a subject.

In vitro methods for detection of the variant proteins associated with VT that are disclosed herein and fragments thereof include, but are not limited to, enzyme linked immunosorbent assays (ELISAs), radioimmunoassays (RIA), Western blots, immunoprecipitations, immunofluorescence, and protein arrays/chips (*e.g.*, arrays of antibodies or aptamers). For further information regarding immunoassays and related protein detection methods, see *Current Protocols in Immunology*, John Wiley & Sons, N.Y., and Hage, “Immunoassays”, *Anal Chem.* 1999 Jun 15;71(12):294R-304R.

Additional analytic methods of detecting amino acid variants include, but are not limited to, altered electrophoretic mobility, altered tryptic peptide digest, altered protein activity in cell-based

or cell-free assay, alteration in ligand or antibody-binding pattern, altered isoelectric point, and direct amino acid sequencing.

Alternatively, variant proteins can be detected *in vivo* in a subject by introducing into the subject a labeled antibody (or other type of detection reagent) specific for a variant protein. For example, the antibody can be labeled with a radioactive marker whose presence and location in a
5 subject can be detected by standard imaging techniques.

Other uses of the variant peptides of the present invention are based on the class or action of the protein. For example, proteins isolated from humans and their mammalian orthologs serve as targets for identifying agents (*e.g.*, small molecule drugs or antibodies) for use in therapeutic
10 applications, particularly for modulating a biological or pathological response in a cell or tissue that expresses the protein. Pharmaceutical agents can be developed that modulate protein activity.

As an alternative to modulating gene expression, therapeutic compounds can be developed that modulate protein function. For example, many SNPs disclosed herein affect the amino acid
15 sequence of the encoded protein (*e.g.*, non-synonymous cSNPs and nonsense mutation-type SNPs). Such alterations in the encoded amino acid sequence may affect protein function, particularly if such amino acid sequence variations occur in functional protein domains, such as catalytic domains, ATP-binding domains, or ligand/substrate binding domains. It is well established in the art that variant proteins having amino acid sequence variations in functional domains can cause or influence
20 pathological conditions. In such instances, compounds (*e.g.*, small molecule drugs or antibodies) can be developed that target the variant protein and modulate (*e.g.*, up- or down-regulate) protein function/activity.

The therapeutic methods of the present invention further include methods that target one or more variant proteins of the present invention. Variant proteins can be targeted using, for
25 example, small molecule compounds, antibodies, aptamers, ligands/substrates, other proteins, or other protein-binding agents. Additionally, the skilled artisan will recognize that the novel protein variants (and polymorphic nucleic acid molecules) disclosed in Table 1 may themselves be directly used as therapeutic agents by acting as competitive inhibitors of corresponding art-known proteins (or nucleic acid molecules such as mRNA molecules).

30 The variant proteins of the present invention are particularly useful in drug screening assays, in cell-based or cell-free systems. Cell-based systems can utilize cells that naturally express the protein, a biopsy specimen, or cell cultures. In one embodiment, cell-based assays involve recombinant host cells expressing the variant protein. Cell-free assays can be used to detect the

ability of a compound to directly bind to a variant protein or to the corresponding SNP-containing nucleic acid fragment that encodes the variant protein.

A variant protein of the present invention, as well as appropriate fragments thereof, can be used in high-throughput screening assays to test candidate compounds for the ability to bind and/or modulate the activity of the variant protein. These candidate compounds can be further screened against a protein having normal function (*e.g.*, a wild-type/non-variant protein) to further determine the effect of the compound on the protein activity. Furthermore, these compounds can be tested in animal or invertebrate systems to determine *in vivo* activity/effectiveness. Compounds can be identified that activate (agonists) or inactivate (antagonists) the variant protein, and different compounds can be identified that cause various degrees of activation or inactivation of the variant protein.

Further, the variant proteins can be used to screen a compound for the ability to stimulate or inhibit interaction between the variant protein and a target molecule that normally interacts with the protein. The target can be a ligand, a substrate or a binding partner that the protein normally interacts with (for example, epinephrine or norepinephrine). Such assays typically include the steps of combining the variant protein with a candidate compound under conditions that allow the variant protein, or fragment thereof, to interact with the target molecule, and to detect the formation of a complex between the protein and the target or to detect the biochemical consequence of the interaction with the variant protein and the target, such as any of the associated effects of signal transduction.

Candidate compounds include, for example, 1) peptides such as soluble peptides, including Ig-tailed fusion peptides and members of random peptide libraries (see, *e.g.*, Lam *et al.*, *Nature* 354:82-84 (1991); Houghten *et al.*, *Nature* 354:84-86 (1991)) and combinatorial chemistry-derived molecular libraries made of D- and/or L- configuration amino acids; 2) phosphopeptides (*e.g.*, members of random and partially degenerate, directed phosphopeptide libraries, see, *e.g.*, Songyang *et al.*, *Cell* 72:767-778 (1993)); 3) antibodies (*e.g.*, polyclonal, monoclonal, humanized, anti-idiotypic, chimeric, and single chain antibodies as well as Fab, F(ab')₂, Fab expression library fragments, and epitope-binding fragments of antibodies); and 4) small organic and inorganic molecules (*e.g.*, molecules obtained from combinatorial and natural product libraries).

One candidate compound is a soluble fragment of the variant protein that competes for ligand binding. Other candidate compounds include mutant proteins or appropriate fragments containing mutations that affect variant protein function and thus compete for ligand. Accordingly, a fragment that competes for ligand, for example with a higher affinity, or a fragment that binds ligand but does not allow release, is encompassed by the invention.

The invention further includes other end point assays to identify compounds that modulate (stimulate or inhibit) variant protein activity. The assays typically involve an assay of events in the signal transduction pathway that indicate protein activity. Thus, the expression of genes that are up or down-regulated in response to the variant protein dependent signal cascade can be assayed. In one embodiment, the regulatory region of such genes can be operably linked to a marker that is easily detectable, such as luciferase. Alternatively, phosphorylation of the variant protein, or a variant protein target, could also be measured. Any of the biological or biochemical functions mediated by the variant protein can be used as an endpoint assay. These include all of the biochemical or biological events described herein, and other functions known to those of ordinary skill in the art.

Binding and/or activating compounds can also be screened by using chimeric variant proteins in which an amino terminal extracellular domain or parts thereof, an entire transmembrane domain or subregions, and/or the carboxyl terminal intracellular domain or parts thereof, can be replaced by heterologous domains or subregions. For example, a substrate-binding region can be used that interacts with a different substrate than that which is normally recognized by a variant protein. Accordingly, a different set of signal transduction components is available as an end-point assay for activation. This allows for assays to be performed in other than the specific host cell from which the variant protein is derived.

The variant proteins are also useful in competition binding assays in methods designed to discover compounds that interact with the variant protein. Thus, a compound can be exposed to a variant protein under conditions that allow the compound to bind or to otherwise interact with the variant protein. A binding partner, such as ligand, that normally interacts with the variant protein is also added to the mixture. If the test compound interacts with the variant protein or its binding partner, it decreases the amount of complex formed or activity from the variant protein. This type of assay is particularly useful in screening for compounds that interact with specific regions of the variant protein (Hodgson, *Bio/technology*, 1992, Sept 10(9), 973-80).

To perform cell-free drug screening assays, it is sometimes desirable to immobilize either the variant protein or a fragment thereof, or its target molecule, to facilitate separation of complexes from uncomplexed forms of one or both of the proteins, as well as to accommodate automation of the assay. Any method for immobilizing proteins on matrices can be used in drug screening assays. In one embodiment, a fusion protein containing an added domain allows the protein to be bound to a matrix. For example, glutathione-S-transferase/¹²⁵I fusion proteins can be adsorbed onto glutathione sepharose™ beads (Sigma Chemical, St. Louis, MO) or glutathione derivatized microtitre plates,

which are then combined with the cell lysates (*e.g.*, ³⁵S-labeled) and a candidate compound, such as a drug candidate, and the mixture incubated under conditions conducive to complex formation (*e.g.*, at physiological conditions for salt and pH). Following incubation, the beads can be washed to remove any unbound label, and the matrix immobilized and radiolabel determined directly, or in the supernatant after the complexes are dissociated. Alternatively, the complexes can be dissociated from the matrix, separated by SDS-PAGE, and the level of bound material found in the bead fraction quantitated from the gel using standard electrophoretic techniques.

Either the variant protein or its target molecule can be immobilized utilizing conjugation of biotin and streptavidin. Alternatively, antibodies reactive with the variant protein but which do not interfere with binding of the variant protein to its target molecule can be derivatized to the wells of the plate, and the variant protein trapped in the wells by antibody conjugation. Preparations of the target molecule and a candidate compound are incubated in the variant protein-presenting wells and the amount of complex trapped in the well can be quantitated. Methods for detecting such complexes, in addition to those described above for the GST-immobilized complexes, include immunodetection of complexes using antibodies reactive with the protein target molecule, or which are reactive with variant protein and compete with the target molecule, and enzyme-linked assays that rely on detecting an enzymatic activity associated with the target molecule.

Modulators of variant protein activity identified according to these drug screening assays can be used to treat a subject with a disorder mediated by the protein pathway, such as VT. These methods of treatment typically include the steps of administering the modulators of protein activity in a pharmaceutical composition to a subject in need of such treatment.

The variant proteins, or fragments thereof, disclosed herein can themselves be directly used to treat a disorder characterized by an absence of, inappropriate, or unwanted expression or activity of the variant protein. Accordingly, methods for treatment include the use of a variant protein disclosed herein or fragments thereof.

In yet another aspect of the invention, variant proteins can be used as "bait proteins" in a two-hybrid assay or three-hybrid assay (see, *e.g.*, U.S. Patent No. 5,283,317; Zervos *et al.* (1993) *Cell* 72:223-232; Madura *et al.* (1993) *J. Biol. Chem.* 268:12046-12054; Bartel *et al.* (1993) *Biotechniques* 14:920-924; Iwabuchi *et al.* (1993) *Oncogene* 8:1693-1696; and Brent WO94/10300) to identify other proteins that bind to or interact with the variant protein and are involved in variant protein activity. Such variant protein-binding proteins are also likely to be involved in the propagation of signals by the variant proteins or variant protein targets as, for example, elements of a protein-mediated signaling pathway. Alternatively, such variant protein-binding proteins are inhibitors of the variant protein.

The two-hybrid system is based on the modular nature of most transcription factors, which typically consist of separable DNA-binding and activation domains. Briefly, the assay typically utilizes two different DNA constructs. In one construct, the gene that codes for a variant protein is fused to a gene encoding the DNA binding domain of a known transcription factor (e.g., GAL-4). In the other construct, a DNA sequence, from a library of DNA sequences, that encodes an unidentified protein ("prey" or "sample") is fused to a gene that codes for the activation domain of the known transcription factor. If the "bait" and the "prey" proteins are able to interact, *in vivo*, forming a variant protein-dependent complex, the DNA-binding and activation domains of the transcription factor are brought into close proximity. This proximity allows transcription of a reporter gene (e.g., LacZ) that is operably linked to a transcriptional regulatory site responsive to the transcription factor. Expression of the reporter gene can be detected, and cell colonies containing the functional transcription factor can be isolated and used to obtain the cloned gene that encodes the protein that interacts with the variant protein.

15 Antibodies Directed to Variant Proteins

The present invention also provides antibodies that selectively bind to the variant proteins disclosed herein and fragments thereof. Such antibodies may be used to quantitatively or qualitatively detect the variant proteins of the present invention. As used herein, an antibody selectively binds a target variant protein when it binds the variant protein and does not significantly bind to non-variant proteins, *i.e.*, the antibody does not significantly bind to normal, wild-type, or art-known proteins that do not contain a variant amino acid sequence due to one or more SNPs of the present invention (variant amino acid sequences may be due to, for example, nonsynonymous cSNPs, nonsense SNPs that create a stop codon, thereby causing a truncation of a polypeptide or SNPs that cause read-through mutations resulting in an extension of a polypeptide).

25 As used herein, an antibody is defined in terms consistent with that recognized in the art: they are multi-subunit proteins produced by an organism in response to an antigen challenge. The antibodies of the present invention include both monoclonal antibodies and polyclonal antibodies, as well as antigen-reactive proteolytic fragments of such antibodies, such as Fab, F(ab)₂, and Fv fragments. In addition, an antibody of the present invention further includes any of a variety of engineered antigen-binding molecules such as a chimeric antibody (U.S. Patent Nos. 4,816,567 and 4,816,397; Morrison *et al.*, *Proc. Natl. Acad. Sci. USA*, 81:6851, 1984; Neuberger *et al.*, *Nature* 312:604, 1984), a humanized antibody (U.S. Patent Nos. 5,693,762; 5,585,089; and 5,565,332), a single-chain Fv (U.S. Patent No. 4,946,778; Ward *et al.*, *Nature* 334:544, 1989), a bispecific antibody with two binding specificities (Segal *et al.*, *J. Immunol. Methods* 248:1, 2001; Carter, J.

Immunol. Methods 248:7, 2001), a diabody, a triabody, and a tetrabody (Todorovska *et al.*, *J. Immunol. Methods*, 248:47, 2001), as well as a Fab conjugate (dimer or trimer), and a minibody.

5 Many methods are known in the art for generating and/or identifying antibodies to a given target antigen (Harlow, *Antibodies*, Cold Spring Harbor Press, (1989)). In general, an isolated peptide (*e.g.*, a variant protein of the present invention) is used as an immunogen and is administered to a mammalian organism, such as a rat, rabbit, hamster or mouse. Either a full-length protein, an antigenic peptide fragment (*e.g.*, a peptide fragment containing a region that varies between a variant protein and a corresponding wild-type protein), or a fusion protein can be used. A protein used as an immunogen may be naturally-occurring, synthetic or recombinantly produced, and may be administered in combination with an adjuvant, including but not limited to, Freund's 10 (complete and incomplete), mineral gels such as aluminum hydroxide, surface active substance such as lysolecithin, pluronic polyols, polyanions, peptides, oil emulsions, keyhole limpet hemocyanin, dinitrophenol, and the like.

15 Monoclonal antibodies can be produced by hybridoma technology (Kohler and Milstein, *Nature*, 256:495, 1975), which immortalizes cells secreting a specific monoclonal antibody. The immortalized cell lines can be created *in vitro* by fusing two different cell types, typically lymphocytes, and tumor cells. The hybridoma cells may be cultivated *in vitro* or *in vivo*. Additionally, fully human antibodies can be generated by transgenic animals (He *et al.*, *J. Immunol.*, 169:595, 2002). Fd phage and Fd phagemid technologies may be used to generate and select recombinant antibodies *in vitro* (Hoogenboom and Chames, *Immunol. Today* 21:371, 2000; 20 Liu *et al.*, *J. Mol. Biol.* 315:1063, 2002). The complementarity-determining regions of an antibody can be identified, and synthetic peptides corresponding to such regions may be used to mediate antigen binding (U.S. Patent No. 5,637,677).

25 Antibodies are preferably prepared against regions or discrete fragments of a variant protein containing a variant amino acid sequence as compared to the corresponding wild-type protein (*e.g.*, a region of a variant protein that includes an amino acid encoded by a nonsynonymous cSNP, a region affected by truncation caused by a nonsense SNP that creates a stop codon, or a region resulting from the destruction of a stop codon due to read-through mutation caused by a SNP). Furthermore, preferred regions will include those involved in 30 function/activity and/or protein/binding partner interaction. Such fragments can be selected on a physical property, such as fragments corresponding to regions that are located on the surface of the protein, *e.g.*, hydrophilic regions, or can be selected based on sequence uniqueness, or based on the position of the variant amino acid residue(s) encoded by the SNPs provided by the present invention. An antigenic fragment will typically comprise at least about 8-10 contiguous amino acid

residues in which at least one of the amino acid residues is an amino acid affected by a SNP disclosed herein. The antigenic peptide can comprise, however, at least 12, 14, 16, 20, 25, 50, 100 (or any other number in-between) or more amino acid residues, provided that at least one amino acid is affected by a SNP disclosed herein.

5 Detection of an antibody of the present invention can be facilitated by coupling (*i.e.*, physically linking) the antibody or an antigen-reactive fragment thereof to a detectable substance. Detectable substances include, but are not limited to, various enzymes, prosthetic groups, fluorescent materials, luminescent materials, bioluminescent materials, and radioactive materials. Examples of suitable enzymes include horseradish peroxidase, alkaline phosphatase, β -
10 galactosidase, or acetylcholinesterase; examples of suitable prosthetic group complexes include streptavidin/biotin and avidin/biotin; examples of suitable fluorescent materials include umbelliferone, fluorescein, fluorescein isothiocyanate, rhodamine, dichlorotriazinylamine fluorescein, dansyl chloride or phycoerythrin; an example of a luminescent material includes luminol; examples of bioluminescent materials include luciferase, luciferin, and aequorin, and
15 examples of suitable radioactive material include ^{125}I , ^{131}I , ^{35}S or ^3H .

Antibodies, particularly the use of antibodies as therapeutic agents, are reviewed in: Morgan, "Antibody therapy for Alzheimer's disease," *Expert Rev Vaccines*. 2003 Feb; 2(1):53-9; Ross *et al.*, "Anticancer antibodies," *Am J Clin Pathol*. 2003 Apr; 119(4):472-85; Goldenberg, "Advancing role of radiolabeled antibodies in the therapy of cancer," *Cancer Immunol Immunother*. 2003 May;
20 52(5):281-96. Epub 2003 Mar 11; Ross *et al.*, "Antibody-based therapeutics in oncology," *Expert Rev Anticancer Ther*. 2003 Feb; 3(1):107-21; Cao *et al.*, "Bispecific antibody conjugates in therapeutics," *Adv Drug Deliv Rev*. 2003 Feb 10; 55(2):171-97; von Mehren *et al.*, "Monoclonal antibody therapy for cancer," *Annu Rev Med*. 2003; 54:343-69. Epub 2001 Dec 03; Hudson *et al.*, "Engineered antibodies," *Nat Med*. 2003 Jan; 9(1):129-34; Brekke *et al.*, "Therapeutic antibodies for
25 human diseases at the dawn of the twenty-first century," *Nat Rev Drug Discov*. 2003 Jan; 2(1):52-62 (Erratum in: *Nat Rev Drug Discov*. 2003 Mar; 2(3):240); Houdebine, "Antibody manufacture in transgenic animals and comparisons with other systems," *Curr Opin Biotechnol*. 2002 Dec; 13(6):625-9; Andreakos *et al.*, "Monoclonal antibodies in immune and inflammatory diseases," *Curr Opin Biotechnol*. 2002 Dec; 13(6):615-20; Kellermann *et al.*, "Antibody discovery: the use of
30 transgenic mice to generate human monoclonal antibodies for therapeutics," *Curr Opin Biotechnol*. 2002 Dec; 13(6):593-7; Pini *et al.*, "Phage display and colony filter screening for high-throughput selection of antibody libraries," *Comb Chem High Throughput Screen*. 2002 Nov; 5(7):503-10; Batra *et al.*, "Pharmacokinetics and biodistribution of genetically engineered antibodies," *Curr Opin*

Biotechnol. 2002 Dec; 13(6):603-8; and Tangri *et al.*, "Rationally engineered proteins or antibodies with absent or reduced immunogenicity," *Curr Med Chem.* 2002 Dec; 9(24):2191-9.

Uses of Antibodies

5 Antibodies can be used to isolate the variant proteins of the present invention from a natural cell source or from recombinant host cells by standard techniques, such as affinity chromatography or immunoprecipitation. In addition, antibodies are useful for detecting the presence of a variant protein of the present invention in cells or tissues to determine the pattern of expression of the variant protein among various tissues in an organism and over the course of normal development or
10 disease progression. Further, antibodies can be used to detect variant protein *in situ*, *in vitro*, in a bodily fluid, or in a cell lysate or supernatant in order to evaluate the amount and pattern of expression. Also, antibodies can be used to assess abnormal tissue distribution, abnormal expression during development, or expression in an abnormal condition, such as VT. Additionally, antibody detection of circulating fragments of the full-length variant protein can be used to identify turnover.

15 Antibodies to the variant proteins of the present invention are also useful in pharmacogenomic analysis. Thus, antibodies against variant proteins encoded by alternative SNP alleles can be used to identify individuals that require modified treatment modalities.

 Further, antibodies can be used to assess expression of the variant protein in disease states such as in active stages of the disease or in an individual with a predisposition to a disease related to
20 the protein's function, particularly VT. Antibodies specific for a variant protein encoded by a SNP-containing nucleic acid molecule of the present invention can be used to assay for the presence of the variant protein, such as to screen for predisposition to VT as indicated by the presence of the variant protein.

 Antibodies are also useful as diagnostic tools for evaluating the variant proteins in
25 conjunction with analysis by electrophoretic mobility, isoelectric point, tryptic peptide digest, and other physical assays well known in the art.

 Antibodies are also useful for tissue typing. Thus, where a specific variant protein has been correlated with expression in a specific tissue, antibodies that are specific for this protein can be used to identify a tissue type.

30 Antibodies can also be used to assess aberrant subcellular localization of a variant protein in cells in various tissues. The diagnostic uses can be applied, not only in genetic testing, but also in monitoring a treatment modality. Accordingly, where treatment is ultimately aimed at correcting the expression level or the presence of variant protein or aberrant tissue distribution or developmental

expression of a variant protein, antibodies directed against the variant protein or relevant fragments can be used to monitor therapeutic efficacy.

The antibodies are also useful for inhibiting variant protein function, for example, by blocking the binding of a variant protein to a binding partner. These uses can also be applied in a therapeutic context in which treatment involves inhibiting a variant protein's function. An antibody can be used, for example, to block or competitively inhibit binding, thus modulating (agonizing or antagonizing) the activity of a variant protein. Antibodies can be prepared against specific variant protein fragments containing sites required for function or against an intact variant protein that is associated with a cell or cell membrane. For *in vivo* administration, an antibody may be linked with an additional therapeutic payload such as a radionuclide, an enzyme, an immunogenic epitope, or a cytotoxic agent. Suitable cytotoxic agents include, but are not limited to, bacterial toxin such as diphtheria, and plant toxin such as ricin. The *in vivo* half-life of an antibody or a fragment thereof may be lengthened by pegylation through conjugation to polyethylene glycol (Leong *et al.*, *Cytokine* 16:106, 2001).

The invention also encompasses kits for using antibodies, such as kits for detecting the presence of a variant protein in a test sample. An exemplary kit can comprise antibodies such as a labeled or labelable antibody and a compound or agent for detecting variant proteins in a biological sample; means for determining the amount, or presence/absence of variant protein in the sample; means for comparing the amount of variant protein in the sample with a standard; and instructions for use.

Vectors and Host Cells

The present invention also provides vectors containing the SNP-containing nucleic acid molecules described herein. The term "vector" refers to a vehicle, preferably a nucleic acid molecule, which can transport a SNP-containing nucleic acid molecule. When the vector is a nucleic acid molecule, the SNP-containing nucleic acid molecule can be covalently linked to the vector nucleic acid. Such vectors include, but are not limited to, a plasmid, single or double stranded phage, a single or double stranded RNA or DNA viral vector, or artificial chromosome, such as a BAC, PAC, YAC, or MAC.

A vector can be maintained in a host cell as an extrachromosomal element where it replicates and produces additional copies of the SNP-containing nucleic acid molecules. Alternatively, the vector may integrate into the host cell genome and produce additional copies of the SNP-containing nucleic acid molecules when the host cell replicates.

The invention provides vectors for the maintenance (cloning vectors) or vectors for expression (expression vectors) of the SNP-containing nucleic acid molecules. The vectors can function in prokaryotic or eukaryotic cells or in both (shuttle vectors).

5 Expression vectors typically contain cis-acting regulatory regions that are operably linked in the vector to the SNP-containing nucleic acid molecules such that transcription of the SNP-containing nucleic acid molecules is allowed in a host cell. The SNP-containing nucleic acid molecules can also be introduced into the host cell with a separate nucleic acid molecule capable of affecting transcription. Thus, the second nucleic acid molecule may provide a trans-acting factor interacting with the cis-regulatory control region to allow transcription of the SNP-containing
10 nucleic acid molecules from the vector. Alternatively, a trans-acting factor may be supplied by the host cell. Finally, a trans-acting factor can be produced from the vector itself. It is understood, however, that in some embodiments, transcription and/or translation of the nucleic acid molecules can occur in a cell-free system.

The regulatory sequences to which the SNP-containing nucleic acid molecules described
15 herein can be operably linked include promoters for directing mRNA transcription. These include, but are not limited to, the left promoter from bacteriophage λ , the lac, TRP, and TAC promoters from *E. coli*, the early and late promoters from SV40, the CMV immediate early promoter, the adenovirus early and late promoters, and retrovirus long-terminal repeats.

In addition to control regions that promote transcription, expression vectors may also include
20 regions that modulate transcription, such as repressor binding sites and enhancers. Examples include the SV40 enhancer, the cytomegalovirus immediate early enhancer, polyoma enhancer, adenovirus enhancers, and retrovirus LTR enhancers.

In addition to containing sites for transcription initiation and control, expression vectors can also contain sequences necessary for transcription termination and, in the transcribed region, a
25 ribosome-binding site for translation. Other regulatory control elements for expression include initiation and termination codons as well as polyadenylation signals. A person of ordinary skill in the art would be aware of the numerous regulatory sequences that are useful in expression vectors (see, *e.g.*, Sambrook and Russell, 2000, *Molecular Cloning: A Laboratory Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY).

30 A variety of expression vectors can be used to express a SNP-containing nucleic acid molecule. Such vectors include chromosomal, episomal, and virus-derived vectors, for example, vectors derived from bacterial plasmids, from bacteriophage, from yeast episomes, from yeast chromosomal elements, including yeast artificial chromosomes, from viruses such as baculoviruses, papovaviruses such as SV40, Vaccinia viruses, adenoviruses, poxviruses, pseudorabies viruses, and

retroviruses. Vectors can also be derived from combinations of these sources such as those derived from plasmid and bacteriophage genetic elements, *e.g.*, cosmids and phagemids. Appropriate cloning and expression vectors for prokaryotic and eukaryotic hosts are described in Sambrook and Russell, 2000, *Molecular Cloning: A Laboratory Manual*, Cold Spring Harbor Laboratory Press, 5 Cold Spring Harbor, NY.

The regulatory sequence in a vector may provide constitutive expression in one or more host cells (*e.g.*, tissue specific expression) or may provide for inducible expression in one or more cell types such as by temperature, nutrient additive, or exogenous factor, *e.g.*, a hormone or other ligand. A variety of vectors that provide constitutive or inducible expression of a nucleic acid sequence in 10 prokaryotic and eukaryotic host cells are well known to those of ordinary skill in the art.

A SNP-containing nucleic acid molecule can be inserted into the vector by methodology well-known in the art. Generally, the SNP-containing nucleic acid molecule that will ultimately be expressed is joined to an expression vector by cleaving the SNP-containing nucleic acid molecule and the expression vector with one or more restriction enzymes and then ligating the fragments 15 together. Procedures for restriction enzyme digestion and ligation are well known to those of ordinary skill in the art.

The vector containing the appropriate nucleic acid molecule can be introduced into an appropriate host cell for propagation or expression using well-known techniques. Bacterial host cells include, but are not limited to, *E. coli*, *Streptomyces*, and *Salmonella typhimurium*. Eukaryotic 20 host cells include, but are not limited to, yeast, insect cells such as *Drosophila*, animal cells such as COS and CHO cells, and plant cells.

As described herein, it may be desirable to express the variant peptide as a fusion protein. Accordingly, the invention provides fusion vectors that allow for the production of the variant peptides. Fusion vectors can, for example, increase the expression of a recombinant protein, 25 increase the solubility of the recombinant protein, and aid in the purification of the protein by acting, for example, as a ligand for affinity purification. A proteolytic cleavage site may be introduced at the junction of the fusion moiety so that the desired variant peptide can ultimately be separated from the fusion moiety. Proteolytic enzymes suitable for such use include, but are not limited to, factor Xa, thrombin, and enterokinase. Typical fusion expression vectors include pGEX (Smith *et al.*, 30 *Gene* 67:31-40 (1988)), pMAL (New England Biolabs, Beverly, MA) and pRIT5 (Pharmacia, Piscataway, NJ) which fuse glutathione S-transferase (GST), maltose E binding protein, or protein A, respectively, to the target recombinant protein. Examples of suitable inducible non-fusion *E. coli* expression vectors include pTrc (Amann *et al.*, *Gene* 69:301-315 (1988)) and pET 11d (Studier *et al.*, *Gene Expression Technology: Methods in Enzymology* 185:60-89 (1990)).

Recombinant protein expression can be maximized in a bacterial host by providing a genetic background wherein the host cell has an impaired capacity to proteolytically cleave the recombinant protein (Gottesman, S., *Gene Expression Technology: Methods in Enzymology* 185, Academic Press, San Diego, California (1990) 119-128). Alternatively, the sequence of the SNP-containing nucleic acid molecule of interest can be altered to provide preferential codon usage for a specific host cell, for example, *E. coli* (Wada *et al.*, *Nucleic Acids Res.* 20:2111-2118 (1992)).

The SNP-containing nucleic acid molecules can also be expressed by expression vectors that are operative in yeast. Examples of vectors for expression in yeast (*e.g.*, *S. cerevisiae*) include pYepSec1 (Baldari, *et al.*, *EMBO J.* 6:229-234 (1987)), pMFa (Kurjan *et al.*, *Cell* 30:933-943(1982)), pJRY88 (Schultz *et al.*, *Gene* 54:113-123 (1987)), and pYES2 (Invitrogen Corporation, San Diego, CA).

The SNP-containing nucleic acid molecules can also be expressed in insect cells using, for example, baculovirus expression vectors. Baculovirus vectors available for expression of proteins in cultured insect cells (*e.g.*, Sf 9 cells) include the pAc series (Smith *et al.*, *Mol. Cell Biol.* 3:2156-2165 (1983)) and the pVL series (Lucklow *et al.*, *Virology* 170:31-39 (1989)).

In certain embodiments of the invention, the SNP-containing nucleic acid molecules described herein are expressed in mammalian cells using mammalian expression vectors. Examples of mammalian expression vectors include pCDM8 (Seed, B. *Nature* 329:840(1987)) and pMT2PC (Kaufman *et al.*, *EMBO J.* 6:187-195 (1987)).

The invention also encompasses vectors in which the SNP-containing nucleic acid molecules described herein are cloned into the vector in reverse orientation, but operably linked to a regulatory sequence that permits transcription of antisense RNA. Thus, an antisense transcript can be produced to the SNP-containing nucleic acid sequences described herein, including both coding and non-coding regions. Expression of this antisense RNA is subject to each of the parameters described above in relation to expression of the sense RNA (regulatory sequences, constitutive or inducible expression, tissue-specific expression).

The invention also relates to recombinant host cells containing the vectors described herein. Host cells therefore include, for example, prokaryotic cells, lower eukaryotic cells such as yeast, other eukaryotic cells such as insect cells, and higher eukaryotic cells such as mammalian cells.

The recombinant host cells can be prepared by introducing the vector constructs described herein into the cells by techniques readily available to persons of ordinary skill in the art. These include, but are not limited to, calcium phosphate transfection, DEAE-dextran-mediated transfection, cationic lipid-mediated transfection, electroporation, transduction, infection, lipofection, and other techniques such as those described in Sambrook and Russell, 2000, *Molecular*

Cloning: A Laboratory Manual, Cold Spring Harbor Laboratory, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY).

5 Host cells can contain more than one vector. Thus, different SNP-containing nucleotide sequences can be introduced in different vectors into the same cell. Similarly, the SNP-containing nucleic acid molecules can be introduced either alone or with other nucleic acid molecules that are not related to the SNP-containing nucleic acid molecules, such as those providing trans-acting factors for expression vectors. When more than one vector is introduced into a cell, the vectors can be introduced independently, co-introduced, or joined to the nucleic acid molecule vector.

10 In the case of bacteriophage and viral vectors, these can be introduced into cells as packaged or encapsulated virus by standard procedures for infection and transduction. Viral vectors can be replication-competent or replication-defective. In the case in which viral replication is defective, replication can occur in host cells that provide functions that complement the defects.

15 Vectors generally include selectable markers that enable the selection of the subpopulation of cells that contain the recombinant vector constructs. The marker can be inserted in the same vector that contains the SNP-containing nucleic acid molecules described herein or may be in a separate vector. Markers include, for example, tetracycline or ampicillin-resistance genes for prokaryotic host cells, and dihydrofolate reductase or neomycin resistance genes for eukaryotic host cells. However, any marker that provides selection for a phenotypic trait can be effective.

20 While the mature variant proteins can be produced in bacteria, yeast, mammalian cells, and other cells under the control of the appropriate regulatory sequences, cell-free transcription and translation systems can also be used to produce these variant proteins using RNA derived from the DNA constructs described herein.

25 Where secretion of the variant protein is desired, which is difficult to achieve with multi-transmembrane domain containing proteins such as G-protein-coupled receptors (GPCRs), appropriate secretion signals can be incorporated into the vector. The signal sequence can be endogenous to the peptides or heterologous to these peptides.

30 Where the variant protein is not secreted into the medium, the protein can be isolated from the host cell by standard disruption procedures, including freeze/thaw, sonication, mechanical disruption, use of lysing agents, and the like. The variant protein can then be recovered and purified by well-known purification methods including, for example, ammonium sulfate precipitation, acid extraction, anion or cationic exchange chromatography, phosphocellulose chromatography, hydrophobic-interaction chromatography, affinity chromatography, hydroxylapatite chromatography, lectin chromatography, or high performance liquid chromatography.

It is also understood that, depending upon the host cell in which recombinant production of the variant proteins described herein occurs, they can have various glycosylation patterns, or may be non-glycosylated, as when produced in bacteria. In addition, the variant proteins may include an initial modified methionine in some cases as a result of a host-mediated process.

5 For further information regarding vectors and host cells, see *Current Protocols in Molecular Biology*, John Wiley & Sons, N.Y.

Uses of Vectors and Host Cells, and Transgenic Animals

10 Recombinant host cells that express the variant proteins described herein have a variety of uses. For example, the cells are useful for producing a variant protein that can be further purified into a preparation of desired amounts of the variant protein or fragments thereof. Thus, host cells containing expression vectors are useful for variant protein production.

15 Host cells are also useful for conducting cell-based assays involving the variant protein or variant protein fragments, such as those described above as well as other formats known in the art. Thus, a recombinant host cell expressing a variant protein is useful for assaying compounds that stimulate or inhibit variant protein function. Such an ability of a compound to modulate variant protein function may not be apparent from assays of the compound on the native/wild-type protein, or from cell-free assays of the compound. Recombinant host cells are also useful for assaying functional alterations in the variant proteins as compared with a known function.

20 Genetically-engineered host cells can be further used to produce non-human transgenic animals. A transgenic animal is preferably a non-human mammal, for example, a rodent, such as a rat or mouse, in which one or more of the cells of the animal include a transgene. A transgene is exogenous DNA containing a SNP of the present invention which is integrated into the genome of a cell from which a transgenic animal develops and which remains in the genome of the mature
25 animal in one or more of its cell types or tissues. Such animals are useful for studying the function of a variant protein *in vivo*, and identifying and evaluating modulators of variant protein activity. Other examples of transgenic animals include, but are not limited to, non-human primates, sheep, dogs, cows, goats, chickens, and amphibians. Transgenic non-human mammals such as cows and goats can be used to produce variant proteins which can be secreted in the animal's milk and then
30 recovered.

A transgenic animal can be produced by introducing a SNP-containing nucleic acid molecule into the male pronuclei of a fertilized oocyte, *e.g.*, by microinjection or retroviral infection, and allowing the oocyte to develop in a pseudopregnant female foster animal. Any nucleic acid

molecules that contain one or more SNPs of the present invention can potentially be introduced as a transgene into the genome of a non-human animal.

Any of the regulatory or other sequences useful in expression vectors can form part of the transgenic sequence. This includes intronic sequences and polyadenylation signals, if not already
5 included. A tissue-specific regulatory sequence(s) can be operably linked to the transgene to direct expression of the variant protein in particular cells or tissues.

Methods for generating transgenic animals via embryo manipulation and microinjection, particularly animals such as mice, have become conventional in the art and are described in, for example, U.S. Patent Nos. 4,736,866 and 4,870,009, both by Leder *et al.*, U.S. Patent No. 4,873,191
10 by Wagner *et al.*, and in Hogan, B., *Manipulating the Mouse Embryo*, (Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y., 1986). Similar methods are used for production of other transgenic animals. A transgenic founder animal can be identified based upon the presence of the transgene in its genome and/or expression of transgenic mRNA in tissues or cells of the animals. A transgenic founder animal can then be used to breed additional animals carrying the transgene.
15 Moreover, transgenic animals carrying a transgene can further be bred to other transgenic animals carrying other transgenes. A transgenic animal also includes a non-human animal in which the entire animal or tissues in the animal have been produced using the homologously recombinant host cells described herein.

In another embodiment, transgenic non-human animals can be produced which contain
20 selected systems that allow for regulated expression of the transgene. One example of such a system is the cre/loxP recombinase system of bacteriophage P1 (Lakso *et al. PNAS* 89:6232-6236 (1992)). Another example of a recombinase system is the FLP recombinase system of *S. cerevisiae* (O'Gorman *et al. Science* 251:1351-1355 (1991)). If a cre/loxP recombinase system is used to regulate expression of the transgene, animals containing transgenes encoding both the Cre
25 recombinase and a selected protein are generally needed. Such animals can be provided through the construction of "double" transgenic animals, *e.g.*, by mating two transgenic animals, one containing a transgene encoding a selected variant protein and the other containing a transgene encoding a recombinase.

Clones of the non-human transgenic animals described herein can also be produced
30 according to the methods described in, for example, Wilmut, I. *et al. Nature* 385:810-813 (1997) and PCT International Publication Nos. WO 97/07668 and WO 97/07669. In brief, a cell (*e.g.*, a somatic cell) from the transgenic animal can be isolated and induced to exit the growth cycle and enter G₀ phase. The quiescent cell can then be fused, *e.g.*, through the use of electrical pulses, to an enucleated oocyte from an animal of the same species from which the quiescent cell is isolated. The

reconstructed oocyte is then cultured such that it develops to morula or blastocyst and then transferred to pseudopregnant female foster animal. The offspring born of this female foster animal will be a clone of the animal from which the cell (*e.g.*, a somatic cell) is isolated.

Transgenic animals containing recombinant cells that express the variant proteins described herein are useful for conducting the assays described herein in an *in vivo* context. Accordingly, the various physiological factors that are present *in vivo* and that could influence ligand or substrate binding, variant protein activation, signal transduction, or other processes or interactions, may not be evident from *in vitro* cell-free or cell-based assays. Thus, non-human transgenic animals of the present invention may be used to assay *in vivo* variant protein function as well as the activities of a therapeutic agent or compound that modulates variant protein function/activity or expression. Such animals are also suitable for assessing the effects of null mutations (*i.e.*, mutations that substantially or completely eliminate one or more variant protein functions).

For further information regarding transgenic animals, see Houdebine, "Antibody manufacture in transgenic animals and comparisons with other systems," *Curr Opin Biotechnol.* 2002 Dec; 13(6):625-9; Petters *et al.*, "Transgenic animals as models for human disease," *Transgenic Res.* 2000; 9(4-5):347-51; discussion 345-6; Wolf *et al.*, "Use of transgenic animals in understanding molecular mechanisms of toxicity," *J Pharm Pharmacol.* 1998 Jun; 50(6):567-74; Echelard, "Recombinant protein production in transgenic animals," *Curr Opin Biotechnol.* 1996 Oct; 7(5):536-40; Houdebine, "Transgenic animal bioreactors," *Transgenic Res.* 2000; 9(4-5):305-20; Purity *et al.*, "Embryonic stem cells, creating transgenic animals," *Methods Cell Biol.* 1998; 57:279-93; and Robl *et al.*, "Artificial chromosome vectors and expression of complex proteins in transgenic animals," *Theriogenology.* 2003 Jan 1; 59(1):107-13.

EXAMPLES

25

The following examples are offered to illustrate, but not to limit, the claimed invention.

Example 1: Statistical Analysis of SNP Allelic Association with VT Risk

A case-control genetic study was performed to determine the association of SNPs in the human genome with VT, and in particular deep vein thrombosis (DVT), using DNA extracted from two independently obtained sample sets. One set was obtained from the Leiden Thrombophilia Study (LETS), as conducted at the Hemostasis and Thrombophilia Research Center, Academic Hospital Leiden, the Netherlands. LETS was a large case-control study designed to investigate risk factors for the development of VT, comprised of nearly 1000 cases and controls. LETS cases were patients with the diagnosis of a confirmed DVT and without a

30

known malignant disorder; controls had no history of venous thromboembolism or malignant disorders. Cases and controls in LETS were matched pairwise by sex and age (+/- five years), and no case-control pair shared a biologic relationship. (FJM van der Meer, FR Rosendaal *et al.*, *Thrombosis and Haemostasis*, 78(1):631-635, 1997).

5 A second sample set was obtained from the Multiple Environmental and Genetic Assessment (MEGA) study. MEGA was a study of over 5,000 cases and controls, also of the Netherlands, designed to investigate risk factors for VT. Cases were patients who had been diagnosed with DVT or PE. Controls were without DVT or PE, and were matched for age and sex with cases. (JW Blom, FR Rosendaal *et al.*, *JAMA* 293(6):715-22, 2005).

10 The inclusion criteria for genotyping in this analysis of SNP association with VT was that cases of the LETS sample set had a confirmed DVT, and cases of the MEGA set had an arm or leg thrombotic event. Excluded from MEGA cases were those patients who had had only pulmonary embolism or cancer. Additionally excluded from the MEGA case set was any patient where it was unknown whether his/her genotype contained the Factor V Leiden mutation (FV
15 1691 C) or prothrombin mutation (FII 2021 C).

DNA was extracted from the LETS and MEGA blood samples using conventional DNA extraction methods or commercially available kits, such as the QIA-amp kit from Qiagen (Valencia, CA), according to the manufacturer's suggestions. SNP markers in the extracted DNA were analyzed by genotyping. Initially, pooling studies were performed in which DNA samples
20 from typically 50 cases or controls from either sample set were pooled, and the allele frequencies for specific markers were obtained using a PRISM® 7900HT Sequence Detection System (Applied Biosystems, Foster City, CA) by allele-specific PCR, similar to the method described by Germer *et al.* *Genome Research* 10:258-266, 2000. (PCR primers used in this example are shown in Table 3, with their corresponding SNP markers by hCV.) In the LETS set, DNA
25 samples from 440 cases and 426 controls were genotyped in the initial pooling studies. In the MEGA set, 1568 cases and 1866 controls were genotyped in the initial pooling studies. Those SNPs with an observed P value < 0.05 (by the Fisher Exact test) and odds ratio (OR) > 1.2 for association with VT in the pooled samples were selected as candidates for validation by individual genotyping.

30 For validation of SNP association with VT based on the results of the pooling studies, samples from LETS and MEGA were individually genotyped for those SNP markers that the risk alleles had shown association with VT by pooling studies (P value < 0.05, OR > 1.2), by a method similar to that described above for pooling studies, using allele-specific real-time PCR as previously described. (S. Germer *et al.*, *Genome Research* 10:258-266, 2000). From the LETS

sample set, 440 cases and 426 controls were individually genotyped. Markers in which the risk alleles had been shown to be associated with VT (P value < 0.05, OR > 1.2) were selected to be individually genotyped in the MEGA set. From the MEGA set, 1568 cases and 1866 controls were individually genotyped.

5 Allele frequencies of SNPs were analyzed for association with DVT as a clinical endpoint in the individually genotyped LETS and MEGA sample sets. Allele frequencies for the tested SNPs were compared between cases and controls to determine VT risk association. Results of this analysis are reported in Tables 4 and 5, for eight SNPs showing association with risk of developing VT, specifically DVT, based on individual genotyping. The magnitude of the allelic effect on risk association (effect size) was estimated by an odds ratio (OR). An allele may be under- or overrepresented in cases. An allele that is overrepresented in cases indicates that the reported allele is a risk factor for disease. An allele that is underrepresented in cases indicates that this allele is protective against disease.

15 A SNP was considered a marker for VT risk if the association analyses in the two individually genotyped sample sets (LETS and MEGA) showed the same risk allele, and P values <0.05 in both sample sets. Allelic P values were calculated using the Fisher Exact test. Table 4 shows those SNP markers associated with VT in patients unstratified by sex. Also shown are the modes (column "Mode") used to calculate allelic association: dominant/recessive (Rec) or homozygous vs. reference (Hom). In the dominant mode, an association with disease is seen when individuals homozygous or heterozygous for the risk allele are compared to those homozygous for the non-risk allele. In the recessive mode, an association is seen only in individuals homozygous for the risk allele, when compared to those heterozygous or homozygous for the non-risk allele. In the homozygous vs. reference mode, an association with VT is seen when individuals homozygous for the risk allele are compared to individuals homozygous for the non-risk allele. In the allelic mode, Allelic indicates a Fisher exact calculation that compares the counts of the risk allele in cases and controls to counts of the non-risk allele in cases and controls.

30 In the case of the one marker where the allele frequencies of cases vs. controls indicated a VT risk association in the female population, this is indicated in the "Stratum" column of Table 5. The sum of eight SNP markers shown by this analysis to be associated with risk of VT are demonstrated in Table 4. The column "Risk Allele" indicates the allele associated with VT risk.

One example of a SNP marker is associated with risk of VT in the homozygous vs. reference mode is hCV11541681, a marker in the LOC200420 gene (Table 4). Samples from LETS and MEGA were individually genotyped for hCV11541681. Association results are

shown in Table 4. Individuals with two copies of the C allele at this SNP position (i.e., homozygous for the risk allele) showed a significant association with increased risk of VT in the LETS and MEGA sample sets (P values of 0.033 and 0.024, OR of 1.57 and 1.27, respectively) when compared to individuals carrying two copies of the non-risk allele (homozygous for the reference allele), in case and control populations unstratified by sex, in both LETS and MEGA sets.

An example of a SNP marker associated with increased risk of VT in the recessive mode is hCV2403368, a marker in the C1QTNF6 gene. Samples from LETS and MEGA were individually genotyped for hCV2403368. Association results are shown in Table 4. Individuals with two copies of the G allele at this SNP position (i.e., homozygous for the risk allele) showed a significant association with increased risk of VT in the LETS and MEGA sample sets (P values of 0.022 and 0.034, OR of 1.37 and 1.16, respectively) when compared to individuals carrying one or no copy of the risk allele (heterozygous for the risk allele or homozygous for the non-risk allele), in case and control populations unstratified by sex.

An example of a SNP marker associated with VT risk only in females is hCV12092542 (Table 5). The same number of samples and types of cases from LETS and MEGA were individually genotyped for this SNP as for SNP hCV11541681, described above. An analysis of the allele frequencies obtained from cases vs. controls showed an association of this SNP with VT when results were stratified by sex; see Table 5. Females with two copies of the T allele at SNP hCV12092542 position showed a significant association with increased risk of VT, relative to females with one or no copy of the risk allele (i.e., recessive mode) in both LETS and MEGA sets (P values of 0.036 and 0.032, OR of 1.52 and 1.25, respectively). By comparison, in samples unstratified by sex (Stratum "All" in Table 5) or from males (Stratum "Male"), this allele did not demonstrate a statistically significant association with risk of VT (P values of 0.383 for "All," 0.192 for "Male").

This and all other SNP markers described herein may also be found in Tables 1 and 2. In Table 1, context sequence information is provided regarding the transcript in which each SNP is found, if it resides in a transcript. For example this SNP marker, Celera SNP ID hCV11541681, is found in the transcript sequence of SEQ ID (Table 1). Also provided in Table 1 is other information regarding the transcript in which the SNP is located: the gene symbol, LOC200420; position of the SNP in the transcript; chromosome, 2; the Public SNP ID (i.e. the rs, or RefSNP, number from the National Center for Biotechnology Information SNP database if known), rs2001490; SNP type, 3' UTR, etc. In Table 2 is provided genomic sequence information for all SNP markers described herein. For example Celera SNP ID hCV11541681 is found in the

genomic sequence of SEQ ID NO: 67 (Table 2). In the event that there are SNPs calculated to be in LD with the interrogated SNP, that information is also provided in Tables 1 and 2, as a "Related Interrogated SNP." E.g., see SEQ ID NO: 67 in Table 2, where the Celera SNP ID for the interrogated SNP is hCV11541681, and the related LD SNP directly following is hCV26996690 (power of 0.51).

The SNPs presented in Tables 4-5 are shown to be associated with risk for VT, specifically DVT. Other SNPs (such as those presented in Tables 1 and 2), including LD SNPs in the genes listed in Table 4 would also be expected to be useful for the diagnosis, prognosis, etc., of VT, and particularly DVT. All such SNPs associated with VT may also be useful for predicting a patient's response to therapeutic agents such as anticoagulant agents.

Example 2: Additional SNPs in LD with VT-Associated Interrogated SNP Markers

An investigation was conducted to identify SNP markers in linkage disequilibrium (LD) with SNPs which have been found to be associated with VT, as shown in Tables 4-5. Briefly, the power threshold (T) was set at 51% for detecting disease association using LD markers. This power threshold is based on equation (31) above, which incorporates allele frequency data from previous disease association studies, the predicted error rate for not detecting truly disease-associated markers, and a significance level of 0.05. Using this power calculation and the sample size, for each interrogated SNP (Table 4) a threshold level of LD, or r^2 value, was derived (r_T^2 , equations (32) and (33)). The threshold value r_T^2 is the minimum value of linkage disequilibrium between the interrogated SNP and its LD SNPs possible such that the non-interrogated SNP still retains a power greater or equal to T for detecting disease-association.

Based on the above methodology, LD SNPs were found for all interrogated SNPs shown in Tables 4-5. LD SNPs are listed in Table 6, each associated with its respective interrogated SNP. Also shown are the public SNP IDs (rs numbers) for interrogated and LD SNPs, the threshold r^2 value and the power used to determine this, and the r^2 value of linkage disequilibrium between the interrogated SNP and its matching LD SNP. As an example, in Table 6, VT-associated interrogated SNP hCV263841 was calculated to be in LD with hCV105917 at a r_T^2 value of 0.46, using 51% power, thus making SNP hCV105917 a marker associated with VT because the r^2 value of .56 is larger than r_T^2 of .46.

Modifications and variations of the described compositions, methods and systems of the invention will be apparent to those skilled in the art without departing from the scope of the invention. Although the invention has been described in connection with specific preferred embodiments and

certain working examples, it should be understood that the invention as claimed should not be unduly limited to such specific embodiments. Indeed, various modifications of the above-described modes for carrying out the invention that are obvious to those skilled in the field of molecular biology, genetics and related fields are intended to be within the scope of the invention.

Table 1. Transcript SNP info and associated gene/protein information

	Gene Number:	1
	Celera Gene:	hCG14694 - 208000030293584
5	Celera Transcript:	hCT2345280 - 208000030293589
	Public Transcript Accession:	
	Celera Protein:	hCP1910568 - 208000030293573
	Public Protein Accession:	
	Gene Symbol:	F9
10	Protein Name:	coagulation factor IX (plasma thromboplastic component, Christmas disease, hemophilia B)
	Celera Genomic Axis:	GA_x5YUV32W21H (1443186..1495938)
	Chromosome:	X
	OMIM NUMBER:	306900
15	OMIM Information:	Hemophilia B (3); Warfarin sensitivity (3)
	Transcript Sequence (SEQ ID NO: 1):	
	Protein Sequence (SEQ ID NO: 21):	
20	SNP Information	
	Context (SEQ ID NO: 41):	
		CCTTTTACCCTCCATGGTCGTTAAAGGAGAGATGGGGAGCATCATTCTGTTAT ACTTCTGTACACAGTTATAACATGTCTATCAAACCCAGACTTGCTTCC R TAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAG TTTGGGGGAAAAGTTTCTTTCAGAGAGTTAAGTTATTTTATATATATA
	Celera SNP ID:	hCV2288124
	Public SNP ID:	rs440051
	SNP in Transcript Sequence	SEQ ID NO: 1
30	SNP Position Transcript:	2028
	Related Interrogated SNP:	hCV596331 (Power=.51)
	SNP Source:	dbSNP; HapMap; ABI_Val; HGBASE;
	Population (Allele,Count):	caucasian (G,73 A,17)
	SNP Type:	UTR3

35

Gene Number: 1
 Celera Gene: hCG14694 - 208000030293584
 Celera Transcript: hCT2345281 - 208000030293579 .
 Public Transcript Accession:
 5 Celera Protein: hCP1910567 - 208000030293571
 Public Protein Accession:
 Gene Symbol: F9
 Protein Name: coagulation factor IX (plasma thromboplastic component,
 Christmas disease, hemophilia B)
 10 Celera Genomic Axis: GA_x5YUV32W21H (1443186..1495938)
 Chromosome: X
 OMIM NUMBER: 306900
 OMIM Information: Hemophilia B (3); Warfarin sensitivity (3)
 15 Transcript Sequence (SEQ ID NO: 2):
 Protein Sequence (SEQ ID NO: 22):
 SNP Information
 Context (SEQ ID NO: 42):
 20 ATGCGAGCAGTTTTGTAAAAATAGTGCTGATAACAAGGTGGTTTGCTCCTGTA
 CTGAGGGATATCGACTTTCACAAACTTCTAAGCTCACCCGTGCTGAG
 R
 CTGTTTTTCCTGATGTGGACTATGTAAATTCTACTGAAGCTGAAACCATTTTG
 GATAACATCACTCAAAGCACCCAATCATTTAATGACTTCACTCGGGT
 25 Celera SNP ID: hCV596331
 Public SNP ID: rs6048
 SNP in Transcript Sequence SEQ ID NO: 2
 SNP Position Transcript: 578
 SNP Source: dbSNP; HapMap;
 30 Population (Allele,Count): caucasian (A,61|G,29)
 SNP Type: Missense Mutation
 Protein Coding: SEQ ID NO: 22, at position 174,(T,ACT) (A,GCT)
 Context (SEQ ID NO: 43):

CCTTTTACCCTCCATGGTCGTTAAAGGAGAGATGGGGAGCATCATTCTGTTAT
 ACTTCTGTACACAGTTATAACATGTCTATCAAACCCAGACTTGCTTCC

R

TAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAG
 5 TTTGGGGGAAAAGTTTCTTTCAGAGAGTTAAGTTATTTTATATATATA

Celera SNP ID: hCV2288124

Public SNP ID: rs440051

SNP in Transcript Sequence SEQ ID NO: 2

SNP Position Transcript: 2071

10 Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE;

Population (Allele,Count): caucasian (G,73|A,17)

SNP Type: UTR3

15 Gene Number: 1

Celera Gene: hCG14694 - 208000030293584

Celera Transcript: hCT2345282 - 208000030293590

Public Transcript Accession:

Celera Protein: hCP1910569 - 208000030293568

20 Public Protein Accession:

Gene Symbol: F9

Protein Name: coagulation factor IX (plasma thromboplastic component,
 Christmas disease, hemophilia B)

Celera Genomic Axis: GA_x5YUV32W21H (1443186..1495938)

25 Chromosome: X

OMIM NUMBER: 306900

OMIM Information: Hemophilia B (3); Warfarin sensitivity (3)

Transcript Sequence (SEQ ID NO: 3):

30 Protein Sequence (SEQ ID NO: 23):

SNP Information

Context (SEQ ID NO: 44):

ATATCGACTTGCAGAAAACCAGAAGTCCTGTGAACCAGCAGTGCCATTTCCA
TGTGGAAGAGTTTCTGTTTCACAACTTCTAAGCTCACCCGTGCTGAG

R

CTGTTTTTCCTGATGTGGACTATGTAAATTCTACTGAAGCTGAAACCATTTTG
5 GATAACATCACTCAAAGCACCCAATCATTTAATGACTTCACTCGGGT

Celera SNP ID: hCV596331

Public SNP ID: rs6048

SNP in Transcript Sequence SEQ ID NO: 3

SNP Position Transcript: 584

10 SNP Source: dbSNP; HapMap;

Population (Allele,Count): caucasian (A,61|G,29)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 23, at position 176,(T,ACT) (A,GCT)

15 Context (SEQ ID NO: 45):

CCTTTTACCCTCCATGGTCGTTAAAGGAGAGATGGGGAGCATCATTCTGTTAT
ACTTCTGTACACAGTTATACATGTCTATCAAACCCAGACTTGCTTCC

R

TAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAG
20 TTTGGGGGAAAAGTTTCTTTCAGAGAGTTAAGTTATTTTATATATATA

Celera SNP ID: hCV2288124

Public SNP ID: rs440051

SNP in Transcript Sequence SEQ ID NO: 3

SNP Position Transcript: 2077

25 Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE;

Population (Allele,Count): caucasian (G,73|A,17)

SNP Type: UTR3

30 Gene Number: 1

Celera Gene: hCG14694 - 208000030293584

Celera Transcript: hCT5715 - 208000030293602

Public Transcript Accession: NM_000133

Celera Protein: hCP35448 - 208000030293575

Public Protein Accession: NP_000124

Gene Symbol: F9

Protein Name: coagulation factor IX (plasma thromboplastic component,
Christmas disease, hemophilia B)

5 Celera Genomic Axis: GA_x5YUV32W21H (1443186..1495938)

Chromosome: X

OMIM NUMBER: 306900

OMIM Information: Hemophilia B (3); Warfarin sensitivity (3)

10 Transcript Sequence (SEQ ID NO: 4):

Protein Sequence (SEQ ID NO: 24):

SNP Information

Context (SEQ ID NO: 46):

15 ATATCGACTTGCAGAAAACCAGAAGTCCTGTGAACCAGCAGTGCCATTCCA

TGTGGAAGAGTTTCTGTTTCACAACTTCTAAGCTCACCCGTGCTGAG

R

CTGTTTTTCCTGATGTGGACTATGTAAATTCTACTGAAGCTGAAACCATTTTG

GATAACATCACTCAAAGCACCCAATCATTTAATGACTTCACTCGGGT

20 Celera SNP ID: hCV596331

Public SNP ID: rs6048

SNP in Transcript Sequence SEQ ID NO: 4

SNP Position Transcript: 638

SNP Source: dbSNP; HapMap;

25 Population (Allele,Count): caucasian (A,61|G,29)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 24, at position 194,(T,ACT) (A,GCT)

Context (SEQ ID NO: 47):

30 CCTTTTACCCTCCATGGTCGTTAAAGGAGAGATGGGGAGCATCATTCTGTTAT

ACTTCTGTACACAGTTATACATGTCTATCAAACCCAGACTTGCTTCC

R

TAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAG

TTTGGGGGAAAAGTTTCTTTCAGAGAGTTAAGTTATTTTATATATATA

Celera SNP ID: hCV2288124
 Public SNP ID: rs440051
 SNP in Transcript Sequence SEQ ID NO: 4
 SNP Position Transcript: 2131

5 Related Interrogated SNP: hCV596331 (Power=.51)
 SNP Source: dbSNP; HapMap; ABI_Val; HGBASE;
 Population (Allele,Count): caucasian (G,73|A,17)
 SNP Type: UTR3

10 Gene Number: 2
 Celera Gene: hCG1776858 - 104000116701164
 Celera Transcript: hCT1815574 - 104000116701165
 Public Transcript Accession:
 Celera Protein: hCP1721024 - 197000064957675

15 Public Protein Accession:
 Gene Symbol: CML2
 Protein Name: putative N-acetyltransferase Camello 2
 Celera Genomic Axis: GA_x5YUV32W5TR (13100168..13120998)
 Chromosome: 2

20 OMIM NUMBER:
 OMIM Information:

Transcript Sequence (SEQ ID NO: 5):
 Protein Sequence (SEQ ID NO: 25):

25

SNP Information

Context (SEQ ID NO: 48):

AAAAAACCCTGGACGCGGTATGTAGACATAGCATTGCGCACAGACATGTCTG
 ACATCACCAAATCCTACCTGAGTGAGTGTGGCTCCTGCTTCTGGGTGG

30

S

TGAATCTGAAGAGAAGGTGGTGGGCACAGTAGGAGCTCTGCCCGTTGATGAT
 CCCACCTTGAGGGAGAAGCGGTTGCAGCTGTTTCATCTCTCTGTGGAC

Celera SNP ID: hCV11541681
 Public SNP ID: rs2001490

SNP in Transcript Sequence SEQ ID NO: 5
 SNP Position Transcript: 370
 SNP Source: Applera
 Population (Allele,Count): caucasian (G,10|C,10) african american (G,18|C,14) total
 5 (G,28|C,24)
 SNP Type: Missense Mutation
 Protein Coding: SEQ ID NO: 25, at position 88,(G,GGT) (A,GCT)
 SNP Source: dbSNP; Celera; HGBASE
 Population (Allele,Count): caucasian (G,40|C,80)
 10 SNP Type: Missense Mutation
 Protein Coding: SEQ ID NO: 25, at position 88,(G,GGT) (A,GCT)
 Gene Number: 3
 Celera Gene: hCG1776867 - 208000043935214
 15 Celera Transcript: hCT1815586 - 208000043935215
 Public Transcript Accession:
 Celera Protein: hCP1720929 - 208000043935207
 Public Protein Accession:
 Gene Symbol: LOC200420
 20 Protein Name: LOC200420
 Celera Genomic Axis: GA_x5YUV32W5TR (13100528..13176592)
 Chromosome: 2
 OMIM NUMBER:
 OMIM Information:
 25
 Transcript Sequence (SEQ ID NO: 6):
 Protein Sequence (SEQ ID NO: 26):
 SNP Information
 30 Context (SEQ ID NO: 49):
 GTCCACAGAGAGATGAAACAGCTGCAACCGCTTCTCCCTCAAGGTGGGATCA
 TCAACGGGCAGAGCTCCTACTGTGCCACACCTTCTCTTCAGATTCA
 S
 CCACCCAGA

Celera SNP ID: hCV11541681
 Public SNP ID: rs2001490
 SNP in Transcript Sequence SEQ ID NO: 6
 SNP Position Transcript: 1660
 5 SNP Source: Applera
 Population (Allele,Count): caucasian (C,10|G,10) african american (C,18|G,14) total
 (C,28|G,24)
 SNP Type: UTR3
 SNP Source: dbSNP; Celera; HGBASE
 10 Population (Allele,Count): caucasian (C,40|G,80)
 SNP Type: UTR3

 Gene Number: 4
 Celera Gene: hCG1787816 - 30000070905589
 15 Celera Transcript: hCT1826884 - 30000070905590
 Public Transcript Accession: NM_207352
 Celera Protein: hCP1743092 - 30000070905586
 Public Protein Accession: NP_997235
 Gene Symbol: CYP4V2
 20 Protein Name: cytochrome P450, family 4, subfamily v, polypeptide 2
 Celera Genomic Axis: GA_x5YUV32W706 (3771943..3813704)
 Chromosome: 4
 OMIM NUMBER: 608614
 OMIM Information: Bietti crystalline corneoretinal dystrophy, 210370 (3)
 25
 Transcript Sequence (SEQ ID NO: 7):
 Protein Sequence (SEQ ID NO: 27):

 SNP Information
 30 Context (SEQ ID NO: 50):
 AATGAGTGAGATGATATTTTCGAAGAATAAAGATGCCCTGGCTTTGGCTTGATC
 TCTGGTACCTTATGTTTAAAGAAGGATGGGAACACAAAAAGAGCCTT
 M

AGATCCTACATACTTTTACCAACAGTGTCATCGCGGAACGGGCCAATGAAAT
GAACGCCAATGAAGACTGTAGAGGTGATGGCAGGGGCTCTGCCCCCTC

Celera SNP ID: hCV25990131

Public SNP ID: rs13146272

5 SNP in Transcript Sequence SEQ ID NO: 7

SNP Position Transcript: 814

SNP Source: Applera

Population (Allele,Count): caucasian (A,22|C,16) african american (A,7|C,11) total
(A,29|C,27)

10 SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 27, at position 259,(K,AAG) (Q,CAG)

SNP Source: dbSNP;

Population (Allele,Count): caucasian (C,45|A,75)

SNP Type: Missense Mutation

15 Protein Coding: SEQ ID NO: 27, at position 259,(K,AAG) (Q,CAG)

Context (SEQ ID NO: 51):

CCTGGCTTTGGCTTGATCTCTGGTACCTTATGTTTAAAGAAGGATGGGAACAC
AAAAAGAGCCTTAAGATCCTACATACTTTTACCAACAGTGTCATCGC

20 K

GAACGGGCCAATGAAATGAACGCCAATGAAGACTGTAGAGGTGATGGCAGG
GGCTCTGCCCCCTCCAAAAATAAACGCAGGGCCTTTCTTGACTTGCTTT

Celera SNP ID: hCV3230097

Public SNP ID: rs3736455

25 SNP in Transcript Sequence SEQ ID NO: 7

SNP Position Transcript: 849

Related Interrogated SNP: hCV25990131 (Power=.7)

SNP Source: Applera

30 Population (Allele,Count): caucasian (G,24|T,16) african american (G,7|T,9) total
(G,31|T,25)

SNP Type: Silent Mutation

Protein Coding: SEQ ID NO: 27, at position 270,(A,GCG) (A,GCT)

SNP Source: dbSNP; Celera; HapMap; HGBASE;

Population (Allele,Count): caucasian (T,42|G,78)

SNP Type: Silent Mutation
 Protein Coding: SEQ ID NO: 27, at position 270,(A,GCG) (A,GCT)

5 Gene Number: 5
 Celera Gene: hCG21777 - 207000006881214
 Celera Transcript: hCT1968270 - 207000006881244
 Public Transcript Accession: NM_003889
 Celera Protein: hCP1782551 - 207000006881210
 Public Protein Accession: NP_003880
 10 Gene Symbol: NR1I2
 Protein Name: nuclear receptor subfamily 1, group I, member 2
 Celera Genomic Axis: GA_x5YUV32VYQG (75803402..75861414)
 Chromosome: 3
 OMIM NUMBER: 603065
 15 OMIM Information:

Transcript Sequence (SEQ ID NO: 8):

Protein Sequence (SEQ ID NO: 28):

20 SNP Information
 Context (SEQ ID NO: 52):
 AAGCACTGCCTTTACTTCAGTGGGAATCTCGGCCTCAGCCTGCAAGCCAAGTG
 TTCACAGTGAGAAAAGCAAGAGAATAAGCTAATACTCCTGTCCTGAA
 M

25 AAGGCAGCGGCTCCTTGGTAAAGCTACTCCTTGATCGATCCTTTGCACCGGAT
 TGTTCAAAGTGGACCCCAGGGGAGAAGTCGGAGCAAAGAACTTACCA

Celera SNP ID: hCV263841

Public SNP ID: rs1523127

SNP in Transcript Sequence SEQ ID NO: 8

30 SNP Position Transcript: 1709

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,49|A,71)

SNP Type: UTR5

Context (SEQ ID NO: 53):

GCAGGCCCCAGATATAGCCCATGCTGTCCTCCTACCCCAGAGCACACTGTTC
AGGCTACTTCCACTGGTACTGAAATCCAGTATTTCACTTACTCTTTT

Y

5 CTTTCCAATATCCTCATGACATTCAATATTTCACTTACTCTAGGTCCTCCCTGC
CTAAGGCCCAAGTCAACTTTCTGTCCAGTGGGATTTGTAATCCAAT

Celera SNP ID: hCV9152783

Public SNP ID: rs1523130

SNP in Transcript Sequence SEQ ID NO: 8

10 SNP Position Transcript: 177

Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (T,49|C,71)

SNP Type: UTR5

15

Context (SEQ ID NO: 54):

GTA CTTC AAAATAATAACA ACTTAAGTCAATAAATAAATGTAAGGAAGTCCA
AATGTTACCTGAAGACA ACTGTGGTCATTTTTTGGCAATCCCAGGTT

Y

20 TCTTTTCTACCTGTTTGCTCAATCGTGGTCTCCCTCTCCCTCTCTTGTTGGGGCC
CATGCCCTGCTTTACTGTTGCCAGAGGCTTGTA CTTGTTTGCCT

Celera SNP ID: hCV27504984

Public SNP ID: rs3814055

SNP in Transcript Sequence SEQ ID NO: 8

25 SNP Position Transcript: 705

Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (C,72|T,48)

SNP Type: UTR5

30

Gene Number: 5

Celera Gene: hCG21777 - 207000006881214

Celera Transcript: hCT2319597 - 207000006881215

Public Transcript Accession: NM_033013

Celera Protein: hCP1786359 - 207000006881211

Public Protein Accession: NP_148934

Gene Symbol: NR1I2

Protein Name: nuclear receptor subfamily 1, group I, member 2

5 Celera Genomic Axis: GA_x5YUV32VYQG (75803402..75861414)

Chromosome: 3

OMIM NUMBER: 603065

OMIM Information:

10 Transcript Sequence (SEQ ID NO: 9):

Protein Sequence (SEQ ID NO: 29):

SNP Information

Context (SEQ ID NO: 55):

15 AAGCACTGCCTTTACTTCAGTGGGAATCTCGGCCTCAGCCTGCAAGCCAAGTG

TTCACAGTGAGAAAAGCAAGAGAATAAGCTAATACTCCTGTCCTGAA

M

AAGGCAGCGGCTCCTTGGTAAAGCTACTCCTTGATCGATCCTTTGCACCGGAT

TGTTCAAAGTGGACCCCAGGGGAGAAGTCGGAGCAAAGAACTTACCA

20 Celera SNP ID: hCV263841

Public SNP ID: rs1523127

SNP in Transcript Sequence SEQ ID NO: 9

SNP Position Transcript: 1709

SNP Source: dbSNP; Celera; HapMap; HGBASE

25 Population (Allele,Count): caucasian (C,49|A,71)

SNP Type: UTR5

Context (SEQ ID NO: 56):

GCAGGCCCCAGATATAGCCCCATGCTGTCCTCCTACCCCAGAGCACACTGTTC

30 AGGCTACTTCCACTGGTACTGAAATCCAGTATTTCACTTACTCTTTT

Y

CTTTCCAATATCCTCATGACATTCAATATTTCACTTACTCTAGGTCCTCCCTGC

CTAAGGCCCAAGTCAACTTTCTGTCCAGTGGGATTTGTAATCCAAT

Celera SNP ID: hCV9152783

Public SNP ID: rs1523130
 SNP in Transcript Sequence SEQ ID NO: 9
 SNP Position Transcript: 177
 Related Interrogated SNP: hCV263841 (Power=.8)
 5 SNP Source: dbSNP; Celera; HapMap; HGBASE
 Population (Allele,Count): caucasian (T,49|C,71)
 SNP Type: UTR5

Context (SEQ ID NO: 57):

10 GTACTTCAAATAATAACA ACTTAAGTCAATAAATAAATGTAAGGAAGTCCA
 AATGTTACCTGAAGACA ACTGTGGTCATTTTTTGGCAATCCCAGGTT
 Y
 TCTTTTCTACCTGTTTGCTCAATCGTGGTCTCCCTCTCCCTCTCTTGTTGGGGCC
 CATGCCCCTGCTTTACTGTTGCCAGAGGCTTGTACTTGTTTGCCT

15 Celera SNP ID: hCV27504984
 Public SNP ID: rs3814055
 SNP in Transcript Sequence SEQ ID NO: 9
 SNP Position Transcript: 705
 Related Interrogated SNP: hCV263841 (Power=.8)
 20 SNP Source: dbSNP; HapMap; ABI_Val; HGBASE
 Population (Allele,Count): caucasian (C,72|T,48)
 SNP Type: UTR5

Gene Number: 6
 25 Celera Gene: hCG21779 - 104000117367418
 Celera Transcript: hCT12872 - 104000117367419
 Public Transcript Accession:
 Celera Protein: hCP39298 - 197000064919027
 Public Protein Accession:
 30 Gene Symbol: C3orf15
 Protein Name: AAT1-alpha
 Celera Genomic Axis: GA_x5YUV32VYQG (75880995..75939455)
 Chromosome: 3
 OMIM NUMBER:

OMIM Information:

Transcript Sequence (SEQ ID NO: 10):

Protein Sequence (SEQ ID NO: 30):

5

SNP Information

Context (SEQ ID NO: 58):

ATACACAGGCAAATATCCAAGCTACCCTGATTCGCAGCAGACTGAGAAAAGT
TCCCAGGTTTAAAACCATGTTCAGTAACCTGATCCATTATCCAAGATA

10

Y

TCTCTATATTGGAGCAAGTCAGATCCTGTCCCACCATTTATCAGTCGGGAATG
GAAGGGACATAAGGAGAAACACAGAGAAGCCCTCCGGCAGCTCACCA

Celera SNP ID: hCV134278

Public SNP ID: rs9848716

15

SNP in Transcript Sequence SEQ ID NO: 10

SNP Position Transcript: 342

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (T,40|C,80)

20

SNP Type: Silent Mutation

Protein Coding: SEQ ID NO: 30, at position 86,(Y,TAC) (Y,TAT)

Gene Number: 6

Celera Gene: hCG21779 - 104000117367418

25

Celera Transcript: hCT2319577 - 104000117367445

Public Transcript Accession:

Celera Protein: hCP1786388 - 197000064919029

Public Protein Accession:

Gene Symbol: C3orf15

30

Protein Name: AAT1-alpha

Celera Genomic Axis: GA_x5YUV32VYQG (75880995..75939455)

Chromosome: 3

OMIM NUMBER:

OMIM Information:

Transcript Sequence (SEQ ID NO: 11):

Protein Sequence (SEQ ID NO: 31):

5 SNP Information

Context (SEQ ID NO: 59):

ATACACAGGCAAATATCCAAGCTACCCTGATTCGCAGCAGACTGAGAAAAGT
TCCCAGGTTTAAAACCATGTTTCAGTAACCTGATCCATTATCCAAGATA
Y

10 TCTCTATATTGGAGCAAGTCAGATCCTGTCCCACCATTTATCAGTCGGGAATG
GAAGGGACATAAGGAGAAACACAGAGAAGCCCTCCGGCAGCTCACCA

Celera SNP ID: hCV134278

Public SNP ID: rs9848716

SNP in Transcript Sequence SEQ ID NO: 11

15 SNP Position Transcript: 342

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (T,40|C,80)

SNP Type: Silent Mutation

20 Protein Coding: SEQ ID NO: 31, at position 86,(Y,TAC) (Y,TAT)

Gene Number: 7

Celera Gene: hCG32403 - 206000045254807

Celera Transcript: hCT2333142 - 206000045254809

25 Public Transcript Accession:

Celera Protein: hCP1877078 - 206000045254796

Public Protein Accession:

Gene Symbol: CASP8AP2

Protein Name: CASP8 associated protein 2

30 Celera Genomic Axis: GA_x54KRFTF0F9 (80703289..80769434)

Chromosome: 6

OMIM NUMBER: 606880

OMIM Information:

Transcript Sequence (SEQ ID NO: 12):

Protein Sequence (SEQ ID NO: 32):

SNP Information

5 Context (SEQ ID NO: 60):

AGATCCTGACTGAAGAAGGAACTGCAAAGGAGGCAACATATAATGATTTGCA
AGTAGAATATGGAAAATGTCAACTACAAATGAAAGAGCTGATGAAAAA

R

TTTAAAGAAATACAGACACAGAATTTTCAGCTTAATAAACGAAAACCAGTCTC

10 TTAAGAAGAATATTTTCAGCACTTATCAAAACTGCCAGAGTGGAAATAA

Celera SNP ID: hCV2744023

Public SNP ID: rs369328

SNP in Transcript Sequence SEQ ID NO: 12

SNP Position Transcript: 475

15 SNP Source: dbSNP; Celera; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,50|G,70)

SNP Type: Silent Mutation

Protein Coding: SEQ ID NO: 32, at position 93,(K,AAA) (K,AAG)

20 Gene Number: 7

Celera Gene: hCG32403 - 206000045254807

Celera Transcript: hCT2348015 - 206000045254845

Public Transcript Accession: NM_012115

Celera Protein: hCP1913265 - 206000045254798

25 Public Protein Accession: NP_036247

Gene Symbol: CASP8AP2

Protein Name: CASP8 associated protein 2

Celera Genomic Axis: GA_x54KRFTF0F9 (80703289..80769434)

Chromosome: 6

30 OMIM NUMBER: 606880

OMIM Information:

Transcript Sequence (SEQ ID NO: 13):

Protein Sequence (SEQ ID NO: 33):

SNP Information

Context (SEQ ID NO: 61):

AGATCCTGACTGAAGAAGGAACTGCAAAGGAGGCAACATATAATGATTTGCA
 5 AGTAGAATATGGAAAATGTCAACTACAAATGAAAGAGCTGATGAAAAA
 R
 TTAAAGAAATACAGACACAGAATTTTCAGCTTAATAAACGAAAACCAGTCTC
 TTAAGAAGAATATTTTCAGCACTTATCAAAACTGCCAGAGTGGAAATAA

Celera SNP ID: hCV2744023

10 Public SNP ID: rs369328

SNP in Transcript Sequence SEQ ID NO: 13

SNP Position Transcript: 472

SNP Source: dbSNP; Celera; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,50|G,70)

15 SNP Type: Silent Mutation

Protein Coding: SEQ ID NO: 33, at position 93,(K,AAA) (K,AAG)

Gene Number: 7

Celera Gene: hCG32403 - 206000045254807

20 Celera Transcript: hCT2348016 - 206000045254795

Public Transcript Accession:

Celera Protein: hCP1913266 - 206000045254794

Public Protein Accession:

Gene Symbol: CASP8AP2

25 Protein Name: CASP8 associated protein 2

Celera Genomic Axis: GA_x54KRFTF0F9 (80703289..80769434)

Chromosome: 6

OMIM NUMBER: 606880

OMIM Information:

30

Transcript Sequence (SEQ ID NO: 14):

Protein Sequence (SEQ ID NO: 34):

SNP Information

Context (SEQ ID NO: 62):

AGATCCTGACTGAAGAAGGAACTGCAAAGGAGGCAACATATAATGATTTGCA
AGTAGAATATGGAAAATGTCAACTACAAATGAAAGAGCTGATGAAAAA

R

5 TTAAAGAAATACAGACACAGAATTCAGCTTAATAAACGAAAACCAGTCTC
TTAAGAAGAATATTCAGCACTTATCAAAACTGCCAGAGTGGAAATAA

Celera SNP ID: hCV2744023

Public SNP ID: rs369328

SNP in Transcript Sequence SEQ ID NO: 14

10 SNP Position Transcript: 1716

SNP Source: dbSNP; Celera; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,50|G,70)

SNP Type: Silent Mutation

Protein Coding: SEQ ID NO: 34, at position 93,(K,AAA) (K,AAG)

15

Gene Number: 7

Celera Gene: hCG32403 - 206000045254807

Celera Transcript: hCT23591 - 206000045254808

Public Transcript Accession: NM_012115

20 Celera Protein: hCP46427 - 206000045254792

Public Protein Accession: NP_036247

Gene Symbol: CASP8AP2

Protein Name: CASP8 associated protein 2

Celera Genomic Axis: GA_x54KRFTF0F9 (80703289..80769434)

25 Chromosome: 6

OMIM NUMBER: 606880

OMIM Information:

Transcript Sequence (SEQ ID NO: 15):

30 Protein Sequence (SEQ ID NO: 35):

SNP Information

Context (SEQ ID NO: 63):

AGATCCTGACTGAAGAAGGAACTGCAAAGGAGGCAACATATAATGATTTGCA
AGTAGAATATGGAAAATGTCAACTACAAATGAAAGAGCTGATGAAAAA

R

TTTAAAGAAATACAGACACAGAATTCAGCTTAATAAACGAAAACCAGTCTC
5 TTAAGAAGAATATTCAGCACTTATCAAAACTGCCAGAGTGGAAATAA

Celera SNP ID: hCV2744023

Public SNP ID: rs369328

SNP in Transcript Sequence SEQ ID NO: 15

SNP Position Transcript: 475

10 SNP Source: dbSNP; Celera; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,50|G,70)

SNP Type: Silent Mutation

Protein Coding: SEQ ID NO: 35, at position 93,(K,AAA) (K,AAG)

15 Gene Number: 8

Celera Gene: hCG40122 - 104000117872607

Celera Transcript: hCT2334641 - 104000117872706

Public Transcript Accession: NM_004347

Celera Protein: hCP1851623 - 197000069364757

20 Public Protein Accession: NP_004338

Gene Symbol: CASP5

Protein Name: caspase 5, apoptosis-related cysteine protease

Celera Genomic Axis: GA_x5YUV32VVY5 (14887755..14936958)

Chromosome: 11

25 OMIM NUMBER: 602665

OMIM Information:

Transcript Sequence (SEQ ID NO: 16):

Protein Sequence (SEQ ID NO: 36):

30

SNP Information

Context (SEQ ID NO: 64):

GGAATACCTGGGCAAAGATGTTCTTCATGGTGTTTTTAATTATTTGGCAAAC
ACGATGTTCTGACATTGAAGGAAGAGGAAAAGAAAAAATATTATGAT

R

CCAAAATTGAAGACAAGGCCCTGATCTTGGTAGACTCTTTGCGAAAGAATCG
CGTGGCTCATCAAATGTTTACCCAAACACTTCTCAATATGGACCAAAA

Celera SNP ID: hCV12092542

5 Public SNP ID: rs507879

SNP in Transcript Sequence SEQ ID NO: 16

SNP Position Transcript: 302

SNP Source: Applera

10 Population (Allele,Count): caucasian (G,3|A,5) african american (G,4|A,4) total
(G,7|A,9)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 36, at position 90,(A,GCC) (T,ACC)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (A,65|G,53)

15 SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 36, at position 90,(A,GCC) (T,ACC)

Gene Number: 8

Celera Gene: hCG40122 - 104000117872607

20 Celera Transcript: hCT2334643 - 104000117872608

Public Transcript Accession: NM_004347

Celera Protein: hCP1851615 - 197000069364749

Public Protein Accession: NP_004338

Gene Symbol: CASP5

25 Protein Name: caspase 5, apoptosis-related cysteine protease

Celera Genomic Axis: GA_x5YUV32VVY5 (14887755..14936958)

Chromosome: 11

OMIM NUMBER: 602665

OMIM Information:

30

Transcript Sequence (SEQ ID NO: 17):

Protein Sequence (SEQ ID NO: 37):

SNP Information

Context (SEQ ID NO: 65):

GGAATACCTGGGCAAAGATGTTCTTCATGGTGTTTTTAATTATTTGGCAAAC
ACGATGTTCTGACATTGAAGGAAGAGGAAAAGAAAAAATATTATGAT

R

5 CCAAATTGAAGACAAGGCCCTGATCTTGGTAGACTCTTTGCGAAAGAATCG
CGTGGCTCATCAAATGTTTACCCAAACACTTCTCAATATGGACCAAAA

Celera SNP ID: hCV12092542

Public SNP ID: rs507879

SNP in Transcript Sequence SEQ ID NO: 17

10 SNP Position Transcript: 302

SNP Source: Applera

Population (Allele,Count): caucasian (G,3|A,5) african american (G,4|A,4) total
(G,7|A,9)

SNP Type: Missense Mutation

15 Protein Coding: SEQ ID NO: 37, at position 90,(A,GCC) (T,ACC)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (A,65|G,53)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 37, at position 90,(A,GCC) (T,ACC)

20

Gene Number: 8

Celera Gene: hCG40122 - 104000117872607

Celera Transcript: hCT31376 - 104000117872684

Public Transcript Accession:

25 Celera Protein: hCP49901 - 197000069364755

Public Protein Accession:

Gene Symbol: CASP5

Protein Name: caspase 5, apoptosis-related cysteine protease

Celera Genomic Axis: GA_x5YUV32VVY5 (14887755..14936958)

30 Chromosome: 11

OMIM NUMBER: 602665

OMIM Information:

Transcript Sequence (SEQ ID NO: 18):

Protein Sequence (SEQ ID NO: 38):

SNP Information

Context (SEQ ID NO: 66):

5 GGAATACCTGGGCAAAGATGTTCTTCATGGTGTTTTTAATTATTTGGCAAAC
 ACGATGTTCTGACATTGAAGGAAGAGGAAAAGAAAAAATATTATGAT
 R
 CCAA AATTGAAGACAAGGCCCTGATCTTGGTAGACTCTTTGCGAAAGAATCG
 CGTGGCTCATCAAATGTTTACCCAAACACTTCTCAATATGGACCAAAA

10 Celera SNP ID: hCV12092542

Public SNP ID: rs507879

SNP in Transcript Sequence SEQ ID NO: 18

SNP Position Transcript: 174

SNP Source: Applera

15 Population (Allele,Count): caucasian (G,3|A,5) african american (G,4|A,4) total
 (G,7|A,9)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 38, at position 48,(A,GCC) (T,ACC)

SNP Source: dbSNP; Celera; HapMap; HGBASE

20 Population (Allele,Count): caucasian (A,65|G,53)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 38, at position 48,(A,GCC) (T,ACC)

Gene Number: 9

25 Celera Gene: hCG2039683 - 209000071858797

Celera Transcript: hCT2301261 - 209000071858813

Public Transcript Accession: NM_031910

Celera Protein: hCP1797817 - 209000071858790

Public Protein Accession: NP_114116

30 Gene Symbol: C1QTNF6

Protein Name: C1q and tumor necrosis factor related protein 6

Celera Genomic Axis: GA_x5YUV32VU5C (15633444..15661565)

Chromosome: 22

OMIM NUMBER:

OMIM Information:

Transcript Sequence (SEQ ID NO: 19):

Protein Sequence (SEQ ID NO: 39):

5

SNP Information

Context (SEQ ID NO: 67):

TCGCCTGGGGAGGCCACAGGACACAGGGTCACCATGGGGACAGCCGCCCTGG
 GTCCCGTCTGGGCAGCGCTCCTGCTCTTTCTCCTGATGTGTGAGATCC

10

S

TATGGTGGAGCTCACCTTTGACAGAGCTGTGGCCAGCGGCTGCCAACGGTGC
 TGTGACTCTGAGGACCCCCTGGATCCTGCCCATGTATCCTCAGCCTCT

Celera SNP ID: hCV2403368

Public SNP ID: rs229526

15

SNP in Transcript Sequence SEQ ID NO: 19

SNP Position Transcript: 201

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (G,95|C,25)

SNP Type: Missense Mutation

20

Protein Coding: SEQ ID NO: 39, at position 42,(P,CCT) (R,CGT)

Gene Number: 9

Celera Gene: hCG2039683 - 209000071858797

Celera Transcript: hCT2301264 - 209000071858835

25

Public Transcript Accession: NM_031910

Celera Protein: hCP1797819 - 209000071858789

Public Protein Accession: NP_114116

Gene Symbol: C1QTNF6

Protein Name: C1q and tumor necrosis factor related protein 6

30

Celera Genomic Axis: GA_x5YUV32VU5C (15633444..15661565)

Chromosome: 22

OMIM NUMBER:

OMIM Information:

Transcript Sequence (SEQ ID NO: 20):

Protein Sequence (SEQ ID NO: 40):

SNP Information

5 Context (SEQ ID NO: 68):

TCGCCTGGGGAGGCCACAGGACACAGGGTCACCATGGGGACAGCCGCCCTGG
GTCCCGTCTGGGCAGCGCTCCTGCTCTTTCTCCTGATGTGTGAGATCC

S

TATGGTGGAGCTCACCTTTGACAGAGCTGTGGCCAGCGGCTGCCAACGGTGC
10 TGTGACTCTGAGGACCCCCTGGATCCTGCCCATGTATCCTCAGCCTCT

Celera SNP ID: hCV2403368

Public SNP ID: rs229526

SNP in Transcript Sequence SEQ ID NO: 20

SNP Position Transcript: 201

15 SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (G,95|C,25)

SNP Type: Missense Mutation

Protein Coding: SEQ ID NO: 40, at position 42,(P,CCT) (R,CGT)

Table 2. Genomic SNP info and associated gene information

	Gene Number:	1
5	Celera Gene:	hCG14694 - 208000030293584
	Gene Symbol:	F9
	Protein Name:	coagulation factor IX (plasma thromboplastic component, Christmas disease, hemophilia B)
	Celera Genomic Axis:	GA_x5YUV32W21H (1443186..1495938)
10	Chromosome:	X
	OMIM NUMBER:	306900
	OMIM Information:	Hemophilia B (3); Warfarin sensitivity (3)
	Genomic Sequence (SEQ ID NO: 69):	
15	SNP Information	
	Context	(SEQ ID NO: 82):
	TGAGAAATATCAGGTTACTAATTTTTCTTCTATTTTTCTAGTGCCATTTCATG	
	TGGAAGAGTTTCTGTTTCACAACTTCTAAGCTCACCCGTGCTGAG	
20	R	
	CTGTTTTTCCTGATGTGGACTATGTAAATTCTACTGAAGCTGAAACCATTTTG	
	GATAACATCACTCAAAGCACCCAATCATTTAATGACTTCACTCGGGT	
	Celera SNP ID:	hCV596331
	Public SNP ID:	rs6048
25	SNP in Genomic Sequence:	SEQ ID NO: 69
	SNP Position Genomic:	30398
	SNP Source:	dbSNP; HapMap;
	Population (Allele,Count):	caucasian (A,61 G,29)
	SNP Type:	MISSENSE MUTATION; HUMAN-MOUSE SYNTENIC
30	REGION; INTRON	
	Context	(SEQ ID NO: 83):
	CTTAGAACCTAATGAAAGTTTGCATTTCCTCAGTAAAATCAGAGACTGCTGATT	
	GACTTAAATGTTTATAGCTTCAAAGTCCTCCTCATTATCATGGCCCA	
35	S	

AAGCCCTTCCATGATTGTCCTTCCCCACCCTCCCCATTACCCTTCTTGCCTCCT
CTGCTACTTCTCTCCTCGCACACTGGGCTCCAGCCACCCTGGCCTT

Celera SNP ID: hCV596323

Public SNP ID: rs438601

5 SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 17052

Related Interrogated SNP: hCV596331 (Power=.6)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (C,71|G,19)

10 SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 84):

CATATAATGGGAATTCTCCACATGTACAAACCACTTCATATGCTAAACTTGTT
GACAACATTCAAAGCTCATCCCTGAATTTGACTATATTGATTACATC

15 R

AAAATGTTACATAGCAACCTTAGAATCCTTGTGTACCTTTTCTTCTCAAAGCC
TAGATTATTTCTTTTTCCGACGTTTCAGTAATTGGAGCAGTAAACCC

Celera SNP ID: hCV596326

Public SNP ID: rs398101

20 SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 21090

Related Interrogated SNP: hCV596331 (Power=.8)

SNP Source: dbSNP; HGBASE

Population (Allele,Count): caucasian (A,69|G,31)

25 SNP Type: INTRON

Context (SEQ ID NO: 85):

CTAGAAGGCCTTTTAGTCTGCAAAGAAACCTTCTTAATCATAAGCAGCAGA
AGTCCCATTACCAAATTGGAAAGTTAAAGTTACAAAGCATCAATCAT

30 M

AGACTTCCATTACAGGGATGGCAATTGGGAGTAAGACTTTTTAGTAAAGAAAC
TAAACACAAAGTCATTAGACTCTGTAAAAGTCTTACCAAATTTGATTC

Celera SNP ID: hCV596330

Public SNP ID: rs422187

SNP in Genomic Sequence: SEQ ID NO: 69
 SNP Position Genomic: 29992
 Related Interrogated SNP: hCV596331 (Power=.9)
 SNP Source: dbSNP; HapMap; HGBASE
 5 Population (Allele,Count): caucasian (A,61|C,29)
 SNP Type: INTRON

Context (SEQ ID NO: 86):

10 GCCAGAGATCAGAGCAGGCTAAGGGACTGCTGGGATCCTGTCCAGCTTTGAG
 ACCCTACAGAGCCATGTTACCTAGCACGTATCCCGTCTGCGGTCACG
 S
 TCATTTCTTACCTTATTCCAGGGCTTTCACCTCAGCTTGCCAGGCTGGAGCCA
 AGGGCCAACGCAGCCGCGCCTTGTTTCGCGATGGTAGCTTCCCAGGAG

Celera SNP ID: hCV596335

15 Public SNP ID: rs413957

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 34649

Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

20 Population (Allele,Count): caucasian (C,73|G,17)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

Context (SEQ ID NO: 87):

25 TACAGAAAATGTCCAGGGAAATGGTCTATTTCTTATTCTATTTTTGACCTAAA
 GAAAATCTTTAAAATGTCTTAGCATTTCCTCCAGTCTCCATCCACTT
 Y
 CCTCAGCTTTGGCCTGAAGCTATCTTTAAAGGTACCCTGTACAGCTCTTGCCC
 TGTACAGCTAGCTACAGAGATTCAATCCTTTCTGTTCGATTAGGACA

Celera SNP ID: hCV596336

30 Public SNP ID: rs110583

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 35664

Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; HGBASE

Population (Allele,Count): caucasian (C,17|T,91)

SNP Type: INTRON

Context (SEQ ID NO: 88):

5 ATGGTTAAGAGAGAGAGTGGAAGAATGAATGAGCCCTGCTATTCCTCACTGCC
TGGATGGCTATAAGCACAGCCCTTATGGAGGCCTTAGGTCTTGCTTCA
Y
AATATTCCAGTTTGAAAAGGGTTTGAAAAGACCTCCTAGAAAAATCAGTAGT
TTTTCTCTTTGAGTAACATGTAGCAAAAAAATTCATCATGTAGGT

10 Celera SNP ID: hCV596339

Public SNP ID: rs370713

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 39526

Related Interrogated SNP: hCV596331 (Power=.51)

15 SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (T,73|C,17)

SNP Type: INTRON

Context (SEQ ID NO: 89):

20 AAAGGGTTTGAAAAGACCTCCTAGAAAAATCAGTAGTTTTCTCTTTGAGTA
ACATGTAGCAAAAAAATTCATCATGTAGGTACAGGGAACACCCTA
R
TAACTATTAATCTCAAGGAGTCAAGCCAGTGTGTTTCCTAATGTATCTGCTGT
ATCCCCATGAAGCAAATTTGCCATCAGAGAACTGACTCATGGGGA

25 Celera SNP ID: hCV596340

Public SNP ID: rs413536

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 39642

Related Interrogated SNP: hCV596331 (Power=.51)

30 SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,73|G,17)

SNP Type: INTRON

Context (SEQ ID NO: 90):

GCTAGACAGTAGTTGCTCAATAATTGTTAGCTGAATCAGAATCCATGTTTATC
CCAGAGTAGCAATTAGTCTTGCATCGAGTATCGTGAAAGAAGGCCAC

R

CTTAAATAAGAATAATGCCTGGGGTTTAGGTTTTATGAAAAAATGAAAGGAA
5 ATTAGTTCTGCTTTTGTGACTAAAGGAAGGGAAGAGAGAAGAGACTA

Celera SNP ID: hCV596344

Public SNP ID: rs445691

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 43879

10 Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; HapMap; HGBASE

Population (Allele,Count): caucasian (G,73|A,17)

SNP Type: INTERGENIC; UNKNOWN

15 Context (SEQ ID NO: 91):

TAAAGTGAACAGCTGCAATGAAAATAAGGGAAGAAAGTTTAGTTCATCTCCG
TTTCTTTCCTTTCCTTTTACTTTCCTTTCCTTTCCTTTTTGGAGTTA

R

TCAGGAAGTAGTCCCAAATACCCAGAAAGTTCATCTTATAAGCCCTTGGTCC
20 TCTTGAGATGGTATCAGATATATTGCTAGACCCTTGAAGAAAGGAAC

Celera SNP ID: hCV596669

Public SNP ID: rs376165

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 14866

25 Related Interrogated SNP: hCV596331 (Power=.8)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,55|G,35)

SNP Type: INTRON

30 Context (SEQ ID NO: 92):

CCTTTTACCCTCCATGGTCGTTAAAGGAGAGATGGGGAGCATCATTCTGTTAT
ACTTCTGTACACAGTTATACATGTCTATCAAACCCAGACTTGCTTCC

R

TAGTGGAGACTTGCTTTTCAGAACATAGGGATGAAGTAAGGTGCCTGAAAAG
TTTGGGGGAAAAGTTTCTTTCAGAGAGTTAAGTTATTTTATATATATA

Celera SNP ID: hCV2288124

Public SNP ID: rs440051

5 SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 42050

Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE;

Population (Allele,Count): caucasian (G,73|A,17)

10 SNP Type: UTR3

Context (SEQ ID NO: 93):

ACAGTGGTCTGAATCCACCTGAGACAGAATTGGGTCTAACTAACTGTGAGTA
TGGCCTTCAATAAGTCACTCTCCATTTGGGAATTTGATTTCTCCACTT

15 S

TATAATGAGAGTATTTGACAGGATGCTCTCCCAAATCCCTTGCAATTTTGTTA
GTCTGTGATTTTCATGTTTTTATTTTATTCCTTCATCCAACAAATAG

Celera SNP ID: hCV2969899

Public SNP ID: rs434144

20 SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 43558

Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (C,71|G,19)

25 SNP Type: INTERGENIC; UNKNOWN; REPEATS

Context (SEQ ID NO: 94):

CAGGATGCTCTCCCAAATCCCTTGCAATTTTGTTAGTCTGTGATTTTCATGTTTT
TATTTTATTCCTTCATCCAACAAATAGTCAAGGAGTAATTGCTGT

30 S

TGCCAAATACCAACAGTATTCATTAAATTGTAATTCAGATTTTATATATATAT
AAATAATGTATAATGTGTATAAATTGCTTTGTGAGTGCCTACTACAC

Celera SNP ID: hCV2969900

Public SNP ID: rs434447

SNP in Genomic Sequence: SEQ ID NO: 69
 SNP Position Genomic: 43677
 Related Interrogated SNP: hCV596331 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap; HGBASE
 5 Population (Allele,Count): caucasian (C,73|G,17)
 SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 95):

10 AGAACTATTTCAAACCTGGCCAGGTCATTCCACTCTAATAGGAGAGCTATC
 CTTCTATTCTCTTGGTTAAGAGAAAGCTGGAAATAGAAACAGGGATA
 Y
 TCCAGGCCAGACACAATGGCTCACGCCTGTAATCCCAACACTTTGGGAGGCC
 GAGGCAGGCAGATCACTTGAGGTCAGGAGTTCAAGATTAGTCTGGCCA

Celera SNP ID: hCV2986569

15 Public SNP ID: rs11095801

SNP in Genomic Sequence: SEQ ID NO: 69
 SNP Position Genomic: 45953
 Related Interrogated SNP: hCV596331 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap
 20 Population (Allele,Count): caucasian (T,73|C,17)
 SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 96):

25 GACCTCATAAAGATAAAGAGTTCCCATGATTTACACATTACGTGCTTTCATAA
 ATATCTATATATAAAAGCCTATTTTCCTCTTGGACTATATTACAAAA
 R
 TAAGTATGCATTTTCATAAGATTCAAACCCAGCTCTAAATGTAAGAAGCCAA
 ATTAAGAATGTAGCAATGTATGAATGAAAAGGAAGGAAAAAAGCCATT

Celera SNP ID: hCV2986570

30 Public SNP ID: rs3117458

SNP in Genomic Sequence: SEQ ID NO: 69
 SNP Position Genomic: 46612
 Related Interrogated SNP: hCV596331 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (G,72|A,18)

SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 97):

5 ATACTTGATTCAAACCTATTTCTGTCTGATCTGATTCTAAAGTCTGTTTTTTCA
CTCAACCACACTGTACAGTCAGCTCTCCTTGTGAGTTCCACAGCCA

M

AGATTCAATTAAGTGCAGATCAAAAATATTCAAGAAAAAATGGATGGTTGC
ATCTCTACTGAACATGTACAGACTCTTTTATCTTTCATTATTCCCTAA

10 Celera SNP ID: hCV2986572

Public SNP ID: rs4149670

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 13775

Related Interrogated SNP: hCV596331 (Power=.51)

15 SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (A,31|C,69)

SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 98):

20 TCTCCTTGTGAGTTCCACAGCCACAGATTCAATTAAGTGCAGATCAAAAATAT
TCAAGAAAAAATGGATGGTTGCATCTCTACTGAACATGTACAGACT

Y

TTTTATCTTTCATTATTCCCTAAACAATACAGCATAACAACACTATTACATAGCA
TTTACATTGTATTAGCTATTAAGAGAAACCTAGAGATGATTTAAAG

25 Celera SNP ID: hCV2986574

Public SNP ID: rs4149672

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 13852

Related Interrogated SNP: hCV596331 (Power=.7)

30 SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (C,64|T,32)

SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 99):

ACTGAACATGTACAGACTCTTTTATCTTTCATTATTCCCTAAACAATACAGCA
TAACAAC TATTACATAGCATTACATTGTATTAGCTATTAAGAGAA

W

CCTAGAGATGATTTAAAGTACAAAGGAGGATGTGTTTAGGTTATATGCAAAT
5 AGTAAGCCATTTTATATCGGAGACTTGAGCATCCACAGATCTTGATA

Celera SNP ID: hCV2986575

Public SNP ID: rs4149674

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 13934

10 Related Interrogated SNP: hCV596331 (Power=.7)

SNP Source: dbSNP; Celera; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,56|T,34)

SNP Type: INTRON; REPEATS

15 Context (SEQ ID NO: 100):

ATGTTGTGTAAACTGTGGAACTATTGAAAGAATCACAGCAGGCAAAGGACT
GTGGGACCCCTGCTCTTTTCAATCAATCCAGGCCCAAAATCACTCTA

S

TCATTTTTTCCTACGGTAGTTTCGAGGCAAATCTTTTTCATCCAGTCTTTGGGG
20 CTATGGACTGCCTTGAGATTTCTGAGTCAATCTCAACTTCTAATAT

Celera SNP ID: hCV26225376

Public SNP ID: rs3117074

SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 45644

25 Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (G,73|C,17)

SNP Type: INTERGENIC; UNKNOWN

30 Context (SEQ ID NO: 101):

TGGAAACTATTGAAAGAATCACAGCAGGCAAAGGACTGTGGGACCCCTGCTC
TTTTCAATCAATCCAGGCCCAAAATCACTCTAGTCATTTTTTCCTAC

R

GTAGTTTCGAGGCAAATCTTTTTCATCCAGTCTTTGGGGCTATGGACTGCCTT
GAGATTTCTGAGTCAATCTCAACTTCTAATATAGGTATCAAAAAGTGA

Celera SNP ID: hCV26225377

Public SNP ID: rs12008759

5 SNP in Genomic Sequence: SEQ ID NO: 69

SNP Position Genomic: 45659

Related Interrogated SNP: hCV596331 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,73|A,17)

10 SNP Type: INTERGENIC; UNKNOWN

Gene Number: 2

Celera Gene: hCG1776858 - 104000116701164

Gene Symbol: CML2

15 Protein Name: putative N-acetyltransferase Camello 2

Celera Genomic Axis: GA_x5YUV32W5TR (13100168..13120998)

Chromosome: 2

OMIM NUMBER:

OMIM Information:

20

Genomic Sequence (SEQ ID NO: 70):

SNP Information

Context (SEQ ID NO: 102):

25 GTCCACAGAGAGATGAAACAGCTGCAACCGCTTCTCCCTCAAGGTGGGATCA

TCAACGGGCAGAGCTCCTACTGTGCCACACCTTCTCTTCAGATTCA

S

CCACCCAGAAGCAGGAGCCCACTCACTCAGGTAGGATTTGGTGATGTCAGA

CATGTCTGTGCGCAATGCTATGTCTACATAACGCGTCCAGGGTTTTTT

30

Celera SNP ID: hCV11541681

Public SNP ID: rs2001490

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 10459

SNP Source: Applera

Population (Allele,Count): caucasian (C,10|G,10) african american (C,18|G,14) total (C,28|G,24)

SNP Type: MISSENSE MUTATION; ESS; TRANSCRIPTION FACTOR BINDING SITE; HUMAN-MOUSE SYNTENIC REGION; UTR3

5 SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (C,40|G,80)

SNP Type: MISSENSE MUTATION; ESS; TRANSCRIPTION FACTOR BINDING SITE; HUMAN-MOUSE SYNTENIC REGION; UTR3

10 Context (SEQ ID NO: 103):

TGAGTGACAGTGAGCACTGCCACAGAGGGGGTGCCCTGCAGTCAGGAGGGG
TGGGAAACGGGTGCAGATCGGCAGCCAACAGAGGCCACCCACATGAG

Y

ACAGGTGGAGTTGCTGGTGCGATGAAGAGTGTGCTGGTTGGGAGTGGAGGGC

15 TGGGGGCCGGTAGCCGGTTTGGGAAACCCAAGCTGAGCGAAAGGTGTG

Celera SNP ID: hCV95670

Public SNP ID: rs4852975

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 14848

20 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,40|T,80)

SNP Type: INTERGENIC; UNKNOWN

25 Context (SEQ ID NO: 104):

CTTGTTGAAAAGCCACCTAAAATTGTCATAAGTAAGAGAGGGTGATCTAATGG
CCCATGAAACTGCCAACTCCAGGAAGACAGACAGGGTGACCTGAAAGT

R

GGTTTGACCTGTTGGTTTTTCTGCCATTCTGTGGGTGGCATCAGCTTCCTGCCA

30 AGACTGGCACTGCATAGGATGTCGGGTGTTTGCCAGCCCAGTTGAA

Celera SNP ID: hCV95671

Public SNP ID: rs11126414

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 14427

Related Interrogated SNP: hCV11541681 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap; ABI_Val
 Population (Allele,Count): caucasian (G,40|A,80)
 SNP Type: INTERGENIC; UNKNOWN

5

Context (SEQ ID NO: 105):

CAGCCGACCTGGGGCTCCACAGCTCTTGCTGCCGCTGGAAGTTGGGCCAAGG
 CTGGTGGCACACACAGGTGTGTTCTGCTGTTGGGATCACCAGACAAG
 K

10

CTGCTTCGCCTTGGGCCGATCTTTTGCCCCGGACCACTCAGCCCTAGGTACTTT
 AGGCCAATCTGCTGATCAAGGCCATAAGTGCCTCATAACCAATCAG

Celera SNP ID: hCV11541694
 Public SNP ID: rs12619258
 SNP in Genomic Sequence: SEQ ID NO: 70

15

SNP Position Genomic: 18301
 Related Interrogated SNP: hCV11541681 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (G,40|T,80)
 SNP Type: INTERGENIC; UNKNOWN

20

Context (SEQ ID NO: 106):

CTCAGCTTCTTGGTGCATTGAGAGAGTACCTACCTCAGGCAGGCCACATGCCC
 AACCCCTTGGTGCTTCTCCAGGGACCATAACATCGCAGAAATCCTGTG
 W

25

CCAAGCACAAGTTTCAGGAAAGGCTCCCTCCAGCTTTATTGTAAGTTTTTAGG
 TTGCTGAGGTGGCATAAGTGCCAAGCATAAAAAACAGCCCAGCCTCC

Celera SNP ID: hCV11941453
 Public SNP ID: rs2001436
 SNP in Genomic Sequence: SEQ ID NO: 70

30

SNP Position Genomic: 11326
 Related Interrogated SNP: hCV11541681 (Power=.51)
 SNP Source: dbSNP; Celera; HGBASE
 Population (Allele,Count): caucasian (T,40|A,76)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTERGENIC;
UNKNOWN

Context (SEQ ID NO: 107):

5 TAAACTTTACTATTGTTGAACTGCCTATGTTGCCTGAACTTTGACTGGTTTAAG
TTGGTTTCATTTATTTTCATAAGAACTCTTGTTAAGGGGTGCACCTT

K

GTATTTTGATATTATTCTCTTGATAGTCATAATAATAGTTCCTGGTGTGCTGC
ATCCTCTCAAGTCATAAATGTTTTGTATGCAGCCATACATTGAGAA

10 Celera SNP ID: hCV26996674

Public SNP ID: rs13006448

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 6785

Related Interrogated SNP: hCV11541681 (Power=.51)

15 SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,40|T,80)

SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 108):

20 ACAATTCCGCTAATGCCAATGCCTCTCCATGTTGCCAATGAGAGTCACCCCAA
AATTCCAAGGGAAGGATAGAAAGCTATCCTTGATATTCTGAACACCA

Y

TGCTACTTGTTTTCTGCACTCACTAAAAGCAACACTTTAACTTTTCCTTTTTT
TTTTTTTTTCTGAGGCAGTCTCACTCTGTTGCCTAGGTTGGAGTGC

25 Celera SNP ID: hCV26996679

Public SNP ID: rs6732812

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 4836

Related Interrogated SNP: hCV11541681 (Power=.51)

30 SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (T,40|C,80)

SNP Type: INTRON

Context (SEQ ID NO: 109):

ATGATTTTGGCATGATAGGGAAACTCTAAGGGTTTCTCACTTCTAGTGGGATC
 ACCATCAAAAATGGACTCCAGTTAGATAAACTCCTAGTTTGCTGTTA

K

TGATCGAAGCTCATACTGCAGAACTGAACCTGAATATCAAGGGAATGCTTA
 5 GCAGATATTTATGCTAAATCAGCTAGTACTGAAACTGTTTCAGATATGC

Celera SNP ID: hCV26996688

Public SNP ID: rs13015885

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 2839

10 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,39|T,79)

SNP Type: INTRON

15 Context (SEQ ID NO: 110):

CATGCCACTGTATTACTGCTGCCTGATGATGGTGAAGACCACAATTGCATAAG
 TGTAGCATCAGAAATAGTGGCCCCTCATGTTAATTTACAAGTTAGTC

S

TTTGGACAATCCTGAGTTAATACTTTTGTTGATGGGTCCTATGCCAAAAGCTC
 20 AGAAGGAAAATATCAGCTAGGATATGCTGTTACCAAAAATGAGTTAAT

Celera SNP ID: hCV26996689

Public SNP ID: rs13014700

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 2478

25 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: Celera; HapMap

Population (Allele,Count): caucasian (C,39|G,79)

SNP Type: INTRON

30 Context (SEQ ID NO: 111):

CCTAAAAGTTAAAACCTATTCAAATTTTCCTTGGCCTGCAACCAAAAAGATA
 ATTAATAGTTTTCTTGGACTTGCAGGATATGTTTTTAGACTTGCAG

R

ATATAAAATTCCTGGGTTCTGAATTTTTCCTTAATAGCCTCACCATTTCATGAG
CTCCCTAAAAAATGCTGTACCAGAGCCTTTATCTTGGGATGATAGT

Celera SNP ID: hCV26996690

Public SNP ID: rs2421575

5 SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 1879

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (G,40|A,80)

10 SNP Type: INTRON

Context (SEQ ID NO: 112):

TCTACTGTGACTTCTGCATTTTGTAGCTTTCCTGATAAAAGCCTGGGCTTCC
TTGTCAGACCTGAAAGATTCATTTATTCATTCAACAAACAGTTCCC

15 R

AGGACCTGAGGTGTGCCGGGCCTGGACTTGGCATGAAGGCACCAGATGTTGG
AAGCGAGGCTGCCGCCAGGAGGACACACCTGATGGGGCTCTGGGAAC

Celera SNP ID: hCV31840149

Public SNP ID: rs12233112

20 SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 14147

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (G,40|A,80)

25 SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTERGENIC;
UNKNOWN; REPEATS

Context (SEQ ID NO: 113):

AGATCCCTCAGAGTTAAAAAGGCTATTTCTGTGGTTGGGGTTTCAAATCAAAT
30 TCAAGAGGTTCCCATATCTGAACCCGTCCAATTGACTTTGGGGCACT

Y

TTCAGAAAATCACACTTTTTTACTGTATGATAGTGCTCCAGTAAATTTGCTAG
AGGGAGATTTACTTTCAAAGCTGAAAGGGCATAGAAGACTACCTATT

Celera SNP ID: hDV69785784

Public SNP ID: rs13000788
 SNP in Genomic Sequence: SEQ ID NO: 70
 SNP Position Genomic: 669
 Related Interrogated SNP: hCV11541681 (Power=.51)
 5 SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (C,39|T,79)
 SNP Type: INTRON

Context (SEQ ID NO: 114):

10 CTGCCCAACTTGCAAAATATTTATTTTTCTTGCTCAATTACTCTGTGAAATTTA
 ATTTGTCTAATGGATTTCTTTTAACAGAACTTATAATTAATAATGGA

R

AGCGGGTGTAATATTTATAAACTTGCAGCCTGGCCATGTGGTAGACAAGG
 AGGAATCCAAGCAGCCTGCTGAGCAACTACTTGCTGGAGACATTAGCA

15 Celera SNP ID: hCV31840159
 Public SNP ID: rs13013228
 SNP in Genomic Sequence: SEQ ID NO: 70
 SNP Position Genomic: 7577
 Related Interrogated SNP: hCV11541681 (Power=.51)
 20 SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (G,40|A,80)
 SNP Type: INTRON

Context (SEQ ID NO: 115):

25 CTCCCAAAGTTCTGGGATTACAGGCGTGAGCCACCATGCCTGGCCTTAACTGT
 TTCAATTAATAACCTGATCAATACCAAAATATAGAAAACCAGTCCAA

R

TGATGCCTGCAAAAGATATATTGTGCTTTCAGGCATCATGCACTCAAGTTACG
 GGAATTACCTATGTGAGTACAAGTAATTGCTTGTATAATATCACCAG

30 Celera SNP ID: hDV70942181
 Public SNP ID: rs17350056
 SNP in Genomic Sequence: SEQ ID NO: 70
 SNP Position Genomic: 5238
 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (G,40|A,80)
 SNP Type: INTRON

5 Context (SEQ ID NO: 116):
 CTGTTTCAATTAATAACCTGATCAATACCAAATATAGAAAACCAGTCCAAGT
 GATGCCTGCAAAGATATATTGTGCTTTCAGGCATCATGCACTCAAG
 W

TACGGGAATTACCTATGTGAGTACAAGTAATTGCTTGTATAATATCACCAGAT
 10 TAAATTCAACAGGGTCTCTTTTTACTAAATGTGTTATACACCCTTAC

Celera SNP ID: hDV70953030

Public SNP ID: rs17434634

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 5287

15 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (T,40|A,80)

SNP Type: INTRON

20 Context (SEQ ID NO: 117):
 CAGGGTCTCTTTTTACTAAATGTGTTATACACCCTTACAACATACTACAGGGC
 GAGCACAAAAGGTCAATTTCCCAGTGGATTTTGTTCAGGAGCTGTG
 Y

AGGTTTATTAATGAACAAATCTAACTGACCCCTGCTCAAATGCAACTATGTGG
 25 CCCTGTTTTTCAACTCCCAAGGGCCTATACTGGGTCTGTGGATAATC

Celera SNP ID: hDV70953035

Public SNP ID: rs17434655

SNP in Genomic Sequence: SEQ ID NO: 70

SNP Position Genomic: 5450

30 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (T,40|C,80)

SNP Type: INTRON

Gene Number: 3
 Celera Gene: hCG1776867 - 208000043935214
 Gene Symbol: LOC200420
 Protein Name: LOC200420
 5 Celera Genomic Axis: GA_x5YUV32W5TR (13100528..13176592)
 Chromosome: 2
 OMIM NUMBER:
 OMIM Information:

10 Genomic Sequence (SEQ ID NO: 71):

SNP Information

Context (SEQ ID NO: 118):

GTCCACAGAGAGATGAAACAGCTGCAACCGCTTCTCCCTCAAGGTGGGATCA
 15 TCAACGGGCAGAGCTCCTACTGTGCCACACCTTCTCTTCAGATTCA
 S
 CCACCCAGAAGCAGGAGCCCACTCACTCAGGTAGGATTTGGTGATGTCAGA
 CATGTCTGTGCGCAATGCTATGTCTACATAACCGCGTCCAGGGTTTTTT

Celera SNP ID: hCV11541681

20 Public SNP ID: rs2001490

SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 66053

SNP Source: Applera

Population (Allele,Count): caucasian (C,10|G,10) african american (C,18|G,14) total
 25 (C,28|G,24)

SNP Type: MISSENSE MUTATION; ESS; TRANSCRIPTION FACTOR
 BINDING SITE; HUMAN-MOUSE SYNTENIC REGION; UTR3

SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (C,40|G,80)

30 SNP Type: MISSENSE MUTATION; ESS; TRANSCRIPTION FACTOR
 BINDING SITE; HUMAN-MOUSE SYNTENIC REGION; UTR3

Context (SEQ ID NO: 119):

TGAGTGACAGTGAGCACTGCCACAGAGGGGGTGCCCTGCAGTCAGGAGGGG
TGGGAAACGGGTGCAGATCGGCAGCCAACAGAGGCCACCCACATGAG

Y

ACAGGTGGAGTTGCTGGTGCATGAAGAGTGTGCTGGTTGGGAGTGGAGGGC
5 TGGGGGCCGGTAGCCGGTTTGGGAAACCCAAGCTGAGCGAAAGGTGTG

Celera SNP ID: hCV95670

Public SNP ID: rs4852975

SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 70442

10 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,40|T,80)

SNP Type: INTERGENIC; UNKNOWN

15 Context (SEQ ID NO: 120):

CTTGTTGAAAAGCCACCTAAAATTGTCATAAGTAAGAGAGGTGATCTAATGG
CCCATGAAACTGCCAACTCCAGGAAGACAGACAGGGTGACCTGAAAGT

R

GGTTTGACCTGTTGGTTTTTCTGCCATTCTGTGGGTGGCATCAGCTTCCTGCCA
20 AGACTGGCACTGCATAGGATGTCGGGTGTTTGCCAGCCCAGTTGAA

Celera SNP ID: hCV95671

Public SNP ID: rs11126414

SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 70021

25 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (G,40|A,80)

SNP Type: INTERGENIC; UNKNOWN

30 Context (SEQ ID NO: 121):

CAGCCGACCTGGGGCTCCACAGCTCTTGCTGCCGCTGGAAGTTGGGCCAAGG
CTGGTGGCACACACAGGTGTGTTTCCTGCTGTTGGGATCACCAGACAAG

K

CTGCTTCGCCTTGGGCCGATCTTTTGCCCCGGACCACTCAGCCCTAGGTACTTT
AGGCCAATCTGCTGATCAAGGCCATAAGTGCCTCATAACCAATCAG

Celera SNP ID: hCV11541694

Public SNP ID: rs12619258

5 SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 73895

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,40|T,80)

10 SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 122):

CTCAGCTTCTTGGTGCATTGAGAGAGTACCTACCTCAGGCAGGCCACATGCCC
AACCCCTTTGGTGCTTCTCCAGGGACCATAACATCGCAGAAATCCTGTG

15 W

CCAAGCACAAGTTTCAGGAAAGGCTCCCTCCAGCTTTATTGTAAGTTTTTAGG
TTGCTGAGGTGGCATAAGTGCCAAGCATAAAAAACAGCCCAGCCTCC

Celera SNP ID: hCV11941453

Public SNP ID: rs2001436

20 SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 66920

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (T,40|A,76)

25 SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTERGENIC;
UNKNOWN

Context (SEQ ID NO: 123):

30 TAAACTTTACTATTGTTGAACTGCCTATGTTGCCTGAACTTTGACTGGTTTAAG
TTGGTTTCATTTATTTTCATAAGAACTCTTGTTAAGGGGTGCACCTT

K

GTATTTTGATATTATTCTCTTGATAGTCATAATAATAGTTCCCTGGTGTGCTGC
ATCCTCTCAAGTCATAAATGTTTTGTATGCAGCCATACATTGAGAA

Celera SNP ID: hCV26996674

Public SNP ID: rs13006448
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 62379
 Related Interrogated SNP: hCV11541681 (Power=.51)
 5 SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (G,40|T,80)
 SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 124):

10 ACAATTCCGCTAATGCCAATGCCTCTCCATGTTGCCAATGAGAGTCACCCCAA
 AATTCCAAGGGAAGGATAGAAAGCTATCCTTGATATTCTGAACACCA
 Y
 TGCTACTTGTTTTCTGCACTCACTAAAAGCAACACTTTAACTTTTCCTTTTTT
 TTTTTTTTTCTGAGGCAGTCTCACTCTGTTGCCTAGGTTGGAGTGC

15 Celera SNP ID: hCV26996679
 Public SNP ID: rs6732812
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 60430
 Related Interrogated SNP: hCV11541681 (Power=.51)
 20 SNP Source: dbSNP; Celera; HapMap; ABI_Val
 Population (Allele,Count): caucasian (T,40|C,80)
 SNP Type: INTRON

Context (SEQ ID NO: 125):

25 ATGATTTTGGCATGATAGGGAAACTCTAAGGGTTTCTCACTTCTAGTGGGATC
 ACCATCAAAAATGGACTCCAGTTAGATAAACTCCTAGTTTGCTGTTA
 K
 TGATCGAAGCTCATACTGCAGAACTGAACCTGAATATCAAGGGAATGCTTA
 GCAGATATTTATGCTAAATCAGCTAGTACTGAAACTGTTTCAGATATGC

30 Celera SNP ID: hCV26996688
 Public SNP ID: rs13015885
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 58433
 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (G,39|T,79)
 SNP Type: INTRON

5 Context (SEQ ID NO: 126):
 CATGCCACTGTATTACTGCTGCCTGATGATGGTGAAGACCACAATTGCATAAG
 TGTAGCATCAGAAATAGTGGCCCCTCATGTTAATTTACAAGTTAGTC
 S

TTTGGACAATCCTGAGTTAATACTTTTGTGATGGGTCCTATGCCAAAAGCTC
 10 AGAAGGAAAATATCAGCTAGGATATGCTGTTACCAAATGAGTTAAT

Celera SNP ID: hCV26996689
 Public SNP ID: rs13014700
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 58072

15 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: Celera; HapMap
 Population (Allele,Count): caucasian (C,39|G,79)
 SNP Type: INTRON

20 Context (SEQ ID NO: 127):
 CCTAAAAAGTTAAAACTATTCAAATTTTCCTTGGCCTGCAACCAAAGATA
 ATTAAATAGTTTTCTTGGACTTGCAGGATATGTTTTTAGACTTGCAG
 R

ATATAAAATTCCTGGGTTCTGAATTTTCCTTAATAGCCTCACCATTTTCATGAG
 25 CTCCCTAAAAAATGCTGTACCAGAGCCTTTATCTTGGGATGATAGT

Celera SNP ID: hCV26996690
 Public SNP ID: rs2421575
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 57473

30 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val
 Population (Allele,Count): caucasian (G,40|A,80)
 SNP Type: INTRON

Context (SEQ ID NO: 128):

TTTTTTTTTTTTTTTTTTTTGATAATGGGTTTGTCTCAATCCAAAGATTCTGGGAC
TCTACCTTCTGGGACTCCATCTAATTTGATATGTAAAATTATGGA

Y

5 CCAGAATATGTGCATTTTTAGAAAAATAGATTAACCTTACTAGAGAAAAGAT
GGCCACAATGGGGAAGTTTAAATTTGTATAAAAATTGTTTATTTGCATG

Celera SNP ID: hCV26996697

Public SNP ID: rs12611487

SNP in Genomic Sequence: SEQ ID NO: 71

10 SNP Position Genomic: 54404

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (C,34|T,62)

SNP Type: INTRON

15

Context (SEQ ID NO: 129):

TCTACTGTGACTTCTGCATTTTTGTAGCTTTCCTGATAAAAGCCTGGGCTTCC
TTGTCAGACCTGAAAGATTCATTTATTCATTCAACAAACAGTTCCC

R

20 AGGACCTGAGGTGTGCCGGGCCTGGACTTGGCATGAAGGCACCAGATGTTGG
AAGCGAGGCTGCCGCCAGGAGGACACACCTGATGGGGCTCTGGGAAC

Celera SNP ID: hCV31840149

Public SNP ID: rs12233112

SNP in Genomic Sequence: SEQ ID NO: 71

25 SNP Position Genomic: 69741

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (G,40|A,80)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTERGENIC;

30 UNKNOWN; REPEATS

Context (SEQ ID NO: 130):

AGATCCCTCAGAGTTAAAAAGGCTATTTCTGTGGTTGGGGTTTCAAATCAAAT
TCAAGAGGTTCCCATATCTGAACCCGTCCAATTGACTTTGGGGCACT

Y

TTCAGAAAATCACACTTTTTTACTGTATGATAGTGCTCCAGTAAATTTGCTAG
 AGGGAGATTTACTTTCAAAGCTGAAAGGGCATAGAAGACTACCTATT

Celera SNP ID: hDV69785784

5 Public SNP ID: rs13000788

SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 56263

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap

10 Population (Allele,Count): caucasian (C,39|T,79)

SNP Type: INTRON

Context (SEQ ID NO: 131):

15 CTGCCCAACTTGCAAAATATTTATTTTTCTTGCTCAATTACTCTGTGAAATTTA
 ATTTGTCTAATGGATTTCTTTTAACAGAACTTATAATTTAAAATGGA

R

AGCGGGTGTAATAATTTATAAACTTGCAGCCTGGCCATGTGGTAGACAAGG
 AGGAATCCAAGCAGCCTGCTGAGCAACTACTTGCTGGAGACATTAGCA

Celera SNP ID: hCV31840159

20 Public SNP ID: rs13013228

SNP in Genomic Sequence: SEQ ID NO: 71

SNP Position Genomic: 63171

Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap

25 Population (Allele,Count): caucasian (G,40|A,80)

SNP Type: INTRON

Context (SEQ ID NO: 132):

30 CTCCCAAAGTTCTGGGATTACAGGCGTGAGCCACCATGCCTGGCCTTAACTGT
 TTCAATTAATAACCTGATCAATACCAAAATATAGAAAACCAGTCCAA

R

TGATGCCTGCAAAAGATATATTGTGCTTTCAGGCATCATGCACTCAAGTTACG
 GGAATTACCTATGTGAGTACAAGTAATTGCTTGTATAATATCACCAG

Celera SNP ID: hDV70942181

Public SNP ID: rs17350056
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 60832
 Related Interrogated SNP: hCV11541681 (Power=.51)
 5 SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (G,40|A,80)
 SNP Type: INTRON

Context (SEQ ID NO: 133):

10 CTGTTTCAATTAATAACCTGATCAATACCAAATATAGAAAACCAGTCCAAGT
 GATGCCTGCAAAGATATATTGTGCTTTCAGGCATCATGCACTCAAG
 W
 TACGGGAATTACCTATGTGAGTACAAGTAATTGCTTGTATAATATCACCAGAT
 TAAATTCAACAGGGTCTCTTTTACTAAATGTGTTATACACCCTTAC

15 Celera SNP ID: hDV70953030
 Public SNP ID: rs17434634
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 60881
 Related Interrogated SNP: hCV11541681 (Power=.51)
 20 SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (T,40|A,80)
 SNP Type: INTRON

Context (SEQ ID NO: 134):

25 CAGGGTCTCTTTTACTAAATGTGTTATACACCCTTACAACATACTACAGGGC
 GAGCACAAAAGGTCAATTTCCCAGTGGATTTTGTTCAGGAGCTGTG
 Y
 AGGTTTATTAATGAACAAATCTAACTGACCCCTGCTCAAATGCAACTATGTGG
 CCCTGTTTTTCAACTCCCAAGGGCCTATACTGGGTCTGTGGATAATC

30 Celera SNP ID: hDV70953035
 Public SNP ID: rs17434655
 SNP in Genomic Sequence: SEQ ID NO: 71
 SNP Position Genomic: 61044
 Related Interrogated SNP: hCV11541681 (Power=.51)

SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (T,40|C,80)
 SNP Type: INTRON

5 Gene Number: 4
 Celera Gene: hCG1787816 - 30000070905589
 Gene Symbol: CYP4V2
 Protein Name: cytochrome P450, family 4, subfamily v, polypeptide 2
 Celera Genomic Axis: GA_x5YUV32W706 (3771943..3813704)
 10 Chromosome: 4
 OMIM NUMBER: 608614
 OMIM Information: Bietti crystalline corneoretinal dystrophy, 210370 (3)

Genomic Sequence (SEQ ID NO: 72):

15

SNP Information

Context (SEQ ID NO: 135):

AATGAGTGAGATGATATTTTCGAAGAATAAAGATGCCCTGGCTTTGGCTTGATC
 TCTGGTACCTTATGTTTAAAGAAGGATGGGAACACAAAAAGAGCCTT

20

M

AGATCCTACATACTTTTACCAACAGTGTAAGTCCCTGACTTTTACAATTGTGG
 TAAAATAGACATAACATAAAAATTTCCCTTTATAACCATTTTAACTGT

Celera SNP ID: hCV25990131

Public SNP ID: rs13146272

25

SNP in Genomic Sequence: SEQ ID NO: 72

SNP Position Genomic: 17271

SNP Source: Applera

Population (Allele,Count): caucasian (A,22|C,16) african american (A,7|C,11) total
 (A,29|C,27)

30

SNP Type: NONSENSE MUTATION; MISSENSE MUTATION; HUMAN-
 MOUSE SYNTENIC REGION

SNP Source: dbSNP;

Population (Allele,Count): caucasian (C,45|A,75)

SNP Type: NONSENSE MUTATION; MISSENSE MUTATION; HUMAN-MOUSE SYNTENIC REGION

Context (SEQ ID NO: 136):

5 GAAAGAACTAGCATATTTTATAAGAAAATGTGTTAACTAGGGTGCATCCAA
GTCCAAACAGAAGCATGTGATTATCATTCAAATCATACAGGTCATCGC

K

GAACGGGCCAATGAAATGAACGCCAATGAAGACTGTAGAGGTGATGGCAGG
GGCTCTGCCCCCTCCAAAAATAAACGCAGGGCCTTTCTTGACTTGCTTT

10 Celera SNP ID: hCV3230097

Public SNP ID: rs3736455

SNP in Genomic Sequence: SEQ ID NO: 72

SNP Position Genomic: 19377

Related Interrogated SNP: hCV25990131 (Power=.7)

15 SNP Source: Applera

Population (Allele,Count): caucasian (G,24|T,16) african american (G,7|T,9) total
(G,31|T,25)

SNP Type: MISSENSE MUTATION; HUMAN-MOUSE SYNTENIC
REGION; SILENT MUTATION

20 SNP Source: dbSNP; Celera; HapMap; HGBASE;

Population (Allele,Count): caucasian (T,42|G,78)

SNP Type: MISSENSE MUTATION; HUMAN-MOUSE SYNTENIC
REGION; SILENT MUTATION

25 Context (SEQ ID NO: 137):

AGAGGTGATCTATCAACACATAATTACAACATGTGATATGAGCTATGAACAC
TTATGAACAAACAGGGTGCTGTGTAAGAATAAAGGAACAAAGATCT

R

30 TGTATAGGAGTTTTCTGGAAAATGTTTGGATTTCGGCAGTCATTTTCAAAGGCA
GAGGGCATTGATAGCAGTATCTTAACATGGAAAACATTAATAACTAAC

Celera SNP ID: hCV15968026

Public SNP ID: rs2292426

SNP in Genomic Sequence: SEQ ID NO: 72

SNP Position Genomic: 17811

Related Interrogated SNP: hCV25990131 (Power=.7)

SNP Source: dbSNP; HGBASE

Population (Allele,Count): caucasian (A,42|G,78)

SNP Type: INTRON

5

Gene Number: 5

Celera Gene: hCG21777 - 207000006881214

Gene Symbol: NR1I2

Protein Name: nuclear receptor subfamily 1, group I, member 2

10 Celera Genomic Axis: GA_x5YUV32VYQG (75803402..75861414)

Chromosome: 3

OMIM NUMBER: 603065

OMIM Information:

15 Genomic Sequence (SEQ ID NO: 73):

SNP Information

Context (SEQ ID NO: 138):

20 AAGCACTGCCTTTACTTCAGTGGGAATCTCGGCCTCAGCCTGCAAGCCAAGTG
TTCACAGTGAGAAAAGCAAGAGAATAAGCTAATACTCCTGTCCTGAA
M

AAGGCAGCGGCTCCTTGGTAAAGCTACTCCTTGATCGATCCTTTGCACCGGAT
TGTTCAAAGTGGACCCCAGGGGAGAAGTCGGAGCAAAGAACTTACCA

Celera SNP ID: hCV263841

25 Public SNP ID: rs1523127

SNP in Genomic Sequence: SEQ ID NO: 73

SNP Position Genomic: 11707

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,49|A,71)

30 SNP Type: UTR5; PUTATIVE UTR5

Context (SEQ ID NO: 139):

ATATCTACCCCCACCACAATGCTATTTAATACTGTATTAATCGTTCCAATAA
TGCAATAAGGAAAGAAAAAAGCATAAAGATCAGAAAAGAAGAAAAC

W

GTCTTTCTTTGCAGAAAACATAATTATTACATTGAACATCCTCAGTAATACT
AAGGAATAACTACTAGAACTATTTAATAAAGTCACAGTTATATCTCA

Celera SNP ID: hCV192027

5 Public SNP ID: rs9821892

SNP in Genomic Sequence: SEQ ID NO: 73

SNP Position Genomic: 2603

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

10 Population (Allele,Count): caucasian (A,84|T,36)

SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 140):

15 GATACAGTATAGGCAGATTGTCAGAAGACAAAGAGAATCTTGAAAACAGTGA
GAGAGAGGCAACTCATCCTGTACAAGGGAGCCTTAATAAGATTAACAG

Y

TGATTTTCGCATCAGAAACCATGGAGGCCTGAAAACAGGAGGCCTGATGCCAC
TTCACACCTACTAGGATGGATATAATTTTTTTTAATGGAAAATAACAA

Celera SNP ID: hCV1833991

20 Public SNP ID: rs11926554

SNP in Genomic Sequence: SEQ ID NO: 73

SNP Position Genomic: 6219

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

25 Population (Allele,Count): caucasian (T,73|C,27)

SNP Type: INTERGENIC; UNKNOWN; REPEATS

Context (SEQ ID NO: 141):

30 ACATATCCAGGAAGCTCAACAAAGCCCAAGCAAATAAACACGAAGAGATA
CACCAAGATACAGTATAGGCAGATTGTCAGAAGACAAAGAGAATCTTGA

R

AACAGTGAGAGAGAGGCAACTCATCCTGTACAAGGGAGCCTTAATAAGATTA
ACAGTTGATTTTCGCATCAGAAACCATGGAGGCCTGAAAACAGGAGGCC

Celera SNP ID: hCV1834240

Public SNP ID: rs1581451
 SNP in Genomic Sequence: SEQ ID NO: 73
 SNP Position Genomic: 6162
 Related Interrogated SNP: hCV263841 (Power=.8)
 5 SNP Source: dbSNP; Celera; HapMap; HGBASE
 Population (Allele,Count): caucasian (G,48|A,72)
 SNP Type: INTERGENIC; UNKNOWN; REPEATS

Context (SEQ ID NO: 142):

10 ATTGCAGGCCTGTTGTTAAGAATTTCAAAGGAGACTTTTCTTTTTGCTCTACCC
 AGAGCCAAGGCTGAGAAGGCCTGGGATAGACCTGTCTCCCTCCATG
 R
 TTTCTACTGAGGATGTCACCCTTAGGGGATCTTGAGTTTATAGTCTATGATCT
 GGTTAGGCTCCAGGCTTTGTCTCCTATTCTACCATTAAAATCCAGGC

15 Celera SNP ID: hCV1834252
 Public SNP ID: rs10934498
 SNP in Genomic Sequence: SEQ ID NO: 73
 SNP Position Genomic: 15049
 Related Interrogated SNP: hCV263841 (Power=.8)
 20 SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (G,49|A,71)
 SNP Type: INTRON

Context (SEQ ID NO: 143):

25 GCAGGCCCCAGATATAGCCCCATGCTGTCCTCCTACCCCAGAGCACACTGTTC
 AGGCTACTTCCACTGGTACTGAAATCCAGTATTTCACTTACTCTTTT
 Y
 CTTTCCAATATCCTCATGACATTCAATATTTCACTTACTCTAGGTCCTCCCTGC
 CTAAGGCCCAAGTCAACTTTCTGTCCAGTGGGATTTGTAATCCAAT

30 Celera SNP ID: hCV9152783
 Public SNP ID: rs1523130
 SNP in Genomic Sequence: SEQ ID NO: 73
 SNP Position Genomic: 10175
 Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (T,49|C,71)

SNP Type: UTR5; PUTATIVE UTR5

5 Context (SEQ ID NO: 144):

CTTATAAAAATAAAAATAAATACATAAAAATGTGACTAGCAATTTTAAGACACT
TAAATATTTGAATTAATTTATTGAATACCTGCTCTGGGAATTAGTGTT

R

TGTTTCCATTTACAGAGACAAAACAGGTCTAATTGATTGACTAAACAATAGAT

10 GTAAATGTGGTCAGTGTTTGGATGTCTATTCATAGAATAGGAACAAT

Celera SNP ID: hCV11230788

Public SNP ID: rs7643038

SNP in Genomic Sequence: SEQ ID NO: 73

SNP Position Genomic: 7708

15 Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (G,48|A,72)

SNP Type: INTERGENIC; UNKNOWN

20 Context (SEQ ID NO: 145):

AAGATGTGTGTGAACACAAATATACCTTCTGTTTGAGGTCAGCATCATAGTGG
GTCGTGAATCATGTTGGCCTTGCTGCTGTCTCCTCATTCTAGGGTG

R

AAAAAAAAAAGCATGAAAACAATCACTTAATGTTGAGCCCCATTACTGATGC

25 TCTCTGGTCCTGCACTAGCCTCCTAGAAAAATCACCACAGGAGAAGCC

Celera SNP ID: hCV15882316

Public SNP ID: rs2276706

SNP in Genomic Sequence: SEQ ID NO: 73

SNP Position Genomic: 11975

30 Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; HapMap; HGBASE

Population (Allele,Count): caucasian (G,72|A,48)

SNP Type: INTRON

Context (SEQ ID NO: 146):

GTACTTCAAATAATAACAACCTTAAGTCAATAAATAAATGTAAGGAAGTCCA
AATGTTACCTGAAGACAACCTGTGGTCATTTTTTGGCAATCCCAGGTT

Y

5 TCTTTTCTACCTGTTTGCTCAATCGTGGTCTCCCTCTCCCTCTCTTGTTGGGGCC
CATGCCCTGCTTTACTGTTGCCAGAGGCTTGTACTTGTTCCT

Celera SNP ID: hCV27504984

Public SNP ID: rs3814055

SNP in Genomic Sequence: SEQ ID NO: 73

10 SNP Position Genomic: 10703

Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (C,72|T,48)

SNP Type: UTR5; PUTATIVE UTR5

15

Context (SEQ ID NO: 147):

AAAGGATGCTAAGGGAGGCCTCAGAGGCTTATTTAATCAGTTGGTTAATGGG
AAAAAGTATTTTGTGGTACACAAATACGCTTTAAAAATATTTTAGTTA

Y

20 GTCATAAAAGTCTAAACCTGGGGTTGACAAACTATGGCTCACAGGCTAACTC
TGGCCCACCAGCTGTTTTTATAAATGAAGCTTTATTGGAACACAGCCA

Celera SNP ID: hCV30747432

Public SNP ID: rs12488820

SNP in Genomic Sequence: SEQ ID NO: 73

25 SNP Position Genomic: 12737

Related Interrogated SNP: hCV263841 (Power=.8)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (T,41|C,65)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

30

Gene Number: 6

Celera Gene: hCG21779 - 104000117367418

Gene Symbol: C3orf15

Protein Name: AAT1-alpha

Celera Genomic Axis: GA_x5YUV32VYQG (75880995..75939455)

Chromosome: 3

OMIM NUMBER:

OMIM Information:

5

Genomic Sequence (SEQ ID NO: 74):

SNP Information

Context (SEQ ID NO: 148):

10

TGGGACCCTGGGTCAGCAGAGCTACAGGGACAGGACACTTACCTAGAGCCAT
GGGGATGGGGCATTGTAGAGTCAAAGGGGCAGAGCCATGGGAATGATG
Y

TGCCACCCCAGTGGGCTTAGAAAGCAAGACACTAAGCCAAAGAAGACTATAC
TCAAGAATTAAGATCCAGTGGAATTTGCCTTGCCAGGTCGTGGACTTG

15

Celera SNP ID: hCV105917

Public SNP ID: rs9289134

SNP in Genomic Sequence: SEQ ID NO: 74

SNP Position Genomic: 62414

Related Interrogated SNP: hCV263841 (Power=.51)

20

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (T,84|C,36)

SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 149):

25

TTTTTTAAAATAACAACCCTCAAATTCTAATGTGTAATTTGCTAGTTTATTAT
GTTTGTTGTTTGTGCTCTCTCTCATGAAGTCAAAGATTTTTGTC
R

TTTTGTTGATGTATTCTAAGCACCTAGATAGAGTAGCCAGTCAATAAACATGT
GTTGAATGAAGGAGCAAAAAGAGCCTGAGCCTCTTGCAGCTTTTCAG

30

Celera SNP ID: hCV134275

Public SNP ID: rs9847068

SNP in Genomic Sequence: SEQ ID NO: 74

SNP Position Genomic: 12752

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val
 Population (Allele,Count): caucasian (G,40|A,80)
 SNP Type: INTRON

5 Context (SEQ ID NO: 150):
 CTGATGGTATAATTTTGTCTTTGCCTCCCTTCTCTAACATTAGAGAAAAGTTC
 CCAGGTTTAAAACCATGTTCAGTAACCTGATCCATTATCCAAGATA
 Y
 TCTCTATATTGGAGCAAGTCAGATCCTGTCCCACCATTATCAGTCGGGAATG
 10 GAAGGGACATAAGGAGAAACACAGAGAAGCCCTCCGGCAGCTCACCA

Celera SNP ID: hCV134278
 Public SNP ID: rs9848716
 SNP in Genomic Sequence: SEQ ID NO: 74
 SNP Position Genomic: 14446

15 Related Interrogated SNP: hCV263841 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (T,40|C,80)
 SNP Type: HUMAN-MOUSE SYNTENIC REGION; SILENT MUTATION

20 Context (SEQ ID NO: 151):
 CAAGAGTCTCACTTCACAGTTTAAGGACTCACATAGACTGAAAGTGAAAGGA
 TGAAAAAAGATACTCAATGCAAATGGAAACCAAATAGAACAAGTGTA
 R
 CTAAACTTACATCAGATAAAAATAGGTTTTAAGTAAAAAACTATAAAATGAGA
 25 CCAAAAAGGTCATTGTGTGGTGATAAAGGGGTCGATTCATCAAGAGGC

Celera SNP ID: hCV178227
 Public SNP ID: rs13070374
 SNP in Genomic Sequence: SEQ ID NO: 74
 SNP Position Genomic: 67774

30 Related Interrogated SNP: hCV263841 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (G,66|A,52)
 SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 152):

GTGTGTTTCATGCAACAGCCATAGAGATGGTGCTATTGAAAGGCAATCTGCCC
ATCCTGTGTATTGACAAGCATCAGGATTCCAGGAGGATTCTTCCTGAT

K

5 CAAAGACTGTCAAATAGGAAAAAAGAAAACAGTAAATGAACAAAAACATCA
ACATGGAGTCCAAGTAATGAATTTATAGAGGAAAATGTATGCTACTCTG

Celera SNP ID: hCV255886

Public SNP ID: rs10511394

SNP in Genomic Sequence: SEQ ID NO: 74

10 SNP Position Genomic: 73296

Related Interrogated SNP: hCV263841 (Power=.6)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (T,85|G,35)

SNP Type: UTR3; INTERGENIC; UNKNOWN

15

Context (SEQ ID NO: 153):

ATATTGTCTATAAAAATGTAAATTCAAGGACCAAAGCATTGGCTACCTTATTA
TTGTGATAAAAAGCAAGCCTTAAGGTTGGGGTAGTCATTTTTCTGACA

Y

20 GCAAATTTAATAAGATTAATAATTTACACAATTTGCATTGAATGCATCTGCCT
TTRACTTCTGATTATGGTACCTACTCTATGACATCCATAGTGACCTGT

Celera SNP ID: hCV278948

Public SNP ID: rs1464599

SNP in Genomic Sequence: SEQ ID NO: 74

25 SNP Position Genomic: 60168

Related Interrogated SNP: hCV263841 (Power=.6)

SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (T,81|C,35)

SNP Type: INTRON

30

Context (SEQ ID NO: 154):

AGCATCATATCATATGAAAATATATAAAGTTGTAAATTAGGCTAACATACTGT
TATATAAAAATTTATTCTACCCACCTACTTTCACAAACATATCTTTGA

R

ATGATAACATTGAAAAACGGGTAAGCAGCTATGGGCTTTATGTAATTGAAAG
 ATGGCAAGAACCAAAGATGGTAAGCTTAGTTATTATTGTACATGTCTG

Celera SNP ID: hCV1834242

Public SNP ID: rs11712308

5 SNP in Genomic Sequence: SEQ ID NO: 74

SNP Position Genomic: 71962

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,64|A,56)

10 SNP Type: INTRON

Context (SEQ ID NO: 155):

GGTATTTGCAAGACCCAGGTCAAGAGGATGAATTGAAGCATCATATCATATG
 AAAATATATAAAGTTGTAAATTAGGCTAACATACTGTTATATAAAATT

15 K

ATTCTACCCACCTACTTTCACAAACATATCTTTGAGATGATAACATTGAAAAA
 CGGGTAAGCAGCTATGGGCTTTATGTAATTGAAAGATGGCAAGAACC

Celera SNP ID: hCV1834243

Public SNP ID: rs9682652

20 SNP in Genomic Sequence: SEQ ID NO: 74

SNP Position Genomic: 71926

Related Interrogated SNP: hCV263841 (Power=.6)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (T,85|G,35)

25 SNP Type: INTRON

Context (SEQ ID NO: 156):

CTAAGGCTCAAAGAGAAGTATAAGACTAGGTCTTTATATTTGTAAGTTACTTA
 CACTCTAGATATCAAACGAAATGTAACACTACATGAAGTGCTTCAATG

30 Y

TGAATAGGTTGGTACAAACTCAAGTGCGGTAAGTGGTAAATCTGGGCTGATG
 TAGTCAGAGGAGACTTCCTGGAGGATCTTGGACATGAGCTTCCAGGT

Celera SNP ID: hCV1834256

Public SNP ID: rs2472662

SNP in Genomic Sequence: SEQ ID NO: 74
 SNP Position Genomic: 58134
 Related Interrogated SNP: hCV263841 (Power=.51)
 SNP Source: dbSNP; Celera; HGBASE
 5 Population (Allele,Count): caucasian (C,64|T,54)
 SNP Type: INTRON

Context (SEQ ID NO: 157):

AAATTCAGCTCTGGTCTTCCCGAGTGTAAGCCTTAGCCAGATTATTTGTGTT
 10 GCTGGAGCAGGCACACTTCACTTAAAAAATGGACTGATATAAAGCCC
 Y
 AGTTAATTTTTTTTCCCAGATAGCTATAACTCATAAGAGATTCAGATATGATG
 AGTAAATCATCTGGATGAACA ACTTGTGTTTAGCACTGTTCTTTACT

Celera SNP ID: hCV1834260

15 Public SNP ID: rs4688033

SNP in Genomic Sequence: SEQ ID NO: 74
 SNP Position Genomic: 47546
 Related Interrogated SNP: hCV263841 (Power=.6)
 SNP Source: Celera; HapMap; HGBASE
 20 Population (Allele,Count): caucasian (C,85|T,35)
 SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

Context (SEQ ID NO: 158):

GCCTTGGCACTAGTCCTAGCTCTTACCCACCCCCACCCCGACCCAAGACACT
 25 GGGTTCACAGTCAGGATGCCTTGCGTGCTAGCTAGAGCTGATGGGGC
 M
 TTCCTAGTTAAAAGACAAATTATGGAATTCTTCCATTAAGAGATGAGTATAT
 AACAAATATACATAAATAACAATATGCTGAGGTA CTGTGGTTGGGTT

Celera SNP ID: hCV30699687

30 Public SNP ID: rs11711386

SNP in Genomic Sequence: SEQ ID NO: 74
 SNP Position Genomic: 33806
 Related Interrogated SNP: hCV263841 (Power=.51)
 SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (A,66|C,54)

SNP Type: INTRON

Context (SEQ ID NO: 159):

5 AGCCGAGATCGTGCCATTGCACTCCAGCCTGGGTGACAAGAGCAAGACTTCA
TCTCAAAAAAAAAAAAAAAAAAGTGCTTGGCATAATAAAAATGCTTAAGTA

R

TGTTAGGTTTTATTATTATAATTCATAGACACATAAATTTGTATATAAAAT
GTTAAGATTGAATAGAAAACCACAAATTCCTTTTTGTGTTTACTTG

10 Celera SNP ID: hCV29841665

Public SNP ID: rs7623217

SNP in Genomic Sequence: SEQ ID NO: 74

SNP Position Genomic: 31138

Related Interrogated SNP: hCV263841 (Power=.51)

15 SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (A,56|G,62)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON;

REPEATS

20 Context (SEQ ID NO: 160):

TGTA ACTATAAGGTGTTTTATATATTCCTCATGGTAAGTACAAAACAAAAAAC
CTATATTAATACACAAAAGGTGAAGAGTAAGAAATCAAGGCATGGC

R

25 CTAGAGAAAGTCACCTAATTACAAAGGAATACAGCAAGAGAGGAAGAAAGG
AATAAGGACCTACAAAACAACCATAAAACAATGAACCAGATGATAGTAG

Celera SNP ID: hCV30562884

Public SNP ID: rs9815093

SNP in Genomic Sequence: SEQ ID NO: 74

SNP Position Genomic: 67458

30 Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; HapMap; ABI_Val

Population (Allele,Count): caucasian (G,84|A,36)

SNP Type: INTRON; REPEATS

Gene Number: 7
 Celera Gene: hCG32403 - 206000045254807
 Gene Symbol: CASP8AP2
 Protein Name: CASP8 associated protein 2
 5 Celera Genomic Axis: GA_x54KRFTF0F9 (80703289..80769434)
 Chromosome: 6
 OMIM NUMBER: 606880
 OMIM Information:

10 Genomic Sequence (SEQ ID NO: 75):

SNP Information

Context (SEQ ID NO: 161):

AAAATAATTTTGAATGATACTTCAATTAATCCTTTTATATTTTATAGTTGCAAG
 15 TAGAATATGGAAAATGTCAACTACAAATGAAAGAGCTGATGAAAAA
 R

TTTAAAGAAATACAGACACAGGTAGAGTATAAATGAAAACATAAAAAACAA
 ATTACCAGGCACAGTGACTCACACCAGTAATTCAGCACTTTGAGAGGC

Celera SNP ID: hCV2744023

20 Public SNP ID: rs369328

SNP in Genomic Sequence: SEQ ID NO: 75

SNP Position Genomic: 35612

SNP Source: dbSNP; Celera; ABI_Val; HGBASE

Population (Allele,Count): caucasian (A,50|G,70)

25 SNP Type: MISSENSE MUTATION; HUMAN-MOUSE SYNTENIC
 REGION; SILENT MUTATION

Context (SEQ ID NO: 162):

TCCAGTTTTCTTTTGAACCAGTAGTATCTTTCGGTTATCATATTGTTTGTGAAC
 30 TCACTTACCTGGTGGTATGAGATCTTCTGAGTGAGAGTATGCACTG
 S

AGAAGGACTGATGGCCAACCCTCTAGTTGGTATTAAGGGTTCAAAGAGAGGG
 AAGGGGACAAGCCCCATGTTGGAGAGCCCTAGCACTACAAAACCACTG

Celera SNP ID: hDV71164887

Public SNP ID: rs2585018
 SNP in Genomic Sequence: SEQ ID NO: 75
 SNP Position Genomic: 56712
 Related Interrogated SNP: hCV2744023 (Power=.6)
 5 SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (C,50|G,70)
 SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 163):

10 TTCTTTCAGTTGTCATTACCCTTTTAGTTTTGCTGCTTTAGGTTGTTCTCTAG
 CTCTTTAGATAGCTTCACATATTTTATTTAGTGTCAACAGTAGTT
 R
 CCTGCAGGAGGATTGGTTTCATAGAACTACTCTGCTATAACTGGAATTGGAA
 CTCTTTCAGTAGAATATATATATATGAGACAGGATCTCGCACTGTCA

15 Celera SNP ID: hCV26081278
 Public SNP ID: rs2585008
 SNP in Genomic Sequence: SEQ ID NO: 75
 SNP Position Genomic: 32250
 Related Interrogated SNP: hCV2744023 (Power=.6)
 20 SNP Source: dbSNP; HapMap; HGBASE
 Population (Allele,Count): caucasian (G,47|A,69)
 SNP Type: INTRON; REPEATS

Gene Number: 8
 25 Celera Gene: hCG40122 - 104000117872607
 Gene Symbol: CASP5
 Protein Name: caspase 5, apoptosis-related cysteine protease
 Celera Genomic Axis: GA_x5YUV32VVY5 (14887755..14936958)
 Chromosome: 11
 30 OMIM NUMBER: 602665
 OMIM Information:

Genomic Sequence (SEQ ID NO: 76):

SNP Information

Context (SEQ ID NO: 164):

TTTTGGTCCATATTGAGAAGTGTTTGGGTAACATTTGATGAGCCACGCGATT
 CTTTCGCAAAGAGTCTACCAAGATCAGGGCCTTGTCTTCAATTTTGG

5 Y

ATCATAATATTTTTCTTTTCCTCTTCCTTCAATGTCAGAACATCGTGTTTTGCC
 AAATAATTAACACCATGAAGAACATCTTTGCCAGGTATTCC

Celera SNP ID: hCV12092542

Public SNP ID: rs507879

10 SNP in Genomic Sequence: SEQ ID NO: 76

SNP Position Genomic: 23268

SNP Source: Applera

Population (Allele,Count): caucasian (C,3|T,5) african american (C,4|T,4) total
 (C,7|T,9)

15 SNP Type: MISSENSE MUTATION; TRANSCRIPTION FACTOR
 BINDING SITE

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (T,65|C,53)

SNP Type: MISSENSE MUTATION; TRANSCRIPTION FACTOR

20 BINDING SITE

Gene Number: 9

Celera Gene: hCG2039683 - 209000071858797

Gene Symbol: C1QTNF6

25 Protein Name: C1q and tumor necrosis factor related protein 6

Celera Genomic Axis: GA_x5YUV32VU5C (15633444..15661565)

Chromosome: 22

OMIM NUMBER:

OMIM Information:

30

Genomic Sequence (SEQ ID NO: 77):

SNP Information

Context (SEQ ID NO: 165):

AGAGGCTGAGGATACATGGGCAGGATCCAGGGGGTCCTCAGAGTCACAGCAC
CGTTGGCAGCCGCTGGCCACAGCTCTGTCAAAGGTGAGCTCCACCATA

S

5 GGATCTCACACATCAGGAGAAAGAGCAGGAGCGCTGCCAGACGGGACCCA
GGGCGGCTGTCCCATGGTGACCTGGAACAAGGAAGGAGGGACAGGAAC

Celera SNP ID: hCV2403368

Public SNP ID: rs229526

SNP in Genomic Sequence: SEQ ID NO: 77

10 SNP Position Genomic: 15215

SNP Source: dbSNP; HapMap; ABI_Val; HGBASE

Population (Allele,Count): caucasian (C,95|G,25)

SNP Type: MISSENSE MUTATION

15 Gene Number: 10

Celera Gene: hCG2036531 - 209000073898476

Gene Symbol: MDN1

Protein Name: MDN1, midasin homolog (yeast)

Celera Genomic Axis: GA_x54KRFTF0F9 (80759613..80956812)

20 Chromosome: 6

OMIM NUMBER:

OMIM Information:

Genomic Sequence (SEQ ID NO: 78):

25

SNP Information

Context (SEQ ID NO: 166):

30 TTTCTTACACCCTACTCAGCTTTATATGTTTTCACTTATCACAACCTAAAATTA
CAGTACCTGACTCCTCTAGAATTAAGCTCAATGAGAAGAACCAGG

W

TTTTTGTCTGCCTTACTCATAGCTATAACCACAGCACCAGGGCACTCAATAAG
TATTTGTATTAAGGAAGAGAAAAAATACACAAACACACAGTAAAACA

Celera SNP ID: hCV2012681

Public SNP ID: rs6918308
 SNP in Genomic Sequence: SEQ ID NO: 78
 SNP Position Genomic: 122998
 Related Interrogated SNP: hCV2744023 (Power=.51)
 5 SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (A,71|T,49)
 SNP Type: INTRON

Context (SEQ ID NO: 167):

10 TTTGTTAGTATTGGTCTAGAAAGGCCTTTGAGGGCTGATTTAGCCCTAGAGTG
 TGGCCCTTCTGGGTTCCACACTAATTATGCCAGTATTCAACTTGGTC
 K
 TTCCACTCCAAGTAGTAGGAGTTGAACATCTCCTACACAGAAAACTCCAGTA
 GTTATTCTTTGCATAGCTCACGAAATTTCCCTAAGAAAACACAGCTT

15 Celera SNP ID: hCV3180835
 Public SNP ID: rs292222
 SNP in Genomic Sequence: SEQ ID NO: 78
 SNP Position Genomic: 175944
 Related Interrogated SNP: hCV2744023 (Power=.6)
 20 SNP Source: dbSNP; Celera; HapMap; ABI_Val; HGBASE
 Population (Allele,Count): caucasian (G,50|T,70)
 SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

Context (SEQ ID NO: 168):

25 TCTAGCCTGGGCAACAGAGTGAGACTCTGTCTCAAAAACACAAAACAAAACA
 AAAAACAATTAACAAAAACAGAAATGTGTTGTCTCACAATTCTAAAGT
 Y
 TGGAAGTCCAAGATCAGGGTGTGTCAGCAAGGCCATATTCCTTTGAAACCTATA
 GGGGAGAATCCTTCCTTGCCTCTTCCTAGCTTCTGGTTGTCATCATCA

30 Celera SNP ID: hCV27392140
 Public SNP ID: rs12661693
 SNP in Genomic Sequence: SEQ ID NO: 78
 SNP Position Genomic: 177979
 Related Interrogated SNP: hCV2744023 (Power=.6)

SNP Source: dbSNP; Celera; HapMap
 Population (Allele,Count): caucasian (T,50|C,70)
 SNP Type: INTRON

5 Context (SEQ ID NO: 169):
 GTGTTTTAAACTGGTACTTATTTTAAAGTGGGGGGGGGACCTCTTATTTTATA
 AATTATCTTCATTAATAATAATCAGGAATGAGAAGGTTAGACCTGT
 Y
 ATCACTTCCAACCTTGCTTGAAGTAGACCCAGAGAACTCTATTACAGGA
 10 CCTCTGGCTCAGCATGCAGACTTACGAGTAAGTGGCTCCAGAGTGAA

Celera SNP ID: hCV28980044
 Public SNP ID: rs6923604
 SNP in Genomic Sequence: SEQ ID NO: 78
 SNP Position Genomic: 155203

15 Related Interrogated SNP: hCV2744023 (Power=.51)
 SNP Source: dbSNP; HapMap
 Population (Allele,Count): caucasian (T,71|C,49)
 SNP Type: INTRON

20 Context (SEQ ID NO: 170):
 CAATGGGAACCAACCTGACACCTCAGAGAACTGAGGATATAAACTACAAACA
 ACCAAACA ACTAAAGAGCAAGATTTTCAAACCGTCTTCTGCAGAAAGG
 Y
 ACCTTGGGAGCCACTTAAATGCCAAGAAGGTCCATAACACCTGCCTTCATGG
 25 AAAATAACTTTGCTTATATCTGTTACATACTGGAACCACATGTAAG

Celera SNP ID: hCV30238224
 Public SNP ID: rs6910277
 SNP in Genomic Sequence: SEQ ID NO: 78
 SNP Position Genomic: 134034

30 Related Interrogated SNP: hCV2744023 (Power=.51)
 SNP Source: dbSNP
 Population (Allele,Count): caucasian (C,69|T,47)
 SNP Type: INTRON

Context (SEQ ID NO: 171):

TGGTGCACGTCTGTAGTCCCAGTTACTCGGGAGGCTGAGGCAAGAGGATCAC
ATGAGCCTAGGAGGTACAGGCTGCAGCAAACCTATGAGGGCACCATGGT

M

5 CTCCAGGCTGGGCAACAGAGTGAGACCCTGTCCCCCCCCAAAAAAGAGTTA
AAATTAAATAAATAAATAAAAATCCAAGACCAGCTACAGTGCCTATAAT

Celera SNP ID: hCV29986476

Public SNP ID: rs9451269

SNP in Genomic Sequence: SEQ ID NO: 78

10 SNP Position Genomic: 136179

Related Interrogated SNP: hCV2744023 (Power=.51)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (A,70|C,50)

SNP Type: INTRON; REPEATS

15

Gene Number: 11

Celera Gene: hCG1810876 - 65000099284460

Gene Symbol: RGS7

Protein Name: regulator of G-protein signalling 7

20 Celera Genomic Axis: GA_x5YUV32VWMC (16766653..17366707)

Chromosome: 1

OMIM NUMBER:

OMIM Information:

25 Genomic Sequence (SEQ ID NO: 79):

SNP Information

Context (SEQ ID NO: 172):

30 TAAGCCATTATTTCTTCAAGTACTTTTTTTTTTTAAGTTCAGCTTTTTATTGAAC
ACATTATAAAAGAGGTTTCGTCAAAAAGACCAAAGCCCATGTCAC

Y

ATCAGACTTCTCGGATTCTTCTTTCTTTGCTTCCACTTTCTTCTCCTCAGCTGGA
GCAGCAGCAGCAGAGGCGGGCAGAAGCTCCTGCTGGTACAGATAG

Celera SNP ID: hCV916107

Public SNP ID: rs670659
 SNP in Genomic Sequence: SEQ ID NO: 79
 SNP Position Genomic: 232752
 SNP Source: dbSNP; HapMap
 5 Population (Allele,Count): caucasian (C,78|T,40)
 SNP Type: INTRON

Context (SEQ ID NO: 173):

ACCATAAAATTTATATAAATGTTAAAAGCACAGTTCAAACCTGAATTATAATTT
 10 TTCAGGAGGTTTGCTTGCGTGCCATATTAATAATATCCCTGATGAAAC
 R
 TATTTCTTAAAACCACTAAGAAAAAAGTTCTTGCAATGGGAAAATGATTTATT
 TACCTGGATCTTCTATAGTTAAGTTCTTTATCAACCATTGAACAATG

Celera SNP ID: hCV25653735

15 Public SNP ID: rs7520707
 SNP in Genomic Sequence: SEQ ID NO: 79
 SNP Position Genomic: 217612
 Related Interrogated SNP: hCV916107 (Power=.51)
 SNP Source: Applera

20 Population (Allele,Count): caucasian (A,9|G,11) african american (A,15|G,3) total
 (A,24|G,14)

SNP Type: INTRON

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (A,71|G,47)

25 SNP Type: INTRON

Context (SEQ ID NO: 174):

TGTCTGCAGTTTTCAAATGTTTATTAAGTGGCTTGGTATGCATTTTCGTTGGGCT
 TATTTATCCTATTTGGGCTTCTCTCAGCTTCTTCAACTGAAGGTTT
 30 R
 TATCTTTGCTATGCTTGATGTTTTTAAGCCATTATTTCTTCAAGTACTTTTTT
 TTTTAAAGTTCAGCTTTTTATTGAACACATTATAAAAAGAGGTTTC

Celera SNP ID: hCV916106

Public SNP ID: rs575226

SNP in Genomic Sequence: SEQ ID NO: 79
 SNP Position Genomic: 232625
 Related Interrogated SNP: hCV916107 (Power=.7)
 SNP Source: dbSNP; HapMap; HGBASE
 5 Population (Allele,Count): caucasian (A,77|G,41)
 SNP Type: UTR3; INTRON; REPEATS

Context (SEQ ID NO: 175):

10 ATTTGAAAATCGTATCAGATGTTACAATTGAGGTTTTAGACTCATTATTTTA
 AAAGGGTGACCTTTTACTGTACAGCAGGAACCTGTCTTGTGTGGTT

W

GAGACAAATGACAATGAAGAGATGGATATGTTCCCTAGGATCATATTTAGGCT
 ATAAAACATCTATCAGTGCCCTTATCACACTTTTTGTACTGCTGTAGT

15 Celera SNP ID: hCV26887401
 Public SNP ID: rs10802916
 SNP in Genomic Sequence: SEQ ID NO: 79
 SNP Position Genomic: 203302
 Related Interrogated SNP: hCV916107 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap
 20 Population (Allele,Count): caucasian (T,68|A,48)
 SNP Type: INTRON

Context (SEQ ID NO: 176):

25 GAGAGATGAACTCCAGCCCAGCTACAGACTATGTTCTCTAACCCCTTGATCTA
 AGAAATCTCAGCAATGGCAAAAATCTGGATATCACCAATTAATGGCC

R

GTGTCCTAGTTCTGTCTGCCTTTGGCCAAAAGCGGGTACTGAAAAAAGAGAA
 AAGCAGCCACTGGCTTCCAGGAACTGGCTGGGCACTCACAGCTATGTC

30 Celera SNP ID: hCV26887441
 Public SNP ID: rs9786932
 SNP in Genomic Sequence: SEQ ID NO: 79
 SNP Position Genomic: 223059
 Related Interrogated SNP: hCV916107 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,74|A,46)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

Context (SEQ ID NO: 177):

5 GGAAAATCTTTTATTTGAAATAAGAGTAAGAAGTCAATTGAGAGAATTCCAA
 AAGGTGCCCCAAAGCAAAAACAAAAGATTGAGGTTTAAATGAACTGTA
 Y
 GGTATAAAAAGGCAACAATTTTAAAAGTAAGATAGTAAGGATAAATAGTAA
 AATAGTAAGGCTTTCTAATTAAGATACAAAACCTAGCAGCCATAAAAC

10 Celera SNP ID: hCV26887461

Public SNP ID: rs4660023

SNP in Genomic Sequence: SEQ ID NO: 79

SNP Position Genomic: 235644

Related Interrogated SNP: hCV916107 (Power=.6)

15 SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,74|T,46)

SNP Type: INTRON

Context (SEQ ID NO: 178):

20 AGTGGGGAAGCCCAGCAACGCCTGTCCCTCTAGATTCTTCTGGCCTCTCTGAG
 CAGCATTCTTCTTCTGGGTATGGGGCAGAGCTCTCTCTGGACTGG
 R
 GGTCTTAGGACCTACAGTCCAACAAGGCAGGTCAGATAACTTCTTTATGACTA
 CTTTTACGGCTGGCTTTGAGGAGAAGGGGTTCTGATTTCATGACCC

25 Celera SNP ID: hCV26887463

Public SNP ID: rs6680767

SNP in Genomic Sequence: SEQ ID NO: 79

SNP Position Genomic: 239051

Related Interrogated SNP: hCV916107 (Power=.6)

30 SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,76|A,42)

SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 179):

GAGCAGCATTTCCTTCTGGGTATGGGGCAGAGCTCTCTCTGGACTGGGGG
TCTTAGGACCTACAGTCCAACAAGGCAGGTCAGATAACTTCTTTATG

R

CTACTTTTACGGCTGGCTTTGAGGAGAAGGGGTTCTGATTTCCATGACCCGCC
5 TCGGGCCAGAGGGATTCTAGTCTCTGGGGAGAATGGAACAGGTCAGG

Celera SNP ID: hCV26887464

Public SNP ID: rs6669640

SNP in Genomic Sequence: SEQ ID NO: 79

SNP Position Genomic: 239101

10 Related Interrogated SNP: hCV916107 (Power=.51)

SNP Source: dbSNP; Celera

Population (Allele,Count): caucasian (A,72|G,48)

SNP Type: INTRON; REPEATS

15 Context (SEQ ID NO: 180):

CATTCTTTATGGTGACTGTTCCAATCTTTCACCATGGTTCAATTTCTCAATCAC
TTCTATGAAGTTGCTCAAAACAAAAGCAATATACAGAAGTTTAACT

R

AAAGGTATAGCCCTATTTGACAGAGGGATAATTTATACATACAGTTTTGCACA
20 GATTATATGCACAGAATTTTAGTCCAAAAGTTGCAAAGTTCTGACAA

Celera SNP ID: hCV26887465

Public SNP ID: rs10802919

SNP in Genomic Sequence: SEQ ID NO: 79

SNP Position Genomic: 239506

25 Related Interrogated SNP: hCV916107 (Power=.6)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (A,75|G,45)

SNP Type: INTRON

30 Context (SEQ ID NO: 181):

GACCTTTCAGTACATTCTAGGTTTGGGGACTTTGTTGTGCAGAGCAATTGGAT
AATTATACAAGTCAATGGAAACAACAGTCACTCGTTGCATTACAGAA

Y

GCAGTATTTAAATGCAAACGCAACACCAAAAACAAATAAGGTGGGGAGTAA
AAGAAAAGAGTGTGGCTAGCTTCTAGGTAAATAAGACAGTTGGTCCCA

Celera SNP ID: hCV31714447

Public SNP ID: rs10926387

5 SNP in Genomic Sequence: SEQ ID NO: 79

SNP Position Genomic: 205022

Related Interrogated SNP: hCV916107 (Power=.51)

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (T,70|C,48)

10 SNP Type: INTRON

Gene Number: 12

Celera Gene: hCG39884 - 65000099284565

Gene Symbol:

15 Protein Name:

Celera Genomic Axis: GA_x5YUV32VWMC (16989259..17009788)

Chromosome: 1

OMIM NUMBER:

OMIM Information:

20

Genomic Sequence (SEQ ID NO: 80):

SNP Information

Context (SEQ ID NO: 182):

25 TAAGCCATTATTTCTTCAAGTACTTTTTTTTTTTAAGTTCAGCTTTTTATTGAAC

ACATTATAAAAGAGGTTTCGTCAAAAAGACCAAAGCCCATGTCAC

Y

ATCAGACTTCTCGGATTCTTCTTTCTTTGCTTCCACTTTCTTCTCCTCAGCTGGA

GCAGCAGCAGCAGAGGCGGGCAGAAGCTCCTGCTGGTACAGATAG

30

Celera SNP ID: hCV916107

Public SNP ID: rs670659

SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 10146

SNP Source: dbSNP; HapMap

Population (Allele,Count): caucasian (C,78|T,40)

SNP Type: INTRON

Context (SEQ ID NO: 183):

5 TGTCTGCAGTTTTCAAATGTTTATTAAGTGGCTTGGTATGCATTTTCGTTGGGCT
TATTTATCCTATTTGGGCTTCTCTCAGCTTCTTCAACTGAAGGTTT

R

TATCTTTGCTATGCTTGTGATGTTTTTAAGCCATTATTTCTTCAAGTACTTTTTT
TTTTTAAGTTCAGCTTTTTATTGAACACATTATAAAAGAGGTTTC

10 Celera SNP ID: hCV916106

Public SNP ID: rs575226

SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 10019

Related Interrogated SNP: hCV916107 (Power=.7)

15 SNP Source: dbSNP; HapMap; HGBASE

Population (Allele,Count): caucasian (A,77|G,41)

SNP Type: UTR3; INTRON; REPEATS

Context (SEQ ID NO: 184):

20 GAGAGATGAACTCCAGCCCAGCTACAGACTATGTTCTCTAACCCCTTGATCTA
AGAAATCTCAGCAATGGCAAAAATCTGGATATCACCAATTAATGGCC

R

GTGTCCTAGTTCTGTCTGCCTTTGGCCAAAAGCGGGTACTGAAAAAAGAGAA
AAGCAGCCACTGGCTTCCAGGAACTGGCTGGGCACTCACAGCTATGTC

25 Celera SNP ID: hCV26887441

Public SNP ID: rs9786932

SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 453

Related Interrogated SNP: hCV916107 (Power=.51)

30 SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,74|A,46)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

Context (SEQ ID NO: 185):

GGAAAATCTTTTATTTGAAATAAGAGTAAGAAGTCAATTGAGAGAATTCCAA
AAGGTGCCCAAAGCAAAAACAAAAGATTGAGGTTTAAATGAACTGTA

Y

GGTATAAAAAGGCAACAATTTTAAAAGTAAGATAGTAAGGATAAATAGTAA
5 AATAGTAAGGCTTTCTAATTAAGATACAAAACCTAGCAGCCATAAAAC

Celera SNP ID: hCV26887461

Public SNP ID: rs4660023

SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 13038

10 Related Interrogated SNP: hCV916107 (Power=.6)

SNP Source: dbSNP; Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,74|T,46)

SNP Type: INTRON

15 Context (SEQ ID NO: 186):

AGTGGGGAAGCCCAGCAACGCCTGTCCCTCTAGATTCTTCTGGCCTCTCTGAG
CAGCATTCTTCTTCTGGGTATGGGGCAGAGCTCTCTCTGGACTGG

R

GGTCTTAGGACCTACAGTCCAACAAGGCAGGTCAGATAACTTCTTTATGACTA
20 CTTTACGGCTGGCTTTGAGGAGAAGGGGTTCTGATTTCCATGACCC

Celera SNP ID: hCV26887463

Public SNP ID: rs6680767

SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 16445

25 Related Interrogated SNP: hCV916107 (Power=.6)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,76|A,42)

SNP Type: INTRON; REPEATS

30 Context (SEQ ID NO: 187):

GAGCAGCATTCTTCTTCTGGGTATGGGGCAGAGCTCTCTCTGGACTGGGGG
TCTTAGGACCTACAGTCCAACAAGGCAGGTCAGATAACTTCTTTATG

R

CTACTTTTACGGCTGGCTTTGAGGAGAAGGGGTTCTGATTTCCATGACCCGCC
TCGGGCCAGAGGGATTCTAGTCTCTGGGGAGAATGGAACAGGTCAGG

Celera SNP ID: hCV26887464

Public SNP ID: rs6669640

5 SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 16495

Related Interrogated SNP: hCV916107 (Power=.51)

SNP Source: dbSNP; Celera

Population (Allele,Count): caucasian (A,72|G,48)

10 SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 188):

CATTCTTTATGGTGACTGTTCCAATCTTTCACCATGGTTCAATTTCTCAATCAC
TTCTATGAAGTTGCTCAAAACAAAAGCAATATACAGAAGTTTAACT

15 R

AAAGGTATAGCCCTATTTGACAGAGGGATAATTTATACATACAGTTTTGCACA
GATTATATGCACAGAATTTTAGTCCAAAAGTTGCAAAGTTCTGACAA

Celera SNP ID: hCV26887465

Public SNP ID: rs10802919

20 SNP in Genomic Sequence: SEQ ID NO: 80

SNP Position Genomic: 16900

Related Interrogated SNP: hCV916107 (Power=.6)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (A,75|G,45)

25 SNP Type: INTRON

Gene Number: 13

Celera Gene: hCG2022817 - 104000117367435

Gene Symbol:

30 Protein Name:

Celera Genomic Axis: GA_x5YUV32VYQG (75856160..75894694)

Chromosome: 3

OMIM NUMBER:

OMIM Information:

Genomic Sequence (SEQ ID NO: 81):

SNP Information

5 Context (SEQ ID NO: 189):

TGGGACCCTGGGTCAGCAGAGCTACAGGGACAGGACACTTACCTAGAGCCAT
GGGGATGGGGCATTGTAGAGTCAAAGGGGCAGAGCCATGGGAATGATG

Y

10 TGCCACCCCAGTGGGCTTAGAAAGCAAGACACTAAGCCAAAGAAGACTATAC
TCAAGAATTAAGATCCAGTGGAATTTGCCTTGCCAGGTCGTGGACTTG

Celera SNP ID: hCV105917

Public SNP ID: rs9289134

SNP in Genomic Sequence: SEQ ID NO: 81

SNP Position Genomic: 17653

15 Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (T,84|C,36)

SNP Type: INTRON; REPEATS

20 Context (SEQ ID NO: 190):

CAAGAGTCTCACTTCACAGTTTAAGGACTCACATAGACTGAAAGTGAAAGGA
TGAAAAAAGATACTCAATGCAAATGGAAACCAAATAGAACAAGTGTA

R

25 CTAAACTTACATCAGATAAAAATAGGTTTTAAGTAAAAAACTATAAAAATGAGA
CCAAAAAGGTCATTGTGTGGTGATAAAGGGGTCGATTCATCAAGAGGC

Celera SNP ID: hCV178227

Public SNP ID: rs13070374

SNP in Genomic Sequence: SEQ ID NO: 81

SNP Position Genomic: 23013

30 Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (G,66|A,52)

SNP Type: INTRON; REPEATS

Context (SEQ ID NO: 191):

ATATCTACCCCCACCACAATGCTATTTAATACTGTATTAATCGTTCCAACAACTAA
TGCAATAAGGAAAGAAAAAAGCATAAAGATCAGAAAAGAAGAAAAC

W

5 GTCTTTCTTTGCAGAAAACATAATTATTACATTGAACATCCTCAGTAATACT
AAGGAATAACTACTAGAACTATTTAATAAAGTCACAGTTATATCTCA

Celera SNP ID: hCV192027

Public SNP ID: rs9821892

SNP in Genomic Sequence: SEQ ID NO: 81

10 SNP Position Genomic: 35883

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap; ABI_Val

Population (Allele,Count): caucasian (A,84|T,36)

SNP Type: INTERGENIC; UNKNOWN

15

Context (SEQ ID NO: 192):

GTGTGTTTCATGCAACAGCCATAGAGATGGTGCTATTGAAAGGCAATCTGCCC
ATCCTGTGTATTGACAAGCATCAGGATTCCAGGAGGATTCTTCCTGAT

K

20 CAAAGACTGTCAAATAGGAAAAAAGAAAACAGTAAATGAACAAAAACATCA
ACATGGAGTCCAAGTAATGAATTTATAGAGGAAAATGTATGCTACTCTG

Celera SNP ID: hCV255886

Public SNP ID: rs10511394

SNP in Genomic Sequence: SEQ ID NO: 81

25 SNP Position Genomic: 28535

Related Interrogated SNP: hCV263841 (Power=.6)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (T,85|G,35)

SNP Type: UTR3; INTERGENIC; UNKNOWN

30

Context (SEQ ID NO: 193):

ATATTGTCTATAAAAATGTAAATTCAAGGACCAAAGCATTGGCTACCTTATTA
TTGTGATAAAAAGCAAGCCTTAAGGTTGGGGTAGTCATTTTTCTGACA

Y

GCAAATTTAATAAGATTAATAATTTACACAATTTGCATTGAATGCATCTGCCT
TTACTTCTGATTATGGTACCTACTCTATGACATCCATAGTGACCTGT

Celera SNP ID: hCV278948

Public SNP ID: rs1464599

5 SNP in Genomic Sequence: SEQ ID NO: 81

SNP Position Genomic: 15407

Related Interrogated SNP: hCV263841 (Power=.6)

SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (T,81|C,35)

10 SNP Type: INTRON

Context (SEQ ID NO: 194):

TTACTTATAATGGCTTGAGTTAAAATGCTCCAACCTTTTCTAATTTTTGCTTTTG
AAAATAAAACACCTATAGTGACAATAAAAACACACATTTAAAGCCA

15 W

GTGGTCATAGTGACTGGAGCCATGTAGTGAATTGGGGCTAGACTAGTACGAT
CAATAGCGATACAAAATCACAATGAACTGTCTGAGGCAGAAAACCTTTT

Celera SNP ID: hCV1834237

Public SNP ID: rs9865270

20 SNP in Genomic Sequence: SEQ ID NO: 81

SNP Position Genomic: 29883

Related Interrogated SNP: hCV263841 (Power=.51)

SNP Source: dbSNP; Celera; HapMap

Population (Allele,Count): caucasian (A,84|T,36)

25 SNP Type: INTERGENIC; UNKNOWN

Context (SEQ ID NO: 195):

AGCATCATATCATATGAAAATATATAAAGTTGTAAATTAGGCTAACATACTGT
TATATAAAATTTATTCTACCCACCTACTTTCACAAACATATCTTTGA

30 R

ATGATAACATTGAAAAACGGGTAAGCAGCTATGGGCTTTATGTAATTGAAAG
ATGGCAAGAACCAAAGATGGTAAGCTTAGTTATTATTGTACATGTCTG

Celera SNP ID: hCV1834242

Public SNP ID: rs11712308

SNP in Genomic Sequence: SEQ ID NO: 81
 SNP Position Genomic: 27201
 Related Interrogated SNP: hCV263841 (Power=.51)
 SNP Source: dbSNP; Celera; HapMap
 5 Population (Allele,Count): caucasian (G,64|A,56)
 SNP Type: INTRON

Context (SEQ ID NO: 196):

GGTATTTGCAAGACCCAGGTCAAGAGGATGAATTGAAGCATCATATCATATG
 10 AAAATATATAAAGTTGTAAATTAGGCTAACATACTGTTATATAAAATT
 K
 ATTCTACCCACCTACTTTCACAAACATATCTTTGAGATGATAACATTGAAAAA
 CGGGTAAGCAGCTATGGGCTTTATGTAATTGAAAGATGGCAAGAACC

Celera SNP ID: hCV1834243
 15 Public SNP ID: rs9682652
 SNP in Genomic Sequence: SEQ ID NO: 81
 SNP Position Genomic: 27165
 Related Interrogated SNP: hCV263841 (Power=.6)
 SNP Source: dbSNP; Celera; HapMap; ABI_Val
 20 Population (Allele,Count): caucasian (T,85|G,35)
 SNP Type: INTRON

Context (SEQ ID NO: 197):

CTAAGGCTCAAAGAGAAGTATAAGACTAGGTCTTTATATTTGTAAGTTACTTA
 25 CACTCTAGATATCAAACGAAATGTA ACTCACATGAAGTGCTTCAATG
 Y
 TGAATAGGTTGGTACAACTCAAGTGCGGTAAGTGGTAAATCTGGGCTGATG
 TAGTCAGAGGAGACTTCCTGGAGGATCTTGGACATGAGCTTTCCAGGT

Celera SNP ID: hCV1834256
 30 Public SNP ID: rs2472662
 SNP in Genomic Sequence: SEQ ID NO: 81
 SNP Position Genomic: 13373
 Related Interrogated SNP: hCV263841 (Power=.51)
 SNP Source: dbSNP; Celera; HGBASE

Population (Allele,Count): caucasian (C,64|T,54)

SNP Type: INTRON

Context (SEQ ID NO: 198): \

5 AAATTCAGCTCTGGTCTTCCCGAGTGTAAGCCTTAGCCAGATTATTTGTGTT
GCTGGAGCAGGCACACTTCACTTAAAAAATGGACTGATATAAAGCCC

Y

AGTTAATTTTTTTTCCCAGATAGCTATAACTCATAAGAGATTCAGATATGATG
AGTAAATCATCTGGATGAACA ACTTGTGTTTAGCACTGTTCTTTACT

10 Celera SNP ID: hCV1834260

Public SNP ID: rs4688033

SNP in Genomic Sequence: SEQ ID NO: 81

SNP Position Genomic: 2785

Related Interrogated SNP: hCV263841 (Power=.6)

15 SNP Source: Celera; HapMap; HGBASE

Population (Allele,Count): caucasian (C,85|T,35)

SNP Type: TRANSCRIPTION FACTOR BINDING SITE; INTRON

Context (SEQ ID NO: 199):

20 TGTA ACTATAAGGTGTTTTATATATTCCTCATGGTAAGTACAAAACAAAAAAC
CTATATTAATAACACAAAAGGTGAAGAGTAAGAAATCAAGGCATGGC

R

CTAGAGAAAGTCACCTAATTACAAAGGAATACAGCAAGAGAGGAAGAAAGG
AATAAGGACCTACAAAACAACCATAAAACAATGAACCAGATGATAGTAG

25 Celera SNP ID: hCV30562884

Public SNP ID: rs9815093

SNP in Genomic Sequence: SEQ ID NO: 81

SNP Position Genomic: 22697

Related Interrogated SNP: hCV263841 (Power=.51)

30 SNP Source: dbSNP; HapMap; ABI_Val

Population (Allele,Count): caucasian (G,84|A,36)

SNP Type: INTRON; REPEATS

Table 3, Primer Sequences.

Marker	Alleles	Sequence A (Allele-specific Primer)	Sequence B (Allele-specific Primer)	Sequence C (Common Primer)
hCV11541681	C/G	CCACCTTCTCTTCAGATTAC (SEQ ID NO:200)	CCACCTTCTCTTCAGATTAC (SEQ ID NO:201)	CTGGACGGGTATGTAGAC (SEQ ID NO:202)
hCV12092542	T/C	GCCTTGCTCTCAATTTTGGT (SEQ ID NO:203)	CCTTGCTTCAATTTTGGC (SEQ ID NO:204)	AACAGTTAAGATGTTGGAATACCT (SEQ ID NO:205)
hCV2403368	C/G	AAAGGTGAGCTCCACCATAC (SEQ ID NO:206)	AAAGGTGAGCTCCACCATAG (SEQ ID NO:207)	GGCATCGGCACATAGTAGA (SEQ ID NO:208)
hCV25990131	A/C	TGTTGGTAAAGTATGTAGGATC TT (SEQ ID NO:209)	TGTTGGTAAAGTATGTAGGATCTG (SEQ ID NO:210)	CCATTTATAGAATGAGTGAGATGATAT (SEQ ID NO:211)
hCV263841	A/C	AGCTAATACTCCTGTCCTGAAA (SEQ ID NO:212)	AAGCTAATACTCCTGTCCTGAAC (SEQ ID NO:213)	AAGTTCCTTGTCCGACTTCT (SEQ ID NO:214)
hCV2744023	A/G	CTACCTGTGTCTGTATTTCTTTAA AT (SEQ ID NO:215)	CTACCTGTGTCTGTATTTCTTTAAAC (SEQ ID NO:216)	ATCTGCCTTCTTAGTGAATTGAT (SEQ ID NO:217)
hCV596331	A/G	AGTCCACATCAGGAAAATAGT (SEQ ID NO:218)	GTCCACATCAGGAAAATAGC (SEQ ID NO:219)	CCATTTGCCAATGAGAAAATATCAGGTT ACT (SEQ ID NO:220)
hCV916107	C/T	AGAATCCGAGAAGTCTGATG (SEQ ID NO:221)	GAAGAATCCGAGAAGTCTGATA (SEQ ID NO:222)	TTCAGCTTTTATTGAACACATTATA (SEQ ID NO:223)

Table 4. Replication of association of the 8 polymorphisms in unstratified analysis of LETS and MEGA genotyping

Marker	Gene Symbol	RS	Non Risk			Case			Control			Discovery (LETS)			Replication (MEGA)		
			Allele	Risk Allele	Risk AF	Risk AF	Risk AF	Risk AF	p Val	OR	p Val	OR	Mode	p Val	OR	Mode	
hCV11541681	LOC200420	rs2001490	G	C	43.1%	38.4%	0.033	1.57	0.024	1.27	Hom vs. Ref	0.024	1.27	Hom vs. Ref			
hCV12092542	CASP5*	rs507879	C	T	56.1%	54.6%	0.383	1.14	0.657	1.04	Recessive	0.657	1.04	Recessive			
hCV2403368	C1QTNF6	rs229526	C	G	77.9%	74.1%	0.022	1.37	0.034	1.16	Recessive	0.034	1.16	Recessive			
hCV25990131	CYP4V2	rs13146272	C	A	69.1%	64.9%	0.004	1.97	0.033	1.28	Hom vs. Ref	0.033	1.28	Hom vs. Ref			
hCV263841	NR112	rs1523127	A	C	42.1%	33.1%	0.0001	2.23	0.007	1.32	Hom vs. Ref	0.007	1.32	Hom vs. Ref			
hCV2744023	CASP8AP2	rs369328	G	A	52.0%	44.9%	0.003	1.81	0.024	1.25	Hom vs. Ref	0.024	1.25	Hom vs. Ref			
hCV596331	F9	rs6048	G	A	73.0%	67.3%	0.009	1.31	0.005	1.16	Allelic	0.005	1.16	Allelic			
hCV916107	RGS7	rs670659	T	C	69.8%	64.5%	0.027	1.67	0.011	1.33	Hom vs. Ref	0.011	1.33	Hom vs. Ref			

LETS indicates Leiden Thrombophilia Study; MEGA indicates Multiple Environmental and Genetic Assessment Study

Hom vs. Ref indicates a Fisher exact calculation that compares the counts of the risk homozygote genotype in cases and controls compared to the reference homozygote genotype in cases and controls.

Allelic indicates a Fisher exact calculation that compares the counts of the risk allele in cases and controls to counts of the non-risk allele in cases and controls.

Recessive indicates a Fisher exact calculation that compares the counts of the risk homozygote genotype in cases and controls compared to the sum of the heterozygote and reference homozygote genotype in cases and controls.

*CASP5 was not a hit in unstratified analysis, but was a hit in gender stratified analysis (see Table 5 below)

Table 5. Replication of association of 2 of the 8 polymorphisms in gender stratified analysis of LETS and MEGA genotyping

Marker	Gene Symbol	RS	Non Risk		Risk Allele	Strata	Case		Control		Discovery (LETS)			Replication (MEGA)		
			Allele	Risk Allele			Risk AF	Risk AF	Risk AF	p Val	OR	p Val	OR	p Val	OR	Mode
hCV12092542	CASP5	rs507879	C	T	T	All	56.1%	54.6%	0.383	1.14	0.657	1.04	Recessive			
hCV12092542	CASP5	rs507879	C	T	T	F	58.9%	52.3%	0.036	1.52	0.032	1.25	Recessive			
hCV12092542	CASP5	rs507879	T	C	C	M	47.6%	42.2%	0.192	1.43	0.304	1.15	Recessive			
hCV263841	NR1I2	rs1523127	A	C	C	All	42.1%	33.1%	0.0001	2.23	0.007	1.32	Hom.vs. Ref			
hCV263841	NR1I2	rs1523127	A	C	C	F	41.5%	34.9%	0.048	1.74	0.355	1.15	Hom.vs. Ref			
hCV263841	NR1I2	rs1523127	A	C	C	M	42.9%	30.7%	0.0003	3.17	0.004	1.54	Hom.vs. Ref			

LETS indicates Leiden Thrombophilia Study; MEGA indicates Multiple Environmental and Genetic Assessment Study.

Hom vs. Ref indicates a Fisher exact calculation that compares the counts of the risk homozygote genotype in cases and controls compared to the reference homozygote genotype in cases and controls.

Recessive indicates a Fisher exact calculation that compares the counts of the risk homozygote genotype in cases and controls compared to the sum of the heterozygote and reference homozygote genotype in cases and controls.

Strata: All indicates unstratified, F indicates Female only strata, M indicates male only strata.

Table 6. LD SNPs

Interrogated SNP	Interrogated rs	LD SNP	LD SNP rs	Power	Threshold r_T^2	r^2
hCV11541681	rs2001490	hCV11541694	rs12619258	0.51	0.97	1
hCV11541681	rs2001490	hCV11941453	rs2001436	0.51	0.97	1
hCV11541681	rs2001490	hCV26996674	rs13006448	0.51	0.97	1
hCV11541681	rs2001490	hCV26996679	rs6732812	0.51	0.97	1
hCV11541681	rs2001490	hCV26996688	rs13015885	0.51	0.97	1
hCV11541681	rs2001490	hCV26996689	rs13014700	0.51	0.97	1
hCV11541681	rs2001490	hCV26996690	rs2421575	0.51	0.97	1
hCV11541681	rs2001490	hCV26996697	rs12611487	0.51	0.97	1
hCV11541681	rs2001490	hCV28000363	rs4852972	0.51	0.97	1
hCV11541681	rs2001490	hCV31840149	rs12233112	0.51	0.97	1
hCV11541681	rs2001490	hCV31840159	rs13013228	0.51	0.97	1
hCV11541681	rs2001490	hCV95670	rs4852975	0.51	0.97	1
hCV11541681	rs2001490	hCV95671	rs11126414	0.51	0.97	1
hCV11541681	rs2001490	hDV69785784	rs13000788	0.51	0.97	1
hCV11541681	rs2001490	hDV70942181	rs17350056	0.51	0.97	1
hCV11541681	rs2001490	hDV70953030	rs17434634	0.51	0.97	1
hCV11541681	rs2001490	hDV70953035	rs17434655	0.51	0.97	1
hCV25990131	rs13146272	hCV15968026	rs2292426	0.51	0.478	0.89744
hCV25990131	rs13146272	hCV3230097	rs3736455	0.51	0.478	0.89744
hCV263841	rs1523127	hCV105917	rs9289134	0.51	0.46	0.56104
hCV263841	rs1523127	hCV11230788	rs7643038	0.51	0.46	0.96599
hCV263841	rs1523127	hCV134275	rs9847068	0.51	0.46	0.54527
hCV263841	rs1523127	hCV134278	rs9848716	0.51	0.46	0.54527
hCV263841	rs1523127	hCV15882316	rs2276706	0.51	0.46	0.96599
hCV263841	rs1523127	hCV178227	rs13070374	0.51	0.46	0.55951
hCV263841	rs1523127	hCV1833991	rs11926554	0.51	0.46	0.55479
hCV263841	rs1523127	hCV1834237	rs9865270	0.51	0.46	0.56104
hCV263841	rs1523127	hCV1834240	rs1581451	0.51	0.46	0.96599
hCV263841	rs1523127	hCV1834242	rs11712308	0.51	0.46	0.54861
hCV263841	rs1523127	hCV1834243	rs9682652	0.51	0.46	0.59664
hCV263841	rs1523127	hCV1834252	rs10934498	0.51	0.46	0.93196
hCV263841	rs1523127	hCV1834256	rs2472662	0.51	0.46	0.55854
hCV263841	rs1523127	hCV1834260	rs4688033	0.51	0.46	0.59664
hCV263841	rs1523127	hCV192027	rs9821892	0.51	0.46	0.56104
hCV263841	rs1523127	hCV255886	rs10511394	0.51	0.46	0.59664
hCV263841	rs1523127	hCV27504984	rs3814055	0.51	0.46	0.96599
hCV263841	rs1523127	hCV278948	rs1464599	0.51	0.46	0.59083
hCV263841	rs1523127	hCV29841665	rs7623217	0.51	0.46	0.56331
hCV263841	rs1523127	hCV30562884	rs9815093	0.51	0.46	0.56104
hCV263841	rs1523127	hCV30699687	rs11711386	0.51	0.46	0.56466
hCV263841	rs1523127	hCV30747432	rs12488820	0.51	0.46	0.92167
hCV263841	rs1523127	hCV9152783	rs1523130	0.51	0.46	0.93196
hCV2744023	rs369328	hCV2012681	rs6918308	0.51	0.755	0.89946
hCV2744023	rs369328	hCV26081278	rs2585008	0.51	0.755	1
hCV2744023	rs369328	hCV27392140	rs12661693	0.51	0.755	1
hCV2744023	rs369328	hCV28980044	rs6923604	0.51	0.755	0.89946
hCV2744023	rs369328	hCV29986476	rs9451269	0.51	0.755	0.86693
hCV2744023	rs369328	hCV30238224	rs6910277	0.51	0.755	0.89589
hCV2744023	rs369328	hCV3180835	rs292222	0.51	0.755	1
hCV2744023	rs369328	hDV71164887	rs2585018	0.51	0.755	1
hCV596331	rs6048	hCV2288124	rs440051	0.51	0.34	0.40494

hCV596331	rs6048	hCV26225376	rs3117074	0.51	0.34	0.40494
hCV596331	rs6048	hCV26225377	rs12008759	0.51	0.34	0.40494
hCV596331	rs6048	hCV2969899	rs434144	0.51	0.34	0.34304
hCV596331	rs6048	hCV2969900	rs434447	0.51	0.34	0.40494
hCV596331	rs6048	hCV2986569	rs11095801	0.51	0.34	0.40494
hCV596331	rs6048	hCV2986570	rs3117458	0.51	0.34	0.37257
hCV596331	rs6048	hCV2986572	rs4149670	0.51	0.34	0.39325
hCV596331	rs6048	hCV2986574	rs4149672	0.51	0.34	0.57908
hCV596331	rs6048	hCV2986575	rs4149674	0.51	0.34	0.64136
hCV596331	rs6048	hCV596323	rs438601	0.51	0.34	0.51923
hCV596331	rs6048	hCV596326	rs398101	0.51	0.34	0.82156
hCV596331	rs6048	hCV596330	rs422187	0.51	0.34	1
hCV596331	rs6048	hCV596335	rs413957	0.51	0.34	0.40494
hCV596331	rs6048	hCV596336	rs110583	0.51	0.34	0.35475
hCV596331	rs6048	hCV596339	rs370713	0.51	0.34	0.40494
hCV596331	rs6048	hCV596340	rs413536	0.51	0.34	0.40494
hCV596331	rs6048	hCV596344	rs445691	0.51	0.34	0.40494
hCV596331	rs6048	hCV596669	rs376165	0.51	0.34	0.71131
hCV916107	rs670659	hCV25653735	rs7520707	0.51	0.53	0.54793
hCV916107	rs670659	hCV26887401	rs10802916	0.51	0.53	0.5505
hCV916107	rs670659	hCV26887441	rs9786932	0.51	0.53	0.60409
hCV916107	rs670659	hCV26887450	rs670659	0.51	0.53	1
hCV916107	rs670659	hCV26887461	rs4660023	0.51	0.53	0.72845
hCV916107	rs670659	hCV26887463	rs6680767	0.51	0.53	0.71206
hCV916107	rs670659	hCV26887464	rs6669640	0.51	0.53	0.6131
hCV916107	rs670659	hCV26887465	rs10802919	0.51	0.53	0.69239
hCV916107	rs670659	hCV31714447	rs10926387	0.51	0.53	0.55437
hCV916107	rs670659	hCV916106	rs575226	0.51	0.53	1

This description contains a sequence listing in electronic form in ASCII text format. A copy of the sequence listing in electronic form is available from the Canadian Intellectual Property Office.

WHAT IS CLAIMED IS:

1. A method of determining whether a human has an increased risk for venous thrombosis (VT), the method comprising testing nucleic acid from said human for the presence or absence of a polymorphism in gene *F9* at position 101 of SEQ ID NO:82 or its complement, wherein A at position 101 of SEQ ID NO:82 or T at position 101 of its complement indicates said human has said increased risk for VT.
2. The method of claim 1, wherein said nucleic acid is a nucleic acid extract from a biological sample from said human.
3. The method of claim 2, wherein said biological sample is blood, saliva, or buccal cells.
4. The method of claim 2 or 3, further comprising preparing said nucleic acid extract from said biological sample prior to said testing.
5. The method of any one of claims 1 to 4, wherein said testing comprises nucleic acid amplification.
6. The method of claim 5, wherein said nucleic acid amplification is carried out by polymerase chain reaction (PCR).
7. The method of any one of claims 1 to 6, wherein said testing is performed using sequencing, 5' nuclease digestion, molecular beacon assay, oligonucleotide ligation assay, size analysis, single-stranded conformation polymorphism (SSCP) analysis, or denaturing gradient gel electrophoresis (DGGE).
8. The method of any one of claims 1 to 7, wherein said testing is performed using an allele-specific method.
9. The method of claim 8, wherein said allele-specific method comprises allele-specific probe hybridization, allele-specific primer extension, or allele-specific amplification.

10. The method of claim 8 or 9, wherein said allele-specific method is carried out using an allele-specific primer having a nucleotide sequence comprising SEQ ID NO:218 or SEQ ID NO:219.
11. The method of any one of claims 1 to 10 which is an automated method.
12. The method of any one of claims 1 to 11, wherein said human previously had venous thrombosis.
13. The method of any one of claims 1 to 12, wherein said VT is deep vein thrombosis (DVT).
14. The method of any one of claims 1 to 13, wherein said human is homozygous for said A or said T.
15. The method of any one of claims 1 to 13, wherein said human is heterozygous for said A or said T.
16. The method of any one of claims 1 to 15, further comprising correlating the presence of said A or said T with said increased risk for VT.
17. The method of any of one claims 1 to 13, further comprising correlating the absence of said A or said T with no said increased risk for VT.
18. The method of claim 16 or 17, wherein said correlating is performed by computer software.
19. An allele-specific polynucleotide for use in the method of determining whether a human has an increased risk for venous thrombosis (VT) as defined in any one of claims 1 to 18, wherein said polynucleotide comprises a segment of SEQ ID NO:82 or its complement at least 16 nucleotides in length that includes said position 101.
20. The allele-specific polynucleotide of claim 19, wherein said polynucleotide is detectably labeled.

21. The allele-specific polynucleotide of claim 20, wherein said polynucleotide is labeled with a fluorescent dye.

22. A kit for use in the method of determining whether a human has an increased risk for venous thrombosis (VT) as defined in any one of claims 1 to 18, wherein said kit comprises at least one polynucleotide as defined in claim 19, 20, or 21, and at least one further component, wherein the at least one further component is a buffer, deoxynucleotide triphosphates (dNTPs), an amplification primer pair, an enzyme, or any combination thereof.

23. The kit of claim 22, wherein said enzyme is a polymerase or a ligase.