# VTL – version 1.1
## (Validation & Transformation Language)

## Part 1 - General Description

*(DRAFT FOR PUBLIC REVIEW)*

*October 2016*

30

# Foreword

The Task force for the Validation and Transformation Language (VTL), created in 2012-2013 under the initiative of the SDMX Secretariat, is pleased to present the draft version of VTL 1.1.

The SDMX Secretariat launched the VTL work at the end of 2012, moving on from the consideration that SDMX already had a package for transformations and expressions in its information model, while a specific implementation language was missing. To make this framework operational, a standard language for defining validation and transformation rules (operators, their syntax and semantics) had to be adopted, while appropriate SDMX formats for storing and exchanging rules, and web services to retrieve them, had to be designed. The present VTL 1.1 package is only concerned with the first element, i.e. a formal definition of each operator, together with a general description of VTL, its core assumptions and the information model it is based on.

The VTL task force was set up early in 2013, composed of members of SDMX, DDI and GSIM communities and the work started in summer 2013. The intention was to provide a language usable by statisticians to express logical validation rules and transformations on data, whether described as dimensional tables or as unit-record data. The assumption is that this logical formalization of validation and transformation rules could be converted into specific programming languages for execution (SAS, R, Java, SQL, etc.) but would provide a "neutral" expression at business level of the processing taking place, against which various implementations can be mapped. Experience with existing examples suggests that this goal would be attainable.

An important point that emerged is that several standards are interested in such a language. However, each standard operates on its model artefacts and produces artefacts within the same model (property of closure). To cope with this, VTL has been built upon a very basic information model (VTL IM), taking the common parts of GSIM, SDMX and DDI, mainly using artefacts from GSIM 1.1, somewhat simplified and with some additional detail. This way, existing standards (GSIM, SDMX, DDI, others) may adopt VTL by mapping their information model against the VTL IM. Therefore, although a work-product of SDMX, the VTL language in itself is independent of SDMX and will be usable with other standards as well. Thanks to the possibility of being mapped with the basic part of the IM of other standards, the VTL IM also makes it possible to collect and manage the basic definitions of data represented in different standards.

For the reason described above, The VTL specifications are designed at a logical level, independent of any other standard, including SDMX. The VTL specifications, therefore, are self-standing and can be implemented either on their own or by other standards (such as SDMX). In particular, the work for the SDMX implementation of VTL is taking place in parallel to the work for designing the VTL 1.1 version, and will entail a future update of the SDMX documentation.

The first public consultation on VTL (version 1.0) was held in 2014. Many comments were incorporated in the VTL 1.0 version, published in March 2015. Other suggestions for

72 improving the language, received afterwards, fed the discussion for building the present draft
73 version 1.1, which contains many new features.

74 The VTL 1.1 package, containing the general VTL specifications independent of other
75 standards possible implementations, will include, in its final release:

a) Part 1 – the user manual, highlighting the main characteristics of VTL, its core
assumptions and the information model on which the language is based;
b) Part 2 – the reference manual, containing the full library of operators ordered by
category, including examples; this version will support more validation and
compilation needs compared to VTL 1.0.
c) eBNF notation (extended Backus-Naur Form) which is the technical notation to be
used as a test bed for all the examples.

83 The present document (part 1) contains the general part, highlighting the main characteristics
84 of VTL, its core assumptions and the information model on which VTL is based.

85 The latest version of VTL is freely available online at https://sdmx.org/?page_id=5096

86

# Table of contents

## 155 Introduction

156    This document presents the Validation and Transformation Language (also known as 'VTL').

157    The purpose of VTL is to allow a formal and standard definition of algorithms to validate
158    statistical data and calculate derived data.

159    The first development of VTL aims at enabling, as a priority, the formalisation of data
160    validation algorithms rather than tackling more complex algorithms for data compilation. In
161    fact, the assessment of business cases showed that the majority of the institutions ascribes
162    (prescribes) a higher priority to a standard language for supporting the validation processes
163    and in particular to the possibility of sharing validation rules with the respective data
164    providers, in order to specify the quality requirements and allow validation also before
165    provision.

166    This document is the outcome of a second iteration of the first phase, and therefore still
167    presents a version of VTL primarily oriented to support the data validation. However, as the
168    features needed for validation also include simple calculations, this version of VTL can
169    support basic compilation needs as well. In general, validation is considered as a particular
170    case of transformation; therefore, the term "Transformation" is meant to be more general,
171    including validation as well.

172    The main categories of operators and functions included in this version of the VTL-ML syntax
173    are:

174         General purpose     (e.g. assignment, data access, data storage ...)
175         String     (e.g. substring, concatenation, length ...)
176         Numeric     (e.g. +, -, *, /, round, absolute value ...)
177         Boolean     (e.g. and, or, not ...)
178         Date     (e.g. string from date)
179         Set     (e.g. union, intersection, ...)
180         Statistical     (e.g. aggregate, analytic functions ...)
181         Data validation     (e.g. check ... of value domains, references, figures ...)
182         Time series     (e.g. time shift ...)
183         Conditional     (e.g. if-then-else ...)
184         Clauses     (e.g. keep, calc, attrcalc ...)

185    The VTL-ML includes operators for defining:

186         IM artefacts     (e.g. Dataset, Datastructure ...)

187         Ruleset     (e.g. mapping ...)

188

189    Although VTL is developed under the umbrella of the SDMX governance, DDI and GSIM users
190    may also be highly interested in adopting a language for validation and transformation. In
191    particular, organizations involved in the SDMX, DDI and GSIM communities and in the High-
192    Level Group for the modernisation of statistical production and services (HLG) expressed
193    their wish of having a unique language, usable in SDMX, DDI and GSIM.

194 Accordingly, the task-force working for the VTL development agreed on the objective of
195 adopting a common language, in the hope of avoiding the risk of having diverging variants.

196 As a consequence, VTL is designed as a language relatively independent of the details of
197 SDMX, DDI and GSIM. It is based on an independent information model (IM), made of the very
198 basic artefacts common to these standards. Other models can inherit the VTL language by
199 unequivocally mapping their artefacts to those of the VTL IM.

## 200 Structure of the document

201 The first part of the document is dedicated to the description of the general characteristics of
202 VTL.

203 The following part describes the Information Model on which the language is based. In
204 particular, it describes the model of the data artefacts for which the language is aimed to
205 validate and transform, the model of the variables and value domains used for defining the
206 data artefacts and the model of the transformations.

207 A third part explains the language fundamentals, i.e. the basic characteristics of manipulated
208 objects, operators, expressions, user-defined functions, core and derived parts of the language
209 and so on.

210 The fourth part clarifies some general features of the language (i.e. the core assumptions of
211 the VTL), such as the types of artefacts involved in the transformations, the general behaviour
212 for the operations on the data sets, the methods for referencing the data sets to be operated
213 on, and the general conventions for the grammar of the language.

214 A final part highlights some issues related to the governance of VTL developments and to
215 future work, following a number of comments, suggestions and other requirements which
216 were submitted to the task-force in order to enhance the VTL package.

217 A short annex gives some background information about the BNF (Backus-Naur Form) syntax
218 used for providing a context-free representation of VTL.

219 The Extended BNF (EBNF) representation of the VTL 1.0 package is available at
220 https://sdmx.org/?page_id=5096. The VTL 1.1 representation will be added as soon as it is
221 available.

222

# General characteristics of the VTL

This section lists and briefly illustrates some general high-level characteristics of the validation and transformation language. They have been discussed and shared as requirements for the language in the VTL working group since the beginning of the work and have been taken into consideration for the design of the language.

## User orientation

⇨ The language is designed for users without information technology (IT) skills, who should be able to define calculations and validations independently, without the intervention of IT personnel;

- o The language is based on a "user" perspective and a "user" information model (IM) and not on possible IT perspectives (and IMs)

- o As much as possible, the language is able to manipulate statistical data at an abstract/conceptual level, independently of the IT representation used to store or exchange the data observations (e.g. files, tables, xml tags), so operating on abstract (from IT) model artefacts to produce other abstract (from IT) model artefacts

- o It references IM objects and does not use direct references to IT objects

⇨ The language is intuitive and friendly (users should be able to define and understand validations and transformations as easily as possible), so the syntax is:

- o Designed according to mathematics, which is a universal knowledge;

- o Expressed in English to be shareable in all countries;

- o As simple, intuitive and self-explanatory as possible;

- o Based on common mathematical expressions, which involve "operands" operated on by "operators" to obtain a certain result;

- o Designed with minimal redundancies (e.g. possibly avoiding operators specifying the same operation in different ways without concrete reasons).

⇨ The language is oriented to statistics, and therefore it is capable of operating on statistical objects and envisages the operators needed in the statistical processes and in particular in the data validation phases, for example:

- o Operators for data validations and edit;

- o Operators for aggregation, even according to hierarchies;

- o Operators for dimensional processing (e.g. projection, filter);

- o At a later stage, operators for time series processing (e.g. moving average, seasonal adjustment, correlation) operators for statistics (e.g. aggregation, mean, median, percentiles, variance, indexes, correlation, sampling, inference, estimation);

## Integrated approach

⇨ The language is independent of the statistical domain of the data to be processed;

- VTL has no dependencies on the subject matter (the data content);
- VTL is able to manipulate statistical data in relation to their structure.

⇨ The language is suitable for the various typologies of data of a statistical environment (for example dimensional data, survey data, registers data, micro and macro, quantitative and qualitative) and is supported by an information model (IM) which covers these typologies;

- The IM allows the representation of the various typologies of data of a statistical environment at a conceptual/logical level (in a way abstract from IT and from the physical storage);
- The various typologies of data are described as much as possible in an integrated way, by means of common IM artefacts for their common aspects;
- The principle of the Occam's razor is applied as an heuristic principle in designing the conceptual IM, so keeping everything as simple as possible or, in other words, unifying the model of apparently different things as much as possible.

⇨ The language (and its IM) is independent of the phases of the statistical process and usable in any one of them;

- Operators are designed to be independent of the phases of the process, their syntax does not change in different phases and is not bound to some characteristic restricted to a specific phase (operators' syntax is not aware of the phase of the process);
- In principle, all operators are allowed in any phase of the process (e.g. it is possible to use the operators for data validation not only in the data collection but also, for example, in data compilation for validating the result of a compilation process; similarly it is possible to use the operators for data calculation, like the aggregation, not only in data compilation but also in data validation processes);
- Both collected and calculated data are equally permitted as inputs of a calculation, without changes in the syntax of the operators/expression;
- Collected and calculated data are represented (in the IM) in a homogeneous way with regards to the metadata needed for calculations.

⇨ The language is designed to be applied not only to SDMX but also to other standards;

- VTL, like any consistent language, relies on a specific information model, as it operates on the VTL IM artefacts to produce other VTL IM artefacts. In principle, a language cannot be applied as-is to another information model (e.g. SDMX, DDI, GSIM); this possibility exists only if there is a unambiguous correspondence between the artefacts of those information models and the VTL IM (that is if their artefacts correspond to the same mathematical notion);
- The goal of applying the language to more models/standards is achieved by using a very simple, generic and conceptual Information Model (the VTL IM),

| 301 | and mapping this IM to the models of the different standards (SDMX, DDI, |
| 302 | GSIM, ...); to the extent that the mapping is straightforward and unambiguous, |
| 303 | the language can be inherited by other standards (with the proper |
| 304 | adjustments); |

| 305 | o | To achieve an unambiguous mapping, the VTL IM is deeply inspired by the |
| 306 | | GSIM IM and uses the same artefacts when possible[1]; in fact, GSIM is designed |
| 307 | | to provide a formal description of data at business level against which other |
| 308 | | information models can be mapped; moreover, loose mappings between GSIM |
| 309 | | and SDMX and between GSIM and DDI are already available[2]; a very small |
| 310 | | subset of the GSIM artefacts is used in the VTL IM in order to keep the model |
| 311 | | and the language as simple as possible (Occam's razor principle); these are the |
| 312 | | artefacts strictly needed for describing the data involved in Transformations, |
| 313 | | their structure and the variables and value domains; |

314 o GSIM artefacts are supplemented, when needed, with other artefacts that are
315 necessary for describing calculations; in particular, the SDMX model for
316 Transformations is used;

317 o As mentioned above, the definition of the VTL IM artefacts is based on
318 mathematics and is expressed at an abstract user level.

## 319 Active role for processing

320 ⇨ The language is designed to make it possible to drive in an active way the execution of
321 the calculations (in addition to documenting them)

322 ⇨ For the purpose above, it is possible either to implement a calculation engine that
323 interprets the VTL and operates on the data or to rely on already existing IT tools (this
324 second option requires a translation from the VTL to the language of the IT tool to be
325 used for the calculations)

326 ⇨ The VTL grammar is being described formally using the universally known Backus
327 Naur Form notation (BNF), because this allows the VTL expressions to be easily
328 defined and processed; the formal description allow the expressions:

329 o To be automatically parsed (against the rules of the formal grammar); on the
330 IT level, this requires the implementation of a parser that compiles the
331 expressions and checks their correctness;

332 o To be automatically translated from the VTL to the language of the IT tool to
333 be used for the calculation; on the IT level, this requires the implementation of
334 a proper translator;

335 o To be automatically translated from one VTL version to another, e.g. following
336 an upgrade of the VTL syntax; on the IT level, this requires the implementation
337 of a proper translator also.

---

[1] See the next section (VTL Information Model) and the section "Relations with the GSIM Information model"

[2] See at: http://www1.unece.org/stat/platform/display/gsim/GSIM+and+standards;

338   ⇨ The inputs and the outputs of the calculations and the calculations themselves are
339      artefacts of the IM

340      o  This is a basic property of any robust language because it allows calculated
341         data to be operands of further calculations;

342      o  If the artefacts are persistently stored, their definition is persistent as well; if
343         the artefacts are non-persistently stored (used only during the calculation
344         process like input from other systems, intermediate results, external outputs)
345         their definition can be non-persistent;

346      o  Because the definition of calculations is based on the data structure definition
347         of its input artefacts, the latter must be available when the calculation is
348         defined;

349      o  The VTL is designed to make the data structure of the output of a calculation
350         deducible from the calculation algorithm and from the data structure of the
351         operands (this feature ensures that the calculated data can be defined
352         according to the IM and can be used as operands of further calculations);

353      o  In the IT implementation, it is advisable to automate (as much as possible) the
354         structural definition of the output of a calculation, in order to enforce the
355         consistency of the definitions and avoid unnecessary overheads for the
356         definers.

357   ⇨ The VTL and its information model make it possible to check automatically the overall
358      consistency of the definition of the calculations, including with respect to the artefact
359      of the IM, and in particular to check:

360      o  the correctness of the expressions with respect to the syntax of the language

361      o  the integrity of the expressions with respect to their input and output artefacts
362         and the corresponding structures and properties (for example, the input
363         artefacts must exist, their structure components referenced in the expression
364         must exist, qualitative data cannot be manipulated through quantitative
365         operators, and so on)

366      o  the consistency of the overall graph of the calculations (for example, in order
367         to avoid that the result of a calculation goes as input to the same calculation
368         there should not be cycles in the sequence of calculations, thus eliminating the
369         risk of producing unpredictable and erroneous results);

## Independence of IT implementation

371   ⇨ According to the "user orientation" above, the language is designed so that users are
372      not required to be aware of the IT solution;

373      o  To use the language, the users need to know only the abstract view of the data
374         and calculations and do not need to know the aspects of the IT
375         implementation, like the storage structures, the calculation tools and so on.

376   ⇨ The language is not oriented to a specific IT implementation and permits many
377      possible different implementations (this property is particularly important in order to
378      allow different institutions to rely on different IT environments and solutions);

| 379 | | o | On the technical level, the connection between the user layer and the IT layer |
| 380 | | | is left to the specific IT implementations; |

| 381 | | o | The VTL approach favours effective IT implementations that decouple the user |
| 382 | | | layer and the IT layer. |

⇨ The language does not require the awareness of the physical data structure; the operations on the data are specified according to the conceptual/logical structure, and so are independent of the physical structure; this ensures that the physical structure may change without necessarily affecting the conceptual structure and the user expressions;

- o Data having the same conceptual/logical structure may be accessed using the same statements, even if they have different IT structures;

- o The VTL provides for commands for data storage and retrieval at a conceptual/logical level; the mapping and the conversion between the conceptual and the physical structures of the data is left to the IT implementation (and users need not be aware of it);

- o By mapping the user and the IT data structures, the IT implementations can make it possible to store/retrieve data in/from different IT data stores (e.g. relational databases, dimensional databases, xml files, spread-sheets, traditional files);

⇨ The language does not require the awareness of the IT tools used for the calculations (e.g. routines in a programming language, statistical packages like R, SAS, Matlab, relational databases (SQL), dimensional databases (MDX), XML tools,…);

- o The syntax of the VTL is independent of existing IT calculation tools;

- o On the IT level, this may require a translation from the VTL to the language of the IT tool to be used for the calculation;

- o By implementing the proper translations at the IT level, institutions can use different IT tools to execute the same algorithms; moreover, it is possible for the same institution to use different IT tools within an integrated solution (e.g. to exploit different abilities of different tools);

- o VTL instructions do not change if the IT solution changes (for example following the adoption of another IT tool), so avoiding impacts on users as much as possible;

## Extensibility, customizability

⇨ It is possible to build and extend the language gradually, enriching the available operators according to the evolution of the business needs, so progressively making the language more powerful;

⇨ In addition, it is possible to call external routines of other languages/tools, provided that they are compatible with the IM; this requisite is aimed to fulfil specific calculation needs without modifying the operators of the language, so exploiting the power of the other languages/tools if necessary for specific purposes

| | |
|---|---|
| 419 | o  The external routines should be compatible with, and relate back to, the |
| 420 | conceptual IM of the calculations as for its inputs and outputs, so that the |
| 421 | integrity of the definitions is ensured |
| 422 | o  The external routines are not part of the language, so their use might be |
| 423 | subject to some limitations (e.g. it might be impossible to parse them as if they |
| 424 | were operators of the language) |
| 425 | o  The use of external routines has some drawbacks, because it may obviously |
| 426 | compromise the IT implementation independence, the abstraction and the |
| 427 | user orientation; therefore external routines should be used only for specific |
| 428 | needs and in limited cases, whereas widespread and generic needs should be |
| 429 | fulfilled through the operators of the language; |
| 430 | ⇨ Whilst an Organisation adopting VTL can extend it by defining customized parts, on its |
| 431 | own total responsibility, in order to improve the standard language for specific |
| 432 | purposes (e.g. for supporting possible algorithms not permitted by the standard part), |
| 433 | it is important that the customized parts remain compliant with the VTL IM and the |
| 434 | VTL core assumptions. Adopting Organizations are totally in charge of any possible |
| 435 | maintenance activity deriving from VTL modifications. Such extensions, however, are |
| 436 | not recommended because they can compromise the exchange of validation rules and |
| 437 | the use of common tools. |

## Language effectiveness

| | |
|---|---|
| 439 | ⇨ The language is oriented to give full support to the various typologies of data of a |
| 440 | statistical environment (for example dimensional data, survey data, registers data, |
| 441 | micro and macro, quantitative and qualitative, …) described as much as possible in a |
| 442 | coherent way, by means of common IM artefacts for their common aspects, and |
| 443 | relying on mathematical notions, as mentioned above. The various types of statistical |
| 444 | data are considered as mathematical functions, having independent variables |
| 445 | (Identifiers) and dependent variables (Measures, Attributes[3]), whose extensions can |
| 446 | be thought as logical tables (DataSets) made of rows (Data Points) and columns |
| 447 | (Identifiers, Measures, Attributes). |
| 448 | ⇨ The language supports operations on the Data Sets (i.e. mathematical functions) in |
| 449 | order to calculate new Data Sets from the existing ones, on the structure components |
| 450 | of the Data Sets (Identifiers, Measures, Attributes), on the Data Points. |
| 451 | ⇨ The algorithms are specified by means of mathematical expressions which compose |
| 452 | the operands (Data Sets, Components …) by means of operators (e.g. +,-,*,/,>,<) to |
| 453 | obtain a certain result (Data Sets, Components …); |
| 454 | ⇨ The validation is considered as a kind of calculation having as an operand the Data Set |
| 455 | to be validated and producing a Data Set containing the outcome of the validation |
| 456 | (typically having values "true" and "false" in the measure, respectively for successful |
| 457 | and unsuccessful validation); being a Data Set, the result of the validation can be |
| 458 | further processed (it can be input of further calculations); |

---

[3] The Measures bear information about the real world and the Attributes about the Data Set or some part of it.

459 ⇨ Calculations on multiple measures are supported, as well as calculations on the
460    attributes of the Data Sets and calculations involving missing values;

461 ⇨ The operations are intended to be consistent with the historical changes of the
462    artefacts (e.g. of the code lists, of the hierarchies …), so allowing a proper behaviour
463    for each reference period; however, because different standards may represent
464    historical changes in different ways, the implementation of this aspect is left to the
465    standards adopting the VTL (e.g. SDMX, DDI …) and therefore at the moment the VTL
466    specification does not prescribe any specific methodology for representing historical
467    changes of the artefacts (e.g. versioning, qualification of time validity);

468 ⇨ The language is ready to allow different algorithms for different reference times
469    (feature to be implemented at a later stage);

470 ⇨ the VTL operators are generally "modular", meaning that it is possible to compose
471    multiple operators in a single expression; in other words, an operator can have an
472    expression as operand, so obtaining a new expression, and this can be made
473    recursively;

474 ⇨ The final and the intermediate results of a calculation can be permanently stored (or
475    not) according to the needs;

476 ⇨ Multiple results may be calculated by means of multiple expressions.

477

# Evolution of VTL 1.1 in respect to VTL 1.0

Important contributions gave origin to the work that brought to this VTL 1.1 version.

Firstly, it was not possible to acknowledge immediately - in VTL 1.0 - all of the remarks received during the public review. Secondly, the publication of VTL 1.0 triggered the launch of reviews and proofs of concepts, by several institutions and organizations, aimed at assessing the ability of VTL of supporting properly their real use cases.

The suggestions coming from these activities had a fundamental role in designing the new version of the language.

The main improvements are described below.

## The Information Model

The VTL  Information Model describes the artefacts that VTL manipulates (i.e. it provides generic models for defining Data and their structures, Variables, Value Domains and so on) and how the VTL is used to define validations and transformations (i.e. a generic model for Transformations).

In VTL 1.1, some mistakes have been corrected and new kinds of artefacts have been added in order to make the representation more complete.

## The artefacts Definition Language

VTL 1.0 was initially intended to work on top of an existing standard, like SDMX, DDI or other, and therefore the definition of the artefacts to be manipulated (Data and their structures, Variables, Value Domains and so on) was assumed to be made using the implementing standards and not VTL itself. In other words, VTL 1.0 was not intended to define its artefacts and therefore only contains a manipulation language.

During the work for VTL 1.1, it was acknowledged as being very recommendable and useful to have a complete definition language in VTL, able to define all of the artefacts that VTL can manipulate. This is useful, first, to express structural and reusable definitions directly in VTL (even independently of other standards); second, to facilitate the use of VTL on top of other standards (through a proper mapping, the structural definitions of other standards could be translated into VTL definitions and vice-versa); third, to make it possible to check at parsing time the coherency of the VTL manipulation expressions against the structure of the artefacts to be manipulated (even defined through VTL).

Therefore, VTL 1.1 is also equipped with a definition language for VTL artefacts. In conclusion, in respect to VTL 1.0:

The VTL definition language (VTL-DL) is completely new (there is no definition language in VTL 1.0).

The VTL manipulation language (VTL-ML) has been upgraded (it is the evolution of the VTL 1.0 language).

## Reusable artefacts and rules

The artefacts defined by means of the VTL definition language (e.g. a set of code items) as well as the artefacts defined by means of an existing standard (like SDMX, DDI, or others) are reusable. In fact, the VTL manipulation language can reference these so called "structural" artefacts as many times as needed.

In order to empower the capability of reusing definitions, a main requirement for VTL 1.1 has been the introduction of reusable rules (for example, validation rules defined once and applicable to different cases).

Often, the same algorithm for manipulating data can be obtained by defining and referencing either structural artefacts or reusable rules. Current practices of various organizations show that both approaches are actually used. In order to empower the ability of the organizations of acknowledging and applying transformation/validation rules defined by others, which is one of the main goals of the VTL standard, the VTL structural artefacts and reusable rules are harmonized as much as possible. If needed, it should be feasible to convert the definitions of rules specified according to one approach (e.g. through reusable rules) into the other one (e.g. structural artefacts) and vice-versa.

The reusable artefacts and rules are defined through the VTL definition language and reused through the VTL manipulation language.

## The core language and the standard library

VTL 1.0 contains a flat list of operators, in principle not related to one another. A main suggestion for VTL 1.1 was to identify a core set of primitive operators able to express all of the other operators present in the language. This was done in order to specify more formally the semantics of available operators, avoiding possible ambiguities about their behaviour and fostering coherent implementations.  The distinction between 'core' and 'standard' library is largely of interest of the VTL technical implementers.

The suggestion above has been acknowledged, so that the VTL 1.1 manipulation language is made of a core set of primitive operators and a standard library of derived operators, definable in term of the primitive ones. The standard library contains VTL 1.0 operators (possibly enhanced) and new operators introduced with VTL 1.1.

The VTL core includes a mechanism called join expressions, described in the following sections, which allows the definition of derived dataset operators and their behaviour, including custom operators (not existing in the standard library) for specific purposes of some institutions.

## The functional paradigm

In the VTL Information Model, the various types of statistical data are considered as mathematical functions, having independent variables (Identifiers) and dependent variables (Measures, Attributes), whose extensions can be thought of as logical tables (DataSets) made of rows (Data Points) and columns  (Identifiers, Measures, Attributes). Therefore, the main artefacts to be manipulated using VTL are the logical DataSets, i.e. mathematical functions.

554 Accordingly, VTL uses a functional programming paradigm, meaning a paradigm that treats
555 computations as the evaluation of mathematical functions, avoiding changing-state and
556 mutable data (see also the Language Fundamentals section).

557 It was observed, however, that the functional paradigm is not completely achieved in VTL 1.0
558 and that in limited cases this might cause some problem.

559 Accordingly, some VTL 1.0 operators have been revised in order to enforce their functional
560 behaviour.

## New operators

561

562 VTL 1.1 introduces new operators. As already said, all of the operators of the VTL definition
563 language are completely new. A series of other new operators has been introduced in the VTL
564 manipulation language.

565 The complete list of the VTL 1.1 operators is in the reference manual.

# VTL Information Model

## Introduction

568    The VTL Information Model (IM) describes the artefacts that VTL can manipulate.

569    The knowledge of the artefacts is essential for performing VTL operations correctly.
570    Therefore, it is assumed that the referenced artefacts are defined beforehand.

571    The results of VTL expressions must be defined as well, because it must always be possible to
572    take these results as operands of further expressions to build a chain of transformations as
573    complex as needed. In other words, VTL is meant to be "closed", meaning that operands and
574    results of the VTL expressions are always artefacts of the VTL IM.

575    VTL can manage persistent or temporary artefacts, the former stored persistently in the
576    information system, the latter only used temporarily.

577    As already mentioned, VTL is designed to be used either on its own or on top of other
578    standards. It provides a formal description of data at business level against which the
579    information models of other standards can be mapped, so that through these possible
580    mappings to the definitions of VTL, artefacts can be obtained from the definitions of the
581    corresponding artefacts of the other standards and vice-versa.

582    This is the same purpose as the Generic Statistical Information Model (GSIM) and,
583    consequently, the VTL Information Model uses GSIM artefacts as much as possible (GSIM 1.1
584    version) [4]. Besides, GSIM already provides a first mapping with SDMX and DDI that can be
585    used for the technical implementation [5]. Note that the description of the GSIM 1.1 classes and
586    relevant definitions can be consulted in the "Clickable GSIM" of the UNECE site [6]. However, the
587    detailed mapping between the VTL IM and the IMs of the other standards is out of the scope of
588    this document and is left to the competent bodies of the other standards.

589    The VTL IM is illustrated in the following sections.

590    The first section describes the generic model for defining the statistical data and their
591    structures, which are the fundamental artefacts to be transformed. In fact, the ultimate goal of
592    the VTL is to act on statistical data to produce other statistical data.

593    In turn, the data are composed of variables, value domains, code items and similar artefacts.
594    These are the basic bricks that compose the data structures, fundamental for understanding
595    the meaning of the data and ensuring harmonization of various data when needed. The
596    second section presents the generic model for these kinds of artefacts.

---

[4] See also the section "Relations with the GSIM Information model"

[5] For the GSIM – DDI and GSIM – SDMX mappings, see also the relationships between GSIM and other standards at the UNECE site http://www1.unece.org/stat/platform/display/gsim/GSIM+and+standards. About the mapping with SDMX, however, note that here it is assumed that the SDMX artefacts Data Set and Data Structure Definition may represent both dimensional and unit data (not only dimensional data) and may be mapped respectively to the VTL artefacts Data Set and Data Structure.

[6] Hyperlink "http://www1.unece.org/stat/platform/display/GSIMclick/Clickable+GSIM"

Finally, the VTL transformations, written in the form of mathematical expressions, apply the operators of the language to proper operands in order to obtain the needed results. The third section depicts the generic model of the transformations.

## Generic Model for Data and their structures

This Section provides a formal model for the structure of data as operated on by the Validation and Transformation Language (VTL).

As already said, GSIM artefacts are used as much as possible. Some differences between this model and GSIM are due to the fact that, in the VTL IM, both unit and dimensional data are considered as mathematical functions having independent and dependent variables and are treated in the same way.

For each Unit (e.g. a person) or Group of Units of a Population (e.g. groups of persons of a certain age and civil status), identified by means of the values of the independent variables (e.g. either the "person id" or the age and the civil status), a mathematical function provides for the values of the dependent variables, which are the properties to be known (e.g. the revenue, the expenses …).

A mathematical function can be seen as a **logical table made of rows and columns**. Each column holds the values of a variable (either independent or dependent); each row holds the association between the values of the independent variables and the values of the dependent variables (in other words, each row is a single "point" of the function).

In this way, the manipulation of any kind of data (unit and dimensional) is brought back to the manipulation of very simple and well-known objects, which can be easily understood and managed by users. According to these assumptions, there would be no longer be a need to distinguish between unit and dimensional data; nevertheless, such a distinction is illustrated here in order to make it easier to map the VTL IM to the GSIM IM and, through GSIM, to the DDI and SDMX models.

Starting from this assumption, each mathematical function (logical table) may be defined as a GSIM Data Set and its structure as a GSIM Data Structure, having Identifier, Measure and Attribute Components. The Identifier components are the independent variables of the function, the Measures and Attribute Components are the dependent variables. Obviously, the GSIM artefacts "Data Set" and "Data Set Structure" have to be strictly interpreted as **logical artefacts** on a mathematical level, not necessarily corresponding to physical data sets and physical data structures.

Please note that the distinction between Dimensional and Unit Data is not used at all by VTL and is not part of the VTL IM. This distinction is present in this document just for clarifying the basic mapping between the VTL IMs and the GSIM and DDI IMs.

In order to avoid any possible misunderstanding with respect to SDMX, also take note that the VTL Data Set in general does not correspond to the SDMX Dataset. In fact, a SDMX dataset is a physical set of data (the data exchanged in a single interaction), while the VTL DataSet is a logical set of data, in principle independent of its possible handling (exchange, calculation and so on). The right mapping is between the VTL Data Set and the SDMX Dataflow.

## Data model diagram



White box:        same artefact as in GSIM 1.1
Light grey box:   similar to GSIM 1.1

## Explanation of the Diagram

**Data Set**: a mathematical function (logical table) that describes some properties of some groups of units of a population. In general, the groups of units may be composed of one or more units. For unit data, each group is composed of a single unit. For dimensional data, each group may be composed of any number of units. A VTL Data Set is considered as a logical set of observations (Data Points) having the same structure and the same general meaning, independently of the possible physical representation or storage. Between the VTL Data Sets and the physical datasets, there can be relationships of any cardinality: for example, a VTL Data Set may be stored either in one or in many physical data sets, as well as many VTL Data Sets may be stored in the same physical datasets (or database tables). The VTL Data Set is similar to the GSIM Data Set, the relationship between them is described in the following section.

**Data Point**: a single value of the function, i.e. a single association between the values of the independent variables and the values of the dependent variables. A Data Point corresponds to a row of the logical table that describes the function. A set of Data Points form the extension of the function (Data Set). The single Data Points do not need to be individually defined, because their definition is the definition of the function (i.e. the Data Set definition).  This artefact is the same as the GSIM Data Point.

**Data Structure:** the structure of a mathematical function, having independent and dependent variables. The independent variables are called "Identifier components", the dependent variables are called either "Measure Components" or "Attribute Components". The distinction between Measure and Attribute components is based on their meaning: the Measure Components give information about the real world, while the Attribute components give information about the function itself. The VTL Data Structure is similar to the GSIM Data Structure, the relationship between them is described in the following section.

**Data Structure Component**: any component of the data structure, which can be either an Identifier, or a Measure, or an Attribute Component. This artefact is the same as in GSIM.

> **Identifier Component** (or simply Identifier): a component of the data structure that is an independent variable of the function. This artefact is the same as in GSIM. In respect to SDMX, an Identifier Component may be either a **Unit Identifier**, which correspond to a SDMX Dimension, or a **Measure Identifier**, which corresponds to a SDMX Measure Dimension. The former is an identifier which contributes to the identification of the Units or Groups of Units, the latter is an identifier which contributes, when needed, to the identification of the Measure[7].

> **Measure Component** (or simply Measure): a component of the data structure that is a dependent variable of the function and gives information about the real world. This artefact is the same as in GSIM[8].

> **Attribute Component** (or simply Attribute): a component of the data structure that is a dependent variable of the function and gives information about the function itself. This artefact is the same as in GSIM.

> Note that the VTL manages Measure and Attribute Components in different ways, as explained in the section "The general behaviour of operations on datasets" below, therefore the distinction between Measures and Attributes is significant for the VTL.

## Relationships between VTL and GSIM

As mentioned earlier, the VTL Data Set and Data Structure artefacts are similar to the corresponding GSIM artefact. VTL, however, does not make a distinction between Unit and Dimensional Data Sets and Data Structures.

In order to explain the relationships between VTL and GSIM, the distinction between Unit and Dimensional Data Sets can be introduced virtually even in the VTL artefacts. In particular, the

---

[7] There can be from 0 to N Identifiers in a Data Structure. The particular case of 0 Identifiers and 1 Measure denotes scalar values, while the particular case of 0 Identifiers and N Measures denote vectors of scalar values.

[8] There can be from 0 to N Measures in a Data Structure. The particular case of 0 Measures denotes a "pure" relationship between the Identifiers (i.e. a relationship that does not have properties). For example, the relationship between the "students" and the "courses" that they follow (without any other information): the corresponding Data Set has StudentId and CourseId as Identifiers and do not have any explicit measure. However, as the existing combination of identifiers are implicitly considered as "TRUE", it can be thought that there is an implicit Boolean measure having the constant value "TRUE".

708  GSIM Data Set may be a GSIM Dimensional Data Set or a GSIM Unit Data Set, while a VTL Data
709  Set may (virtually) be:

710      either a (virtual) **VTL Dimensional Data Set**: a kind of (Logical) Data Set describing
711      groups of units of a population that may be composed of many units. This (virtual)
712      artefact would be the same as the GSIM Dimensional Data Set;

713      or a (virtual) **VTL Unit Data Set**: a kind of (Logical) Data Set describing single units of
714      a population. This (virtual) artefact would be the same as the Unit Data Record in
715      GSIM, which has its own structure and can be thought of as a mathematical function.
716      The difference is that the VTL Unit Data Set would not correspond to the GSIM Unit
717      Data Set, because the latter cannot be considered as a mathematical function: in fact it
718      can have many GSIM Unit Data Records with different structures.

719  A similar relationship exists between VTL and GSIM Data Structures. In particular, introducing
720  in VTL the virtual distinction between Unit and Dimensional Data Structures, while a GSIM
721  Data Structure may be a GSIM Dimensional Data Structure or a GSIM Unit Data Structure, a
722  VTL Data Structure may (virtually) be:

723      either a (virtual) **VTL Dimensional Data Structure**: the structure of (0..n)
724      Dimensional Data Sets. This artefact would be the same as in GSIM;

725      or a (virtual) **VTL Unit Data Structure**: the structure of (0..n) Unit Data Sets. This
726      artefact would be the same as the Logical Record in GSIM, which corresponds to a
727      single structure and can be thought as the structure of a mathematical function. The
728      difference is that the VTL Unit Data Structure would not correspond to the GSIM Unit
729      Data Structure, because the latter cannot be considered as the structure of a
730      mathematical function: in fact, it can have many Logical Records with different
731      structures.

732  GSIM – VTL mapping diagram:

733

734

735          GSIM                    mappings        Virtual artefacts
736          Unit DataRecord                         VTL Unit Data
                                                     Set
                                                                         VTL
737          GSIM Dimens.                                                 Data Set
738          Data Set                                VTL Dimens.
                                                     Data Set            0..N
739
                                                                         structured by
740          GSIM                                    VTL Unit Data
741          Logical Record                          Structure           1..1
742                                                                      VTL
            GSIM Dimens.                                                 Data Structure
743         Data Structure                          Dimensional
                                                    Data Structure

744

745  The distinction between Dimensional and Unit Data Set and Data Structure is not used by the
746  VTL language and is not part of the VTL IM. This virtual distinction is highlighted here just for
747  clarifying the mapping of the VTL IM with GSIM and DDI.

**Examples**

As a first simple example of Data Sets seen as mathematical functions, let us consider the following table:

*Production of the American Countries*

| Ref.Date | Country | Meas.Name | Meas.Value | Status |
|----------|---------|-----------|------------|--------|
| 2013 | Canada | Population | 50 | Final |
| 2013 | Canada | GNP | 600 | Final |
| 2013 | USA | Population | 250 | Temporary |
| 2013 | USA | GNP | 2400 | Final |
| … | … | … | … | … |
| 2014 | Canada | Population | 51 | Unavailable |
| 2014 | Canada | GNP | 620 | Temporary |
| … | … | … | … | … |

This table is equivalent to a proper mathematical function: in fact, its rows have the same structure (in term of columns). The Table can be defined as a Data Set, whose name can be "Production of the American Countries". Each row of the table is a Data Point belonging to the Data Set. The Data Structure of this Data Set has five Data Structure Components:

- Reference Date     (Identifier Component)
- Country     (Identifier Component)
- Measure Name     (Identifier Component - Measure Identifier)
- Measure Value     (Measure Component)
- Status     (Attribute Component)

As a second example, let us consider the following physical table, in which the symbol "###" denotes cells that are not allowed to contain a value.

*Institutional Unit Data*

| Row Type | I.U. ID | Ref.Date | I.U. Name | I.U. Sector | Assets | Liabilities |
|----------|---------|----------|-----------|-------------|--------|-------------|
| I | A | ### | AAAAA | Private | ### | ### |
| II | A | 2013 | ### | ### | 1000 | 800 |
| II | A | 2014 | ### | ### | 1050 | 750 |
| I | B | ### | BBBBB | Public | ### | ### |
| II | B | 2013 | ### | ### | 1200 | 900 |

| | | | | | | |
|---|---|---|---|---|---|---|
| II | B | 2014 | ### | ### | 1300 | 950 |
| I | C | ### | CCCCC | Private | ### | ### |
| II | C | 2013 | ### | ### | 750 | 900 |
| II | C | 2014 | ### | ### | 800 | 850 |
| … | … | … | … | … | … | … |

This table, as a whole, is not equivalent to a proper mathematical function because its rows (i.e. the Data Points) have different structures (in term of allowed columns). However, it is easy to recognize that there exist two possible functional structures (corresponding to the Row Types I and II), so that the original table can be split in the following ones:

*Row Type I - Institutional Unit register*

| I.U. ID | I.U. Name | I.U. Sector |
|---|---|---|
| A | AAAAA | Private |
| B | BBBBB | Public |
| C | CCCCC | Private |
| … | … | … |

*Row Type II - Institutional Unit Assets and Liabilities*

| I.U. ID | Ref.Date | Assets | Liabilities |
|---|---|---|---|
| A | 2013 | 1000 | 800 |
| A | 2014 | 1050 | 750 |
| B | 2013 | 1200 | 900 |
| B | 2014 | 1300 | 950 |
| C | 2013 | 750 | 900 |
| C | 2014 | 800 | 850 |
| … | … | … | … |

Each one of these two tables corresponds to a mathematical function and can be represented like in the first example above. Therefore, these would be 2 distinct Data Sets according to the VTL IM, even if stored in the same physical table.

In correspondence to one physical table (the former) there are two logical tables (the latter), so that the definitions will be the following ones:

810    **Data Set 1**:    *Record type I - Institutional Units register*

811    Data Structure 1:
812      • I.U. ID                 (Identifier Component)
813      • I.U. Name               (Measure Component)
814      • I.U. Sector             (Measure Component)
815

816    **Data Set 2**:    *Record type II - Institutional Units Assets and Liabilities*

817    Data Structure 2:
818      • I.U. ID                 (Identifier Component)
819      • Reference Date          (Identifier Component)
820      • Assets                  (Measure Component)
821      • Liabilities             (Measure Component)

822

823    ### The data artefacts

824    The list of the VTL artefacts for the definition of the data is given here, together with the
825    information that the definer have to provide. For the sake of simplicity, we may omit the parts
826    of the names shown between parentheses.

827

828    ***Data Set***

829        *DataSetId*              *mandatory*

830        *DataSetDescr*           *optional*

831        *DataStructureId*        *mandatory    [this is the reference to the data structure of*
832                                 *the Data Set]*

833        *IsCollected*            *mandatory    [YES if the Data Set is collected, NO if it is.*
834                                 *result of a Transformation (i.e. calculated)]*

835

836    ***Data Structure***

837        *DataStructureId*        *mandatory*

838        *DataStructureDescr*     *optional*

839

840    *(Data Structure)* ***Component***

841        *DataStructureId*        *mandatory    [this is part of the identifier of the*
842                                 *Component: the data structure which the Component*
843                                 *belongs to]*

844        *VariableId*             *mandatory    [this is part of the identifier of the*
845                                 *Component: the Represented Variable which defines the*
846                                 *Component (see also hereinafter]*

847        *ComponentRole*          *mandatory    [IDENTIFIER | MEASURE | ATTRIBUTE]*

848     *(Sub)SetId*                          *optional       [possible reference to the (sub)Set containing*
849                                           *the  allowed values for the Component, see the section about*
850                                           *the generic model for Variables and Value Domains]*

851

852     The Data Points have the same structure of the Data Sets they belong to; VTL does not require
853     to define them beforehand.

854     The Validation and Transformation Definition Language introduces the operators for defining
855     the artefacts above (see the VTL reference manual).

856

## Generic Model for Variables and Value Domains

This Section provides a formal model for the Variables, the Value Domains, their Values and the possible (Sub)Sets of Values. These artefacts can be referenced in the definition of the VTL Data Structures and as parameters of some VTL Operators.

**Variable and Value Domain model diagram**

White box:       same as in GSIM 1.1
Light grey:      similar to GSIM 1.1
Dark grey        additional detail (in respect to GSIM 1.1)

**Explanation of the Diagram**

Even in this case, the GSIM artefacts are used as much as possible.  The slight differences are mainly due to the fact that GSIM does not distinguish explicitly between Value Domains and

895 their (Sub)Sets, while in the VTL IM this is made more explicit in order to allow different Data
896 Structure Components relevant to the same aspect of the reality (e.g. the geographic area) to
897 share the same Value Domain and, at the same time, to take values in different Subsets of it.
898 This is essential for VTL for several operations and in particular for validation purposes.  For
899 example, it may happen that the same Variable, say the "place of birth", in a Data Structure
900 takes values in the Set of the European Counties, in another one takes values in the set of the
901 African countries, and so on, even at different levels of details (e.g. the regions, the cities). The
902 definition of the exact Set of Values that a Variable can take may be very important for VTL, in
903 particular for validation purposes.

904 **Data Structure Component**: a component of the data structure (see the explanation already
905 given above, in the data model section). A Data Structure Component is defined by a
906 Represented Variable (see below) and takes values in a subset of its Value Domain (this
907 subset of allowed values may either coincide with the set of all the values belonging to the
908 Value Domain or be a proper subset of it).

909 **Represented Variable**: a characteristic of a statistical population (e.g. the country of birth)
910 represented in a specific way (e.g. through the ISO code). This artefact is the same as in GSIM.
911 A represented variable may define any number of Data Structure Components and takes value
912 in one Value Domain.

913 **Value Domain**: the domain of allowed values for one or more variables. This artefact is very
914 similar to the corresponding artefact in GSIM. Because of the distinction between Value
915 Domain and its Value Domain Subsets, a Value Domain is the wider set of values that can be of
916 interest for representing a certain aspect of the reality (like the time, the geographical area,
917 the economic sector and so on). As for the mathematical meaning, a Value Domain is meant to
918 be the representation of a "space of events" with the meaning of the probability theory[9].
919 Therefore, a single Value of a Value Domain is a representation of a single "event" belonging to
920 this space of events[10].

921 An important characteristic of the Value Domain is the data type (e.g. String, Number,
922 Boolean, Date), which is the type that any Value of the Value Domain must correspond to.

923 **Described Value Domain**: a Value Domain defined by a criterion (e.g. the domain of
924 the positive integers). This artefact is the same as in GSIM.

925 **Enumerated Value Domain**: a Value Domain defined by enumeration of the allowed
926 values (e.g. domain of ISO codes of the countries). This artefact is the same as in GSIM.

927 For completeness, consider that in general a Value Domain can be represented also in a multi-
928 dimensional Cartesian space, therefore a 1-dim Value Domain is a Value Domain defined in a

---

[9] According to the probability theory, a random experiment is a procedure that returns a result belonging a
predefined set of possible results (for example, the determination of the "geographic location" may be
considered as a random experiment that returns a point of the Earth surface as a result). The "space of results" is
the space of all the possible results.

[10] An "event" is a set of results (going back to the example of the geographic location, the event "Europe" is the
set of points of the European territory, more in general an "event" correspond to a "geographical area").  The
"space of events" is the space of all the possible "events" (in the example, the space of the geographical areas).

929 1-dimensional Cartesian space, while a N-dim Value Domain is a Value Domain defined in a N-
930 dimensional Cartesian space and therefore composed by 1-dim Value Domains.

931 The following artefacts are aimed at representing possible subsets of the Value Domains. This
932 is needed for validation purposes, because very often not all the values of the Value Domain
933 are allowed in a Data Structure Component, but only a subset of them (e.g. not all the
934 countries but only the European countries). This is needed also for transformation purposes,
935 for example to filter the Data Points according to a subset of Values of a certain Data Structure
936 Component (e.g. extract only the European Countries from some data relevant to the World
937 Countries) . Although this detail does not exist in GSIM, these artefacts are compliant with the
938 GSIM artefacts described above, representing Value Domains:

939 **Value Domain Subset** (or simply **Set**): a subset of Values of a Value Domain. This
940 artefact does not exist in GSIM, however it is compliant with the GSIM Value Domain. A
941 Value Domain Subset has the same data type as its Value Domain and the same
942 dimensionality. Hereinafter a Value Domain Subset is simply called **Set**, in fact a Value
943 Domain subset can be any set of Values belonging to the Value Domain (even the set of
944 all the values of the Value Domain).

945 **Described Value Domain Subset** (or simply **Described Set**): a described
946 (defined by a criterion) subset of Values of a Value Domain (e.g. the countries
947 having more than 100 million inhabitants, the integers between 1 and 100).
948 This artefact does not exist in GSIM, however it is compliant with the GSIM
949 Described Value Domain.

950 **Enumerated Value Domain Subset** (or simply **Enumerated Set**): an
951 enumerated subset of a Value Domain (e.g. the enumeration of the European
952 countries). This artefact does not exist in GSIM, however it is compliant with the
953 GSIM Enumerated Value Domain.

954 **Value**: an allowed value of a Value Domain. Please note that on a logical /
955 mathematical level, both the Described and the Enumerated Value Domains contain
956 Values, the only difference is that the Values of the Enumerated Value Domains are
957 explicitly represented by enumeration, while the Values of the Described Value
958 Domains are implicitly represented through a criterion.

959 **Code Item**: an allowed item of an enumerated Value Domain. A Code Item is the
960 association of a Value with the relevant meaning (called "category" in GSIM). An
961 example of Code Item is a single countries' ISO code (the Value) associated to the name
962 of the country it represents (the category). As for the mathematical meaning, a Code
963 Item is the representation of an "event" of a space of events (i.e. the relevant Value
964 Domain), according to the notions of "event" and "space of events" of the probability
965 theory (see also the note above).

966 **Code List**: the list of Code Items belonging to an enumerated Value Domain. This
967 artefact is the same as in GSIM except for the multiplicity of the relationship with the
968 Value Domain. Because of the distinction between Value Domain and Value Domain
969 Subsets and because the Value Domain is meant to be the representation of a space of
970 events, a Code List is assumed to contain all the possible Values of interest of the
971 relevant Value Domain (e.g. all the possible GeoAreas of interest), therefore in the VTL
972 IM each enumerated Value Domain has just one Code List.

973    **Set List**: the list of the Code Items belonging to an enumerated Set (e.g. the list of the
974    ISO codes of the European countries). This artefact does not exist in GSIM. However, it
975    has the same role than the Code List in GSIM. The Set List refers only to the Values
976    contained in the list (e.g. the country codes), without the associated categories (e.g. the
977    names of the countries), because the latter are already maintained in the Code List of
978    the relevant Value Domain (which contains all the possible Values with the associated
979    categories).

## Relations and operations between Code Items

981    The VTL allows the representation of logical relations between Code Items, considered as
982    events of the probability theory.

983    As already explained, each Code Item is the representation of an event, according to the
984    notions of "event" and "space of events" of the probability theory. The relations between Code
985    Items aim at expressing the logical implications between the events of a space of events (i.e. in
986    a Value Domain). The occurrence of an event, in fact, may imply the occurrence or the non-
987    occurrence of other events. For example:

988    • The event UnitedKingdom implies the event Europe (e.g. if a person lives in UK he/she
989      also lives in Europe), meaning that the occurrence of the former implies the occurrence
990      of the latter. In other words, the geo-area of UK is included in the geo-area of the
991      Europe.
992    • The events Belgium, Luxembourg, Netherlands are mutually exclusive (e.g. if a person
993      lives in one of these countries he/she does not live in the other ones), meaning that the
994      occurrence of one of them implies the non-occurrence of the other ones (Belgium AND
995      Luxembourg = impossible event; Belgium AND Netherlands = impossible event;
996      Luxembourg and Netherlands = impossible event). In other words, these three geo-
997      areas do not overlap.
998    • The occurrence of one of the events Belgium, Netherlands or Luxembourg (i.e. Belgium
999      OR Netherlands OR Luxembourg) implies the occurrence of the event Benelux (e.g. if a
1000     person lives in one of these countries he/she also lives in Benelux) and vice-versa (e.g.
1001     if a person lives in Benelux, he/she lives at least in one of these countries). In other
1002     words, the union of these three geo-areas coincides with the geo-area of the Benelux.

1003   The logical relationships between Code Items are very useful for validation and
1004   transformation purposes. Considering for example some positive and additive data, like for
1005   example the population, from the relationships above it can be deduced that:

1006   • The population of United Kingdom should be lower than the population of Europe.
1007   • There is no overlapping between the populations of Belgium, Netherlands and
1008     Luxembourg, so that these populations can be added in order to obtain aggregates.
1009   • The sum of the populations of Belgium, Netherlands and Luxembourg gives the
1010     population of Benelux.

1011   A **Code Item Relation** is composed by two members, a 1st (left) and a 2nd (right) member. The
1012   envisaged types of relations are: "is equal to" (=), "implies" (<), "implies or is equal to" (<=),
1013   "is implied by" (>), and "is implied by or is equal to" (>=). "Is equal to" means also "implies
1014   and is implied".  For example:

1015   UnitedKingdom < Europe          means (UnitedKingdom implies Europe)

| 1016 | In other words, this means that if a point of space belongs to United Kingdom it also |
| 1017 | belongs to Europe. |

1018 The left members of a Relation is a single Code Item. The right member can be either a single
1019 Code Item, like in the example above, or a logical composition of Code Items giving another
1020 Code Item as result: these are the **Code Item Relation Operands**. The logical composition can
1021 be defined by means of Operators, whose goal is to compose some Code Items (events) in
1022 order to obtain another Code Item (event) as a result. In this simple algebra, two operators
1023 are envisaged:

1024 • the logical OR of mutually exclusive Code Items, denoted "+", for example:

1025 Benelux = Belgium + Luxembourg + Netherlands

1026 This means that if a point of space belongs to Belgium OR Luxembourg OR Netherlands
1027 then it also belongs to Benelux and that if a point of space belongs to Benelux then it
1028 also belongs either to Belgium OR to Luxembourg OR to Netherlands (disjunction). In
1029 other words, the statement above says that territories of Belgium, Netherland and
1030 Luxembourg are non-overlapping and their union is the territory of Benelux.
1031 Consequently, as for the additive measures (and being equal the other possible
1032 Identifiers), the sum of the measure values referred to Belgium, Luxembourg and
1033 Netherlands is equal to the measure value of Benelux.

1034 • the logical complement of an implying Code Item in respect to another Code Item
1035 implied by it, denoted "-", for example:

1036 EUwithoutUK = EuropeanUnion - UnitedKingdom

1037 In simple words, this means that if a point of space belongs to the European Union and
1038 does not belong to the United Kingdom, then it belongs to EUwithoutUK and that if a
1039 point of space belongs to EUwithoutUK then it belongs to the European Union and not
1040 to the United Kingdom. In other words, the statement above says that territory of the
1041 United Kingdom is contained in the territory of the European Union and its
1042 complement is the territory of EUwithoutUK. As a consequence, considering a positive
1043 and additive measure (and being equal the other possible Identifiers), the difference of
1044 the measure values referred to EuropeanUnion and UnitedKingdom is equal to the
1045 measure value of EUwithoutUK.

1046 Please note that the symbols "+" and "-" do not denote the usual operations of sum and
1047 subtraction, but logical operations between Code Items seen as events of the probability
1048 theory. In other words, two or more Code Items cannot be summed or subtracted to obtain
1049 another Code Item, because they are events (and not numbers), and therefore they can be
1050 manipulated only through logical operations like "OR" and "Complement".

1051 Note also that the "+" also acts as a declaration that all the Code Items denoted by "+" are
1052 mutually exclusive (i.e. the corresponding events cannot happen at the same time), as well as
1053 the "-" acts as a declaration that all the Code Items denoted by "-" are mutually exclusive.
1054 Furthermore, the "-" acts also as a declaration that the relevant Code item implies the result of
1055 the composition of all the Code Items denoted by the "+".

1056 At intuitive level, the symbol "+" means "*with*" (Benelux = Belgium *with* Luxembourg *with*
1057 Netherland) while the symbol "-" means "*without*" (EUwithoutUK = EuropeanUnion *without*
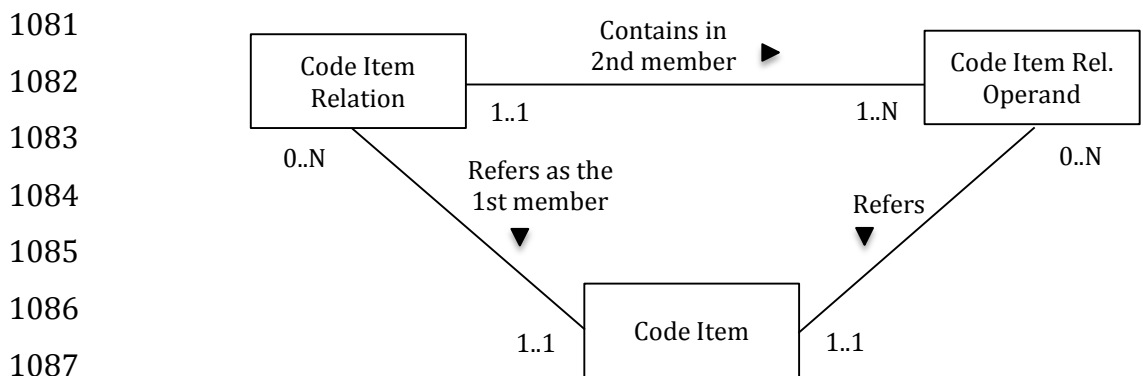1058 UnitedKingdom).

When these relations are applied to additive numeric measures (e.g. the population relevant to geographical areas), they allow the measure values to be obtained from the compound Code Items (i.e. the population of Benelux and EUwithoutUK) by summing or subtracting the measure values relevant to the component Code Items (i.e. the population of Belgium, Luxembourg and Netherland in the former case, EuropeanUnion and UnitedKingdom in the latter). This is why these logical operations are denoted in VTL through the same symbols as the usual sum and subtraction. Please note also that this is valid whichever the Data Set and the additive measure are (provided that possible other dimensions have the same values).

These relations occur between Code Items (events) belonging to the same Value Domain (space of events). They are typically aimed at defining aggregation hierarchies, either structured in levels (classifications), or without levels (chains of free aggregations) or a combination of these options.

For example, the following relations are aimed at defining the continents and the whole world in terms of individual countries:

- World = Africa + America + Asia + Europe + Oceania
- Africa = Algeria + ... + Zimbabwe
- America = Argentina + ... + Venezuela
- Asia = Afghanistan + ... + Yemen
- Europe = Albania + ... + Vatican City
- Oceania = Australia + ... + Vanuatu

A simple model diagram for the Code Item Relations and Code Item Relation Operands is the following:



### The historical changes

The changes in the real world may induce changes in the artefacts and in the relationships between them, so that some definitions may be considered valid only with reference to certain time values. For example, the birth of a new country as well as the split or the merge of existing countries in the real world would induce changes in the Code Items belonging to the Geo Area Value Domain, in the composition of the relevant Sets, in the relationships between the Code Items and so on.

A correct representation of the historical changes of the artefacts is essential for VTL, because the VTL operations are meant to be consistent with these historical changes, in order to ensure a proper behaviour in relation to each time. With regard to this, VTL must face a

complex environment, because it is intended to work also on top of other standards, whose assumptions for representing historical changes may be heterogeneous. Moreover, institutions and even departments of the same Institutions often use different conventions for representing historical changes. The VTL IM tries to manage this heterogeneity by allowing multiple options when possible and clarifying the relationships between these options.

Please note that there are two main temporal aspects: the so-called validity time and operational time. The validity time is the time during which a definition is true in the real world. The operational time is the time period during which a definition is available and may produce operational effects. In this context only the former is considered, while the latter is left to the concrete implementations of processing systems.

Even the **identification of the artefacts** is related to temporal assumptions. Regard to this aspect, two main options can be considered:

a) The artefacts are assumed to be variable in time and therefore represent a given abstraction of the reality even if it changes. As a consequence, a single artefact may represent the whole history of an abstraction. For example, under this option the same artefact (e.g. EU) may represent the European Union even if its geographic area changes (i.e. even if the participant countries change, like happened many times so far). This option follows the intuitive conceptualization in which abstractions are identified independently of time and may change with time maintaining the same identity.

b) The artefacts are assumed to be invariable in time and therefore represent a given abstraction of the reality only for the period in which this abstraction does not change. As a consequence, more artefacts have to be used to represent the whole history of an abstraction, one for each period in which the abstraction does not change. For example, under this option the European Union can be represented by more artefacts, one for each period during which its geographic area was stable (e.g. EU1, … , EU9). This option is based on the conceptualization in which the artefacts are identified in connection with the time, so that an artefact corresponds to the abstraction of some aspects of the reality (e.g. Geo Area) in association with certain times. VTL conventionally assimilates to this case also the common practice of giving a version to the artefacts for representing time changes (e.g. EUv1, … , EUv9 where v=version), being each version of the artefact assumed as invariable.

The general assumptions of VTL in relation to the representation of the historical changes are the following:

- VTL artefacts are identified and referenced by means of their univocal identifier, therefore, for VTL, in the option a) there would exist one artefact for Europe (e.g. EU) while in the option b) there would exist 9 different artefacts for Europe (e.g. EU1, … , EU9).
- possible versions of the artefacts aimed at managing temporal changes are considered to be part of the univocal artefact identifier, so that different versions are considered as different artefacts like in the option b); the Europe in this case would be represented by many artefacts (e.g. EUv1, … , EUv9). More in general, the univocal identifiers of the artefacts may be composite in the implementations, so that the adopting standards and organizations may use their own identification conventions, provided that the version is considered part of the VTL identifier.

- The characteristics of the invariable artefacts obviously cannot change with time, so they are assumed to be constant and their time validity is not explicitly considered by VTL (if required, a time validity for these artefact can be managed by the implementations).
- The variable artefacts can have characteristics variable with time. There can be many occurrences of these characteristics for the same artefact, but only one of them is valid in a time instant; the same applies to variable relations between artefacts (for example, the United Kingdom may belong to Europe only for a certain time). In these cases, each occurrence is qualified by means of a validity period (start date - end date). As obvious, the validity periods of these different occurrences cannot overlap. Validity periods are considered as "optional", because they would not be needed if the option b) is assumed. If not specified, the validity period is assumed to be "ever".
- VTL does not consider explicitly possible variations with time of the textual descriptions of the artefacts (if required, this can be managed in the implementations).

## The Variables and Value Domains artefacts

The list of the VTL artefacts related to Variables and Value Domains is given here, together with the information that the definer have to provide.

### *(Represented) Variable*

| | | |
|---|---|---|
| *VariableId* | *mandatory* | |
| *VariableDescr* | *optional* | |
| *ValueDomainId* | *mandatory* | *[reference to the Value Domain which measures the Variable, i.e. in which the Variable takes values]* |

### *Value Domain*

| | | |
|---|---|---|
| *ValueDomainId* | *mandatory* | |
| *ValueDomainDescr* | *optional* | |
| *IsEnumerated* | *mandatory* | *[YES if the Domain is Enumerated, NO if it is Described]* |
| *DataType* | *mandatory* | *[this is the data type of the Values of the Value Domain, i.e. one of the allowed VTL data types (see hereinafter)]* |
| *ValueRestriction* | *optional* | *[this is a regular expression which expresses a criterion for restricting the allowed Values if needed, for example by specifying a max length, an upper or/and a lower value, and so on]* |

### *Code List (composition)*      *[mandatory for Enumerated Value Domains]*

| | | | |
|---|---|---|---|
| 1184<br>1185 | *ValueDomainId* | *mandatory* | *[this is part of the identifier of the Value: the Value Domain which the Value belongs to]* |
| 1186<br>1187<br>1188 | *ValueId* | *mandatory* | *[also named Code Item, this is part of the identifier of the Value: i.e. the univocal name of the Value within the Value Domain it belongs to]* |
| 1189<br>1190 | *ValueDescr* | *optional* | *[in GSIM terms, this is the category associated to the Code Item]* |
| 1191<br>1192 | *StartDate* | *optional* | *[needed if a Value belongs to a Value Domain only for a certain period]* |
| 1193<br>1194 | *EndDate* | *optional* | *[needed if a Value belongs to a Value Domain only for a certain period]* |

1195

1196 ### *N-dimensional Value Domain*

1197
1198
1199
1200

*A N-dim Value Domain is a combined space of 1-dim Value Domains. It is not required to define explicitly the N-dim Value Domains, because all the possible combinations of 1-dim Value Domains are considered as defined by default. The Values of a N-dim value domains are combination of Values of the component 1-dim Value Domains.*

1201

1202 *(Value Domain Sub)***Set**

| | | | |
|---|---|---|---|
| 1203<br>1204 | *ValueDomainId* | *mandatory* | *[this is part of the Identifier of the Set: the Value Domain which the set belongs to]* |
| 1205<br>1206<br>1207 | *Set_Id* | *mandatory* | *[this is part of the identifier of the Set: i.e. the univocal name of the Set within the Value Domain it belongs to]* |
| 1208 | *SetDescr* | *optional* | |
| 1209<br>1210 | *IsEnumerated* | *mandatory* | *[YES if the the Set is Enumerated, NO if it is Described]* |
| 1211<br>1212<br>1213 | *SetCriterion* | *mandatory for Described Sets* | *[a regular expression which expresses a criterion for identifying the Values belonging to the Set]* |
| 1214<br>1215 | *StartDate* | *optional* | *[needed if a Set belongs to a Value Domain only for a certain period]* |
| 1216<br>1217 | *EndDate* | *optional* | *[needed if a Set belongs to a Value Domain only for a certain period]* |

1218

1219 **Set List (composition)** *[mandatory for Enumerated Sets]*

| | | | |
|---|---|---|---|
| 1220<br>1221<br>1222 | *ValueDomainId* | *mandatory* | *[this is part of the identifier of the Set List: reference to the Value Domain which the Set and the Value belongs to]* |

| 1223 | SetId | mandatory | [this is part of the identifier of the Set List: reference to the Set which contains the Value] |
| 1224 | | | |
| 1225 | ValueId | mandatory | [this is part of the identifier of the Set List: reference to the Value which belongs to the Set] |
| 1226 | | | |
| 1227 | StartDate | optional | [needed if a Value belongs to a Set only for a certain period] |
| 1228 | | | |
| 1229 | EndDate | optional | [needed if a Value belongs to a Set only for a certain period] |
| 1230 | | | |

1231

1232

### Code Item Relation

| 1234 | 1stMemberDomainId | mandatory | [this is part of the identifier of a Relation: reference to the Value Domain of the first member of the Relation; e.g. Geo_Area] |
| 1235 | | | |
| 1236 | | | |
| 1237 | 1stMemberValueId | mandatory | [this is part of the identifier of a Relation: reference to the Value of the first member of the Relation; e.g. Benelux] |
| 1238 | | | |
| 1239 | | | |
| 1240 | 1stMemberCompositionId | mandatory | [this is part of the identifier of a Relation: conventional name of the composition related with the first member, needed to distinguish possible different compositions related to the same first member Value. It must be univocal within the 1stMemberValueId. Not necessarily it has to be meaningful, it can be simply a progressive number ; e.g. "1"] |
| 1241 | | | |
| 1242 | | | |
| 1243 | | | |
| 1244 | | | |
| 1245 | | | |
| 1246 | | | |
| 1247 | CompositionDescr | optional | [e.g. "Benelux from its countries"] |
| 1248 | Relation Type | mandatory | [relation between the first and the second member, having as possible values =, <, <=, >, >=] |
| 1249 | | | |
| 1250 | StartDate | optional | [needed if a Relation is valid only for a certain period] |
| 1251 | | | |
| 1252 | EndDate | optional | [needed if a Relation is valid only for a certain period] |
| 1253 | | | |

1254

### Code Item Relation Operand

| 1256 | 1stMemberDomainId | mandatory | [this is part of the identifier of a Relation Operand: see its description above; e.g. Geo Area] |
| 1257 | | | |
| 1258 | 1stMemberValueId | mandatory | [this is part of the identifier of a Relation Operand: see its description above; e.g. Benelux] |
| 1259 | | | |
| 1260 | 1stMemberCompositionId | mandatory | [this is part of the identifier of a Relation Operand: see its description above; e.g. "1"] |
| 1261 | | | |

| | | |
|---|---|---|
| 1262 | *2ndMemberValueId* | *mandatory*   *[this is part of the identifier of a Relation Operand: it references the ValueId of an operand; e.g. Belgium]* |
| 1265 | *Operator* | *optional*   *[it specifies the applied operator, its possible values are "+" and "- "; the default is "+"; e.g. "+"]* |
| 1267 | *StartDate* | *optional*   *[needed if an Operand of a Relation is valid only for a certain period]* |
| 1269 | *EndDate* | *optional*   *[needed if an Operand of a Relation is valid only for a certain period]* |

1271

1272

# Generic Model for Transformations

1274 The purpose of this section is to provide a formal model for describing the validation and
1275 transformation of the data.

1276 A Transformation is assumed to be an algorithm to produce a new model artefact (typically a
1277 Data Set) starting from existing ones. It is also assumed that the data validation is a particular
1278 case of transformation, therefore the term "transformation" is meant to be more general and
1279 to include the validation case as well.

1280 This model is essentially derived from the SDMX IM[11], as DDI and GSIM do not have an explicit
1281 transformation model at the moment[12]. In its turn, the SDMX model for Transformations is
1282 similar in scope and content to the Expression metamodel that is part of the Common
1283 Warehouse Metamodel (CWM) [13] developed by the Object Management Group (OMG).

1284 The model represents the user logical view of the definition of algorithms by means of
1285 expressions.  In comparison to the SDMX and CWM models, some more technical details are
1286 omitted for the sake of simplicity, including the way expressions can be decomposed in a tree
1287 of nodes in order to be executed (if needed, this detail can be found in the SDMX and CWM
1288 specifications).

1289 The basic brick of this model is the notion of Transformation.

1290 A Transformation specifies the algorithm to obtain a certain artefact of the VTL information
1291 model, which is the result of the Transformation, starting from other existing artefacts, which
1292 are its operands.

---

[11] The SDMX specification can be found at https://sdmx.org/?page_id=5008  (see Section 2 - Information Model, package 13 - "Transformations and Expressions").

[12] The Transformation model described here is not a model of the processes, like the ones that both SDMX and GSIM have, and has a different scope. The mapping between the VTL Transformation and the Process models is not covered by the present document, and will be addressed in a separate work task with contributions from several standards experts.

[13] This specification can be found at http://www.omg.org/cwm.

1293 Normally the artefact produced through a Transformation is a Data Set (as usual considered
1294 at a logical level as a mathematical function). Therefore, a Transformation is mainly an
1295 algorithm for obtaining a derived Data Set starting from already existing ones.

1296 The general form of a Transformation is the following:

1297      `variable parameter := expression`

1298 ":=" is the assignment operator, meaning that the result of the evaluation of *expression* in the
1299 right-hand side is assigned to the *variable parameter* in the left-hand side, which is the a-
1300 priori unknown output of *expression* (typically a Data Set).

1301 In turn, the *expression* in the right-hand side composes some operands (e.g. some input Data
1302 Sets) by means of some operators (e.g. sum, product ...) to produce the desired results (e.g.
1303 the validation outcome, the calculated data).

1304 For example:        $D_r := D_1 + D_2$       ($D_r$, $D_1$, $D_2$ are assumed to be Data Sets)

1305 In this example the measure values of the Data Set $D_r$ is calculated as the sum of the measure
1306 values of the Data Sets $D_1$ and $D_2$.

1307 A validation is intended to be a kind of Transformation. For example, the simple validation
1308 that $D_1 = D_2$ can be made through an "If" operator, with an expression of the type:

1309     $D_r$     :=     If ($D_1 = D_2$, then TRUE, else FALSE)

1310 In this case, the Data Set $D_r$ would have a Boolean measure containing the value TRUE if the
1311 validation is successful and FALSE if it is unsuccessful.

1312 These are only fictitious examples for explanation purposes. The general rules for the
1313 composition of Data Sets (e.g. rules for matching their Data Points, for composing their
1314 measures ...) are described in the sections below, while the actual Operators of the VTL are
1315 described in the VTL reference manual.

1316 The *expression* in the right-hand side of a Transformation must be written according to a
1317 formal language, which specifies the list of allowed operators (e.g. sum, product ...), their
1318 syntax and semantics, and the rules for composing the expression (e.g. the default order of
1319 execution of the operators, the use of parenthesis to enforce a certain order ...). The Operators
1320 of the language have Parameters[14], which are the a-priori unknown inputs and output of the
1321 operation, characterized by a given role (e.g. dividend, divisor or quotient in a division).

1322 Note that this generic model does not specify the formal language to be used. As a matter of
1323 fact, not only the VTL but also other languages might be compliant with this specification,
1324 provided that they manipulate and produce artefacts of the information model described
1325 above. This is a generic and formal model for defining Transformations of data through
1326 mathematical expressions, which in this case is applied to the VTL, agreed as the standard
1327 language to define and exchange validation and transformation rules among different
1328 organizations

1329 Also note that this generic model does not actually specify the operators to be used in the
1330 language. Therefore, the VTL may evolve and may be enriched and extended without impact
1331 on this generic model.

---

[14] The term is used with the same meaning of "argument", as usual in computer science.

1332 In the practical use of the language, Transformations can be composed one with another to
1333 obtain the desired outcomes. In particular, the result of a Transformation can be an operand
1334 of other Transformations, in order to define a sequence of calculations as complex as needed.
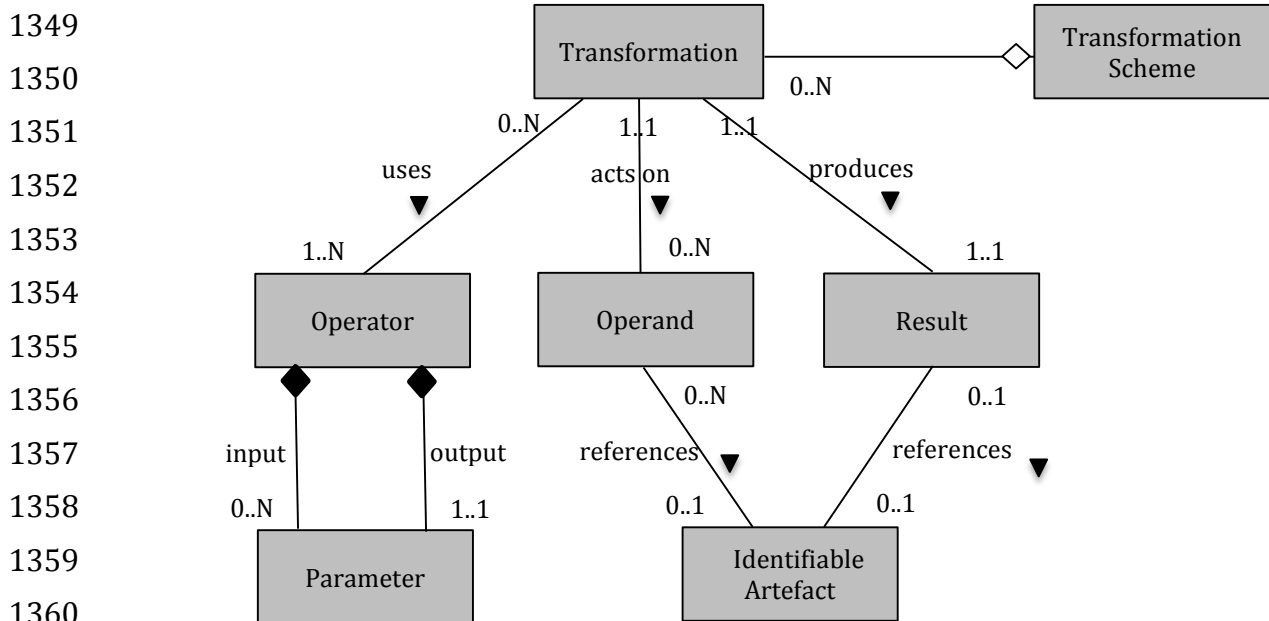
1335 Moreover, the Transformations can be grouped into Transformations Schemes, which are sets
1336 of transformations meaningful to the users. For example, a Transformation Scheme can be the
1337 set of transformations needed to obtain some specific meaningful results, like the validations
1338 of one or more Data Sets.

1339 A set of Transformations takes the structure of a graph, whose nodes are the model artefacts
1340 (usually Data Sets) and whose arcs are the links between the operands and the results of the
1341 single Transformations. This graph is directed because the links are directed from the
1342 operands to the results and is acyclic because it should not contain cycles (like in the
1343 spreadsheets), otherwise the result of the Transformations might become unpredictable.

1344 The ability of generating this graph is a main goal of the VTL, because the graph documents
1345 the operations performed on the data, just like a spreadsheet documents the operations
1346 among its cells.

1347 **Transformations model diagram**

1348



1362   White box:       same as in GSIM 1.1
1363   Dark grey box:   additional detail (in respect to GSIM 1.1)

1364 (These artefacts match the SDMX artefact having the same name; however, the identifiable artefacts are intended
1365 to be the ones of the VTL model)

1366

### Explanation of the diagram

**Transformation**: the basic element of the calculations, which consists of a statement which assigns the outcome of the evaluation of an Expression to an Artefact of the Information model;

**Expression**: a finite combination of symbols that is well-formed according to the syntactical rules of the language. The goal of an Expression is to compose some Operands in a certain order by means of the Operators of the language, in order to obtain the desired result. Therefore, the symbols of the Expression designate Operators, Operands and the order of application of the Operators (e.g. the parenthesis); an expression is defined as a string and is a property of a Transformation, as in the SDMX IM;

**Transformation Scheme**: a set of Transformations aimed at obtaining some meaningful results for the user (like the validation of one or more Data Sets); the Transformation Scheme may also be considered as a VTL program;

**Operator**: the specification of a type of operation to be performed on some Operands (e.g. +, -, *, /);

**Parameter**: a-priori unknown input or output of an Operator, having a definite role in the operation (e.g. dividend, divisor or quotient for the division) and corresponding to a certain type of artefact (e.g. a "Data Set", a "Data Structure Component" …);

**Operand**: a specific Artefact referenced in the expression as an input (e.g. a specific input Data Set); the distinction between Operand and Result is not explicit in the SDMX IM;

**Result**: a specific Artefact to which the result of the expression is assigned (e.g. the calculated Data Set); the distinction between Operand and Result is not explicit in the SDMX IM;

**Identifiable Artefact**: a persistent Identifiable Artefact of the VTL information model (e.g. a persistent Data Set); a persistent artefact can be result of no more than one Transformation;

Note that with regards to the SDMX Transformation and Expression Model, some artefacts are intentionally not shown here, essentially to avoid more technical details (i.e. the decomposition of the operations in the Expression, described in SDMX by means of the ExpressionNode and its sub-types ReferenceNode, ConstantNode, OperatorNode). For this reason, in the diagram above, the Transformation references directly Operators and Artefacts (through its Expression), instead in the SDMX IM the Transformation contains ExpressionNodes which in turn reference Operators and Artefacts. On the technical implementation perspective, however, the model would be the same as the SDMX one (except some details that are specific to the SDMX context).

### Example

Imagine that $D_1$, $D_2$ and $D_3$ are Data Sets containing information on some goods, specifically: $D_1$ the stocks of the previous date, $D_2$ the flows in the last period, $D_3$ the current stocks. Assume that it is desired to check the consistency of the Data Sets using the following statement:

$D_r$ := If $((D_1 + D_2) = D_3$ , then "true", else "false")

In this case:

The Transformation may be called "Consistency check between stocks and flows" and is formally defined through the statement above.

| 1409 | • $D_r$ | is the Result |
|---|---|---|
| 1410 | • $D_1, D_2$ and $D_3$ | are the Operands |
| 1411 | • If $((D_1 + D_2) = D_3$ , then TRUE, else FALSE) | is the Expression |
| 1412 | • ":=", "If",  "+" , "=" | are Operators |

1413      Each operator has some predefined parameters, for example in this case:

| 1414 | • input parameters of "+": | two numeric Data Sets (to be summed) |
|---|---|---|
| 1415 | • output parameters of "+": | a numeric Data Sets (resulting from the sum) |
| 1416 | • input parameters of "=": | two Data Sets (to be compared) |
| 1417 | • output parameter of "=": | a Boolean Data Set (resulting from the comparison) |
| 1418 | • input parameters of "If": | an Expression defining a condition, i.e. $(D_1+D_2)=D_3$ |
| 1419 | • output parameter of "If": | a Data Set (as resulting from the "then", "else" clauses) |

1420

## 1421 Persistency and Identification of the artefacts of the model

1422 The artefacts of the model can be either persistent or non-persistent. An artefact is persistent
1423 if it is permanently stored, and vice-versa.

1424 A persistent artefact exists externally and independently of a VTL program, while a non-
1425 persistent artefact exists only within a VTL program.

1426 The VTL grammar provides for the identification of the non-persistent artefacts (see the
1427 section about the conventions for the grammar of the language) and leaves the accurate
1428 definition of the identification mechanism of the persistent artefacts to the standards
1429 adopting the VTL (e.g. SDMX, DDI …)[15].

1430 However, the VTL aims at promoting international sharing of rules, which should have a clear
1431 identification. Therefore, VTL just gives some minimum requirements about the structure of
1432 this universal identifier, assuming that the standards adopting the VTL will ensure that the
1433 identifier of a persistent artefact is unique.

1434 In practice, the VTL considers that many definers need to operate independently and
1435 simultaneously (e.g. many organizations, units,…), so that they should be made independent
1436 as much as possible in assigning names to the artefacts, making sure that nevertheless the
1437 resulting names are unique.

1438 Therefore, VTL foresees:

1439 • the **Name** of the  artefact (a generic string), which is unique in the environment of the
1440      definer;
1441 • an optional **Namespace** (generic string beginning with an alphabetic character) which
1442      is a supplementary qualifier that identifies the environment in which the artefact
1443      Name is assumed to be unique, to avoid name conflicts.

---

[15] Different standards may have different identification mechanisms.

1444 The Name of the artefact may be composite. For example, in case of versioned artefacts, the
1445 Name is assumed to contain the version as well. It is the responsibility of the definer to ensure
1446 that the artefact Names are unique in the environment.

1447 The Namespace may be composite as well. For example, a composite structure may be useful
1448 to make reference to environments and sub-environments. Notice that VTL does not provide
1449 for a general mechanism to ensure that a Namespace is universally unique, which is left to the
1450 standards implementing the VTL.

1451 When the context is clear, as typically happens in validation, the Namespace can be omitted.
1452 In other words, the Name of the artefact is always mandatory, while the Namespace is
1453 required only for the operands that belong to a different Namespace than the Transformation.

1454 As intuitive, the Namespace may begin with the name of the institution ("maintenance
1455 agency" in SDMX terms). Assuming the dot (".") as separator character between environments
1456 and sub-environments, examples of possible Namespaces are:

1457 - ESCB.analyis&insight
1458 - EuropeanStatisticalSystem.validation
1459 - OECD.Stat
1460 - Unesco
1461 - Bancaditalia.dissemination.public

1462

1463 The artefact identifier as a whole is also a string, composed of the concatenation of the
1464 Namespace – if needed – and the artefact Name, where the slash ("/") symbol is a typical and
1465 recommended choice (e.g. "NAMESPACE/NAME" for explicit Namespace definition or simply
1466 "NAME" for referencing the default Namespace).

1467

1469 VTL 1.1 is a powerful language that allows the user to express complex validation and
1470 transformation operations on one or more datasets in a clear, concise and readable manner,
1471 without the need to program low-level data handling details. The whole language has been
1472 designed to simplify the problem of writing validation and transformation tasks, and to free
1473 the programmer from writing the usual boilerplate code, therefore making the program
1474 maintenance easier and reducing the chance of introducing bugs.

1475 In the Reference Manual chapter on core operators, including the join expressions, we shall
1476 present in detail how VTL allows user to write dataset expressions using the familiar
1477 arithmetic, logic, string, date and other elementary (or scalar) operators, while the language
1478 itself takes care of all low-level details, such as joining and traversing the datasets involved in
1479 such an expression. In order to lay the foundation for such treatment, in this chapter we cover
1480 the preliminaries -- the key language concepts upon which VTL is built. Considering the power
1481 and expressiveness of VTL, there are surprisingly few of them, and the sections that follow
1482 aim at providing a thorough and not too technical overview of each of them.

## Objects and Types

1483

1484 In VTL, an object is any entity that can be processed or computed. This includes elementary
1485 objects as small and simple as numerical, Boolean, string or date scalar values, or as large and
1486 complex as the datasets of the Information Model (IM). Whatever their size and complexity,
1487 objects share some common features:

1488 - All objects in VTL are immutable. This means that VTL programs never change the
1489   content of an input object (e.g., a collection or a dataset), but can, when necessary,
1490   generate a new updated version, which is also immutable. VTL internally uses some
1491   clever tricks to make sure that working with immutable objects does not incur
1492   excessive penalties in terms of computing time and resources.

1493 - Each object has a type. At runtime, each object carries with itself so-called runtime
1494   type information, which describes its structure and can be (and is) inspected by
1495   various VTL operations in order to decide how that object should be processed. But
1496   VTL is also a statically typed system, meaning that before the program is executed, the
1497   compiler uses the information about types of literals, variable parameters, and other
1498   program elements to automatically infer, or at least approximate as much as possible,
1499   the type of more and more complex program constructs. In this way, the compiler can
1500   optimize code and prevent an important and large class of potential type errors that
1501   might otherwise occur at runtime.

1502 Type `any` is the most general type, and includes all possible objects, without telling us
1503 anything about them. On the other extreme, type `()` is the empty type, containing no objects.
1504 Nested between these two extremes are all other types in VTL, organized in the following
1505 main type families:

1506 - *Scalar types* refer to basic numeric, string, Boolean, and date values that can be stored
1507   in a single numeric, string, Boolean, or date-time values in a tabular representation of a
1508   dataset. Type `scalar` is the most generic, denoting any scalar value, and type `null`

contains only the value `null`, denoting a missing, non-applicable, or undefined scalar value. Nested between `scalar` and `null` are all other scalar types, as described in the text that follows. The scalar types include:

- ○ `integer` -- any integer, taking implicitly into account the range of supported values, as described below in the *Basic VTL Assumptions*.

- ○ `integer` [$a$:] -- any integer greater than or equal to some integer constant $a$.

- ○ `integer` [:$b$] -- any integer less than or equal to some integer constant $b$.

- ○ `integer` [$a$:$b$] -- any integer that falls between two integer constants $a$ and $b$, both inclusive  (where $a<b$).

- ○ `integer` {$x_1$, ..., $x_n$} -- one of integers enumerated in  {$x_1$, ..., $x_n$}

- ○ `float` -- any floating-point number compatible with double-precision IEEE 754.

- ○ `number` -- the generalization of `integer` and `float`

- ○ `boolean` -- a  Boolean value, either true or false.

- ○ `date` -- a date-time timestamp

- ○ `string` -- any string of characters from the UNICODE character set

- ○ `string` [$a$] -- any strings consisting of exactly $a$ characters

- ○ `string` [$a$:$b$] -- any string consisting of between $a$ and $b$ characters

- ○ `string` {$s_1$, ..., $s_n$} -- one of strings enumerated in  {$s_1$, ..., $s_n$}; in effect this type describes elements of a code list.

- *Collection types* are lists and sets of elements of the same type:

  - ○ `list<`*t*`>` is a list of elements of type *t*. For instance, `list<integer>` is a list of integers

  - ○ `set<`*t*`>` is a set of elements of type *t*. For instance, `set<string>` is a set of character strings

  - ○ `collection<t>` the generalization of `list<t>` and `set<t>`

- *Dataset types.* Dataset types describe VTL datasets by summarizing the information about their structure (i.e., components) as needed by different functions and procedures operating on datasets, and as seen or inferred at compile-time:

  - ○ `dataset` -- any dataset

  - ○ `dataset` { *role$_1$ name$_1$* as *type$_1$*, *role$_2$ name$_2$* as *type$_2$*, ..., *role$_N$ name$_N$* as *type$_N$*} -- any dataset that has exactly the listed components. Each *role* is either `identifier`, `measure` or `attribute`, each *name* must be distinct, and each *type* is a scalar type.

  - ○ `dataset` { *role$_1$ name$_1$* as *type$_1$*, ..., *role$_N$ name$_N$* as *type$_N$* ...} (with "..." before the closing "}")-- any dataset that has at least the listed components, and possibly some more. Each *role* is either `identifier`, `measure` or `attribute`, each *name* must be distinct, and each *type* is a scalar type.

- *Record types*. These types are analogous to the dataset types, except that they use keyword `record` instead of `dataset`, and refer not to a complete dataset, but to an individual row in it.

- *Product types*. Type $t_1$ * $t_2$ * ... * $t_n$ (where n>1) describes all n-tuples whose components belong to the corresponding types t1, ..., tn. E.g., `integer * string * boolean` is a type of all triples whose first component is an integer, the second component is a string, and the third component is a Boolean. For instance, (105,`"Luxembourg"`, `false`) is a triple that belongs to this type.

- *Function types*. Type of the form *t* -> *T* describes a function that takes an object of type *t* and produces a result of type *T*. For instance, `integer -> string` is the type a function that takes an integer and returns a string. Or, (`integer * string`) `-> boolean` is the type of a function that takes a pair consisting of an integer and string, and returns a Boolean.

One of the objectives of the VTL type system is to encode useful information about the objects that belong to a type. That includes meta-information from the data model. Using enumerated string types, one can effectively encode a code list:

```
type BENELUX = string { "BE", "NL", "LU" }
```

```
type EU12 = string { "BE", "DE", "DK", "ES", "FR", "GR", "IE", "IT", "LU", "NL", "PT", "UK" }
```

This is an example of two user-defined named types.

Another way the compiler can use the type information are integer computations. If the variable parameter *x* is declared as `integer[0:10]`, then the compiler can infer that the expression `y:= 2*x+3` has type `integer [3:23]`, and therefore *y* cannot be negative or zero in looping, branching, or filtering constructs.


## Identifiers and Values

As in many other programming languages, VRL uses identifiers to refer to objects of different kinds. Syntactically, regular identifiers start with a (lowercase or uppercase) English alphabet letter, followed by zero or more letters, decimal digits, or underscores. However, such a regular identifier cannot be the same as a keyword or a reserved word.

Regular identifiers (just like keywords) are not case sensitive. Internally, VTL system may either convert them to uppercase or lowercase. In that sense, Pos, pos, and POS are treated as the same identifier.

Also, a regular identifier cannot start with an underscore, which denotes an argument placeholder in a function, as described below.

However, VTL 1.1 allows us to escape the limitations imposed on regular identifiers by enclosing them in single quotes (apostrophes). For instance, `'1'` is a valid VTL identifier, as well as `'_'`, `'a-2'`. `'a:b:c'`, `'a/b/c'`, or `'?x%'`. Also, `'string'` is a valid quoted identifier, while `string` is not (because it is a keyword). Quoted identifiers also may contain apostrophes, but they have to be doubled. For instance `'a''b'` is an identifier consisting of letter `a`, an apostrophe, and letter `b`. And, unlike the regular identifiers, the quoted identifiers

are case-sensitive: `'Pos'` is different from `'pos'`, and both are different from `'POS'`. Whether unquoted identifier `pos` is the same as `'Pos'`, `'pos'` or `'POS'` is implementation dependent, and users are advised not to depend on any capitalization scheme in order to ensure portability of their VTL code.

VTL 1.1 makes no difference between the regular and the quoted identifiers. That is to say that wherever an identifier is expected, we can freely use one form or another.

One common use of identifiers in VTL is to store results of computations. For instance:

```
D := 0.2*D1 + 0.8*D2
```

is an assignment statement, where the expression `0.2*D1 + 0.8*D2` is computed, and (supposing that *D1* and *D2* are dataset variable parameters) the resulting dataset is stored in the variable parameter *D*. After the assignment, we can use D to refer to the computed value.

We use the word "variable parameter" for historical reasons, because that is the term commonly used in mathematics and programming. Hereinafter, we shorten this term, for sake of simplicity, to simply "variable". Please note that the same term ("variable") is used in the "VTL Information Model" section with a different meaning, i.e. as an abbreviation of Represented Variable, which is a GSIM artefact also used by the VTL IM, denoting a Statistical Variable that has a representation and can be used as a Component of a Data Structure. Hereinafter, instead, the term "variable" is used as an abbreviation of "variable parameter", so meaning an argument, a priori unknown, of an Operator of the language. Speaking about VTL expressions, therefore, variables are synonym of parameters. However, variables in VTL are less like storage locations in computer memory that can be updated at will, but more like logical variables in mathematics: they are immutable. This is to say that once the assignment is executed, we cannot change the value to which *D* refers. We are allowed to write:

```
D := 0.2*D1 + 0.8*D2
```

```
D := 1.2 * D
```

```
/* other code using D */
```

but this is internally translated into:

```
D := 0.2*D1 + 0.8*D2
```

```
U := 1.2 * D   /* U is a "fresh" variable name not appearing in the
original code */
```

```
/* other code using U instead of D */
```

In other words, the second assignment of the form "D := ..." hides the "original" value of D from the rest of the code.

To understand how variables work, we need to understand the concept of scope. A scope is a mapping from a set of identifiers visible at some point in VTL to values or objects to which they refer.

Each assignment statement changes or updates the scope for the statements that follow by associating the assigned variable name to the result of the expression to the right side of ":=". Therefore, when two statements in sequence assign to the same variable name, the first computed value of the variable is visible in the second assignment, but gets overwritten by the second assignment. This creates the illusion of variable update.

1627 It is sometimes useful to limit the scope of variables.  For instance, in formula:

```
1628        D := (D1+D2-1-D3)/(D1+D2+1+D3)
```

1629 it may be useful to isolate `D1+D2` and `1+D3` in an auxiliary variable A and B, which makes the
1630 code more readable:

```
1631        A := D1 + D2
```

```
1632        B := 1 + D3
```

```
1633        D := (A-B)/(A+B)
```

1634 However, we may want to limit the scope of *A* and *B* only to the computation of *D*. This can be
1635 done using a nested assignment block enclosed in curly braces:

```
1636        D := {
1637          A := D1 + D2
1638          B := 1 + D3
1639            (A-B)/(A+B)
1640        }
```

1641 This points to a general rule in VTL: wherever an expression is expected (as, for instance, to
1642 the right of ":="), we can insert a block in curly brackets that introduces local assignments,
1643 whose visibility is limited to the block. The final element of the block must be an expression,
1644 whose result is the result of the entire block.

1645 The whole VTL program can also be seen as one global block, implicitly closed in curly braces.
1646 It may contain zero or more assignments, and may end in a dataset expression which is,
1647 effectively, the result of the program. For compatibility with VTL 1.0 and unlike in normal
1648 blocks, we allow the last statement in the program to be an assignment, in which case the
1649 result of the whole program is the value of the last computed variable.

## 1650 Expressions

1651 Each VTL program is, essentially, an expression that takes some inputs and returns a result,
1652 which on the program level needs to be a dataset.  The same holds for user-defined functions
1653 that we shall mention later: each function is defined as an expression.

1654 Expressions are built from the following ingredients:

1655 • Literals, such as `1` or `-105` (integer) `2.0` or `10.5e-4` (float), `"Abcdef"` (string), `true`
1656   or `false` (Boolean). As a special case, function abstractions (described in the following
1657   subsections) such as `_+_` and `\x,y{x+y}` -- both are functions that take two
1658   arguments and add them together -- can also be considered a special form of "function
1659   literals."

1660 • Variable or column names, given as identifiers (regular or quoted).

1661 • References to dataset components, of the form *D.X*, where *D* is a dataset variable name,
1662   and *X* is an identifier naming the component.

- Qualified names of module or object members, of the form *M::X*, where *M* is the name of the module or object, and *X* is the identifier naming a member of *M* (a value, function, or other object).

- Function calls of the form *name*(*arg*1, ..., *argN*) (where *N*>0), where *name* is the name of a built-in or user-defined function, and $arg_1$, ..., $arg_N$ are the function call arguments.

- Built-in unary (prefix and postifx), binary (infix) and ternary (infix) operators, given in the Reference Manual. These can be used to build (sub-)expressions using the prefix, infix, or postfix operator notation.

- Join expressions, discussed in the chapter on Core Operators.

- Dataset clauses, discussed in the Reference Manual.

As usual, parentheses override binary and unary operator priorities.

Expressions in VTL are interpreted in two possible ways, depending on the context in which they appear:

- *General expressions* are those found in the top-level program assignment statements, and bodies of user-defined functions. In these expressions, identifier *X* is always interpreted as a variable name (used as a parameter in an expression), referring to a program input, function argument, or an assigned variable. General expressions can be of any type.  For instance, in `A:=D1+D2`, *D1* and *D2* are variable names.

- *Component expressions* appear in record-level statements inside the join expression body and in dataset clauses. In them, identifier *X* (not followed by "." or "::") is interpreted as a dataset component name. To use variable *X*, we have to write $X*. Component expressions are always scalar.  For instance, in `D[filter X>0]`, X is not a variable name, but a component name (of dataset D). However, in `D[filter X>$Limit]`, element `$Limit` stands for variable *Limit* (which may be, for instance, a function argument).

## Data Flow Optimization

As we could see in the preceding examples, expressions can be complex and may contain nested blocks that compute temporary variables. For complex block expressions, it is important to understand that in VTL their actual computation may differ from what is usually found in imperative programming languages. In the latter, each assignment is computed sequentially, followed with the computation of the final result.

It is important to understand that from the programmer's perspective, VTL block expressions produce results *as if* they are executed sequentially.  For instance, in the block expression:

```
{
  A := D1 + D2
  B := 1 + D3
  D := (A−B)/(A+B)
  D    /* result */
}
```

1702 we can logically think about the result *D* as being computed gradually: first *A* is computed,
1703 then *B*, and finally *D*. The semantics of VTL complex expressions guarantees that the final
1704 result is going to be the same *as if* such step-by-step computation has taken place. This makes
1705 it easy for the programmer to think about the programming problem and organize and write
1706 code in as clear and correct manner as possible.

1707 However, the VTL compiler may perform data flow analysis to infer the data flow graph in the
1708 program in order to optimize the handling of datasets. For instance, computing *A*, *B* and *D*
1709 sequentially in the previous example would be inefficient, since *A* would require one dataset
1710 join and traversal (*D1* and *D2*), *B* another (*D3*), and *D* the third (*A* and *B*). Instead, the
1711 compiler can transform this into a more efficient single join and traversal of datasets *D1*, *D2*,
1712 and *D3*, where all calculations are done in a single run. The way this optimization is done must
1713 guarantee that the result of the block is the same *as if* the computation is performed
1714 sequentially. But the actual execution strategy used by a VTL implementation can range from
1715 a centralized sequential computation, to translating programs into database or data
1716 warehouse queries, to executing different operations on different interconnected nodes in a
1717 distributed computing system, by routing or streaming data between them. Whatever
1718 execution strategy is actually used, it must be transparent to the programmer.


## 1719 User-Defined Functions

1720 VTL 1.1 adopts many features from the functional programming languages. In particular, each
1721 scalar or dataset operation and operator can be seen as a function that accepts some
1722 arguments and returns a result. This means that most of the processing can be viewed as
1723 application of functions to data. Sometimes, this is explicit in using functional notation, as in
1724 `size(D)`, but even when using infix or prefix operator notation as in `2*X-3`, this is equivalent
1725 to (and can indeed be written as) a function call of the form `'-'('*'(2,X),3)`. That makes
1726 functions one of the fundamental concepts in VTL, along with the join expressions.

1727 There are essentially two ways to define functions in VTL 1.1. Suppose, for instance, a sorting
1728 algorithm that operates on collections of objects of some type *t*, which requires to be supplied
1729 with a function of type (*t*,*t*) -> `boolean`, which takes a pair of objects of type *t* and returns
1730 `true` exactly when the first element is considered to precede the second element (making it
1731 "smaller" in some ordering scheme). The sorting algorithm is neutral with the respect to the
1732 type *t* of collection elements, and it depends on this function to perform comparison.

1733 Now, let us suppose we want to use that algorithm to sort a collection of integers in a
1734 descending order. For that we have to supply a function of type `(integer,integer)-`
1735 `>boolean` which returns `true` for arguments (x,y) exactly when x>y. The classical way to do
1736 that is to write a named function definition of the form:

```
define function compare_integer_descending(x as integer, y as
integer)
returns boolean
as x > y
```

1741 We can normally omit the "`returns boolean`" part, as the return type information can be
1742 inferred by the compiler from the expression "`x>y`".

1743 This is an example of a named function definition. As a result of it, identifier
1744 compare_integer_descending refers, in the scope in which it is defined, to a function object of

1745 type `(integer,integer)->boolean`. We can then pass this function to the sorting
1746 algorithm by name, using identifier compare_integer_descending.

1747 However, for this kind of relatively simple cases, VTL 1.1 allows us to specify a function object
1748 directly, without the need to define/create it separately. This we call the anonymous function,
1749 and in our case it can look like this:

1750    `\x as integer, y as integer { x > y }`

1751 The anonymous function starts with a backslash, followed by arguments (and optionally their
1752 types), followed by a block expression that produces the result. This is also a simple example
1753 of a function in whose body arguments appear only once, and that in the order in which are
1754 listed. When the type of the arguments is unambiguous from the context (i.e., when the
1755 compiler can decide that both arguments must be integer, because it already knows we are
1756 sorting a collection of integers, we can be even terser and write:

1757    `_>_`

1758 Here, we use underscores as placeholders for arguments. When the compiler encounters an
1759 underscore, it converts the expression in which it appears into an anonymous function, and
1760 turns each underscore into a function argument:

1761    `\x,y{x>y}`

1762 An anonymous function that computes an average of three numbers can be written as:

1763    `(_+_+_)/3`

1764 We can even write:

1765    `_ between _ and _`

1766 (noting the spaces surrounding underscores, to prevent underscores to be treated as a part of
1767 identifiers) to denote a function of type (number,number,number)->boolean which takes
1768 three numbers and checks whether the first one falls between the second and the third one.

1769 Anonymous functions specified using underscores have a limitation that each argument can
1770 be used only once. And, by definition, the anonymous functions cannot be recursive, because
1771 they have no way of calling themselves. Therefore, to achieve more general computation
1772 tasks, we need to use the most general way for defining functions, which is using the named
1773 functions.

1774 As seen in the example above, the general template for defining a function is:

1775    **define function** $name(arg_1, ..., arg_N)$

1776    **returns** $t$

1777    **as** $E$

1778 where *name* is the function identifier, and each *arg* is an identifier, optionally followed with
1779 keyword **as** and the argument type. The **returns** part is also optional, and it specifies the type
1780 of the function's result.

1781 The return type, as well as argument types, are optional, because in many case (although not
1782 always) the compiler can infer their type from the context. For instance, function that checks
1783 if a quadratic equation $ax^2+bx+c=0$ has a solution can be defined as:

```
1784        define function has_solution(a, b, c)
1785        as b*b-4*a*c>0
```

1786  The compiler knows that the comparison (>) produces a Boolean result. Also, since the right-
1787  hand side of > is a numeric literal 0, the left-hand side also has to be a number. And, since the
1788  left-hand side produces from variables *a*, *b*, and *c* and arithmetic operators, the compiler is
1789  able to convert the above definition into:

```
1790        define function has_solution(a as number, b as number, c as number)
1791        returns boolean
1792        as b*b-4*a*c>0
```

1793  However, it is advisable to provide argument and return types for the more complex or
1794  externally visible user-defined functions, in order to help their users, and to make the
1795  compiler check that their implementation really produces the result of the desired type.

1796  So far, all function arguments were obligatory. For functions that perform complex
1797  operations, this may lead to a large number of function arguments, most of which have some
1798  sensible default value.Let us take, for instance, a function that computes *n*-degree distance
1799  between measurements in two datasets. For two real numbers *x* and *y*, the distance of *n*-th
1800  degree is *n*-th root of $x^n$-$y^n$.So, the first degree distance is simply *x-y*, the second degree
1801  distance is $sqrt(x^2$-$y^2)$, etc. We can write the function as:

```
1802        define function distance(d1 as dataset,
1803            d2 as dataset,
1804            n as integer := 2)
1805        as
1806          (d1^n - d2^n)^(1/n)
```

1807  Note how we added ":= 2" in the declaration of argument n. This makes it an optional named
1808  argument, which, if not specified in a function call, takes on the default value 2. A call
1809  distance(x,y) is equivalent to distance(x,y,n:2).  The optional named arguments must come
1810  after the non-optional arguments, and in a call their values are preceded with the argument
1811  name followed by a colon, "n: 2". If we have a function with more than one optional named
1812  argument, such as:

```
1813        define function z_transform(x as number,
1814            mu as number := 0,
1815            sigma as number := 1)
1816        as (x-mu)/sigma
```

1817  then we can write both

```
1818        z_transform(x, mu: 50, sigma: 4.3)
```

1819  and

```
1820        z_transform(x, sigma: 4.3, mu: 50)
```

1821  That means that the relative ordering of the optional named arguments in a function call is
1822  not important, since the compiler always looks at the definition to pass the arguments in a
1823  correct sequence. However, as mentioned above, all positional (i.e., not named) arguments
1824  must come first.

# Procedures

Besides functions, VTL supports procedures. Procedures differ from functions in several important respects.

- Procedures are aimed at automating common processing tasks, and can be used as a means for shortening the code by replacing common processing tasks with a procedure call. On the other hand, functions are concerned only with computing results from arguments.

- Procedures may have several input and output arguments, which are passed by reference, while the procedure call has no return value of its own. In contrast, functions defined via a single expression (which may be a complex, block expression), and exhibit so-called referential integrity. That is to say that a function call with same arguments (always passed by value) should always return the same result.

To understand procedures, let us take a simple example of a procedure that computes a quotient and a remainder of a division of measures in a dataset and a number (the same can be easily extended to two datasets):

```
define procedure quot_rem(in ds as dataset, in divisor as number,
     out quot as dataset, out rem as dataset)
as {
  quot := floor(ds / divisor)
  rem := ds - quot*divisor
}
```

We first note that each argument of a procedure is qualified as **in** or **out**. Input arguments, such as ds and divisor in our example, are passed by value, just like function arguments, and we can pass any expression with compatible type when calling the function. However, output parameters, such as quot and rem in our example, must be specified as names of variables that will hold results computed in the procedure body.

For instance, we can call the above procedure like:

```
call quot_rem(PopPerCountry, AvgPop, Multiple, Remainder)
```

and this call is equivalent to the following two assignments:

```
Multiple := floor(PopPerCountry / AvgPop)
Remainder := PopPerCountry - Multiple*AvgPop
```

Note that in our case the body of the procedure is a sequence of assignments enclosed in curly braces. In general, it is always a sequence of assignments or procedure calls. Also, any assignment in the procedure body to a variable that is not marked as output is invisible to the calling code.

Procedures may look a lot like macros, but they are much more powerful. Firstly, the body of a procedure is compiled and type checked, which means that any syntax or semantic errors in a procedure definition are detected at compile time. This extends to the type checking of input and output arguments. Finally, procedures can be stored in modules and reused.

## Language Core

The ability to define user functions and procedures allows development of libraries of reusable and standardized VTL validation and transformation building blocks, which, in turn, adds to the effectiveness and expressiveness of use of VTL 1.1 in normal use case scenarios. But to be useful, these functional and procedural facilities need to rest on a solid foundation directly provided by the language. This includes the two main components:

- Core constructs, which represent the fundamental building blocks into which any dataset processing in VTL 1.1 can be decomposed, and
- Standard library, which contains a large number of utility functions and operators built from the core constructs or other standard library constructs.

Both the core constructs and the standard library are explained in detail in the Reference Manual.

The role of the core constructs is to express the semantics of simple and complex operations in VTL in an unambiguous manner. For instance, using the scalar operators '+' and '*' that add and multiply numbers, and a join expression, we can define the function:

```
define function midway(d1 as dataset {measure x as number, …},
        d2 as dataset {measure x as number, …})
returns dataset {number x as number, …}
as
        [d1 outer join d2] {
                filter d1.x is not null or d2.x is not null
                x := 0.5*d1.x + 0.5*d2.x
        }
```

which takes two dataset arguments *d1* and *d2*, each containing (at least) a numeric measure component named *x*, and returns a dataset with a numeric measure component named *avg* which is the mean of *x* from *d1* and *d2*. Without going here into too much detail, well explained in the Reference Manual, the function body after `as` is a join expression that:

- Performs a *join* of *d1* and *d2*, by matching records from *d1* and *d2* that share the same values of identifier components. The set of identifier components of *d1* must be equal to, a subset of, or a superset of, the set of identifier components of *d2*.
- The type of join is *outer*, which means that if for some record in d1 there is no matching record in *d2* (or vice versa), the join "invents" the latter with all measure and attribute component values set to `null`.
- The body of the join expression is given inside the curly braces, '{' and '}'. Inside the body, *d1* and *d2* refer to the matched records from the corresponding joined datasets.
- The `filter` statement skips the cases where the numeric measure *x* is undefined In both *d1* and *d2.* This is important, because datasets *d1* and *d2* may have more than one measure component,
- For each pair of matched records *d1* and *d2*, the result contains one record that inherits all identifier component values from *d1* and *d2*, and has a numeric measure component *x* which is computed as `0.5*d1.x + 0.5*d2.x`.

For instance, let us suppose we have these two data sets:

*d1* :=

---

| Year | Geo | X |
|------|-----|---:|
| 2011 | LU | 104 |
| 2011 | NL | 812 |
| 2012 | LU | 97 |

1907  and *d2* :=

| Geo | X |
|-----|---:|
| LU | 128 |
| NL | 768 |

1908  Then `midway(d1,d2)` will produce:

| Year | Geo | X |
|------|-----|---:|
| 2011 | LU | 116 |
| 2011 | NL | 790 |
| 2012 | LU | 112.5 |

1909  Incidentally, the same operation can be directly and simply written in VTL as a dataset
1910  expression:

1911       `0.5*d1.x + 0.5*d2.x`

1912  where *d1* and *d2* are two dataset variables. This simple dataset expression is internally
1913  automatically translated by the compiler translated into the same expression as given in the
1914  body of the function given above. Note that in the dataset expression '*' is a mixed
1915  scalar/dataset operator (multiplying a scalar value 0.5 with a dataset), and '+' is a dataset
1916  operator (both operands are datasets). However, the meaning of these two scalar/dataset and
1917  dataset operators and of the entire expression does not need to be separately defined: it is
1918  systematically derived from the core operators and constructs, scalar '+' and '*' and the join,
1919  as described in the corresponding chapter below.

1920  It is important to note that the selection of core operators and constructs is entirely driven by
1921  the language design and the need for semantic soundness. Users need not be concerned
1922  whether they are using a "core" or a "library"operator, function, or another construct. Users
1923  should always try to use the construct which is best suited for their intended purpose.

1924  For the language implementers, the existence of the language core represents a contract that
1925  controls the correct behaviour of their VTL implementation. It does not always necessarily
1926  mean that every implementer needs to use the core constructs as the back-end. While every
1927  VTL construct needs to be expressible in terms of the language core, implementations may
1928  use more efficient backend-specific algorithms and techniques (in R, SAS, SQL, etc.). However,

1929  the implementers must ensure that the user-observable behaviour of their implementations
1930  respect the behaviour required by the contract.


## Compilation Units and Dialect Selection

1931

1932  Programs and modules are two types of compilation units in VTL. By a compilation unit we
1933  here mean a unit of code stored in a single file or transmitted as a message. The main
1934  difference between a VTL program and a VTL module is that the former executes some
1935  particular dataset processing task (some form of validation or transformation), while the
1936  latter creates and packages functions, procedures, values, named types, and other objects so
1937  that they can be used by programs and other modules.

1938  Since VTL comes in several versions, which may use different syntax or may interpret the
1939  same syntactic forms differently. To indicate the version or dialect of VTL used in a
1940  compilation unit, its first line (after leading whitespace and comments) should be the
1941  following directive:

1942  ```
        use syntax "X.Y"
      ```

1943  (optionally followed by a semicolon) where X.Y is the version number of VTL dialect in which
1944  the compilation unit is written. For instance:

1945  ```
        use syntax "1.1"
      ```

1946  indicates that what follows in the file uses the VTL 1.1 syntax.

1947  The version number in **use syntax** directive can be followed by one or more of case
1948  insensitive tags of the form "+tag" where *tag* consists of one or more Latin letters, decimal
1949  digits and underscores. For instance:

1950  ```
        use syntax "1.1+estat+strict"
      ```

1951  may indicate VTL 1.1 syntax with custom Eurostat (ESTAT) tags, and strict type checking
1952  option.

1953  If a VTL system does not support the version indicated in **use syntax** directive, it is obliged to
1954  reject the compilation unit and report an error. However, each VTL implementation can freely
1955  decide which tags to recognize, and should ignore all unecognized tags (possibly issuing a
1956  compile-time warning).


## Program and Module Structure

1957

1958  A module is distinguished from a program by starting with a **module** directive after the
1959  leading whitespace, comments, and the optional **use syntax** directive. If the first thing after
1960  the leading workspace, comments and the optional **use syntax** directive is not a **module**
1961  directive, then the compilation unit is treated as a program, not module.

1962  **Module Declaration**

1963  The simplest form of the **module** directive is:

1964  ```
        module name
      ```

1965 (optionally terminated with a semicolon). *Name* is an identifier giving the module name. This
1966 defines a *transient* module, which is created in memory when the module is loaded by the
1967 compiler because it is used by a program or another module.

1968 Another more complex form of the **module** directive is:

1969     `module` *name* `in` "*AGENCY*:*ENTITY*:*VERSION*"

1970 (optionally terminated with a semicolon). *AGENCY* is a code for the owner of *ENTITY*, which is
1971 a logical name of a persistent entity in the underlying information model used by the VTL
1972 system. For instance, in VTL systems based on SDMX, *ENTITY* refers to a named versionable
1973 artefact, such as a data structure definition or a dataflow. Finally, *VERSION* gives the version of
1974 the *ENTITY* to which the module is associated.

1975 The latter form of the **module** directive creates a persistent module, which the VTL system
1976 associates with *ENTITY*.

## Module Usage

1978 Both VTL programs and modules can depend on other modules. These dependencies are
1979 expressed with **use module** directives. The first form:

1980     `use module` *name*

1981 (optionally followed by a semicolon) expresses a dependency on a transient module with the
1982 given *name* which is locally available, i.e., it is supplied together with the program or module
1983 using it, for instance as a file in the same directory tree or a part or attachment of the same
1984 message.

1985 The second form:

1986     `use module` *name* `in` "*AGENCY:ENTITY:VERSION*"

1987 (optionally followed by a semicolon) expresses a dependency on a persistent module with the
1988 given *name* which is attached to the persistent *ENTITY* owned by *AGENCY*. *VERSION* is either a
1989 version number, or an asterisk (*) that signifies the latest version. In this case, depending on
1990 the underlying concrete information model (such as, for instance, SDMX), the compiler needs
1991 to retrieve the module from a

1992 Module dependencies cannot be circular. One advantage of expressing the module
1993 dependencies with **use module** directives is that the compiler (or any other source code
1994 handling tool, such as a registry) can analyse module dependencies, construct dependency
1995 graphs, and detects any problems (such as missing modules or circular dependencies)
1996 statically, i.e., before the VTL program is deployed and run.

## Definitions

1998 A VTL program or a module can contain zero or more definitions. These include:

1999   ● Type definitions

2000   ● Function and procedure definitions

2001   ● Validation rule / rule set definitions

2002 All definitions introduce a named object (a type, a function, a procedure, a validation rule /
2003 rule set) in the scope of the program or module.

2004 To refer to identifier *x* in module *module*, we use the double column syntax:

2005     *name* :: *x*

2006 which is called a qualified name, in contrast with a simple identifier or simple name.

2007 ### Module-Level Computations

2008 After definitions, modules can contain computations, which take the form of assignments:

2009     $x := E$

2010 where *x* is a variable, and *E* is an expression. Like a definition, each assignment also associates
2011 an object which is the result of *E* with identifier *x* in the module scope, but this time using the
2012 general expression syntax. This is useful, for instance, when the module describes a data
2013 structure, and needs to have a member which is a set of tuples describing constraints on the
2014 dataset component values.

2015 Or, a mathematical module can contain assignment:

2016
```
PI := 4*atan(1.0)
```

2017 Another example where computations come handy is re-exporting a named object from a
2018 used module. In the following example:

2019
```
use syntax "1.1"
2020 module A
2021 use module B
2022 /* definitions */
2023 X := B::X
```

2024 module A uses module B, and can refer from A to member named X in B as B::X. But, by
2025 assigning it to name X in its own scope, module A re-exports B::X as A::X which is accessible
2026 from any module using A (and not necessarily using B).

2027 ### Program-Level Computations

2028 While computations are optional in modules, they are mandatory in programs. In fact,
2029 performing a computation and returning a result is the whole purpose of a program. The
2030 computation statements consist of zero or more assignments or procedure call statements,
2031 followed by an expression which is the result of the whole program. This final expression can
2032 be omitted if the last statement in the program is an assignment; in this case, the result of the
2033 program is the result of the last assignment.

2034 ## Module Instantiation and Incremental Compilation

2035 In the preceding section, we already said that circular dependencies between modules are
2036 forbidden in VTL. In fact, we go one step further by requiring that a module needs to be
2037 *instantiated* before being used in a program or another module.

2038 A module is instantiated when:

2039     ● All modules on which it depends (if any) are (transitively) instantiated

2040     ● All type, function, procedure, rule, etc., definitions in the module have created the
2041        corresponding objects and bound them to the names in the module scope.

2042 ● All module computations have been performed, and all values have been bound to the
2043   corresponding variable names in the module scope.

2044 The instantiated module can be seen simply as a map from module member names to the VTL
2045 objects created from definitions or computed from assignments. Of course, on the technical
2046 level the situation is somewhat more complex, since an instantiated module also needs to
2047 carry additional information about types and module dependencies.

2048 One advantage of this approach is that an instantiated module is not only limited to an in-
2049 memory representation, but can also be written to a persistent store in some appropriate
2050 external format -- for instance by serializing to a file, or populating database tables. Unless the
2051 module source code or some of its dependencies change, the VTL compiler needs to compile
2052 and instantiate the module only once. This may significantly improve the speed of compilation
2053 and execution of VTL programs.

2054 Besides, by requiring that all modules used by a VTL compilation unit need to be previously
2055 instantiated, it becomes natural for the compiler to perform incremental compilation, starting
2056 from the bottom of the module dependency tree and going upwards towards the top-level
2057 target (a program or a module). A recompilation and re-instantiation of a module would be
2058 triggered only when its instantiated form is outdated or missing, or when one or more of its
2059 dependencies change.

## Principle of Introspection

2061 It has already been hinted above that one of the important uses of modules in VTL is to
2062 describe data structures of different datasets that are used in a program. Note that the dataset
2063 structure can be described in several different ways:

2064 ● Using compile-time type information -- we have already seen that the structure of a
2065   dataset can be fully or partially described using **dataset** type. The level of detail and
2066   precision of a dataset type reflects the information put into the code by the
2067   programmer and the characteristics of operations applied to the datasets.

2068 ● Using runtime type information -- each dataset at runtime carries with itself a full and
2069   precise description of its structure, as fed on input or computed in the VTL program.
2070   This information is typically more precise than the type information inferred at
2071   compile time.

2072 ● By explicitly constructing a description of dataset structure at runtime -- this means
2073   constructing VTL objects that represent dataset components, their types, roles, or
2074   constraints.

2075 Each of these approaches has certain advantages and disadvantages. The compile-time type
2076 analysis prevents using objects that are not datasets in dataset operations, or using datasets
2077 that lack the necessary components with the required data types and roles. For instance, if $f$ is
2078 a function and $ds$ a dataset variable, the type system ensures that in the call:

2079 ```
     f(ds)
     ```

2080 $ds$ always meets the minimum of requirements imposed on its structure by $f$. However, the
2081 compile-time type analysis is limited by what is known before a program is run and before it

2082 receives any inputs. Therefore, its characterization of datasets can be sometimes too general
2083 and coarse.

2084 We can also define type that describes a particular dataset structure.  For instance:

```
2085        type population = dataset {
2086          identifier geo as string
2087          identifier year as integer
2088          measure population as float
2089          attribute status as string
2090        }
```

2091 If we define *f* to accept an argument of type *population*, the compiler raises a red flag if we try
2092 to use a dataset that may not be compliant. But what if we want to check if *ds* can be fed to *f*
2093 not in general, but in a particular case of program execution?

2094 At runtime, each input to the program and each result of computation carries with it the
2095 precise description of its structure.  If *ds* is a dataset variable, we can use **is** operator to ask:

```
2096        ds is population
```

2097 Note that this construct allows us to use the runtime type information of *ds* against a statically
2098 defined type *population*. If this test succeeds (returns **true**), we know that passing this
2099 particular *ds* to *f* is safe even if at compile-time we had no information to justify the safety of
2100 passing *ds* to *f* in general.

2101 The "trick" on which this is based is that population on the right-hand side of **is** is *reified,*
2102 which is to say that it is represented as an object at runtime. Thus, **is** takes the run-time type
2103 information of *ds* and the reified type information of *population*, and compares them.

2104 But let us go one step further, and imagine we have an arbitrary dataset *ds* and want to
2105 inspect its structure from within a VTL program.

2106 One drawback on relying on runtime type information is that the objects describing it can be
2107 very complex and unstable in the sense that they can change from one version of the language
2108 to another. This means that if a VTL program wants to look into the structure of a dataset at
2109 runtime, it would need to rely on a very complex internal API, which would likely change as
2110 new features are added to the language.

2111 This seems to suggest that it is better to keep the structure of the runtime type information
2112 representation hidden from the programmer. As an alternative, we can construct a simplified
2113 description of the structure, which faithfully reflects the data type.

2114 Imagine that from the *population* data set type we generate the following module

2115        **use syntax** "1.1"

2116        **module** pop_ds

2117        **type** t = **dataset** {

2118          **identifier** geo **as string**

2119          **identifier** year **as integer**

2120          **measure** population **as float**

```
2121         attribute status as string
2122       }
2123     structure := list(
2124       module {
2125        name := "geo"
2126        role := "identifier"
2127        type t = string
2128       },
2129       module {
2130        name := "year"
2131        role := "identifier"
2132        type t = integer
2133       },
2134       module {
2135        type t = float
2136        name := "population"
2137        role := "measure"
2138       },
2139       module {
2140        name := "status"
2141        role := "attribute"
2142        type t = string
2143       }
2144       }
```

2145 In this module, we have encoded the desired dataset structure in two ways: by defining a type
2146 *t* and by providing the list of objects describing individual components.  Each **module** { ... }
2147 inside list is a component descriptor object.

2148 If we have a module or a program that uses *pop_ds*:

2149
```
use module pop_ds
```

2150 then we can refer to the database type as:

2151
```
pop_ds :: t
```

2152 and if the following test returns true:

2153
```
ds is pop_ds :: t
```

2154   we can inspect the structure of the dataset at runtime by looking at:

2155   ```
       pop_ds :: structure
       ```

2156   In order to allow introspection of dataset structure for arbitrary datasets, we can use built-in
2157   function *get_dataset_structure* which takes an arbitrary dataset and returns a list of
2158   component descriptors whose structure is illustrated our example. In that sense, the dynamic
2159   introspection is still possible, but the API is kept at minimum.

2160   Looking at the *pop_ds* module above, it becomes obvious that this kind of modules can be
2161   automatically generated from the information model. Indeed, in VTL 1.1 each dataset
2162   structure that is identifiable with *AGENCY:NAME:VERSION* coordinates behaves *as if* it has
2163   attached a VTL module describing the dataset structure in the described manner. Of course,
2164   these modules are not written by hand, but are automatically generated from the information
2165   model itself.

2166   For instance, one can write:

2167   ```
       use module pop_ds in "acme:population:*"
       ```

2168   to import the dataset structure description module *pop_ds* for the latest version of *population*
2169   table owned by *acme*, and then use pop_ds::t and pop_ds::structure in the described manner.

2170   This is the principle of automated introspection of dataset structures from the information
2171   model in VTL code.

## 2173 Scalar Core Operators

2174 VTL 1.1 scalar operators are unary and binary operators that accept a scalar argument and
2175 return a scalar value. In this section, we present only the operators that are "natively" scalar,
2176 but can be automatically lifted to the dataset/scalar and dataset levels. There are a number of
2177 other operators that can take scalar values, but are not amenable to the automatic lifting.
2178 They are all presented systematically in the Reference Manual, and below we give only a brief
2179 overview.

2180 Binary scalar operators are always infix, and can be left-associative, right-associative and non-
2181 associative. If operator `@` is left-associative, then `X@Y@Z` is the same as `(X@Y)@Z`, and if it is
2182 right-associative, then `X@Y@Z` is the same as `X@(Y@Z)`. If `@` is non-associative, the form
2183 `X@Y@Z` is syntactically invalid. Unary scalar operators can be prefix and postfix.

### 2184 Scalar arithmetic operators

2185 The next table presents the arithmetic operators, which take number operands and produce a
2186 number result:

| Operator | Usage | Associativity | Description |
|---|---|---|---|
| *Additive operators* | | | |
| + | $E + E'$ | Left | Addition |
| − | $E - E'$ | Left | Subtraction |
| *Multiplicative operators* | | | |
| * | $E * E'$ | Left | Multiplication |
| / | $E / E'$ | Left | Division |
| div | $E \operatorname{div} E'$ | None | Integer division |
| mod | $E \operatorname{mod} E'$ | None | Remainder |
| *Power operators* | | | |
| ^ | $E \char`\^ E'$ | Right | Exponentiation |
| *Unary operators* | | | |
| − | $-E$ | Prefix | Sign inversion |
| + | $+E$ | Prefix | Sign preservation |

2187

2188 The unary operators have the highest priority, then the power operators, then the
2189 multiplicative operators, and finally the additive operators.

2190 The operands to the scalar arithmetic operators can be any `number`. If at least one operand is
2191 `null`, the result is also `null`.

### Scalar string operators

2193 There is a string concatenation operator:

| Operator | Usage | Associativity | Description |
|----------|-------|---------------|-------------|
| `\|\|` | $E$ `\|\|` $E'$ | Left | String concatenation |

2194

2195 VTL does not distinguish between `null` and the empty string `""`.

### Scalar  Boolean operators

2197 Scalar Boolean operators correspond to the logical connectives `or`, `xor`, `and`, and `not`. They
2198 take Boolean operands and return a Boolean value.  Unary `not` has the highest priority, then
2199 the multiplicative operator `and`, and finally the two additive operators `or` and `xor`.

| Operator | Usage | Associativity | Description |
|----------|-------|---------------|-------------|
| *Additive operators* | | | |
| `or` | $E$ `or` $E'$ | Right | Logical disjunction |
| `xor` | $E$ `xor` $E'$ | Right | Logical exclusive disjunction |
| *Multiplicative operators* | | | |
| `and` | $E$ `and` $E'$ | Left | Logical conjunction |
| *Unary operators* | | | |
| `not` | `not` $E$ | Prefix | Logical negation |

2200

2201

2202 The treatment of nulls is the following:

| X | `not` X | Y | X `and` Y | X `or` Y | X `xor` Y |
|---|---------|---|-----------|----------|-----------|
| `true` | `false` | `true` | `true` | `true` | `false` |
| | | `false` | `false` | `true` | `true` |
| | | `null` | `null` | `true` | `null` |
| `false` | `true` | `true` | `false` | `true` | `true` |

| | | | | | |
|---|---|---|---|---|---|
| | | false | false | false | false |
| | | null | false | null | null |
| null | null | true | null | true | null |
| | | false | false | null | null |
| | | null | null | null | null |

2203

## Scalar relational and test operators

2204

| Operator | Usage | Associativity | Description |
|---|---|---|---|
| | | *Binary operators* | |
| = | $E = E'$ | None | Value equality. $E$ and $E'$ are the same |
| <> | $E <> E'$ | None | $E$ and $E'$ are not the same |
| < | $E < E'$ | None | $E$ is smaller than $E'$ |
| <= | $E <= E'$ | None | $E$ is smaller than or equal to $E'$ |
| > | $E > E'$ | None | $E$ is greater than $E'$ |
| >= | $E >= E'$ | None | $E$ is greater than or equal to $E'$ |
| not = | $E$ not $= E'$ | None | Equivalent to $E <> E'$ |
| not <> | $E$ not $<> E'$ | None | Equivalent to $E = E'$ |
| not < | $E$ not $< E'$ | None | Equivalent to $E >= E'$ |
| not <= | $E$ not $<= E'$ | None | Equivalent to $E > E'$ |
| not > | $E$ not $> E'$ | None | Equivalent to $E <= E'$ |
| not >= | $E$ not $>= E'$ | None | Equivalent to $E < E'$ |
| | | *Ternary operators* | |
| between | $E$ between $E'$ and $E''$ | None | Equivalent to $(E'<=E$ and $E<=E'')$ |
| not between | $E$ not between $E'$ and $E''$ | None | Equivalent to $(E'>E$ or $E>E'')$ |
| | | *Unary operators* | |
| is null | $E$ is null | Postfix | Returns true iff $E$ is null. Does not distinguish between empty strings and |

| | | | nulls. |
|---|---|---|---|
| is not null | $E$ is not null | Postfix | Equivalent to not ($E$ is null) |

2205

2206

2207 The equality and inequality operators (=, <>, and their negated variants) can take any scalar
2208 values as operands. Scalar relational operators (<, <=, >, >=, between and their negated
2209 variants) only take numeric operands. If at least one operand to a relational operator is null,
2210 the result is also null.

2211 Unary test operators is null and is not null test whether the operand is (or is not
2212 null) and return the corresponding Boolean value as a result.

### Scalar Functions

2214 In VTL 1.1, scalar functions (i.e., functions whose arguments are only scalar and that return
2215 scalar as a result) can also be automatically lifted to dataset/scalar and dataset levels,
2216 similarly to the unary and binary operators. For instance, pow(X,N) computes N-th power of
2217 number X, and log(X) computes the natural logarithm of X. When one or more arguments
2218 to such a function are datasets, they get automatically lifted. For instance,
2219 pow(D1.X,D2.N) joins D1 and D2, and then for each matched row computes the scalar
2220 power, taking the measure X from D1 as the base and measure N from D2 as the exponent, and
2221 the result is a joint dataset with an additional column holding the result.

## Join Expressions

2223 VTL 1.1 introduces the join expressions as the base mechanism for combining and
2224 manipulating datasets, including the lifting of the scalar operators and functions to the
2225 dataset/scalar and dataset levels.. The general join expression syntax has the form:

2226 $$[JOIN] \ \{BODY\}$$

2227 where *JOIN* is one of several join specifications described below, and *BODY* is a list of zero or
2228 more join expression statements that perform data filtering, computation, manipulation,
2229 grouping and ordering, also described in more detail in the text that follows. The start of the
2230 join expression is distinguished by the open square bracket (" [").

### Join Specifications

2232 The join specification is one of the following:

2233 • *d* – a single dataset variable. In this case we have **dataset traversal** (no join is
2234    performed). *BODY* is executed for each record in *d*. Inside *BODY*, *d* refers to the current
2235    record in dataset *d*.

2236 • $d_1, d_2, …, d_n$ – where *n*>1, performs an **outer join** of datasets held in dataset variables
2237    $d_1,d_2,…,d_n$. These datasets must be joinable: for some index *j*, the set of identifier
2238    components in $d_j$ must include identifiers from all other datasets; $d_j$ is called the pivot
2239    dataset. Then, $d_j$ is joined using a full outer join with each of datasets $d_i$ (*i*<>*j*) on shared

2240     identifier components. Inside *BODY*, each of $d_1,d_2,...,d_n$ refers to the matched record
2241     from the respective dataset.

2242   •  $d_1$ `outer` $d_2, ..., d_n$ – where $n>1$, is synonymous to the previous case $d_1, d_2, ..., d_n$.

2243   •  $d_1$ `inner` $d_2, ..., d_n$ – where $n>1$, performs an **inner join** of datasets held in dataset
2244     variables $d_1,d_2,...,d_n$. As in the outer join case, the datasets must be joinable: for some
2245     index $j$, the set of identifier components in $d_j$ must include identifiers from all other
2246     datasets; $d_j$ is called the pivot dataset.. Then, $d_j$ is joined using an inner join with each of
2247     datasets $d_i$ ($i<>j$) on shared identifier components. Inside *BODY*, each of $d_1,d_2,...,d_n$
2248     refers to the matched record from the respective dataset.

2249   •  $d_1$ `cross` $d_2, ..., d_n$ – where $n>1$, performs a **cross join** (or a Cartesian product) of
2250     datasets held in dataset variables $d_1,d_2,...,d_n$. All combinations of records are processed
2251     in *BODY*, and each of $d_1,d_2,...,d_n$ refers to the matched record from the respective
2252     dataset.

2253 The meaning of the inner and the outer join is the same as the meaning of `INNER JOIN`
2254 and `FULL OUTER JOIN` constructs, respectively, in the SQL-92 standard. In the cross join
2255 case, *BODY* of the join expression typically filters out record combinations that do not fit
2256 some logical condition.

2257 It is possible for two or more dataset variables involved in a join to refer to (i.e., act as
2258 aliases for) the same dataset. Inner and outer joins recognize dataset aliases, and
2259 automatically simplify the join structure to ensure that each dataset variable refers to a
2260 distinct dataset, while the aliases can still be used in the join body and refer to the same
2261 matched record from the original dataset. This is an automatic process that is transparent
2262 to the user. Indeed, aliases can be safely removed in an inner or outer join because joining
2263 a dataset with itself on the same set of identifier components always matches each record
2264 with itself.

2265 However, in a cross join, each dataset variable is used, whether or not two or more of them
2266 refer to a same dataset. This allows matching of two or more records from the same
2267 dataset using custom filter criteria, and is instrumental in implementing multiple-record
2268 (combinatorial, first-order, or "diagonal") validation rules.

2269 **Functional Integrity**

2270 The VTL information model requires of each dataset a functional dependency between the
2271 identifier components and all other components.  If we look at a dataset as a tabular structure
2272 with a finite number of columns (which correspond to components) and rows (which
2273 correspond to individual records), this translates into the following *functional integrity*
2274 requirements:

2275   •  A dataset can have an arbitrary number of identifier, measure and attribute columns.
2276     Each column has a distinct name in the dataset, and a scalar data type.

2277   •  All `null` values in string columns are implicitly converted into the empty string, and
2278     are not seen as `null`s in the points below.

2279   •  If a dataset has no identifier columns, but it has at least one measure or attribute
2280     column, it must have exactly one row. A dataset that has no columns whatsoever

2281    cannot have any rows. The points below apply only to datasets with one or more
2282    identifier components.

2283    • No identifier column can have a `null` value in any dataset row.

2284    • The combination of identifier column values in a dataset row is called the key. Two or
2285    more rows in the same dataset cannot have the same key.

2286    • When a measure or attribute column has value `null` in a dataset row, it is considered
2287    undefined for that row's key.

2288    The join expressions not only expect the input datasets to be functionally integral, but are
2289    engineered in a way that ensures functional integrity of the result. The key to this is the
2290    behaviour of join clauses and elements of *BODY*, explained below. Therefore, any construct
2291    built with the join expressions, including the lifting of the scalar operators and functions to
2292    the dataset/scalar and dataset levels, respects functional integrity by construction.

### Successive Dataset Transformations

2294    To explain the meaning of the join expressions, we can logically view it as a series of
2295    successive dataset transformations:

2296    First, the join specification that starts a join expression (traversal, inner, outer, or cross join)
2297    creates by itself the initial "joined" dataset:

2298    o  For a ***dataset traversal***, the initial dataset is identical to the traversed dataset.
2299    • For For ***inner and outer joins***, the initial working record consists of the identifier
2300    components from the pivot dataset matching record.

2301    • For ***cross join***, the initial working record consists of identifier components from all
2302    input datasets: identifier component $X$ from input $d_i$ appears under name $d_i\_X$
2303    (name of the dataset variable $d_i$ plus an underscore, plus the name of the
2304    component $X$). To avoid possible ambiguities, in the cross join case the names of
2305    input dataset variables cannot contain an underscore.

2306    Second, the first join expression statement in *BODY* (if any) operates on this initial dataset and
2307    produces a resulting dataset. which is fed as input to the next statement in *BODY*, etc. The
2308    dataset which is the result of the last statement in *BODY* is the result of the entire join
2309    expression.

2310    It should be noted that this is a logical view on the semantics of the join specification and the
2311    statements in *BODY*, which makes it easy to explain and understand. In reality, having each
2312    statement making a separate pass through its input dataset would not be efficient. Indeed, it is
2313    often the case that all *BODY* statements can be executed in a single pass (e.g., a single SQL
2314    query) through the joined datasets.

### Kinds of Body Statements

2316    The element *BODY* in a join expression consists of zero or more ***join expression statements***
2317    that define the processing steps applied to the (joined) input datasets inside the join
2318    expression.  These statements can be divided in two main groups:

- **_Record-level statements_** process each individual record of the statement's input dataset, by adding or updating columns, computing temporary values (i.e., local variables), or deciding whether to keep or discard a record based on a filter condition.

- **_Transposition statements_**, which unfold an identifier component (a measure dimension) from several records from its input dataset into a single output record, or perform a symmetric folding operation. The measure dimension breakdown for folding and unfolding is either given explicitly as a part of the transposition statement, or by reference to an externally defined hierarchy.

## Record-Level Statements

Several record-level statements use *scalar expressions in the column mode*. These are expressions that evaluate to a scalar value, but differ from normal scalar expressions (in the general mode) in the interpretation of identifiers. In the column mode expressions, the identifiers (that are not followed by an open parenthesis or a `.`) refer to components in the working record which is the input to the statement, and not to variables. To refer to a variable, one has to prefix its name with a dollar sign.

### Explicit component computations

These statements compute the value of a component in the working record.

| Form | Description |
|------|-------------|
| $X := E$ | Computing new/updated measure |
| `measure` $X := E$ | Same as the previous |
| `attribute` $X := E$ | Computing new/updated attribute |
| `identifier` $X := E$ | Computing new/updated identifier |

In the above table, $X$ is a component name (an identifier) for the newly computed component, and $E$ is a scalar expression in the column mode. By default, if an explicit role keyword (`measure`, `attribute`, or `identifier`) is omitted, role `measure` is assumed.

An explicit component computation adds to the working record a component named $X$ with a given role and value specified by $E$. The working record may already contain a measure or attribute component named $X$, which can be used in $E$, but is replaced with the newly computed $X$ (which may have a different role and/or type). An error is raised if the working record has an identifier component named $X$.

The type of component $X$ in the resulting working record is the type of expression $E$.

$E$ is not a string expression and it evaluates to `null`.Example 1:

```
[D] {
  Total := Men + Women
  WomenRatio := Women / Total
  MenRatio := 1.0 – WomenRatio
  attribute ObsStatus := ObsStatus || "A"
}
```

2353 Example 2:

```
2354 [D] {
2355   Population := Population * 1.01
2356   attribute ObsStatus := ObsStatus || "I"
2357 }
```

2358 Example 3:

```
2359 [D1, D2] {
2360   Population := D1.Population + D2.Population
2361   attribute ObsStatus := D1.ObsStatus || D2.ObsStatus
2362 }
```

2363

2364 Implicit component computations

2365 The implicit component computation statements compute the value of a component if it is not
2366 already present in the working record.

| Form | Description |
|------|-------------|
| implicit $X := E$ | Computing implicit measure |
| implicit measure $X := E$ | Same as the previous |
| implicit attribute $X := E$ | Computing implicit attribute |
| implicit identifier $X := E$ | Computing implicit identifier |

2367 In the above table, *X* is a component name (an identifier), and *E* is a scalar expression in the
2368 column mode. By default, if an explicit role keyword (measure, attribute, or
2369 identifier) is omitted, role measure is assumed.

2370 The implicit component computation statements behave similarly like their explicit
2371 counterparts (without keyword implicit), but they are executed only if the working record
2372 does not already have a component named *X*. An error is raised if there is already a
2373 component named *X*, but with a different role.

2374 The type of component *X* in the resulting working record is the type of expression *E*.

2375 *E* is a non-string expression that evaluates to null. Example 1:

```
2376 [D] {
2377   implicit attribute ObsStatus := ""
2378 }
```

2379 Example 2:

```
2380 [D1, D2] {
2381   Population := D1.Population + D2.Population
2382   attribute ObsStatus := D1.ObsStatus || D2.ObsStatus
2383   implicit identifier RefArea := "EU"
2384 }
```

2385 Computing local variables

2386 Local variables store a value for the remainder of the record-level statements in *BODY*.

---

| Form | Description |
|---|---|
| $\$X := E$ | Computing a local variable |

2387 In the table above, *X* is an identifier, used as a variable name, and *E* is a scalar expression in
2388 the column mode.

2389 This statement is useful for computing a value and storing the result temporarily for easier
2390 reference, without making it appear in the result.

2391 Example:

```
[D] {
  $Total := Men + Women + Children
  WomenRatio := Women / $Total
  MenRatio := Men / $Total
  ChildrenRatio := 1.0 - WomenRatio - MenRatio
}
```

2398 Filtering records

2399 The filtering statement decides whether to keep the working record in the result or to omit it.

| Form | Description |
|---|---|
| `filter` *E* | Permit only records satisfying condition *E* |

2400 In this statement, *E* is a Boolean expression in the column mode.

2401 If at runtime *E* does not evaluate to `true`, no further record-level statements are executed,
2402 and the working record is discarded.

2403 Example 1:

```
[D] {
  $Total := Men + Women + Children
  WomenRatio := Women / $Total
  MenRatio := Men / $Total
  filter MenRatio + WomenRatio >= 0.6 /* Treat only these cases. */
  ChildrenRatio := 1.0 - WomenRatio - MenRatio
}
```

2411 Example 2:

```
[D1, D2] {
  filter D1.Pop is not null or D2.Pop is not null
        /* At least one of D1.Pop and D2.Pop must be defined. */
  Pop := D1.Pop + D2.Pop
}
```

2417 Example 3:

```
[D1 cross D2] {
  filter D1.Pop < D2.Pop    /* Custom join condition. */
  Ratio := D1.Pop / D2.Pop
}
```

2422 Function application to components of the working record

2423 These statements transform components of the working record by applying a function to
2424 them.

| Form | Description |
|---|---|
| apply $F$ | Apply function to measures of the matching type |
| apply $F$ to attributes | Apply function to attributes of the matching type |
| apply $F$ to measures and attributes | Apply function to measures and attributes of the matching type |

Here, $F$ is a function that takes one argument of some scalar type $t$ and returns a result of some scalar type $T$. The first form transforms value of each measure $X$ from the working record whose type is compatible with $t$ to value $F(X)$ of type $T$ in the resulting working record.

The statement forms that include `to attributes' and `to measures and attributes' apply function $F$ to components with the respective roles, not just to measures as in the first form.

Example:

```
[D] {
  apply _*1000           /* Multiplies all numeric measures by 1000. */
  apply _&"x" to attributes /* Adds "x" to all string attributes. */
}
```

Function application to components of the matched input records

These statements combine components from the matched records of the input datasets by applying a function to their values and adding the result to the working record.

| Form | Description |
|---|---|
| apply $F$ to $d_{k1}$, ..., $d_{km}$ | Apply function to measures from $d_{k1}$, ..., $d_{km}$ with same names and matching types |
| apply $F$ to attributes in $d_{k1}$, ..., $d_{km}$ | Apply function to attributes from $d_{k1}$, ..., $d_{km}$ with same names and matching types |
| apply $F$ to measures and attributes in $d_{k1}$, ..., $d_{km}$ | Apply function to measures and attributes from $d_{k1}$, ..., $d_{km}$ with same names and matching types |

Here, $F$ is a function that takes $m>0$ arguments of the corresponding scalar types $t_1,...,t_m$, and returns a scalar result of type $T$. $d_{k1},...,d_{km}$ is a subset of the input dataset variables from *JOIN* that represent the matched records in *BODY*.

The first form of the statement looks for the same-name measure components that appear in each of $d_{k1},...,d_{km}$ and whose respective types are compatible with $t_1,...,t_m$. For each such shared component named $X$, a measure component $X$ of type $T$ is added (or replaced) in the resulting working record, with value $F(d_{k1}.X,...,d_{km}.X)$.

The forms with `to attributes in' and `to measures and attributes in' apply $F$ to the components of the respective role, not just to measures in $d_{k1},...,d_{km}$.

2448    Example:

```
2449  [D1,D2]                                                        {
2450    apply 0.3*_+0.7*_  to D1, D2         /* Weighted sum of numeric measures */
2451    apply _&_  to attributes in D1, D2  /* Concatenating string attributes */
2452    apply _or_  to attributes in D1, D2 /* Disjunction of Boolean attribs */
2453  }
```

2454    <u>Component renaming statements</u>

2455    These statements change names of one or more components in the working record
2456    simultaneously.

| Form | Description |
|---|---|
| `rename` $X_1$ `to` $Y_1$, $X_2$ `to` $Y_2$, ..., $X_n$ `to` $Y_n$ | Simultaneously rename $X$s to $Y$s |
| `rename` $X_1$ -> $Y_1$, $X_2$ -> $Y_2$, ..., $X_n$ -> $Y_n$ | Same as the above |

2457    Each $X_i$ and $Y_i$ (i=1,...,$n$, $n>0$) in the table above an identifier specifying a column name,
2458    optionally preceded with a role (`identifier`, `measure`, or `attribute`). Identifiers in $X_1$,
2459    ..., $X_n$ must be mutually distinct, as well as those in $Y_1$, ..., $Y_n$.

2460    Each $X_i$ must exist in the working record. If $X_i$ does not specify the source role, the actual role
2461    of the component with that name is in the working record is used. If $Y_i$ does not specify the
2462    target role, the source role is used.

2463    The renaming statement (between { }) is performed simultaneously as a whole, which makes
2464    column name and role swapping and cycling possible with a single statement. If the working
2465    record has a measure or attribute whose name is in $Y_1$, ..., $Y_n$, but not in $X_1$, ..., $X_n$, that
2466    component is replaced by the renamed component. However, an error is raised if such
2467    component is an identifier.

2468    It is also an error to change the role of an identifier component using `rename`.

2469    Example 1:

```
2470  [D] {
2471    rename A to B, B to A /* Swap component names */
2472  }
```

2473    Example 2:

```
2474  [D] {
2475    rename identifier Geo to RefArea, /* Rename identifier Geo */
2476           Age to identifier Age,      /* Make Age an identifier */
2477           attribute ObsStatus to measure Status,
2478                                        /* Convert attribute to a measure */
2479           Z to attribute Z            /* Error if Z is an identifier */
2480  }
```

2481    <u>Component filtering statements</u>

2482    These statements keep or drop the specified components in the working record.

| Form | Description |
|---|---|
| `keep` $X_1$, ..., $X_n$ | Keep measures or attributes in the working record |

| | |
|---|---|
| `drop` $X_1, ..., X_n$ | Drop measures or attributes from the working record |

2483 Each $X_i$ (i=1,...,n, n>0) is an identifier giving the column name, optionally preceded with a role
2484 `measure` or `attribute`.

2485 Statement `keep` keeps in the working record only the measures and attributes given by $X_1$, ...,
2486 $X_n$, which must all exist in the working record. Identifiers are not affected.

2487 Statement `drop` drops from the working record those measures and attributes given by $X_1$, ...,
2488 $X_n$ that exist in the working record. An error is raised if any of $X_1$, ..., $X_n$ is an identifier.

2489 Example 1:

```
[D] {
  $Total := Men + Women + Children
  WomenRatio := Women / $Total
  MenRatio := Men / $Total
  ChildrenRatio := 1.0 – WomenRatio – MenRatio
  keep WomenRatio, MenRatio, ChildrenRatio
                    /* Keep only these measures (no attributes kept)  */

}
```

2499 Example 2:

```
[D] {
  $Total := Men + Women + Children
  WomenRatio := Women / $Total
  MenRatio := Men / $Total
  ChildrenRatio := 1.0 – WomenRatio – MenRatio
  drop Women, Men, Children
      /* Keep all measures and attributes except these three  */
}
```

### Transposition Statements

2509 The transposition statements can be used instead of the aggregation statements. These
2510 statements also operate on all records resulting from the join and the record-level statements,
2511 but instead of aggregating, they transpose columns from several input records into a single
2512 output record and back.

| Form | Description |
|---|---|
| `unfold` $X, Y$ `to` $B_1, ..., B_n$ | Unfold identifier $X$ and measure $Y$ into columns $B_1$, ..., $B_n$ (n>0). |
| `unfold` $X, Y$ `using` $H$ | Unfold identifier $X$ and measure $Y$ using hierarchy definition $H$. |
| `fold` $B_1, ..., B_n$ `to` $X, Y$ | Fold columns $B_1, ..., B_n$ (n>0) into a new identifier $X$ and measure $Y$. |
| `fold` `using` $H$ `to` $X, Y$ | Fold a new identifier $X$ using hierarchy definition H and measure $Y$. |

2513 In the above table, $X$ is the name of an identifier column

2514  Each $B_i$ in breakdown $B_1, ..., B_n$ is either a base element (an identifier), or a computed element
2515  of the form $Z=C_1+ ... + C_m$, where $Z$ is an identifier, and $C_1, ..., C_m$ ($m>0$) are other breakdown
2516  elements (base or computed) that go into $Z$. Circular dependencies between computed
2517  breakdown elements are not allowed. Each breakdown element $B_i$ has the base set $U_i$ of base
2518  elements that it "covers". If $B_i$ is a base breakdown element, its $U_i=\{B_i\}$. If $B_i$ is a computed
2519  breakdown element of the form $Z=C_1+ ... + C_m$, its elementary set is the union of the base sets
2520  of $C_1, ..., C_m$.

2521  The breakdown structure $B_1, ..., B_n$ can be specified explicitly in the statement, or it can be
2522  defined in a hierarchy object $H$ defined elsewhere (i.e., in metadata). In the text that follows
2523  we shall assume that in the latter case the actual structure $B_1, ..., B_n$ has been retrieved from $H$.

2524  The `unfold` statement divides the input dataset with a string identifier component $X$ (the
2525  measure dimension) and a numeric measure component $Y$ into groups of records sharing the
2526  value of all identifiers other than $X$.

2527  Each input group is then transformed into a single output record that has:

2528  - A copy of all identifier components from the input group except $X$.

2529  - Numeric measure columns $B_1, ..., B_n$ instead of the single measure column $Y$. For each $B_i$
2530    ($i=1..n$), the value of the measure column named $B_i$ in the output record is the sum of $Y$
2531    in the group records where the value of $X$ belongs to the base set of $B_i$ (as a set of string
2532    literals).

2533  - All other measure and attribute components, whose value is taken as the maximum in
2534    the group.

2535  The `fold` statement works in the opposite direction: for each input record it generates a
2536  group of output records, with one output record for each breakdown element $B_i$ ($i=1..n$)
2537  where the value of component $B_i$ is not `null`, consisting of:

2538  - A copy of all identifier components from the input record.

2539  - A new string identifier component named $X$ with value equal to $B_i$ (as a string literal).

2540  - A new numeric measure component $Y$ with value equal to the value of $B_i$ in the input
2541    record.

2542  - A copy of all attribute and measure components (other than $B_1, ..., B_n$) taken from the
2543    input record.

2544  Example 1:

2545  Suppose *BeNeLuxPop* is the following dataset:

| **Year** | **Geo** | Pop | *Status* |
|---|---|---:|---|
| 2015 | BE | 11,324 | A |
| 2015 | NE | 16,948 | |
| 2015 | LU | 563 | AP |

2546

2547 Then the result of the join expression:

2548 ```
[BeNeLuxPop] {
2549    unfold Geo, Pop to BE, NE, LU, Total = BE + NE + LU
2550 }
```

2551 is:

| Year | BE | NE | LU | Total | *Status* |
|---|---|---|---|---|---|
| 2015 | 11,324 | 16,948 | 563 | 28,835 | AP |

2552 Example 2:

2553 If D is the result of the previous example, then the following join expression:

2554 ```
[D] {
2555    fold BE, NE, LU, Total = BE + NE + LU to Geo, Pop
2556 }
```

2557 gives the result:

| Year | Geo | Pop | *Status* |
|---|---|---|---|
| 2015 | BE | 11,324 | AP |
| 2015 | NE | 16,948 | AP |
| 2015 | LU | 563 | AP |
| 2015 | Total | 28,835 | AP |

2558

2559 Note that this result is very similar to the original input, except for a couple of differences that
2560 illustrate some important aspect of the `fold` and `unfold` statements:

2561 ● The computed breakdown element Total appears in the result, while it was not present
2562    in the original input dataset *BeNeLuxPop*. If this is undesirable, the `fold` statement
2563    should use only the base (not computed) breakdown components BE, NE, and LU.

2564 ● In the `fold` statement, the computed breakdown elements, such as *Total*, are not
2565    computed, but are treated in the same way as the base breakdown elements (*BE*, *NE*,
2566    and *LU*).

2567 ● While the *Status* attribute varies in the original input dataset *BeNeLuxPop*, it is
2568    uniformly equal to "AP" in all result rows. The reason for this is that unfolding entails a
2569    loss of information for attributes like *Status*, where it takes the maximum for the whole
2570    group of records where *Year*=2015. Folding, on the other hand, does not entail any loss
2571    of information (it can, in fact, create additional information, as seen in the previous
2572    point).

# Lifting Scalar Operators and Functions With Join Expressions

We now turn to the issue of lifting the scalar operators and functions to the dataset/scalar and dataset level using the join expressions. This lifting is not something a VTL programmer needs to do manually -- it is done automatically under the hood by the compiler. However, it is important for both the programmers and language implementers to understand clearly how the lifting works in order to ensure the correct behaviour.

## Liftable Expressions

As a preliminary, we need to define what is a "liftable" expression. For an expression to be liftable, it has to satisfy certain structural and typing constraints. The typing constraints are important because the syntactic form of an expression does not provide sufficient information for deciding whether an expression needs to be lifted and how. For instance, `A+B` may be a scalar or a dataset expression, depending on the types of $A$ and $B$. For what we need here, we shall take a simplified look at the type analysis:

- The typing of an expression is decided inductively, or bottom-up: from the operation or function argument types to the type of the operator application or function call.

- After determining that the type of an expression $E$ is $t$, we shall be making simple assertions, such as: "$t$ is a scalar type (i.e., $E$ is a scalar expression)", or "$t$ is a dataset type (i.e., $E$ is a dataset expression)".

Intuitively, we can define a scalar-based expression as an expression that uses only scalar operators and functions on arguments that are scalar variables or literals, datasets and their components, or scalar-based subexpressions. A liftable expression is then a scalar-based expression that returns a dataset, because one or more of the arguments to a scalar operator or function is given as a dataset. Or, in other words, only a scalar-based expression can be liftable, but the property of being liftable is stronger.

More formally, we say that an expression of the form $f(E_1, ..., E_n)$, $n>0$, is a scalar-based expression if:

- $f$ accepts $n$ scalar arguments and returns a scalar result

- Each argument $E_i$ ($i=1..n$) is one of the following:
    - a scalar variable or a numeric, string or Boolean literal [weak argument]
    - a dataset variable [strong argument]
    - an expression of the form $D.X$ where $D$ is a dataset variable, and $X$ is a component identifier [strong argument]
    - a scalar-based expression [strong argument exactly when $E$ is liftable]
- If at least one argument is strong, then the scalar-based expression $f(E_1, ..., E_n)$ is liftable.

We wrote $f(E_1, ..., E_n)$ to denote both a call to function $f$ and an application of an $n$-ary operator (prefix, infix, or postfix) to its arguments.

Example 1:

2611 Expressions $-X$, `log(X)`, and `X*Y`, where *X* and *Y* are scalar variables, are all scalar-based,
2612 but they are not liftable, because they do not use any dataset. However, expression
2613 `D.X*log(D.X)`, where *D* is a dataset variable, is both scalar-based and liftable.

2614 Example 2

2615 Expression:

2616 $$D1^2+2*D1*D2+D2^2$$

2617 where *D1* and *D2* are dataset variables, is liftable, because it uses these two dataset variables
2618 as arguments to basically scalar operators +, *, and ^.

## Component Selection And Lifting Scheme

2620 A liftable expression *E* must contain one or more dataset references of the form *D* or dataset
2621 component references of the form *D.X*, where *D* is a dataset variable. The shape of these
2622 references significantly affects the computation that is performed.

2623 In the sub-sections that follow we cover all dataset and dataset component reference cases
2624 that may occur, and give representative examples of the lifting scheme.

### Operating on All Shared Components

2626 The first case is when *E* contains only dataset references (*D*), but no dataset component
2627 references (*D.X*). In this case, the computation is performed on all shared measure
2628 components, i.e., the measure components with the same name and type that appear in all
2629 referenced datasets. The resulting dataset uses these shared measure components to hold the
2630 result.

2631 Example 1:

2632 As a simple example, `D1+D2`, where *D1* and *D2* are dataset variables with numeric measure
2633 components *A* and *B*, will create a result with measure components *A* and *B* whose value is the
2634 sum of *A*s and *B*s from *D1* and *D2*. The lifting is then done using a join expression and `apply`:

```
2635 [D1,D2] {
2636    apply _+_ to D1, D2
2637 }
```

2638 Example 2:

2639 Expression:

2640 $$D1^2+2*D1*D2+D2^2$$

2641 is lifted with:

```
2642 [D1,D2] {
2643    apply \x,y{x^2+2*x*y+y^2} to D1, D2
2644 }
```

2645 In this example, we had to explicitly name the arguments *x* and *y* in the function, because *D1*
2646 and *D2* appear more than once in the original expression.

### Operating on Single Named Component

2648 The second case is when *E* contains one or more dataset component references of the form
2649 *D.X* where *D* may vary, but *X* is a single component name. In this case, we only operate on that

single component *X* in all referenced datasets, and the result contains a single measure component *X* holding the result. All dataset references of the form *D* in *E* are implicitly rewritten into *D.X*. The fixed component *X* must not be null in at least one referenced dataset.

Example 1:

Expression:

```
                    D1.Pop + D2
```

where *D1* and *D2* are dataset variables, operates on a single component *Pop*. It is therefore equivalent to:

```
                    D1.Pop + D2.Pop
```

And is lifted as:

```
[D1,D2] {
  filter D1.Pop is not null or D2.Pop is not null
  Pop := D1.Pop + D2.Pop
}
```

The result contains a single measure component named *Pop*.

Example 2:

Expression:

```
                    D1.Pop * 1.02
```

also uses the single named component *Pop*. It is lifted as follows:

```
[D1] {
  filter D1.Pop is not null
  Pop := D1.Pop * 1.02
}
```

Note that in this example the join expression traverses a single dataset *D1*, and therefore all other measures and attributes are kept unchanged in the result.

### Operating on Multiple Named Components

Finally, we may have a case where *E* contains two or more dataset component references of the form *D.X* where *X* is not always the same. This case was illegal in VTL 1.0 because of the rule that differently named components from different datasets cannot mix in a computation. The experience indicates that his requirement can sometimes be too strict, and may force the programmer to frequently explicitly rename components in order to be able to compute on them.

That is why VTL 1.1 allows mixing two or more differently named dataset components in a single liftable expression *E*, provided that *E* contains no dataset references of the form *D* (i.e., only contains dataset component references of the form *D.X*). The resulting dataset contains a single measure component named *Value* holding the result of the computation.

Example 1:

Expression:

```
      D1.Pop + D2.Population + D3.Residents + D4.Inhabitants
```

is lifted as follows:

```
2690   [D1, D2, D3, D4] {
2691     filter D1.Pop is not null or D2.Population is not null
2692             or D3.Residents is not null or D4.Inhabitants is not null
2693     Value := D1.Pop * D2.Population + D3.Residents + D4.Inhabitants
2694   }
```

2695   Example 2:

2696   Expression:

```
2697                     D1.Pop between D2.Min and D2.Max
```

2698   is lifted as follows:

```
2699   [D1, D2] {
2700     filter D1.Pop is not null or D2.Min is not null
2701             or D2.Max is not null
2702     Value := D1.Pop between D2.Min and D2.Max
2703   }
```

2704   The resulting measure *Value* is Boolean.

2705   Example 2:

2706   Expression:

```
2707           (D1.Pop between D2.Min and D2.Max)[Value->InRange]
```

2708   is lifted as follows:

```
2709   [D1, D2] {
2710     filter D1.Pop is not null or D2.Min is not null
2711             or D2.Max is not null
2712     Value := D1.Pop between D2.Min and D2.Max
2713     rename Value -> InRange
2714   }
```

2715   The resulting Boolean measure generically named *Value* has been renamed to more domain-
2716   specific *InRange*.

### Allowing Non-Scalar-Based Subexpressions

2718   The approach for lifting expressions built with scalar operators and functions to the
2719   database/scalar and database levels explained above restricts the structure of such
2720   expressions to scalar-based expressions defined above. This limitation can sometimes be too
2721   strict. For instance, expression:

```
2722                         D1.Total + size(D2)
```

2723   where *D1* and *D2* are dataset operations, and *size* is a function that returns the number of
2724   records in a dataset, is not scalar-based (and therefore misses the precondition to be lifted)
2725   because *size* does not take a scalar, but a dataset argument. Therefore, in this expression *D2*
2726   should be treated differently than *D1*: we do not need to join these two datasets, we just first
2727   need to count rows in *D2*, remember the result and then use it in the main expression.

2728   Another example is:

```
2729                         union(D1, D2) * 1.02
```

2730 This is also a valid expression, where we increase all numeric measures in the union of two
2731 datasets *D1* and *D2* by 2%. But it is not a scalar-based expression (and therefore not a liftable
2732 one), because *union* is not a scalar function. Still it is clear that first we have to make a union
2733 of *D1* and *D2* , and then multiply the result with 1.02.

2734 These two examples hint at a general solution: we can often transform a non-scalar-based
2735 expression into a scalar-based one by proceeding step-by-step.

2736 Let us first take $E$ to be an expression that contains some sub-expression $A$. It is clear that $E$ is
2737 equivalent to a VTL block:

2738 {
2739    $V$ := $A$
2740    $E[V/A]$
2741 }

2742 where $V$ is a variable name that does not appear in $E$, and $E[V/A]$ is a copy of $E$ where $V$
2743 replaces $A$.

2744 This scheme can be automatically applied to all scalar or dataset subexpressions $A_1$, ..., $A_n$ of $E$
2745 that are not scalar-based. As a result, we transform $E$ into the form:

2746 {
2747    $V_1$ := $A_1$ /* $V_1$ does not appear in $E$ */
2748    $V_2$ := $A_2$ /* $V_2$ does not appear in $E$ */
2749    ...
2750    $V_n$ := $A_n$ /* $V_n$ does not appear in $E$ */
2751    $E[V/A]$    /* Becomes liftable expression! */
2752 }

2753 This transformation can be automatically done by the compiler.

2754 *Example 1*:

2755 Expression:

2756 ```
D1.Total + size(D2)
```

2757 becomes:

2758 ```
{
2759   V := size(D2)
2760   D1.Total + V /* liftable */
2761 }
```

2762 which after lifting becomes:

2763 ```
{
2764   V := size(D2)
2765   [D1] {
2766     filter D1.Total is not null
2767     Total := D1.Total + V
2768   }
2769 }
```

2770 *Example 2:*

2771 Expression:

```
2772   union(D1, D2) * 1.02
```

2773   becomes:

```
2774   {
2775     V := union(D1, D2)
2776     V * 1.02 /* liftable */
2777   }
```

2778   which after lifting becomes:

```
2779   {
2780     V := union(D1, D2)
2781     [V] {
2782       apply _*1.02
2783     }
2784   }
```

## Expressing Validation Rules With Join Expressions

2786   In the previous sections we have shown how the VTL 1.1 join expressions can be used for
2787   lifting of basically scalar expressions and functions to the dataset/scalar and dataset levels.
2788   This lifting is performed automatically and transparently by the compiler, and provides a
2789   well-defined semantics for the lifted constructs. We can therefore think about the join
2790   expressions as a "core" mechanism for expressing the behaviour of higher-level dataset
2791   operations.

2792   The same approach can be used for expressing the behaviour of some important classes of
2793   validation rules:

- ● *Horizontal rules* -- these rules check validity of individual records (or rows) in a
    dataset.  For the sake of simplicity, let us say that each horizontal rule has a condition
    *SCOPE_COND* that selects records to which the validation rule needs to be applied, a
    condition *VALID_COND* that defines when a row is valid, and a string *RULE_CODE* that
    is inserted in the result column *ERR_CODE* if the validation fails on a record.  The
    validation of a dataset *D* using a horizontal rule is then equivalent to:

  - ● `[D] {`
      `implicit attribute ERR_CODE := ""`
      `filter` *SCOPE_COND*
      `attribute` *RULE* `:=` *VALID_COND*
      `attribute ERR_CODE :=`
        `if` *RULE* `then ERR_CODE else` *RULE_CODE*
    `}`

- ● *Vertical rules* --  these rules apply to values of some measure component *Y* that are
    stacked "vertically" one under another in each group of records, so that each value of *Y*
    corresponds to a particular code of some measurement dimension *X*.  The breakdown
    of *X* to individual codes is typically given explicitly in a vertical rule as $B_1, ..., B_n$.  Again,
    for the sake of simplicity, let us say that each vertical rule has a condition *SCOPE_COND*

that selects groups of records to which it applies, a condition *VALID_COND* that defines when a row is valid, and a string *RULE_CODE* inserted in the result column *ERR_CODE* if the validation fails on a record.  The validation of a dataset *D* using a vertical rule is then equivalent to:

- {
    ```
      U := [D] { unfold X, Y to B₁, …, Bₙ }
      [U] {
        implicit attribute ERR_CODE := ""
        filter SCOPE_COND
        attibute RULE := VALID_COND
        attribute ERR_CODE :=
          if RULE then ERR_CODE else RULE_CODE
      }
    }
    ```

- ***First-order or combination rules*** -- these rules apply to combination of records from two or more datasets $D_1$, ..., $D_n$ (the same dataset variable can be repeated several times). The criteria for matching these records is specified as *MATCH_COND* , and we here take the other (simplified) assumptions about *VALID_COND*, *RULE*, and *RULE_CODE* as in the examples of the horizontal and vertical rules above. Then, the validation of a dataset *D* using this kind of rules is then equivalent to:

  - ```
    [D₁ cross D₂, …, Dₙ] {
        implicit attribute ERR_CODE := ""
        filter MATCH_COND
        attribute RULE := VALID_COND
        attribute ERR_CODE :=
          if RULE then ERR_CODE else RULE_CODE
    }
    ```

The above examples were simplified (among other things) because they refer to a single rule, while VTL 1.1 allows more powerful rule sets to be defined. However, at this point it should be evident that there are ways for expressing rule sets using the same kind of constructs. Suppose, for instance, we have a horizontal rule set consisting of three rules, *RULE1, RULE2* and *RULE3.*  The translation would look like this:

```
[D] {
  implicit attribute ERR_CODE := ""
  filter SCOPE_COND1 or SCOPE_COND2 or SCOPE_COND3
  $ERR_CODE := ERR_CODE
  $RULE1 := not (SCOPE_COND1) or VALID_COND1
  $ERR_CODE :=
    if $RULE1 then $ERR_CODE else paste($ERR_CODE, RULE_CODE1, ",")
  $RULE2 := not (SCOPE_COND2) or VALID_COND2
  $ERR_CODE :=
    if $RULE2 then $ERR_CODE else paste($ERR_CODE, RULE_CODE2, ",")
```

```
2854    $RULE3 := not (SCOPE_COND3) or VALID_COND3
2855    $ERR_CODE :=
2856      if $RULE3 then $ERR_CODE else paste($ERR_CODE, RULE_CODE3, ",")
2857    attribute RULESET := $RULE1 and $RULE2 and $RULE3
2858    attribute ERR_CODE := $ERR_CODE
2859  }
```

2860 This construct would check all three horizontal rules in the rule set in a single traversal of *D*,
2861 and would look only on records where at least one rule is applicable. It would create the
2862 attribute *ERR_CODE* if it did not exist, and would add to it (as a comma-separated list) error
2863 codes of all failed rules. The result would also have an attribute column *RULESET* (named
2864 after the rule set) which holds Boolean true if the `record` has passed all three rules, or
2865 `false` if at least one rule has failed on the record.

2866

2868 In this chapter we present some of the main assumption on which the Validation and
2869 Transformation Language bases the semantics of its Operators. These core assumptions
2870 complement the core language elements presented in the previous chapter, and they specify
2871 the general behaviour of the language, and is by default stable.  The standard library of
2872 operators is presented in detail in the Reference Manual, and presents the built-in
2873 functionality that can be gradually enriched following the evolution of the user needs.
2874 Possible new functions and operators must obviously comply with the core assumptions
2875 presented here.

2876 The main assumptions include:

2877 • Details of operand and result types
2878 • The general behaviour of operations on datasets
2879 • Storage and retrieval of datasets
2880 • The conventions for the grammar of the language

2881 The main assumptions are explained in the following sections.

## 2882 Details of operand and result types

### 2883 The Data types of the VTL

2884 As explained in the previous chapter, the type system of VTL 1.1 presents an outline of a type
2885 system, which is able to characterize all kind of objects that are used as an input, an
2886 intermediate result or auxiliary parameter, or produced as the result of any expression in a
2887 VTL program.

2888 In this section, we are concentrating on a subset of VTL types which we call the data types.
2889 Data types differ from other types in that they have a well-defined external representation,
2890 covered by the VTL Information Model (IM). Obviously, different parts of VTL programs can
2891 use or produce other objects, such as anonymous functions or tuples and collections of
2892 arbitrary objects, which are transient in nature. Such transient objects exist only in memory
2893 during the execution of a VTL program, but cannot be "materialized," i.e., they have no well-
2894 defined representation in the IM.

2895 The VTL data types, on the other hand, correspond to various artefacts represented in the IM.
2896 They include:

2897 • Datasets, composed of identifier, measure, and attribute components; each component
2898   contains a data of the same scalar type.
2899 • Collections of scalar types, or of Cartesian products of scalar types, which are used to
2900   express constraints, i.e., the permissible values for one or more scalar variables.
2901 • Modules representing dataset structure, as well as user-defined functions, types, and
2902   special objects such as validation rules.

### 2903 Basic scalar data types

2904 The **basic (unconstrained) scalar data types** of the language are: *string*, *number* (including
2905 *integer* and *float*), *boolean* and *date.* Their instances written directly in VTL code (i.e. the real

2906    objects of those types) are called *literals*. Thecharacteristics of the base scalar types are
2907    described in the following table.

| Basic scalar data types | |
|---|---|
| string | A sequence of zero or more UNICODE characters enclosed in double quotes ("). Examples of allowed literals for this data type are: "hello", "test", "x", "this is a string" and "" (the empty string). Note that in the VTL syntax the double quotes are intended to be the standard ones ("),  i.e. the same character to open and close the string, even if in this document and in the Part 2 the styled double quotes may be shown.  If a string literal needs to include a double quote in its contents, the quote needs to be doubled: literal "a""b" consists of three charaters: letter *a*, the double quote, and letter *b*. |
| number | Includes both *integer* and  a *float*. |
| float | Floating point numbers, whose precision is compatible with or greater than the IEEE 754 quadruple precision (128 bits encoding). At least the range of floating point numbers (absolute values) between $2^{-16949}$ (approx. $10^{-4965}$) and $2^{16384}-2^{16271}$ (approx. $1.1897*10^{4932}$) with 34 significant decimal digits should be representable. Alternatively, implementations may use arbitrary-precision floating point numbers. The point (.) is used as the decimal separator and must be present in the literal. Examples of allowed literals for this type are: 1.0, 234.56, 456.45; also the scientific notation is allowed: 12.23E+12, 35.2E-150, -2E10+3, 0.0.  The uppercase letter "E" can be written also as the lowercase "e". |
| integer | The basic signed integer type. At least 64 bit in size. Alternatively, implementations may use arbitrary-precision integers. Examples of allowed literals for this type are: 2, 5, 7, 24, -14, 0. |
| boolean | The Boolean data type. The allowed literals are *true* and *false*. |
| date | A point-in-time value. The type stores the year, the month, the day, the hours the minutes and the seconds (after midnight). Date are in 24-hours format: YYYY-MM-DD HH24:MI:SS While the YYYY-MM-DD is mandatory, HH24:MI:SS is optional and, if omitted, 00:00:00 is implied. Examples of allowed literal values are: 2012-09-30, 2013-10-02, 2014-01-01 12:23:35. The format for Date literals is customizable, in the sense that specific supplementary formats may be used in implementations in addition to this one, if properly configured in the system. Alternate literals may also include the ones adopted by commercial systems for compatibility reasons, for example: date'2012-09-30'. |

2908

2909 With reference to the VTL information model, the data type is a characteristic of the Value
2910 Domain. In turn, the data type of the Value Domain is inherited by its Values and its Subsets.

2911 A Represented Variable has the same data type of its Value Domain.

2912 A Structure Component has the same data type of the corresponding Represented Variable
2913 (i.e. the data type of its Value Domain).

2914 Also the Data Set has a data type, which is a "composite" one and corresponds to the set of the
2915 data types of its Structure Components.

2916 A Transformation (Expression) has the data type of its result.

### Type management and checking

2918 The language does not have explicit operators for converting the type (typecasting).

2919 It is envisaged that there will be "implicit upcasting" between the integer and the number data
2920 types. This means that wherever in the language it is possible to use a number, an integer or
2921 float is allowed. Obviously, the opposite is not allowed.

2922 The VTL is strongly typed, in the sense that any operand or parameter in an operation belongs
2923 to one of the possible types.

2924 The various VTL functions and operators have specific constraints in terms of number and
2925 types of parameters (see the corresponding sections in the Part 2).

2926 The type of an expression is computer at compile time.

2927 The function and operator constraints in terms of number and types of their arguments are
2928 statically checked (at compile time) so that type errors are not possible at runtime. Moreover,
2929 only type-safe upcast conversion for integers into num is performed.

2930 Type errors result in **compile time errors** preventing the Transformations from being used
2931 (exchanged, executed …).

2932

# The general behaviour of operations on datasets

### General rules

2935 As already mentioned, normally the model artefact produced through a Transformation is a
2936 Data Set (considered at a logical level as a mathematical function). Therefore, a
2937 *Transformation* is mainly an algorithm for obtaining a derived Data Set starting from already
2938 existing ones. As a matter of fact, the Data Set at the moment is the only type of Parameter
2939 that is possible to store permanently through a command of the language (see the Put section
2940 in the Part 2).

2941 If we assume that $F$ is a Data Set Operator (i.e., an operation that takes some inputs and
2942 produces a dataset), that $D_r$ is its result Data Set and that $D_{i\ (i=1,\dots n)}$ are its input Data Sets, the
2943 general form of a Transformation based on $F$ can be written as follows:

$$D_r := F (D_1, D_2, \dots , Dn)$$

2945 Operator F composes the Data Points of $D_{i\ (i=1,\dots n)}$ to obtain the Data Points of $D_r$.

2946 For computing the result of this operation, F follows a number of default behaviours
2947 described here.

2948 In general the Data Sets $D_{i\ (i=1,\dots n)}$ and consequently their Data Points may have any number of
2949 Identifier, Measure and Attribute Components, nevertheless the VTL Data Set Operators may
2950 require specific constraints on the Data Structure Components of their input Data Sets[16].

2951 The Data Structure Components of the result Data Set $D_r$ will be determined as a function of
2952 the Data Structure Components of the input Data Sets and the semantics of the Operator $F$.

2953 There can exist different cases of application of the Data Set Operators, having specific default
2954 behaviours and constraints.

2955 In particular, as for the number of operands, a **Data Set Operator** is called "**unary**" if it uses
2956 only one Data Set as input operand (e.g. minimum, maximum, absolute value ...) and "**n-ary**" if
2957 it requires more than one Data Set as input operand (e.g. sum, product, merge ...). The **n-ary**
2958 Operators require a preliminary matching between the Data Points of the various input Data
2959 Sets.

2960 **Data Sets** may be also usefully categorized with reference to the number of their Measure
2961 Components. A Data Set is called "**mono-measure**" if it has just one Measure Component and
2962 "**multi-measure**" if it has two or more Measure Components. For the multi-measure Data
2963 Sets it may be necessary to specify which measures should be considered in the operation.

2964 Other cases originate from the possible existence of missing data and Attribute Components.
2965 If there are missing values in the input Data Sets, the operation may generate meaningless
2966 outcomes, so inducing missing values in the result according to certain rules. On the other
2967 hand, there can be the need of producing the values for the Attribute Components of the result
2968 starting from the values of the Attributes of the operands.

2969 ### The Identifier Components and the Data Points default matching

2970 By default, the unary Data Set Operators leave the Identifier Components unchanged, so that
2971 the result has the same identifier components as the operand. The operation applies only on
2972 the Measures and no matching between Data Points is needed.

2973 The "n-ary" VTL Data Set Operators compose more than one input Data Sets. A simple
2974 example is: $D_r\ :=\ D_1\ +\ D_2$

2975 These Operators (i.e. the + ) require a preliminary match between the Data Points of the input
2976 Data Sets (i.e. $D_1$ and $D_2$) in order to compose their measures (e.g. summing them) and obtain
2977 the Data Points of the result (i.e. $D_r$).

2978 For example, let us assume that $D_1$ and $D_2$ contain the population and the gross product of the
2979 United States and the European Union respectively and that they have the same Structure
2980 Components, namely the Reference Date and the Measure Name as Identifier Components,
2981 and the Measure Value as Measure Component:

2982        $D_1$ = United States Data

| Ref.Date | Meas.Name | Meas.Value |
|---|---|---|

---

[16] To adhere to the needed constraints, the identification structure of the Data Sets can be manipulated by means
of appropriate VTL Operators, also described in this document.

| | | |
|---|---|---|
| 2013 | Population | 200 |
| 2013 | Gross Prod. | 800 |
| 2014 | Population | 250 |
| 2014 | Gross Prod. | 1000 |

$D_2$ = European Union Data

| Ref.Date | Meas.Name | Meas.Value |
|---|---|---|
| 2013 | Population | 300 |
| 2013 | Gross Prod. | 900 |
| 2014 | Population | 350 |
| 2014 | Gross Prod. | 1000 |

The desired result of the sum is the following:

$D_r$ = United States + European Union

| Ref.Date | Meas.Name | Meas.Value |
|---|---|---|
| 2013 | Population | 500 |
| 2013 | Gross Prod. | 1700 |
| 2014 | Population | 600 |
| 2014 | Gross Prod. | 2000 |

In this operation, the Data Points having the same values for the Identifier Components are matched, then their Measure Components are combined according to the semantics of the specific Operator (in the example the values are summed).

The example above shows what happens under a **strict constraint**: when the input Data Sets have exactly the same Identifier Components. The result will also have the same Identifier Components as the operands.

However, most of Data Set operations (including the sum) are also be possible also under a more **relaxed constraint**, that is when the Identifier Components of one Data Set are a superset of those of the other Data Set.[17]

For example, let us assume that $D_1$ contains the population of the European countries (by reference date and country) and $D_2$ contains the population of the whole Europe (by reference date):

$D_1$ = European Countries

| Ref.Date | Country | Population |
|---|---|---|

---

[17] This corresponds to the "outer join" form of the join expressions, explained in details in the Reference Manual.

| 2012 | U.K. | 60 |
|---|---|---|
| 2012 | Germany | 80 |
| 2013 | U.K. | 62 |
| 2013 | Germany | 81 |

$D_2$ = Europe

| Ref.Date | Population |
|---|---|
| 2012 | 480 |
| 2013 | 500 |

In order to calculate the percentage of the population of each single country on the total of Europe, the Transformation will be:

$$D_r := D_1 \ / \ D_2 \ * \ 100$$

The Data Points will be matched according to the Identifier Component*s* common to $D_1$ and $D_2$ (in this case only the Ref.Date), then the operation will take place.

The result Data Set will have the Identifier Components of both the operands:

$D_r$ = European Countries / Europe * 100

| Ref.Date | Country | Population |
|---|---|---|
| 2013 | U.K. | 12.5 |
| 2013 | Germany | 16.7 |
| 2014 | U.K. | 12.4 |
| 2014 | Germany | 16.2 |

More formally, let $F$ be a generic n-ary VTL Data Set Operator, $D_r$ the result Data Set and $D_i$ $(i=1,… n)$ the input Data Sets, so that:

$$D_r := F(D_1, D_2, … , D_n)$$

The "strict" constraint requires that the Identifier Components of the $D_i$ $_{(i=1,… n)}$ are the same. The result $D_r$ will also have the same Identifier components.

The "relaxed" constraint requires that at least one input Data Set $D_k$ exists such that for each $D_i$ $_{(i=1,… n)}$ the Identifier Components of $D_i$ are a (possibly improper) subset of those of $D_k$. The output Data Set $D_r$ will have the same Identifier Components of $D_k$.

The n-ary Operator $F$ will produce the Data Points of the result by matching the Data Points of the operands that share the same values for the common Identifier Components and by operating on the values of their Measure Components according to its semantics.

### Behaviour for Measure Components

As already mentioned, given $D_r := F(D_1, D_2, … , D_n)$, the input Data Sets $D_i$ $_{(i=1,… n)}$ may have any number of Measure Components. Therefore, to enforce the desired behaviour it is necessary to understand which Measures the Operator is applied to. This Section shows the general VTL

3050 assumptions about how Measure Components are handled, while the behaviour of the single
3051 operators is described in the Part 2.

3052 The simplest case is the **application of unary Operators to mono-measure Data Sets**,
3053 which does not generate ambiguity; in fact, the Operator is intended to be applied to the only
3054 Measure of the input Data Set. The result Data Set will have the same Measure, whose values
3055 are the result of the operation.

3056 For example, let us assume that $D_1$ contains the salary of the employees (the only Identifier is
3057 the Employee ID and the only Measure is the Salary):

3058 $D_1$ = Salary of Employees

| Employee ID | Salary |
|:---:|:---:|
| A | 1000 |
| B | 1200 |
| C | 800 |
| D | 900 |

3065 The Transformation $D_r := D_1 * 1.10$ applies to the only Measure (the salary)
3066 and calculates a new value increased by 10%, so the result will be:

3067 $D_r$ = Increased Salary of Employees

| Employee ID | Salary |
|:---:|:---:|
| A | 1100 |
| B | 1320 |
| C | 880 |
| D | 990 |

3074 In case of **unary Operators applied to a multi-measure Data Set**, the Operator $F$ is by
3075 default intended to be applied separately to all its Measures, unless differently specified. The
3076 result Data Set will have the same Measures as the operand.

3077 For example, given the import and export by reference date:

3078 $D_1$ = Import & Export

| Ref.Date | Import | Export |
|:---:|:---:|:---:|
| 2011 | 1000 | 1200 |
| 2012 | 1300 | 1100 |
| 2013 | 1200 | 1300 |

3083 The Transformation $D_r := D_1 * 0.80$ applies to all the Measures (e.g. to
3084 both the Import and the Export) and calculates their 80%:

3085 $D_r$ = 80% of Import & Export

| Ref.Date | Import | Export |
|----------|--------|--------|
| 2011     | 800    | 960    |
| 2012     | 1040   | 880    |
| 2013     | 960    | 1040   |

If there is the need to **apply an Operator only to specific Measures**, the dot (.) operator can be used, which allows referencing specific Components within a Data Set. The syntax is: *dataset_name.component_name* (for a better description see the corresponding section in the Part 2).

For example, in the Transformation $D_r := D_1.Import * 0.80$

the operation applies only to the Import (and calculates its 80%):

$D_r$ = 80% of the Import, 100% of the Export

| Ref.Date | Import | Export |
|----------|--------|--------|
| 2011     | 800    | 1200   |
| 2012     | 1040   | 1100   |
| 2013     | 960    | 1300   |

Note that in the example above, the Import is kept and left unchanged. In fact, by default all the Measures are kept in the result, even the ones that are not operated on. If there is the need to keep only some Measures, the "keep" clause can be used (see the Part 2).

In case of **n-ary Operators**, by default **the operation is applied on the Measures of the input Data Sets having the same names**, unless differently specified. To avoid ambiguities and possible errors, the input Data Sets are constrained to have the same Measures and the result will have the same Measures too.

For example, let us assume that $D_1$ and $D_2$ contain the births and the deaths of the United States and the European Union respectively.

$D_1$ = Births & Deaths of the United States

| Ref.Date | Births | Deaths |
|----------|--------|--------|
| 2011     | 1000   | 1200   |
| 2012     | 1300   | 1100   |
| 2013     | 1200   | 1300   |

$D_2$ = Birth & Deaths of the European Union

| Ref.Date | Births | Deaths |
|----------|--------|--------|
| 2011     | 1100   | 1000   |
| 2012     | 1200   | 900    |
| 2013     | 1050   | 1100   |

The Transformation $D_r := D_1 + D_2$ will produce:

$D_r$ = Births & Deaths of United States + European Union

| Ref.Date | Births | Deaths |
|----------|--------|--------|
| 2011 | 2100 | 2200 |
| 2012 | 2500 | 2000 |
| 2013 | 2250 | 2400 |

The Births of the first Data Set have been summed with the Births of the second to calculate the Births of the result (and the same for the Deaths).

If there is the need to **apply an Operator on Measures having different names**, the "rename" clause can be used to make their names equal (for a complete description of the clause see the corresponding section in the Part 2).

For example, given these two Data Sets:

$D_1$   (Residents in the United States)

| Ref.Date | Residents |
|----------|-----------|
| 2011 | 1000 |
| 2012 | 1300 |
| 2013 | 1200 |

$D_2$   (Inhabitants of the European Union)

| Ref.Date | Inhabitants |
|----------|-------------|
| 2011 | 1100 |
| 2012 | 1200 |
| 2013 | 1050 |

A Transformation for calculating the population of United States + European Union is:

$D_r := D_1[Residents \rightarrow Population] + D_2[Inhabitants \rightarrow Population]$

The result will be:

$D_r$   (Population of United States + European Union)

| Ref.Date | Population |
|----------|------------|
| 2011 | 2100 |
| 2012 | 2500 |

3155 Note that the number and the names of the Measure Components of the input Data Sets are
3156 assumed to match (following their renaming if needed), otherwise the Expression is
3157 considered in error.

3158 To avoid a potentially excessive renaming, VTL 1.1 additionally allows operations where each
3159 participating dataset has an explicitly specified component using the dot notation. For
3160 instance,

3161      $D_r := D_1.Residents + D_2.Inhabitants$

3162 creates a result with a single measure component named *Result*, which can then be renamed,
3163 if necessary, at will:

3164      $D_r := (D_1.Residents + D_2.Inhabitants)[Result->Population]$

3165 If there is the need to **apply an Operator only to specific Measures**, the dot (.) operator can
3166 be used as in the case of unary Operators. Even in this case, by default all the Measures are
3167 kept in the result, even the ones that are not operated on; if there is the need to keep only
3168 some Measures, the "keep" clause can be used (see the Part 2).

3169 Finally, note that **each Operator may be applied on Measures of certain data types**,
3170 corresponding to its semantics.  For example, *abs* and *round* will require the Measures to be
3171 numeric, while *substr* will require them to be a string. Expressions which violate this
3172 constraint are obviously considered in error.

3173 For example consider the Transformation:      $D_r := abs (D_1)$

3174 As already described, this expression is assumed to apply the *abs* Operator (i.e. absolute
3175 value) to all the Measures Components of $D_1$. If all these Measures are quantitative the
3176 expression is considered correct, on the contrary, if at least one Measure is of an incompatible
3177 data type, the expression is considered in error. The general description of the VTL data types
3178 is given above while the description of the data types on which each operator can be applied
3179 is given in the Part 2.

3180 **Order of execution**

3181 VTL allows the application of many Operators in a single expression. For example:

3182      $Dr := D_1 + D_2 / (D_3 - D_4 / D_5)$

3183 When the order of execution of the Operators is not explicitly defined (through the use of
3184 parenthesis), a default order of execution applies.

3185 In the case above, according to the VTL precedence rules, the order will be:

3186   I.     $D_4 / D_5$       (default precedence order)
3187   II.    $D_3 - I$        (explicitly defined order)
3188  III.    $D_2 / II$       (default precedence order)
3189   IV.   $D_1 + III$      (default precedence order)

3190 The default order of execution depends on the precedence and associativity order of the VTL
3191 Operators and is described in detail in the Part 2.

## Missing Data

The awareness of missing data is very important for correct VTL operations, because the knowledge of the Data Points of the result depends on the knowledge of the Data Points of the operands. For example, assume $D_r := D_1 + D_2$ and suppose that some Data Points of $D_2$ are unknown, it follows that the corresponding Data Points of $D_r$ cannot be calculated and are unknown too.

Missing data can take up two basic forms.

In the first form, **the lack of information is explicitly represented**. This is the case of Data Points that show a "missing" value for some Measure or Attribute Components, which denotes the absence of a true value for a Component. The "missing" value is not allowed for the Identifier Components, in order to ensure that the Data Points are always identifiable.

In the second form, **the lack of information remains implicit**. This is the case of Data Points that are not present at all in the Data Set. For example, given a Data Set containing the reports to an international organization relevant to different countries and different dates, and having as Identifier Components the Country and the Reference Date, this Data Set may lack the Data Points relevant to some dates (for example the future dates) or some countries (for example the countries that didn't send their data) or some combination of dates and countries.

The handling of missing data in VTL dataset operation can be handled in several ways. One way is to require all participating dataset components used in a computation to be known (corresponding to the notion of "inner join" of dataset components). Another way is to allow some, but not all, components from the participating dataset components to be unknown (corresponding to the notion of "outer join" of components). The mechanics of these approaches is explained in details in the section on the joinexpressions and treatment of NULLs in the Reference Manual.

On the basic level, most of the scalar operations (arithmetic, logical, and others) return `null` when any of their arguments is `null`.

The general properties of the `null` are the following ones:

- **Data type:** `null` value belongs to its own type named null. Type null is subsumed by all scalar types, which is to say that `null` value can (in principle) appear wherever a scalar data is expected; this means that it is an allowed value for any scalar type (string, number, boolean, date). However, complex data types (collections, datasets, records, modules, etc.) do not allow `null` values.
- **Testing**. A built-in Boolean operator **is null** can be used to test if a scalar value is `null`.
- **Comparisons**. Whenever a `null` value is involved in a comparison (>, <, >=, <=, in, not in, between) the result of the comparison is `null`.
- **Arithmetic operations**. Whenever a `null` value is involved in a mathematical operation (+, -, *, /, …), the result is `null`.
- **String operations**. In operations on Strings, `null` is considered an empty String ("").
- **Boolean operations**. VTL adopts 3VL (three-value logic). Therefore the following deduction rules are applied:

        TRUE   *or*   null  →    TRUE

        FALSE  *or*   null  →    null

        TRUE   *and*  null  →    null

$$\text{FALSE } and \text{ null } \rightarrow \text{ FALSE}$$

- **Conditional operations**. The `null` is considered equivalent to FALSE; for example in the control structures of the type (*if (p) -then -else*), the action specified in *–then* is executed if the predicate *p* is TRUE, while the action *-else* is executed if the *p* is FALSE or `null`;
- **Filter clauses**. The `null` is considered equivalent to FALSE; for example in the filter clause [*filter p*], the Data Points for which the predicate *p* is TRUE are selected and returned in the output, while the Data Points for which *p* is FALSE or `null` are discarded.
- **Aggregations**. The aggregations (like *sum*, *avg* and so on) return one Data Point in correspondence to a set of Data Points of the input. In these operations the input Data Points having a `null` value are in general not considered. In the average, for example, they are not considered both in the numerator (the sum) and in the denominator (the count). Specific cases for specific operators are described in the respective sections.
- **Implicit zero**. Arithmetic operators assuming implicit zeros (+,-,*,/) may generate `null` values for the Identifier Components in particular cases (superset-subset relation between the set of the involved Identifier Components). Because `null` values are in general forbidden in the Identifiers, the final outcome of an expression must not contain Identifiers having `null` values. As a momentary exception needed to allow some kinds of calculations, Identifiers having `null` values are accepted in the <u>partial results</u>. To avoid runtime error, possible `null` values of the Identifiers have to be fully eliminated in the final outcome of the expression (through a selection, or other operators), so that the operation of "assignment" (:=) does not encounter them.

If a different behaviour is desired for `null` values, it is possible to **override** them. This can be achieved with the combination of the *calc* clauses and *is null* operators.

For example, suppose that in a specific case the `null` values of the Measure Component $M_1$ of the Data Set $D_1$ have to be considered equivalent to the number 1, the following Transformation can be used to multiply the Data Sets $D_1$ and $D_2$, preliminarily converting `null` values of $D_1.M_1$ into the number 1. For detailed explanations of *calc* and *is null* refer to the specific sections in the Part 2.

$$D_r := D_1 \text{[M1 := if M1 is null then 1 else M1] } * D_2$$

### The Attribute Components

Given as usual $D_r := F(D_1, D_2, \dots , D_n)$ and considering that the input Data Sets $D_i$ *(i=1,… n)* may have any number of Attribute Components, there can be the need of calculating the desired Attribute Components of $D_r$. This Section describes the general VTL assumptions about how Attributes are handled (specific cases are dealt with in description of the single operators in the Part 2).

It should be noted that the Attribute Components of a Data Set are dependent variables of the corresponding mathematical function, just like the Measures. In fact, the difference between Attribute and Measure Components lies only in their meaning: it is intended that the Measures give information about the real world and the Attributes about the Data Set itself (or some part of it, for example about one of its measures).

The VTL has different optional behaviours for Attributes and for Measures.

As specified above, Measures are kept in the result by default, whereas Attributes may be assigned a characteristic called "**virality**", which determines if the Attribute is kept in the result by default or not: a "**viral**" Attribute is kept while a "**non-viral**" Attribute is not kept (the virality is applied when no explicit indication about the keeping of the Attribute is provided in the expression, if the virality is not defined, the Attribute is considered as non-viral).

A second aspect is the "virality" of the Attribute in the result. By default, a viral Attribute is considered viral also in the result.

A third aspect is the operation performed on an Attribute. By default, **the operations which apply to the Measures are not applied to the Attributes**, so that the operations on the Attributes need a dedicated specification. If no operations are explicitly defined on a viral Attribute, a default calculation algorithm is applied in order to determine the Attribute's values in the result. If needed, the VTL default behaviour described here may be overridden by customized default behaviours.

As already mentioned, when the default behaviour is not desired, a different behaviour can be specified by means of the proper use of the *keep*, *calc* and *attrcalc* clauses. In particular, through these clauses, it is possible to override the virality (to keep a *non-viral* Attribute or not to keep a *viral* one), to alter the virality of the Attributes in the result (from *viral* to *non-viral* or vice-versa) and to define a specific calculation algorithm for an Attribute (see the detailed description of these clauses in the Part 2).[18]

Hence, the **default Attribute propagation rule** behaves as follows:

- the non-viral Attributes are not kept in the result and their values are not considered;
- the viral Attributes of the operand are kept and are considered viral also in the result; in other words, if an operand has a viral Attribute V, the result will have V as viral Attribute too;
- The Attributes, like the Measures, are combined according to their names, e.g. the Attributes having the same names in multiple Operands are combined, while the Attributes having different names are considered as different Attributes;
- the values of the Attributes which exist and are viral in only one operand are simply copied (obviously, in the case of unary Operators this applies always);
- the Attributes which exist and are viral in multiple operands (i.e. Attributes having the same names) are combined in one Attribute of the result (having the same name also), whose values are calculated according to the default calculation algorithm explained below;

Extending an example already given for unary Operators, let us assume that $D_1$ contains the salary of the employees of a multinational enterprise (the only Identifier is the Employee ID, the only Measure is the Salary, and there are two other Components defined as viral Attributes, namely the Currency and the Scale of the Salary):

---

[18] In particular the *keep* clause allows the specification of whether or not an attribute is kept in the result while the *calc* and the *attrcalc* clauses make it possible to define calculation formulas for specific attributes. The *calc* can be used both for Measures and for Attributes and is a unary Operator, e.g. it may operate on Components of just one Data Set to obtain new Measures / Attributes, while the *attrcalc* is dedicated to the calculation of the Attributes in the N-ary case

3316  $D_1$ = Salary of Employees

| Employee ID | Salary | Currency | Scale |
|---|---|---|---|
| A | 1000 | U.S. $ | Unit |
| B | 1200 | € | Unit |
| C | 800 | yen | Thousands |
| D | 900 | U.K. Pound | Unit |

The Transformation $D_r := D_1 * 1.10$ applies only to the Measure (the salary) and calculates a new value increased by 10%, the viral Attributes are kept and left unchanged, so the result will be:

$D_r$ = Increased Salary of Employees

| Employee ID | Salary | Currency | Scale |
|---|---|---|---|
| A | 1100 | U.S. $ | Unit |
| B | 1320 | € | Unit |
| C | 880 | yen | Thousands |
| D | 990 | U.K. Pound | Unit |

The Currency and the Scale of $D_r$ will be considered viral too and therefore would be kept also in case $D_r$ becomes operand of other Transformations.

For n-ary operations, the VTL **default Attribute calculation algorithm** produces the values of the Attributes of the result Data Set from those of its operands and is applied by default if no operations on the Attributes are explicitly defined. This algorithm is independent of the Operator applied on the Measures and works as follows:

- Whenever in the evaluation of a VTL expression, two data points $P_i$ and $P_j$ are combined as for their Measures, the Attributes having the same name, if viral, are combined as well (non-viral Attributes are ignored)
- It is assumed that each possible value of an Attribute is associated to a **default weight** (in the IM, this is a type of property of the Value Domain which contains the possible values of the Attribute);
- the result of the combination is **the value having the highest weight**;
- if multiple values have the same weight, the result of the combination is the first in lexicographical order.

Note that the default weight for each possible value of an Attribute can be overridden, if desired. However, this is out of the scope of the language: the specific implementations will provide configuration mechanisms (e.g. a user modifiable text file) to alter such values.

For example, let us assume that $D_1$ and $D_2$ contain the births and the deaths of the United States and the Europe respectively, plus a viral Attribute that qualifies if the Value is estimated (having values True or False).

$D_1$ = Births & Deaths of the United States

| Ref.Date | Births | Deaths | Estimate |
|----------|--------|--------|----------|
| 2011 | 1000 | 1200 | False |
| 2012 | 1300 | 1100 | False |
| 2013 | 1200 | 1300 | True |

$D_2$ = Birth & Deaths of the European Union

| Ref.Date | Births | Deaths | Estimate |
|----------|--------|--------|----------|
| 2011 | 1100 | 1000 | False |
| 2012 | 1200 | 900 | True |
| 2013 | 1050 | 1100 | False |

Assuming the weights 1 for "false" and 2 for "true", the Transformation   Dr := D1 + D2 will produce:

$D_r$ = Births & Deaths of United States + European Union

| Ref.Date | Births | Deaths | Estimate |
|----------|--------|--------|----------|
| 2011 | 2100 | 2200 | False |
| 2012 | 2500 | 2000 | True |
| 2013 | 2250 | 2400 | True |

Note also that:

- if the attribute *Estimate* was non-viral in both the input Data Sets, it would not be kept in the result
- if the attribute *Estimate* was viral only in one Data Set, it would be kept in the result with the same values as in the viral Data Set

The VTL default Attribute propagation rule (here called A) ensures the following properties (in respect to the application of a generic VTL operator "§" on the measures):

**Commutative law (1)**

$A(D_1 \S D_2) = A(D_2 \S D_1)$

The application of *A* produces the same result (in term of Attributes) independently of the ordering of the operands. For example, $A(D_1 + D_2) = A(D_2 + D_1)$. This may seem quite intuitive for "sum", but it is important to point out that it holds for every operator, also for non-commutative operations like difference, division, logarithm and so on; for example $A(D_1 / D_2) = A(D_2 / D_1)$

**Associative law (2)**

$A(D_1 \S A(D_2 \S D_3)) = A(A(D_1 \S D_2) \S D_3)$

Within one operator, the result of *A* (in term of Attributes) is independent of the sequence of processing.

3391     **Reflexive law (3)**

3392     $A(\ §(D_1)) = A(D_1)$

3393     The application of $A$ to an Operator having a single operand gives the same result (in
3394     term of Attributes) that its direct application to the operand (in fact the propagation
3395     rule keeps the viral attributes unchanged).

3396 Having these properties in place, it is always possible to avoid ambiguities and circular
3397 dependencies in the determination of the Attributes' values of the result. Moreover, it is
3398 sufficient without loss of generality to consider only the case of binary operators (i.e. having
3399 two Data Sets as operands), as more complex cases can be easily inferred by applying the VTL
3400 Attribute propagation rule recursively (following the order of execution of the operations in
3401 the VTL expression).

3402 With regard to this last aspect, the VTL assumes that the **order of execution** of the operations
3403 in an expression is determined by the precedence and associativity rules of the Operators
3404 applied on the Measures, as already explained in the relevant section. The operations on the
3405 Attributes are performed in the same order, independently of the application of the default
3406 Attribute propagation rule or user defined operations.

3407 For example, recalling the example already given:

3408     $D_r := D_1 + D_2\ /\ (D_3 - D_4\ /\ D_5)$

3409 The evaluation of the Attributes will follow the order of composition of the Measures:

3410   I.    $A(D_4\ /\ D_5)$          (default precedence order)
3411   II.   $A(D_3 - I)$            (explicitly defined order)
3412   III.  $A(D_2\ /\ II)$          (default precedence order)
3413   IV.  $A(D_1 + III)$        (default precedence order)

## 3414 Storage and retrieval of the Data Sets

3415 **The Storage**

3416 As mentioned, the general form of Transformation can be written as follows:

3417                     $D_r\ := F\ (D_1, D_2, \dots , Dn)$

3418 In practice, the right-hand side is a mathematical expression like the one described above:

3419                     $D_r := D_1 + D_2\ /\ (D_3 - D_4\ /\ D_5)$

3420 As already shown, this expression implies the calculation of many Data Sets in different steps:

3421      I.    $(D_4\ /\ D_5)$
3422      II.   $(D_3 - I)$
3423      III.  $(D_2\ /\ II)$
3424      IV.  $(D_1 + III)$

3425 Calculated Data Sets are assumed to be non-persistent (temporary), as well as $D_r$, to which is
3426 assigned the final result of the expression (step IV).

3427 A temporary result within the expression can only be input of other operators in the same
3428 expression.

3429 Parameter $D_r$ , which the result of the whole expression is assigned to, can be directly
3430 referenced as operand by other Transformations of the same VTL program  (a VTL program is
3431 a  set of Transformations, that is a Transformation Scheme, aimed to obtain some meaningful
3432 results for the users, supposed to be executed in the same run).

3433 The *Put* command is used to specify that a result must be persistent. Any step of the
3434 calculation can be made persistent (including all the steps).

3435 The *Put* has two parameters, the first is the (partial) result of the calculation that has to be
3436 made persistent (a non-persistent parameter of *Dataset* type), the second is the reference to
3437 the persistent Data Set, for example:

3438 $$D_r := Put(D_1 + D_2 \ / (D_3 - D_4 / D_5), \text{``PDS1''})$$

3439 means that the overall result of the expression is stored in the persistent Data Set having
3440 name PDS1. The expression:

3441 $$D_r := Put(D_1 + D_2 \ / Put((D_3 - D_4 / D_5), \text{``PDS1''}), \text{``PDS2''})$$

3442 Specifies that $(D_3 - D_4 / D_5)$ is stored in *PDS1* and the overall result in *PDS2*.

3443 ## The Retrieval

3444 Considering again the general form of Transformation:

3445 $$D_r \ := F \ (D_1, D_2, \dots , Dn)$$

3446 the "n" Data Sets  $D_{i \ (i=1,\dots n)}$  are the operands of the Expression and their values have to be
3447 retrieved.

3448 The generic $D_i$ may be retrieved either as the temporary result of another Transformation (of
3449 the same VTL program) or from a persistent data source. In the former case $D_i$ is the name of
3450 the left-hand parameter ($D_r$) of the other Transformation. In the latter, $D_i$ is the reference to a
3451 persistent Data Set (see the following sections).

3452 A specific Operator (Get) ensures powerful features for accessing persistent data (see the
3453 detail in the Part 2). A direct reference to a persistent Data Set is equivalent to the application
3454 of the Get command.

3455 The Operators Get and Put are also called "commands" because they allow the interaction
3456 with the persistent storage.

3457 ## The references to persistent Data Sets

3458 In defining the Transformations, persistent Data Sets can be retrieved or stored by means of
3459 the Get and Put commands respectively.

3460 As described in the VTL IM, the Data Set is considered as an artefact at a logical level,
3461 equivalent to a mathematical function having independent variables (Identifiers) and
3462 dependent variables (Measures and Attributes). A Data Set is a set of Data Points, which are
3463 the occurrences of the function. Each Data Point is an association between a combination of
3464 values of the independent variables and the corresponding values of the dependent variables.

3465 Therefore, the VTL references the conceptual/logical Data Sets and does not reference the
3466 physical objects where the Data Points are stored. The link between the Data Set at a logical
3467 level and the corresponding physical objects is out of the scope of the VTL and left to the
3468 implementations.

3469 Also the versioning of the artefacts of the information model, including the Data Sets, is out of
3470 the scope of the VTL and left to the implementations.

3471 The VTL allows reference through commands (Get and Put) to any persistent Data Set defined
3472 and identified according the VTL IM. For correct operation, knowledge of the Data Structure of
3473 the input Data Sets is essential, in order to check the correctness of the expression and
3474 determine the Data Structure of the result. For this reason, the VTL requires that at
3475 compilation time the Data Structures of the referenced Data Sets are available.

3476 In addition, to simplify some kind of operations, the VTL makes it possible to reference also
3477 Cartesian subsets of the already defined Data Sets (i.e. sub Data Sets specified as Cartesian
3478 products of Value Domain Subsets of some Identifier Components).

3479 This is consistent with the IM, because any subset of the Data Points of a Data Set may be
3480 considered in its turn a Data Set, and with correct VTL operations, because the Data Structure
3481 of a sub Data Set is deducible from the Data Structure of the original Data Set, once that the
3482 specification of the subset is given.

3483 Note however that it is not possible to reference directly a non-Cartesian sub Data Set (i.e. a
3484 sub Data Set that cannot be obtained as a Cartesian product of Value Domain Subsets). As any
3485 other kind of Data Set, however, non-Cartesian subsets can be obtained through an
3486 Expression, as partial or final results.

3487 For example, in case of unit data, given the Data Set "Legal Entity" having as Identifiers of the
3488 Country, the IssuerOrganization, and the LegalEntityNumber, the VTL allows direct reference
3489 to either the whole Data Set or a sub-Data Set obtained specifying some countries, and/or
3490 issuers, and/or numbers. By specifying a single value for each identifier it is possible to
3491 reference even a single Legal Entity (i.e. a single Data Point).

3492 In case of Dimensional Data Sets, assuming that the Country and the Date are the Identifiers, it
3493 is possible to reference the sub Data Sets corresponding to one or some countries, to one or
3494 some dates, and to a combination of them. If the dates are periodical, the sub Data Set
3495 corresponding to one country is a time-series. The sub Data Set corresponding to a certain
3496 date is a cross-section. The sub Data Set corresponding to one country and one date is a single
3497 Data Point. Therefore, VTL allows direct reference to dimensional data, time-series, cross-
3498 sections, and single observations.

3499 In conclusion, a VTL reference to a persistent (sub)Data Set is composed of two parts:

3500 • The identification of the Data Set (mandatory)
3501 • The specification of a subset of it (optional)

### The Identification of a persistent Data Set

3503 The identification of the persistent Data Sets to read from (Get) or to store into (Put) follows
3504 the general rules of identification of the persistent artefact (see the corresponding section
3505 above).

3506 Therefore, the Data Set identifier is the **Data Set Name**, which is unique in the environment.
3507 As different environments can use the same Data Set Names for their artefacts, the Data Set
3508 Name can optionally be qualified by a **Namespace**, to avoid name conflicts.

3509 In case the Data Set identifier has a Namespace, a separator character can be chosen (and
3510 configured in the system) among the non-alphanumeric ones. A typical, and recommended,

3511 choice is the slash ("/") symbol. If the Data Set identifier does not have a Namespace, the same
3512 namespace as the respective Transformation is assumed.

3513 Examples of good references to Data Sets are:

3514 "NAMESPACE/DS_NAME"    (explicit Namespace definition)

3515 "DS_NAME"                        (the Namespace of the Transformation is assumed)

### The specification of a subset of a persistent Data Set

3517 The VTL allows the retrieval or storage of a subset of a predefined Data Set by filtering the
3518 values of its Identifier Components.

3519 Two basic options are allowed in the grammar of this specification:

- 3520 A **full notation (query string), specifying both the Identifiers and the values** to be
  3521 filtered (e.g. Date= 2014, Country=USA, Sector=Public …); in this case the filtering
  3522 condition is preceded by the "?" symbol.
- 3523 A **short notation (ordered concatenation), specifying only the values** to be filtered
  3524 (e.g. 2014.USA.Public); in this case the filtering condition is preceded by the "/"
  3525 symbol; the values have to be specified following a predefined order of the Identifiers.

3526 The **query string** is a postfix syntax specifying the filter in case the order of the identifiers is
3527 not defined beforehand or not known.

3528 The filter is specified by concatenating the filtering conditions on the Identifiers, expressed in
3529 any order and separated by "&". If a filtering condition is not specified for an Identifier, the
3530 latter is not constrained and all the available values are taken.  For example:

3531   I.    `DS_NAME?DATE=2014&COUNTRY=USA&SECTOR=PUBLIC`

3532 In the example above, **single values** are specified for each filtering condition.

3533 In the same way, it is also possible to specify **multiple values** for some filtering conditions,
3534 separating the values by the "**+**" keyword (list). For example, to take the years 2013 and 2014
3535 and the countries USA and Canada:

3536  II.    `DS_NAME?DATE=2013+2014&COUNTRY=USA+CANADA&SECTOR=PUBLIC`

3537 Finally, where the Values have an order like the one for the "Date" data type, it is possible to
3538 specify ranges of values for some filtering conditions, separating the first and last values of
3539 the range by the "**-**" keyword (range). For example, to take all the years from 2008 to 2014:

3540 III.    `DS_NAME?DATE=2008-2014&COUNTRY=USA+CANADA&SECTOR=PUBLIC`

3541 The **ordered concatenation** is a simplified syntax to specify the filter in case the order of the
3542 identifiers is defined beforehand and known.

3543 The filter is specified by concatenating the filtering conditions in the predefined order of the
3544 Identifiers; the filtering conditions do not require the specification of the name of the
3545 Identifier, which can be deduced by their predefined order, therefore only the values are
3546 specified, separated by "**.**", i.e. a dot. If a value is omitted, the corresponding Identifier is not
3547 constrained and all the available values are taken.  For example, (assuming that the order on
3548 the identifiers is 1-Date, 2-Country, 3-Sector):

3549   I.    `DS_NAME/2014.USA.PUBLIC`

3550        This definition in the query string syntax corresponds to:

3551        `DS_NAME?DATE=2014&COUNTRY=USA&SECTOR=PUBLIC`

3552   II.   `DS_NAME/.USA.PUBLIC`

3553        This definition filters all the available years for the USA and the public sector, and
3554        in the query string syntax corresponds to:

3555        `DS_NAME?COUNTRY=USA&SECTOR=PUBLIC`

3556   III.   `DS_NAME/..PUBLIC`

3557        This definition filters all the available years and countries for the public sector and
3558        in the query string syntax corresponds to:

3559        `DS_NAME?SECTOR=PUBLIC`

3560 If needed, the list ("+") and/or range ("-") keywords can be used to specify lists or range of
3561 values respectively. For example:

3562   IV.   `DS_NAME/2008-2014.USA+CANADA.PUBLIC`

3563        This definition in the query string syntax corresponds to:

3564        `DS_NAME?DATE=2008-2014&COUNTRY=USA+CANADA&SECTOR=PUBLIC`

3565

# 3566 Conventions for the grammar of the language

### 3567 General conventions

3568 A VTL program is a set of Transformations executed in the same run, which is defined as a
3569 Transformation Scheme.

3570 Each Transformation consists in a *statement* that is an assignment of the form:

3571        `variable parameter := expression`

3572 ":=" is the assignment operator, meaning that the result of the evaluation of the *expression* in
3573 the right-hand side is assigned to the *variable parameter* in the left-hand side (which is the
3574 output parameter of the assignment).

3575 Examples of assignments are (assuming that ds_i *(i=1…n)* are Data Sets):

3576   •   `ds_1 := ds_2`
3577   •   `ds_3 := ds_4 + ds_6`

### 3578 Variable Parameter names

3579 The variable parameters are non-persistent (temporary).

3580 The names of the variable parameters are alphanumeric (starting with an alphabetic
3581 character). Also non alphabetic characters ("_","-") are allowed, but not in the first position.
3582 Parameter names are case-sensitive.

3583 Examples of allowed names for the parameters are: par1, p_1, VarPar_ABCD, paraMeterXY.

### 3584 Reserved words

3585 Certain words are reserved **keywords** in the language and cannot be used as parameter
3586 names, they include:

| 3587 | - all the names of the operators / clauses |
|---|---|
| 3588 | - all the symbols used by the language (assignment ":=", parenthesis "(",")","[" ,"]", |
| 3589 | ampersand "&", hash "#" ...) |
| 3590 | - true |
| 3591 | - false |
| 3592 | - all |
| 3593 | - imbalance |
| 3594 | - errorlevel |
| 3595 | - condition |
| 3596 | - msg_code |
| 3597 | - dataset |
| 3598 | - script |

### Comments

3600 VTL allows comments within the statements in order to provide textual explanations of the
3601 operations. Whatever is enclosed between /* and */ shall not be processed by VTL parsers, as
3602 it shall be considered as comment.

3603 For example:

```
3604   /* Set constant for 'π'*/
3605   numpi := 3.14
3606   popA := populationDS + 1 /* Assign temp Dataset popA */
```

### Constraints and errors

3608 VTL supports a number of error types, which can occur in different situations; errors are
3609 divided into three main categories **compile time, runtime, validation.** Each category is
3610 divided in turn in subcategories, containing the specific errors.

3611 An error is identified by the string "VTL-" followed by a four digit code CSEE, where:

3612      - C identifies the category (0: compile time, 1: runtime, 2: validation)
3613      - S identifies the subcategory
3614      - EE identifies the specific error in the subcategory

3615 While the three categories (and subcategories for compile errors) are standardized with
3616 codes reported in the remainder of this section, an encoding for specific errors (identified by
3617 the last two digits, EE) is not enforced here and can be independently defined by the adopting
3618 organization.[19]

3619 A compile time error prevents an expression from being used (exchanged, executed ...) and
3620 results in an exception reporting the error code (VTL-0XXX) and the wrong expression to the
3621 definer.

3622 In contrast, when a runtime error is raised, it can cause:

3623      a) an abnormal termination of the running VTL program, with an exception reporting the
3624         error code (VTL-1XXX) and the wrong expression to the user
3625      b) the current expression to be discarded, without generating any exception

---

[19] However, notice that in a following version of the language, a standardization is foreseen also for subcategories and specific error codes.

3626    c)  only the violating Data Point to be discarded, without generating any exception.

3627    The choice between these three behaviours should be dependent on the runtime system and
3628    is not part of the language, nor linked to the error codes.

3629    Validation errors are errors resulting from data validation (e.g. *check* operator), which can be
3630    stored in Datasets and used for further elaboration. Indeed, validation errors are not VTL
3631    errors and do not influence the use of the expression or the normal execution of a VTL
3632    program.

## Compile Time errors (VTL-0xxx)

3634    The VTL grammar specifies the rules to be followed in writing expressions. The VTL language
3635    allows the detection at compile time of the possible violation of the **correct syntax**, the use of
3636    **wrong types** as parameters for the operators or the **violation of any of the static
3637    constraints of the operators** (with respect to the rules described in the Part 2).

3638    A VTL compiler has to be able to detect all the syntax errors, help the user understand the
3639    reason and recover. Three subcategories are predetermined (see below). The specific error
3640    can be represented by the adopting organization with any code ranging from 00 to 99
3641    (examples are: unclosed literal string; unexpected symbol, etc.)

### Syntax errors (VTL-01xx)

3643    A violation of the VTL syntax with respect to the syntax templates of operators in names of
3644    operators or number of operands.

3645    Examples of syntactically invalid expressions are:

3646    `R := C1 +`                        –the second operand is missing

3647    `R := C1 exist_in_all C2`          -  the correct syntax is "exists_in_all".

3648    `R := if k1>4 then else K3 + 3`    – the "then" operand is missing

### Type errors (VTL-02xx)

3650    A violation of the types of the operands allowed for the operators.

3651    Examples of expressions that are type-invalid are:

3652    `R := C1 + '2'`     –   if C1 has a measure component that is not <String>

3653    `R := C1 + C2`      –   if  C1  has  a  MeasureComponent<String>  and  C2  has  a
3654    MeasureComponent<Numeric>

3655    `R := C1 / 5`       -   if C1 has a MeasureComponent<String>.

3656    `R:= if (K1 > 3 and k1 < 5) then 0 else "hello"`   – the "then" and the "else"
3657    operands must be of the same type

3658    Since the language is strongly typed, all type violations can be reported at compile time.

### Static constraint violation errors (VTL-03xx)

3660    Every operator may have additional constraints. They are reported in the respective
3661    "Constraints" sections in the Part 2. Some of them are static, in the sense that they can be
3662    checked at compile type.

3663    A constraint violation error is the violation of a static VTL constraint .

3664    Examples of expressions that violate static constraints are:

3665    `R := C1 + C2`        – if the IdentifierComponents of C1 and C2 are not the same or
3666    are not contained in the ones of the other operator.

3667    `R := 3 + 5`        – in the plus (+) operator, at least one operand must be a Dataset.

3668

### Runtime errors (VTL-1xxx)

3670    These errors can be detected only at runtime, typically because they are generated by the
3671    data.

3672    Examples are the classical mathematical constraints on operators arguments (negative or
3673    zero logarithm argument, division by zero, etc.).

3674    Particular types of runtime errors are:

3675    • presence of **duplicate** Data Points to be assigned to a Data Set (it is not allowed that
3676       two Data Points in a Data Set have the same values for all the Identifier Components
3677       because the Data Point identification would be impossible)
3678    • presence of a `null` **value** in an Identifier Component of a Data Point.

3679    These two errors result in a runtime exception only if the inconsistent Data Points are
3680    assigned (:=) to a Data Set in the left-hand side of a Transformation or are stored in a
3681    persistent Data Set. In other words, if such Data Points are only partial and temporary results
3682    inside the expression on the right-hand side, no runtime exceptions will be raised provided
3683    that the anomalies (duplications or NULLS) are removed before the execution of the
3684    assignment or the Put command.

3685    Examples of expressions generating runtime errors are:

3686    `R := C1 / C2`        – where C2 is 0 for any observation

3687    `R := substr(A, 2, 5)`   – if A is 1 character long, causing an "out of range"

3688    `R := C1`        – if C1 contains `null` values for some IdentifierComponents.
3689    Notice that the assignment causes the runtime error; the fact that C1 contains a `null` value for
3690    an IdentifierComponent is accepted as partial and temporary result in the right-hand side of
3691    the expression.

3692    `R := C1`        – if C1 contains duplicates on an IdentifierComponent. Also in this
3693    case, notice that the assignment causes the runtime error; the fact that C1 contains a duplicate
3694    is accepted as partial and temporary result in the right-hand side of the expression.

3695    A VTL runtime environment will be able to detect a wide number of runtime errors. The
3696    specific errors can be divided into subcategories by the adopting organization; moreover, the
3697    specific error can be represented with any code ranging from 00 to 99.

3698

### Validation errors (VTL-2xxx)

3700    They represent the outcome of a failed user-defined validation. The code can be used for
3701    further elaboration or to report discrepancies.

3702 Error codes can be associated with the single validations with the *check* operator, whose last
3703 parameter is *errorCode*. This is the code to be used for each Data Point having FALSE for its
3704 MeasureComponent.

3705 For example:

3706 `R := check(C1 >= C2, all, 2601)`

3707 Checks if C1 is greater or equal than C2 and, if not the case, stores the code 2601 in the
3708 *errorCode* attribute.

3709

| C1 | | |
|----|----|----|
| K1 | K2 | M1 |
| 1 | A | 1000 |
| 2 | B | 200 |

3714

| C2 | | | |
|----|----|----|----|
| K1 | K2 | K3 | M1 |
| 1 | A | X | 1000 |
| 2 | B | Y | 350 |
| 2 | B | Z | 150 |

3720 and produces:

3721

| R | | | | |
|----|----|----|-----------|-----------|
| K1 | K2 | K3 | CONDITION | ERRORCODE |
| 1 | A | X | TRUE | |
| 2 | B | Y | FALSE | 2601 |
| 2 | B | Z | TRUE | |

3728 A set of VTL validation rules, will be able to detect a wide number of validation errors. The
3729 specific errors can be divided into subcategories by the adopting organization; moreover, the
3730 specific error can be represented with any code ranging from 00 to 99.

# Governance, other requirements and future work

The SDMX Technical Working Group, as mandated by the SDMX Secretariat, is responsible for ensuring the technical maintenance of the Validation and Transformation Language through a dedicated VTL task-force. The VTL task-force is open to the participation of experts from other standardisation communities, such as DDI and GSIM, as the language is designed to be usable within different standards.

## The governance of the extensions

According to the requirements, it is envisaged that the language can be enriched and made more powerful in future versions according to the evolution of the business needs. For example, new operators and clauses can be added, and the language syntax can be upgraded.

The VTL governance body will take care of the evolution process, collecting and prioritising the requirements, planning and designing the improvements, releasing future VTL versions.

The release of new VTL versions is considered as the preferred method of fulfilling the requirements of the user communities. In this way the possibility of exchanging standard validation and transformation rules would be preserved to the maximum extent possible.

In order to fulfil specific calculation features not yet supported, the VTL provides for a specific operator (Evaluate) whose purpose is to invoke an external calculation function (routine), provided that this is compatible with the VTL IM and data types.

The operator "Evaluate" (also "Eval") allows defining and making customized calculations (also reusing existing routines) without upgrading or extending the language, because the external calculation function is not considered as an additional operator. The expressions containing Eval are standard VTL expressions and can be parsed through a standard parser. For this reason, when it is not possible or convenient to use other VTL operators, Eval is the recommended method of customizing the language operations.

However, as explained in the section "Extensibility and Customizability" of the "General Characteristics of VTL" above, calling external functions has some drawbacks in respect to the use of the proper VTL operators. The transformation rules would be not understandable unless such external functions are properly documented and shared and could become dependent on the IT implementation, less abstract and less user oriented. Moreover, the external functions cannot be parsed (as if they were built through VTL operators) and this could make the expressions more error-prone. External routines should be used only for specific needs and in limited cases, whereas widespread and generic needs should be fulfilled through the operators of the language.

While the "Eval" operator is part of VTL, the invoked external calculation functions are not. Therefore, they are considered as customized parts under the governance, and are responsibility and charge of the organizations which use it.

Another possible form of customization is the extension of VTL by means of non-standard operators/clauses. This kind of extension is deprecated, because it would compromise the possibility of sharing validation rules and using common tools (for example, a standard parser would consider an expression containing non-standard operators as in error).

3771 Organizations possibly extending VTL through non-standard operators/clauses would
3772 operate on their own total risk and responsibility, also for any possible maintenance activity
3773 deriving from VTL modifications.

## Relations with the GSIM Information Model

3775 As explained in the section "VTL Information Model", VTL 1.0 is inspired by GSIM 1.1 as much
3776 as possible, in order to provide a formal model at business level against which other
3777 information models can be mapped, and to facilitate the implementation of VTL with
3778 standards like SDMX, DDI and possibly others.

3779 GSIM faces many aspects that are out of the VTL scope; the latter uses only those GSIM
3780 artefacts which are strictly related to the representation of validations and transformations.
3781 The referenced GSIM artefacts have been assessed against the requirements for VTL and, in
3782 some cases, adapted or improved as necessary, as explained earlier. No assessment was made
3783 about those GSIM artefacts which are out of the VTL scope.

3784 In respect to GSIM, VTL considers both unit and dimensional data as mathematical functions
3785 having a certain structure in term of independent and dependent variables. This leads to a
3786 simplification, as unit and dimensional data can be managed in the same way, but it also
3787 introduces some slight differences in data representation. The aim of the VTL Task Force is to
3788 propose the adoption of this adjustment for the next GSIM versions.

3789 The VTL IM allows defining the Value Domains (as in GSIM) and their subsets (not explicitly
3790 envisaged in GSIM), needed for validation purposes. In order to be compliant, the GSIM
3791 artefacts are used for modelling the Value Domains and a similar structure is used for
3792 modelling their subsets. Even in this case, the VTL task force will propose the explicit
3793 introduction of the Value Domain Subsets in future GSIM versions.

3794 VTL is based on a model for defining mathematical expressions which is called
3795 "Transformation model". GSIM does not have a Transformation model, which is however
3796 available in the SDMX IM.  The VTL IM has been built on the SDMX Transformation model,
3797 with the intention of suggesting its introduction in future GSIM versions.

3798 Some misunderstanding may arise from the fact that GSIM, DDI, SDMX and other standards
3799 also have a Business Process model. The connection between the Transformation model and
3800 the Business Process model has been neither analysed nor modelled in VTL 1.0. One reason is
3801 that the business process models available in GSIM, DDI and SDMX are not yet fully
3802 compatible and univocally mapped.

3803 It is worth nothing that the Transformation and the Business Process models address
3804 different matters. In fact, the former allows defining validation and calculation rules in the
3805 form of mathematical expressions (like in a spreadsheet) while the latter allows defining a
3806 business process, made of tasks to be executed in a certain order.  The two models may
3807 coexist and be used together as complementary. For example, a certain task of a business
3808 process (say the validation of a data set) may require the execution of a certain set of
3809 validation rules, expressed through the Transformation model used in VTL. Further progress
3810 in this reconciliation is a task which needs some parallel work in GSIM, SDMX and DDI, and
3811 could be reflected in a future VTL version.

# Annex 1 - EBNF

The VTL language is also expressed in EBNF (Extended Backus-Naur Form).

EBNF is a standard[20] meta-syntax notation, typically used to describe a Context-Free grammar and represents an extension to BNF (Backus-Naur Form) syntax. Indeed, any language described with BNF notation can also be expressed in EBNF (although expressions are typically lengthier).

Intuitively, the EBNF consists of terminal symbols and non-terminal production rules. Terminal symbols are the alphanumeric characters (but also punctuation marks, whitespace, etc.) that are allowed singularly or in a combined fashion. Production rules are the rules governing how terminal symbols can be combined in order to produce words of the language (i.e. legal sequences).

More details can be found at http://en.wikipedia.org/wiki/Extended_Backus–Naur_Form

## Properties of VTL grammar

VTL can be described in terms of a Context-Free grammar[21], with productions of the form $V \rightarrow w$, where $V$ is a single non-terminal symbol and $w$ is a string of terminal and non-terminal symbols.

VTL grammar aims at being unambiguous. An ambiguous Context-Free grammar is such that there exists a string that can be derived with two different paths of production rules, technically with two different leftmost derivations.

In theoretical computer science, the problem of understanding if a grammar is ambiguous is undecidable. In practice, many languages adopt a number of strategies to cope with ambiguities. This is the approach followed in VTL as well. Examples are the presence of *associativity* and *precedence* rules for infix operators (such as addition and subtraction), and the existence of compulsory *else* branch in *if-then-else* operator.

These devices are reasonably good to guarantee the absence of ambiguity in VTL grammar. Indeed, real parser generators (for instance YACC[22]), can effectively exploit them, in particular using the mentioned associativity and precedence constrains as well as the relative ordering of the productions in the grammar itself, which solves ambiguity by default.

---

[20] ISO/IEC 14977

[21] http://en.wikipedia.org/wiki/Context-free_grammar

[22] http://en.wikipedia.org/wiki/Yacc