

SDMX Technical Working Group

VTL Task Force

Validation & Transformation Language

Part 1 - General Description

Version 1.0

February 2015

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

32 The SDMX Technical Working Group is pleased to present the version 1.0 of the Validation
33 and Transformation Language, in short VTL.

34
35 The work on VTL was launched at the end of 2012 by the SDMX Secretariat. SDMX already has
36 a package for transformations and expressions which is present in the information model,
37 although a specific language does not yet exist. To make this framework fully operational, a
38 standard “language” for defining validation and transformation rules (set of operators, their
39 syntax and semantics) should be adopted, appropriate IT formats for exchanging such rules
40 and related metadata should be introduced, and the web services to store and retrieve them
41 should be designed.

42
43 A task force was put in place, composed of members of SDMX, DDI and GSIM communities and
44 the work started in summer 2013. The intention was to provide a language which is usable by
45 statisticians to express logical validation rules and transformations on data, whether
46 described as dimensional tables or as unit-record data. The assumption is that this logical
47 formalization of validation and transformation rules would be converted into specific
48 programming languages for execution (SAS, R, Java, SQL, etc.) but would provide a “neutral”
49 expression at business level of the processing taking place, against which various
50 implementations can be mapped. Experience with existing examples suggests that this goal
51 would be attainable.

52
53 An important point that emerged is that several standards are interested in such a language.
54 However, each standard operates on its model artefacts and produces artefacts within the
55 same model (property of closure). To cope with this, VTL has been built upon a very basic
56 information model, taking the common parts of GSIM, SDMX and DDI, mainly using artefacts
57 from GSIM 1.1, somewhat simplified and with some additional detail. This way the existing
58 standards (SDMX, DDI, others) may adopt VTL by mapping their information model against
59 the VTL one. Therefore, although a work-product of SDMX, the VTL language will be usable
60 also with other standards.

61
62 The VTL 1.0 package includes:

- 63 a) Part 1, highlighting the main characteristics of VTL, its core assumptions and the
64 information model the language is based on;
- 65 b) Part 2, containing the full library of operators ordered by category, including examples; this
66 first version can support validation and basic compilation needs. Future versions will include
67 more features related to transformation of data.
- 68 c) BNF notation (Backus-Naur Form) which is the technical notation to be used as a test bed
69 for all the examples throughout the document.

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

72

73 The latest version of the VTL is freely available online at www.sdmx.org.

74

75 **Acknowledgements**

76 This publication has been prepared thanks to the collective input of experts from Bank of
77 Italy, Bank for International Settlements (BIS), European Central Bank (ECB), Eurostat, ILO,
78 ISTAT-Italy, OECD, Statistics Netherlands, and UNESCO. Other experts from the SDMX
79 Technical Working Group, the SDMX Statistical Working Group and the DDI initiative were
80 also consulted and participated in reviewing the documentation.

81 The list of contributors and reviewers includes the following experts: Sami Airo, Foteini
82 Andrikopoulou, David Barraclough, Luigi Bellomarini, Marc Bouffard, Maurizio Capaccioli,
83 Vincenzo Del Vecchio, Fabio Di Giovanni, Jens Dossé, Heinrich Ehrmann, Bryan Fitzpatrick,
84 Arofan Gregory, Edgardo Greising, Angelo Linardi, Chris Nelson, Stratos Nikoloutsos, Marco
85 Pellegrino, Michele Romanelli, Juan Alberto Sanchez, Angel Simon Delgado, Daniel Suranyi,
86 Olav ten Bosch, Laura Vignola, Nikolaos Zisimos.

87 Feedback and suggestions for improvement are encouraged and can be sent to the SDMX
88 Technical Working Group (twg@sdmx.org).

89 Table of contents

90		
91	FOREWORD	3
92	TABLE OF CONTENTS	5
93	INTRODUCTION	6
94	STRUCTURE OF THE DOCUMENT	7
95	GENERAL CHARACTERISTICS OF THE VTL	8
96	USER ORIENTATION	8
97	INTEGRATED APPROACH	9
98	ACTIVE ROLE FOR PROCESSING.....	10
99	INDEPENDENCE OF IT IMPLEMENTATION	11
100	EXTENSIBILITY, CUSTOMIZABILITY	12
101	LANGUAGE EFFECTIVENESS.....	13
102	VTL INFORMATION MODEL	15
103	GENERIC MODEL FOR DATA AND THEIR STRUCTURES	15
104	GENERIC MODEL FOR VARIABLES AND VALUE DOMAINS.....	21
105	GENERIC MODEL FOR TRANSFORMATIONS.....	23
106	PERSISTENCY AND IDENTIFICATION OF THE ARTEFACTS OF THE MODEL	27
107	VTL CORE ASSUMPTIONS	29
108	THE TYPES OF OPERANDS AND RESULTS	29
109	THE OPERATIONS ON THE DATA SETS	33
110	STORAGE AND RETRIEVAL OF THE DATA SETS	47
111	CONVENTIONS FOR THE GRAMMAR OF THE LANGUAGE	51
112	GOVERNANCE, OTHER REQUIREMENTS AND FUTURE WORK	57
113	RELATIONS WITH THE GSIM INFORMATION MODEL.....	58
114	FUTURE DIRECTIONS	59
115	ANNEX 1 – EBNF	61
116	PROPERTIES OF VTL GRAMMAR	61
117		

119 This document presents the Validation and Transformation Language (aka VTL).

120 The purpose of the VTL is to allow a formal and standard definition of algorithms to validate
121 statistical data and calculate derived data.

122 The VTL development is organized in a first phase aimed to allow the formalisation of the data
123 validation algorithms and in following phases aimed to tackle more complex algorithms for
124 data compilation. In fact, the assessment of business cases showed that the majority of the
125 institutions ascribes a higher priority to a standard language for supporting the validation
126 processes and in particular to the possibility of sharing validation rules with the respective
127 data providers, in order to specify the quality requirements and allow validation also before
128 provision.

129 This document is the outcome of the first phase and therefore presents a first version of the
130 VTL primarily oriented to support the data validation. However, because the features needed
131 for the validation include simple calculations, this first version of the VTL can also support
132 basic compilation needs. In general, validation is assumed to be a particular case of
133 transformation; therefore, the term "Transformation" is meant to be more general and to
134 include validation as well.

135 The main categories of operators included in this version of the VTL syntax are:

136	General	(e.g. assignment, data access, data storage ...)
137	String	(e.g. substring, concatenation, length ...)
138	Mathematical	(e.g. +, -, *, /, round, absolute value ...)
139	Boolean	(e.g. and, or, not ...)
140	Relational	(e.g. selection, union, intersection, merge ...)
141	Statistical	(e.g. minimum, maximum, aggregation ...)
142	Validation	(e.g. of value domains, references, figures ...)
143	Conditional	(e.g. if-then-else ...)

144

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

150 Accordingly, the working group for the VTL development includes representatives of
151 institutions involved in the DDI and GSIM initiatives and there has been agreement on the
152 objective of adopting a common language, applicable to SDMX as well as to DDI and GSIM, in
153 the hope of avoiding the risk of having diverging variants.

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

158

159 Structure of the document

160 The first part of the document is dedicated to the description of the general characteristics of
161 the VTL.

162 The following part describes the Information Model on which the language is based. In
163 particular, it describes the model of the data artefacts that the language is aimed to validate
164 and transform, the model of the variables and value domains used in the data artefacts and
165 the model of the transformations.

166 A third part clarifies some general features of the language (i.e. the core assumptions of the
167 VTL), such as the types of artefacts involved in the transformations, the general rules for the
168 operations on the data sets, the methods for referencing the data sets to be operated on, and
169 the general conventions for the grammar of the language.

170 A final part highlights some issues related to the governance of VTL developments and to
171 future work, following a number of comments, suggestions and other requirements which
172 were submitted to the task-force in order to enhance the current VTL 1.0 package.

173 A short annex gives some background information about the BNF (Backus-Naur Form) syntax
174 which has been used for providing a context-free representation of VTL. The Extended BNF
175 (EBNF) representation is part of the VTL 1.0 package available at www.sdmx.org.

176

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

182 User orientation

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

186 ○ The language is based on a “user” perspective and a “user” information model
187 (IM) and not on possible IT perspectives (and IMs)

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

193 ○ It references IM objects and does not use direct references to IT objects

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

196 ○ Designed according to mathematics, which is a universal knowledge;

197 ○ Expressed in English to be shareable in all countries;

198 ○ As simple, intuitive and self-explanatory as possible;

199 ○ Based on common mathematical expressions, which involve “operands”
200 operated on by “operators” to obtain a certain result;

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

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

206 ○ Operators for data validations and edit;

207 ○ Operators for aggregation, including according to hierarchies;

208 ○ Operators for dimensional processing (e.g. projection, filter);

209 ○ At a later stage, operators for time series processing (e.g. time shift, change of
210 periodicity, moving average, seasonal adjustment, correlation) operators for
211 statistics (e.g. aggregation, mean, median, percentiles, variance, indexes,
212 correlation, sampling, inference, estimation);

213 Integrated approach

- 214 ⇒ The language is independent of the statistical domain of the data to be processed;
- 215 ○ VTL has no dependencies on the subject matter (the data content);
- 216 ○ VTL is able to manipulate statistical data in relation to their structure.
- 217 ⇒ The language is suitable for the various typologies of data of a statistical environment
218 (for example dimensional data, survey data, registers data, micro and macro,
219 quantitative and qualitative) and is supported by an information model (IM) which
220 covers these typologies;
- 221 ○ The IM allows the representation of the various typologies of data of a
222 statistical environment at a conceptual/logical level (in a way abstract from IT
223 and from the physical storage);
- 224 ○ The various typologies of data are described as much as possible in an
225 integrated way, by means of common IM artefacts for their common aspects;
- 226 ○ The principle of the Occam's razor is applied as an heuristic principle in
227 designing the conceptual IM, so keeping everything as simple as possible or, in
228 other words, unifying the model of apparently different things as much as
229 possible.
- 230 ⇒ The language (and its IM) is independent of the phases of the statistical process and
231 usable in any one of them;
- 232 ○ Operators are designed to be independent of the phases of the process, their
233 syntax does not change in different phases and is not bound to some
234 characteristic restricted to a specific phase (operators' syntax is not aware of
235 the phase of the process);
- 236 ○ In principle, all operators are allowed in any phase of the process (e.g. it is
237 possible to use the operators for data validation not only in the data collection
238 but also, for example, in data compilation for validating the result of a
239 compilation process; similarly it is possible to use the operators for data
240 calculation, like the aggregation, not only in data compilation but also in data
241 validation processes);
- 242 ○ Both collected and calculated data are equally permitted as inputs of a
243 calculation, without changes in the syntax of the operators/expression;
- 244 ○ Collected and calculated data are represented (in the IM) in a homogeneous
245 way with regards to the metadata needed for calculations.
- 246 ⇒ The language is designed to be applied not only to SDMX but also to other standards;
- 247 ○ VTL, like any consistent language, relies on a specific information model, as it
248 operates on the VTL IM artefacts to produce other VTL IM artefacts. In
249 principle, a language cannot be applied as-is to another information model
250 (e.g. SDMX, DDI, GSIM); this possibility exists only if there is a unambiguous
251 correspondence between the artefacts of those information models and the
252 VTL IM (that is if their artefacts correspond to the same mathematical notion);
- 253 ○ The goal of applying the language to more models/standards is achieved by
254 using a very simple, generic and conceptual Information Model (the VTL IM),

255 and mapping this IM to the models of the different standards (SDMX, DDI,
256 GSIM, ...); to the extent that the mapping is straightforward and unambiguous,
257 the language can be inherited by other standards (with the proper
258 adjustments);

- 259 ○ To achieve an unambiguous mapping, the VTL IM is deeply inspired by the
260 GSIM IM and uses the same artefacts when possible¹; in fact, GSIM is designed
261 to provide a formal description of data at business level against which other
262 information models can be mapped; moreover, loose mappings between GSIM
263 and SDMX and between GSIM and DDI are already available²; a very small
264 subset of the GSIM artefacts is used in the VTL IM in order to keep the model
265 and the language as simple as possible (Occam's razor principle); these are the
266 artefacts strictly needed for describing the data involved in Transformations,
267 their structure and the variables and value domains;
- 268 ○ GSIM artefacts are supplemented when needed, with other artefacts that are
269 necessary for describing calculations; in particular, the SDMX model for
270 Transformations is used;
- 271 ○ As mentioned above, the definition of the VTL IM artefacts is based on
272 mathematics and is expressed at an abstract user level.

273 Active role for processing

274 ⇒ The language is designed to possibly drive in an active way the execution of the
275 calculations (in addition to documenting them)

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

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

- 283 ○ To be automatically parsed (against the rules of the formal grammar); on the
284 IT level, this requires the implementation of a parser that compiles the
285 expressions and checks their correctness;
- 286 ○ To be automatically translated from the VTL to the language of the IT tool to
287 be used for the calculation; on the IT level, this requires the implementation of
288 a proper translator;
- 289 ○ To be automatically translated from one VTL version to another, e.g. following
290 an upgrade of the VTL syntax; on the IT level, this requires the implementation
291 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>;

- 292 ⇒ The inputs and the outputs of the calculations and the calculations themselves are
293 artefacts of the IM
- 294 ○ This is a basic property of any robust language because it allows calculated
295 data to be operands of further calculations;
 - 296 ○ If the artefacts are persistently stored, their definition is persistent as well; if
297 the artefacts are non-persistently stored (used only during the calculation
298 process like input from other systems, intermediate results, external outputs)
299 their definition can be non-persistent;
 - 300 ○ Because the definition of a calculation needs the data structure definition of its
301 input artefacts, the latter must be available when the calculation is defined;
 - 302 ○ The VTL is designed to make the data structure of the output of a calculation
303 deducible from the calculation algorithm and from the data structure of the
304 operands (this feature ensures that the calculated data can be defined
305 according to the IM and can be used as operands of further calculations);
 - 306 ○ In the IT implementation, it is advisable to automate (as much as possible) the
307 structural definition of the output of a calculation, in order to enforce the
308 consistency of the definitions and avoid unnecessary overheads for the
309 definers.
- 310 ⇒ The VTL and its information model make it possible to check automatically the overall
311 consistency of the definition of the calculations, including with respect to the artefact
312 of the IM, and in particular to check:
- 313 ○ the correctness of the expressions with respect to the syntax of the language
 - 314 ○ the integrity of the expressions with respect to their input and output artefacts
315 and the corresponding structures and properties (for example, the input
316 artefacts must exist, their structure components referenced in the expression
317 must exist, qualitative data cannot be manipulated through quantitative
318 operators, and so on)
 - 319 ○ the consistency of the overall graph of the calculations (for example, there
320 should not be cycles in the sequence of calculations in order to avoid that the
321 result of a calculation goes as input to the same calculation, so producing
322 unpredictable and erroneous results);

323 Independence of IT implementation

- 324 ⇒ According to the “user orientation” above, the language is designed so that users are
325 not required to be aware of the IT solution;
- 326 ○ To use the language, the users need to know only the abstract view of the data
327 and calculations and do not need to know the aspects of the IT
328 implementation, like the storage structures, the calculation tools and so on.
- 329 ⇒ The language is not oriented to a specific IT implementation and permits many
330 possible different implementations (this property is particularly important in order to
331 allow different institutions to rely on different IT environments and solutions);

- 332 ○ On the technical level, the connection between the user layer and the IT layer
333 is left to the specific IT implementations;
- 334 ○ The VTL approach favours effective IT implementations that decouple the user
335 layer and the IT layer.
- 336 ⇒ The language does not require the awareness of the physical data structure; the
337 operations on the data are specified according to the conceptual/logical structure,
338 and so are independent of the physical structure; this ensures that the physical
339 structure may change without necessarily affecting the conceptual structure and the
340 user expressions;
- 341 ○ Data having the same conceptual/logical structure may be accessed using the
342 same statements, even if they have different IT structures;
- 343 ○ The VTL provides for commands for data store and retrieve at a
344 conceptual/logical level; the mapping and the conversion between the
345 conceptual and the physical structures of the data is left to the IT
346 implementation (and users need not be aware of it);
- 347 ○ By mapping the user and the IT data structures, the IT implementations can
348 make it possible to store/retrieve data in/from different IT data stores (e.g.
349 relational databases, dimensional databases, xml files, spread-sheets,
350 traditional files);
- 351 ⇒ The language does not require the awareness of the IT tools used for the calculations
352 (e.g. routines in a programming language, statistical packages like R, SAS, Matlab,
353 relational databases (SQL), dimensional databases (MDX), XML tools,...);
- 354 ○ The syntax of the VTL is independent of existing IT calculation tools;
- 355 ○ On the IT level, this may require a translation from the VTL to the language of
356 the IT tool to be used for the calculation;
- 357 ○ By implementing the proper translations at the IT level, institutions can use
358 different IT tools to execute the same algorithms; moreover, it is possible for
359 the same institution to use different IT tools within an integrated solution (e.g.
360 to exploit different abilities of different tools);
- 361 ○ VTL instructions do not change if the IT solution changes (for example
362 following the adoption of another IT tool), so avoiding impacts on users as
363 much as possible;

364 Extensibility, customizability

- 365 ⇒ It is possible to build and extend the language gradually, enriching the available
366 operators according to the evolution of the business needs, so progressively making
367 the language more powerful;
- 368 ⇒ In addition, it is possible to call external routines of other languages/tools, provided
369 that they are compatible with the IM; this requisite is aimed to fulfil specific
370 calculation needs without modifying the operators of the language, so exploiting the
371 power of the other languages/tools if necessary for specific purposes

- 372 ○ The external routines should be compatible with, and relate back to, the
373 conceptual IM of the calculations as for its inputs and outputs, so that the
374 integrity of the definitions is ensured
- 375 ○ The external routines are not part of the language, so their use might be
376 subject to some limitations (e.g. it might be impossible to parse them as if they
377 were operators of the language)
- 378 ○ The use of external routines has some drawbacks, because it may obviously
379 compromise the IT implementation independence, the abstraction and the
380 user orientation; therefore external routines should be used only for specific
381 needs and in limited cases, whereas widespread and generic needs should be
382 fulfilled through the operators of the language;
- 383 ⇒ Nothing can prevent the Organizations adopting the VTL from extending it by defining
384 customized parts, on their own total responsibility and charge, in order to improve the
385 standard language for their specific purposes (e.g. for supporting possible algorithms
386 not permitted by the standard part); also the customized parts must be compliant with
387 the VTL IM and the VTL core assumptions (adopting Organizations are totally in
388 charge of any possible maintenance activity deriving from VTL modifications); such
389 extensions however are not recommended because they can compromise the
390 exchange of validation rules and the use of common tools.

391 Language effectiveness

- 392 ⇒ The language is oriented to give full support to the various typologies of data of a
393 statistical environment (for example dimensional data, survey data, registers data,
394 micro and macro, quantitative and qualitative, ...) described as much as possible in a
395 coherent way, by means of common IM artefacts for their common aspects, and
396 relying on mathematical notions, as mentioned above. The various types of statistical
397 data are considered as mathematical functions, having independent variables
398 (Identifiers) and dependent variables (Measures, Attributes³), whose extensions can
399 be thought as logical tables (DataSets) made of rows (Data Points) and columns
400 (Identifiers, Measures, Attributes).
- 401 ⇒ The language supports operations on the Data Sets (i.e. mathematical functions) in
402 order to calculate new Data Sets from the existing ones, on the structure components
403 of the Data Sets (Identifiers, Measures, Attributes), on the Data Points.
- 404 ⇒ The algorithms are specified by means of mathematical expressions which compose
405 the operands (Data Sets, Components ...) by means of operators (e.g. +, -, *, /, >, <) to
406 obtain a certain result (Data Sets, Components ...);
- 407 ⇒ The validation is considered as a kind of calculation having as an operand the Data Set
408 to be validated and producing a Data Set containing the outcome of the validation
409 (typically having values “true” and “false” in the measure, respectively for successful
410 and unsuccessful validation); being a Data Set, the result of the validation can be
411 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.

- 412 ⇒ Calculations on multiple measures are supported, as well as calculations on the
413 attributes of the Data Sets and calculations involving missing values;
- 414 ⇒ The operations are intended to be consistent with the historical changes of the
415 artefacts (e.g. of the code lists, of the hierarchies ...), so allowing a proper behaviour
416 for each reference period; the support to this aspect is left to the standards adopting
417 the VTL (e.g. SDMX, DDI ...) because different standards may represent historical
418 changes in different ways;
- 419 ⇒ The language is ready to allow different algorithms for different reference times
420 (feature to be implemented at a later stage);
- 421 ⇒ the VTL operators are generally “modular”, meaning that it is possible to compose
422 multiple operators in a single expression; in other words, an operator can have an
423 expression as operand, so obtaining a new expression, and this can be made
424 recursively;
- 425 ⇒ The final and the intermediate results of a calculation can be permanently stored (or
426 not) according to the needs;
- 427 ⇒ Multiple results may be calculated by means of multiple expressions.
- 428

430 Generic Model for Data and their structures

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

433 The purpose is to provide a formal description of data at business level against which other
434 information models (IMs) can be mapped, to facilitate the implementation of VTL with other
435 standards like SDMX, DDI and possibly others. This is the same purpose as the Generic
436 Statistical Information Model (GSIM) and, consequently, this formal model uses the GSIM
437 artefacts as much as possible (GSIM 1.1 version)⁴. Besides, GSIM already provides a first
438 mapping with SDMX and DDI that can be used for the technical implementation⁵. Note that the
439 description of the GSIM 1.1 classes and relevant definitions can be consulted in the “Clickable
440 GSIM” of the UNECE site⁶.

441 Some slight differences between this model and GSIM are due to the fact that in the VTL IM
442 both unit and dimensional data are considered as mathematical functions having independent
443 and dependent variables and are treated in the same way.

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

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

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

459 Starting from this assumption, each mathematical function (logical table) may be defined as a
460 GSIM Data Set and its structure as a GSIM Data Structure, having Identifier, Measure and

⁴ See also the section “Relations with the GSIM Information model”

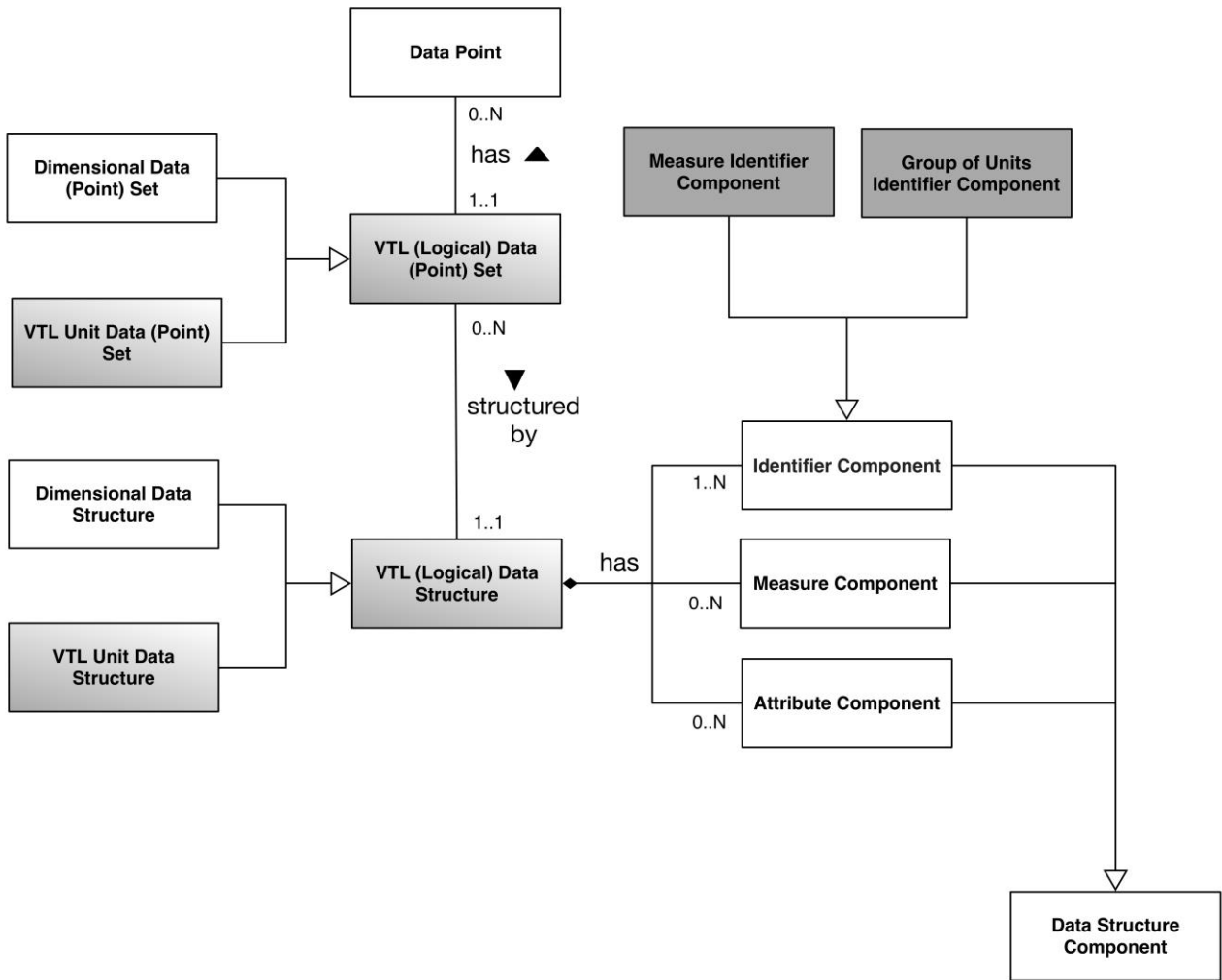
⁵ For the GSIM – DDI and GSIM – SDMX mappings, see also the relationships between GSIM and other standards at the UNECE site <http://www1.unece.org/stat/platform/display/gsim/GSIM+and+standards>. About the mapping with SDMX, however, note that here it is assumed that the SDMX artefacts Data Set and Data 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>”

461 Attribute Components. The Identifier components are the independent variables of the
 462 function, the Measures and Attribute Components are the dependent variables. Obviously the
 463 GSIM artefacts “Data Set” and “Data Set Structure” have to be strictly interpreted as **logical**
 464 **artefacts** on a mathematical level, not necessarily corresponding to physical data sets and
 465 physical data structures.

466 As earlier pointed out, in respect to GSIM this assumption leads to a representation that is
 467 identical for the dimensional data and very similar for the unit data, as described below. The
 468 same names as in GSIM are used for the Artefacts, the “VTL” prefix is applied to the Artefact
 469 that are very similar to the GSIM ones but not exactly corresponding.

470 **ER diagram - Data model**



471

472

473 White box: same artefact as in GSIM 1.1
 474 Light grey box: similar to GSIM 1.1
 475 Dark grey box: additional detail (in respect to GSIM 1.1)
 476

477 **Explanation of the Diagram**

478 **VTL (Logical) Data (Point) Set:** a mathematical function (logical table) that describes some
479 properties of some groups of units of a population. In general, the groups of units may be
480 composed of one or more units. For unit data, each group is composed of a single unit. For
481 dimensional data, each group may be composed of any number of units. A VTL Data Set is
482 considered as a logical set of observations (Data Points) having the same structure and the
483 same general meaning, independently of the possible physical representation or storage. This
484 artefact is similar to the “Data Set” in GSIM. In particular, the GSIM Data Set may be a GSIM
485 Dimensional Data Set or a GSIM Unit Data, while the VTL Data Set may be:

486 **Dimensional Data (Point) Set:** a kind of (Logical) Data Set describing groups of units
487 of a population that may be composed of many units. This artefact is the same as the
488 GSIM Dimensional Data Set.

489 **VTL Unit Data (Point) Set:** a kind of (Logical) Data Set describing single units of a
490 population. This is similar to GSIM because the VTL Unit Data Set is the same as the
491 Unit Data Record in GSIM, which has its own structure and can be thought of as a
492 mathematical function. The difference is that the VTL Unit Data Set takes the place of
493 the GSIM Unit Data Set, which is omitted because it cannot be considered as a
494 mathematical function: in fact it can have many GSIM Unit Data Records with different
495 structures.

496 **Data Point:** a single value of the function, i.e. a single association between the values of the
497 independent variables and the values of the dependent variables. A Data Point corresponds to
498 a row of the table that describes the function. This artefact is the same as the GSIM Data Point.

499 **VTL (Logical) Data Structure:** the structure of a mathematical function, having independent
500 and dependent variables. The independent variables are called “Identifier components”, the
501 dependent variables are called either “Measure Components” or “Attribute Components”. The
502 distinction between Measure and Attribute components is based on their meaning: the
503 Measure Components give information about the real world, while the Attribute components
504 give information about the function itself. This artefact is similar to the Data Structure in
505 GSIM. In particular, the GSIM Data Structure may be a Dimensional Data Structure or a Unit
506 Data Structure, while the VTL Data Structure may be:

507 **Dimensional Data Structure:** the structure of (0..n) Dimensional Data Sets. This
508 artefact is the same as in GSIM.

509 **VTL Unit Data Structure:** the structure of (0..n) Unit Data Sets. This is similar to GSIM
510 because the VTL Unit Data Structure is the same as the Logical Record in GSIM, which
511 corresponds to a single structure. The difference is that the VTL Unit Data Structure
512 takes the place of the GSIM Unit Data Structure, which is omitted because it cannot be
513 considered as the structure of a mathematical function: in fact it can have many Logical
514 Records with different structures.

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

517 **Identifier Component** (or simply Identifier): a component of the data structure that is
518 an independent variable of the function. This artefact is the same as in GSIM. On the
519 other hand, the following distinction is a detail that does not exist in GSIM, needed to

520 distinguish proper Identifier Components and possible Identifiers Components used in
521 some cases to identify the measures:

522 **Group of Units Identifier Component:** a “proper” Identifier Component that
523 contributes to identify the groups of units (composed of either single or many
524 units) that the function describes.

525 **Measure Identifier Component:** an Identifier Component that contributes to
526 identify the measures of the function when more measures are conveyed
527 through the same Measure Component. This artefact corresponds to the SDMX
528 Measure Dimension.

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

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

535 Examples

536 As a first simple example, let us consider the following table:

537 *Production of the American Countries*

<i>Ref.Date</i>	<i>Country</i>	<i>Meas.Name</i>	<i>Meas.Value</i>	<i>Status</i>
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
...

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

- 551 • Reference Date (Identifier Component)
- 552 • Country (Identifier Component)
- 553 • Measure Name (Measure Identifier Component)
- 554 • Measure Value (Measure Component)
- 555 • Status (Attribute Component)

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

558

Institutional Unit Data

559

560

561

562

563

564

565

566

567

568

569

570

<i>Row Type</i>	<i>I.U. ID</i>	<i>Ref.Date</i>	<i>I.U. Name</i>	<i>I.U. Sector</i>	<i>Assets</i>	<i>Liabilities</i>
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
...

571

572

573

574

This table is not equivalent as a whole 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 structures (corresponding to the Row Types I and II), so that the original table can be split in the following ones:

575

Row Type I - Institutional Unit register

576

577

578

579

580

<i>I.U. ID</i>	<i>I.U. Name</i>	<i>I.U. Sector</i>
A	AAAAA	Private
B	BBBBB	Public
C	CCCCC	Private
...

581

Row Type II - Institutional Unit Assets and Liabilities

582

583

584

585

586

587

588

589

<i>I.U. ID</i>	<i>Ref.Date</i>	<i>Assets</i>	<i>Liabilities</i>
A	2013	1000	800
A	2014	1050	750
B	2013	1200	900
B	2014	1300	950
C	2013	750	900
C	2014	800	850
...

590

591

Each one of these two tables corresponds to a mathematical function and can be represented like in the first example above.

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

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

595 Data Structure 1:

- 596 • I.U. ID (Identifier Component)
- 597 • I.U. Name (Measure Component)
- 598 • I.U. Sector (Measure Component)

599

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

601 Data Structure 2:

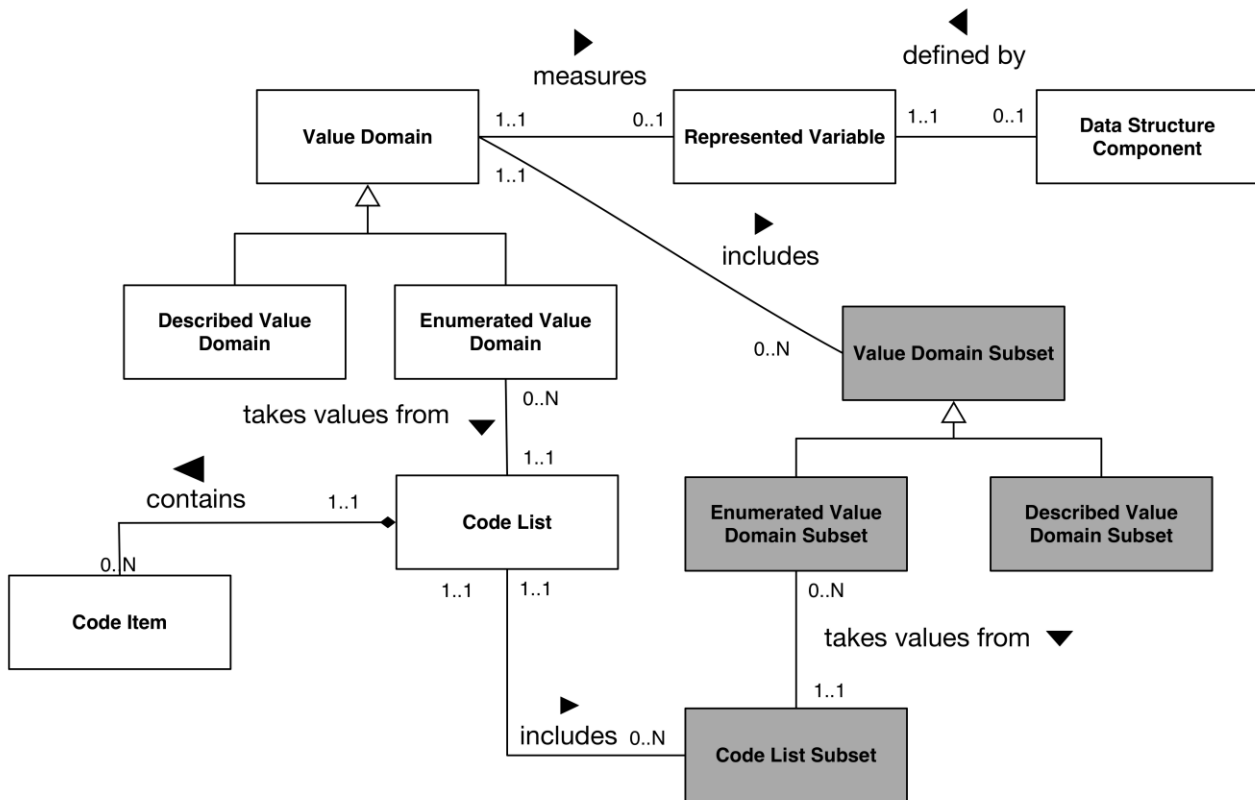
- 602 • I.U. ID (Identifier Component)
- 603 • Reference Date (Identifier Component)
- 604 • Assets (Measure Component)
- 605 • Liabilities (Measure Component)

606

607

608 Generic Model for Variables and Value Domains

609 ER diagram – Variable and Value domain model



610

611 Whitebox: same as in GSIM 1.1

612 Dark grey box: additional detail (in respect to GSIM 1.1)

613

614 Explanation of the Diagram

615 **Data Structure Component:** see the explanation already given above, in the data model
616 section.

617 **Represented Variable:** a characteristic of a statistical population (e.g. the country of birth)
618 represented in a specific way (e.g. through the ISO code). This artefact is the same as in GSIM.

619 **Value Domain:** the domain of the allowed values for a variable. This artefact is the same as in
620 GSIM. An important characteristic of the Value Domain is the data type (e.g. String, Numeric,
621 Integer, Boolean, Date), which is the type that any Value of the Value Domain must
622 correspond to.

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

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

627 **Code List:** a list of allowed codes (values) of an Enumerated Value Domain, with associated
628 categories (e.g. the list of the ISO codes of the countries, each one associated with the name of
629 the country). This artefact is the same as in GSIM.

630 *The following artefacts are aimed to represent possible subsets of the GSIM Value Domains and*
631 *Code Lists. This is needed for validation purposes, because very often not all the values of the*
632 *Value Domain are allowed, but only a subset of them (e.g. not all the countries but only the*
633 *European countries). Although this detail does not exist in GSIM, these artefacts are fully*
634 *compliant with the GSIM artefacts described above, representing Domains and Code Lists:*

635 **Value Domain Subset:** a subset of the domain of the allowed values for a variable. This
636 artefact does not exist in GSIM, however it is compliant with the GSIM Value Domain. A
637 Value Domain Subset has the same data type as its Value Domain.

638 **Described Value Domain Subset:** a described (defined by a criterion) subset of
639 a Value Domain (e.g. the countries having more than 100 million inhabitants,
640 the integers between 1 and 100). This artefact does not exist in GSIM, however
641 it is compliant with the GSIM Described Value Domain.

642 **Enumerated Value Domain Subset:** an enumerated subset of a Value Domain
643 (e.g. the enumeration of the European countries). This artefact does not exist in
644 GSIM, however it is compliant with the GSIM Enumerated Value Domain.

645 **Code List Subset:** the list of the codes of an Enumerated Value Domain Subset (e.g. the
646 list of the ISO codes of the European countries). This artefact does not exist in GSIM,
647 however is consistent with the GSIM Code List. The Code List Subset enumerates only
648 the codes and does not associate the categories (e.g. the names of the countries),
649 because the latter are already maintained in the Code List artefact (which contains all
650 the possible codes with the associated categories).

651

652 Generic Model for Transformations

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

655 A transformation is assumed to be an algorithm to produce a new model artefact (typically a
656 Data Set) starting from existing ones. It is also assumed that the data validation is a particular
657 case of transformation, therefore the term “transformation” is meant to be more general and
658 to include the validation case as well.

659 This model is essentially derived from the SDMX IM⁷, as DDI and GSIM do not have an explicit
660 transformation model at the moment⁸. In its turn, the SDMX model for Transformations is
661 similar in scope and content to the Expression metamodel that is part of the Common
662 Warehouse Metamodel (CWM)⁹ developed by the Object Management Group (OMG).

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

668 The basic brick of this model is the notion of a Transformation.

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

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

675 The general form of a Transformation is the following:

676 `variable parameter := expression`

677 “:=” is the assignment operator, meaning that the result of the evaluation of *expression* in the
678 right-hand side is assigned to the *variable parameter* in the left-hand side, which is the a-
679 priori unknown output of *expression* (typically a Data Set).

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

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

⁷ The SDMX specification can be found at http://sdmx.org/wp-content/uploads/2011/08/SDMX_2-1-1_SECTION_2_InformationModel_201108.pdf (see package 13 - “Transformations and Expressions”).

⁸ The Transformation model described here is not a model of the processes, like the ones that both SDMX and GSIM have. 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.

⁹ This specification can be found at <http://www.omg.org/cwm>.

684 In this example the measure values of the Data Set D_r is calculated as the sum of the measure
685 values of the Data Sets D_1 and D_2 .

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

688 $D_r \quad := \quad \text{If } (D_1 = D_2, \text{ then “true”, else “false”})$

689 In this case, the Data Set D_r would have a Boolean measure containing the value “true” if the
690 validation is successful and “false” if it is unsuccessful.

691 These are only fictitious examples for explanation purposes. The general rules for the
692 composition of Data Sets (e.g. rules for matching their Data Points, for composing their
693 measures ...) are described in the sections below, while the actual Operators of the VTL are
694 described in the Part 2.

695 The *expression* in the right-hand side of a Transformation must be written according to a
696 formal language, which specifies the list of allowed operators (e.g. sum, product ...), their
697 syntax and semantics, and the rules for composing the expression (e.g. the default order of
698 execution of the operators, the use of parenthesis to enforce a certain order ...). The Operators
699 of the language have Parameters¹⁰, which are the a-priori unknown inputs and output of the
700 operation, characterized by a given role (e.g. dividend, divisor or quotient in a division).

701 Note that this generic model does not specify the language to be used. As a matter of fact, not
702 only the VTL but also other languages might be compliant with this specification, provided
703 that they manipulate and produce artefacts of the information model described above.
704 However the VTL has been agreed as the standard language to define and exchange validation
705 and transformation rules among different organizations.

706 Also note that this generic model does not actually specify the operators to be used in the
707 language. Therefore, the VTL may evolve and may be enriched and extended.

708 In the practical use of the language, Transformations can be composed one with another to
709 obtain the desired outcomes. In particular, the result of a Transformation can be an operand
710 of other Transformations, in order to define a sequence of calculations as complex as needed.

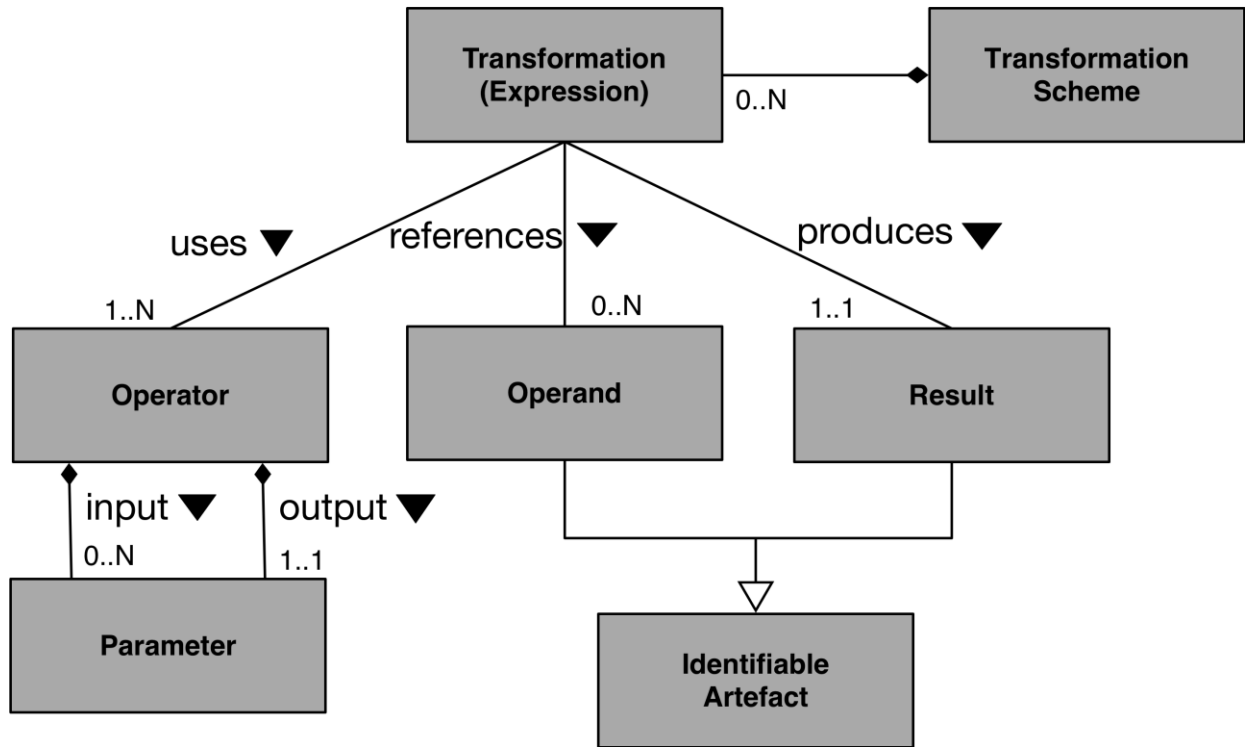
711 Moreover, the Transformations can be grouped into Transformations Schemes, which are sets
712 of transformations meaningful to the users. For example a Transformation Scheme can be the
713 set of transformations needed to obtain some specific meaningful results, like the validations
714 of one or more Data Sets.

715 A set of Transformations takes the structure of a graph, whose nodes are the model artefacts
716 (usually Data Sets) and whose arcs are the links between the operands and the results of the
717 single Transformations. This graph is directed because the links are directed from the
718 operands to the results and is acyclic because it should not contain cycles (like in the spread-
719 sheets), otherwise the result of the Transformations might become unpredictable.

720

¹⁰ The term is used with the same meaning of “argument”, like usual in computer science.

721 **ER Diagram - Transformations**



722

723

724 White box: same as in GSIM 1.1

725 Dark grey box: additional detail (in respect to GSIM 1.1)

726 (All these artefacts match the SDMX artefact having the same name; however the identifiable artefacts are
727 intended to be the ones of the VTL model)

728

729 **Explanation of the diagram**

730 **Transformation:** the basic element of the calculations, which consists in a statement which
731 assigns the outcome of the evaluation of an Expression to an Identifiable Artefact of the
732 Information model; the Transformation artefact is the same as in SDMX;

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

739 **Transformation Scheme:** a set of Transformations aimed to obtain some meaningful results
740 for the user (like the validation of one or more Data Sets); the Transformation Scheme may be
741 also considered as a VTL program; this artefact is the same as in SDMX;

742 **Operator:** the specification of an operation to be performed on some Operands (e.g. +, -, *, /);
743 this artefact is the same as in SDMX;

744 **Parameter:** a-priori unknown input or output of an Operator, having a definite role in the
745 operation (e.g. dividend, divisor or quotient for the division) and corresponding to a certain
746 type of artefact (e.g. a “Data Set”, a “Data Structure Component” ...), the Parameter artefact is
747 the same as in SDMX;

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

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

752 **Identifiable Artefact:** an Identifiable Artefact of the VTL information model (e.g. a Data Set, a
753 Data Structure Component); this artefact is the same as in SDMX;

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

762 Example

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

767 $D_r \quad := \quad \text{If } ((D_1 + D_2) = D_3, \text{ then “true”, else “false”})$

768 In this case:

769 The Transformation may be called “Consistency check between stocks and flows” and is
770 formally defined through the statement above.

- 771
- 772 • D_r is the Result
 - 773 • D_1, D_2 and D_3 are the Operands
 - 774 • $\text{If } ((D_1 + D_2) = D_3, \text{ then “true”, else “false”})$ is the Expression
 - 775 • “:=”, “If”, “+”, “=” are the Operators

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

- 777 • input parameters of “+”: two numeric Data Sets (to be summed)
- 778 • output parameters of “+”: a numeric Data Sets (resulting from the sum)
- 779 • input parameters of “=”: two Data Sets (to be compared)
- 780 • output parameter of “=”: a Data Set (resulting from the comparison)
- 781 • input parameters of “If”: an Expression defining a condition, i.e. $(D_1+D_2)=D_3$
- 782 • output parameter of “If”: a Data Set (as resulting from the “then”, “else” clauses)

783 Persistency and Identification of the artefacts of the model

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

786 A persistent artefact exists externally independently of a VTL program, while a non-persistent
787 artefact exists only within a VTL program.

788 The VTL grammar provides for the identification of the non-persistent artefacts (see the
789 section about the conventions for the grammar of the language) and leaves the accurate
790 definition of the identification mechanism of the persistent artefacts to the standards
791 adopting the VTL (e.g. SDMX, DDI ...) ¹¹.

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

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

800 Therefore, VTL foresees:

- 801 • the **Name** of the artefact (a generic string), which is unique in the environment of the
802 definer;
- 803 • an optional **Namespace** (generic string beginning with an alphabetic character) which
804 is a supplementary qualifier that identifies the environment in which the artefact
805 Name is assumed to be unique, to avoid name conflicts.

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

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

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

816 As intuitive, the Namespace may begin with the name of the institution (“maintenance
817 agency” in SDMX terms). Assuming the dot (“.”) as separator character between environments
818 and sub-environments, examples of possible Namespaces are:

- 819 • ESCB.analysis&insight
- 820 • EuropeanStatisticalSystem.validation
- 821 • OECD.Stat

¹¹ Different standards may have different identification mechanisms.

- 822 • Unesco
- 823 • Bancaditalia.dissemination.public

824

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

829

830 VTL core assumptions

831 The Validation and Transformation Language is based on two parts: the **core assumptions**
832 and the **standard library of Operators**. The former specifies the general behaviour of the
833 language, and is by default stable. The latter contains the standard set of Operators of the
834 language, and can be gradually enriched following the evolution of the user needs. Possible
835 new operators must obviously comply with the core assumptions.

836 The core assumptions include:

- 837 • The types of Operands and Results
- 838 • The operations on the Data Sets
- 839 • Storage and retrieval of the Data Sets
- 840 • The conventions for the grammar of the language

841 The core assumptions are explained in the following sections. The standard library of
842 operators is described in the Part 2.

843 The Types of Operands and Results

844 The Data types of the VTL

845 The VTL assumes that operands and results belong to a data type, which influences the
846 operations that can be applied on the data.

847 The instances of the various data types (i.e. the real objects of those types) are called *literals*.

848 The **basic data types** of the language are five: *String*, *Numeric*, *Integer*, *Boolean* and *Date*. They
849 are described in the following table.

Basic data types	
String	A sequence of one or more characters enclosed in double quotes ("). Examples of allowed literals for this data type are: "hello", "test", "x", "this is a 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.
Numeric	Fixed and floating point numbers, up to 38 digits of precision. At least the following numbers should be representable as numeric in implementations: <ul style="list-style-type: none">- Positive numbers in the range 1×10^{-130} to $9.99\dots9 \times 10^{125}$ with up to 38 significant digits.- Negative numbers from -1×10^{-130} to $9.99\dots99 \times 10^{125}$ with up to 38 significant digits.- Zero (0).- Positive (+Inf) and negative infinity (-Inf). 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.

Integer	The basic signed integer type. At least 16 bit in size, although the actual size may vary by implementation. Examples of allowed literals for this type are: 2, 5, 7, 24, -14, 0.
Boolean	The Boolean data type. The allowed literals are <i>true</i> and <i>false</i> .
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'.

850

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

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

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

856 Also the Data Set has a data type, which is a “composite” one and corresponds to the set of the
857 data types of its Structure Components.

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

859 In conclusion, a data type can be assigned to any artefact of the VTL model (either a basic or a
860 composite data type).

861 **The Parameters of the VTL Operators**

862 As already mentioned, a Parameter is a generic input or output of an Operator and has a
863 definite role in the operation (e.g. dividend, divisor or quotient for the division).

864 A Parameter corresponds either to a certain type of artefacts of the information model or to
865 some kind of constant value (for the sake of simplicity, the constants have not been described
866 in the IM).

867 The parameters corresponding to a type of artefacts of the IM are called variable parameters,
868 because their values are not known beforehand (i.e. when the Expression is written and
869 compiled) and can be considered as the “language variables”. The types of variable parameter
870 are the Data Set type, the Structure Component type (hereinafter simply Component), the
871 Value Domain Subset type and, possibly, other IM artefacts.

872 The parameters corresponding to constant values are called constant parameters, because
873 their values are known beforehand (they are written directly in the expressions).

- 874 The instances of the various types of Parameters (i.e. the real objects of those types, both
 875 variable and constants) are named *literals* (like the instances of the simple data types above).
 876 The following table contains the main types of variable parameters.

Types of variable Parameters	
Dataset<T>	<p>A Data Set, having the composite data type T, which corresponds to the set of the data types of its Structure Components.</p> <p>A Data Set may be persistent or not. A persistent Data Set is permanently stored, i.e. maintained in a storage media and therefore exists also independently of a VTL program. A temporary Data Set is not stored and exists only within a VTL program.</p> <p>These sub-types of Datasets are specified by writing:</p> <ul style="list-style-type: none"> • PersistentDataset<T> • TemporaryDataset<T>
Component<T>	<p>A Structure Component having the data type T.</p> <p>A Structure component has the role of Identifier, Measure or Attribute Component, this role can be specified by writing:</p> <ul style="list-style-type: none"> • IdentifierComponent<T> • MeasureComponent<T> • AttributeComponent<T> <p>Structure Components can be classified according to their Data Type as:</p> <ul style="list-style-type: none"> • String: Component<String> • Numeric: Component<Numeric> • Integer: Component<Integer> • Boolean: Component<Boolean> • Date: Component<Date> <p>Allowed literals are the names of the Structure Components of the Data Sets, as defined in the IM. The membership (#) operator allows referencing specific Components within a Data Set. The syntax is: <i>dataset_name#component_name</i> (for a better description see the corresponding section in the Part 2). For the <i>dataset name</i> an alias can be used.</p>
ValueDomainSubset<T>	<p>A Value Domain Subset of data type T.</p> <p>Value Domain Subsets can be classified according to their Data Type as:</p> <ul style="list-style-type: none"> • String: ValueDomainSubset<String> • Numeric: ValueDomainSubset<Numeric> • Integer: ValueDomainSubset<Integer> • Boolean: ValueDomainSubset<Boolean> • Date: ValueDomainSubset<Date>

- 877
- 878 In addition to the IM artefacts, the Operators can also use constant values of the following
 879 types (they have not been described in the IM for the sake of simplicity):
- 880 • Simple Constants (meaning scalar constants belonging to one of the basic data types)
 - 881 • Sets of Constants (meaning unordered sets of constants having a common data type)
 - 882 • Lists of Constants (meaning ordered sets of constants having a common data type)

883 The following table contains the main types of constant parameters.

Types of constant Parameters	
Constant<T>	A constant value of data type <T>. Constants can be classified according to their Data Type as: <ul style="list-style-type: none">• String: Constant<String>• Numeric: Constant<Numeric>• Integer: Constant<Integer>• Boolean: Constant<Boolean>• Date: Constant<Date>
ConstantSet<T>	An unordered collection, without duplicates, of Constants of the same type T. The round brackets "()" denote that the order is not significant. Examples of allowed literals: ("a","b","c","d"), (1,2,3,4), (1.2, 3.4, 0.0).
ConstantList<T>	An ordered collection of Constants of the same type T, enclosed in square brackets, which denotes that the order is significant. Examples of allowed literals: ["a","b","c","d"], [1,2,3,4], [1.2, 3.4, 0.0].

884
885

886 **Type management and checking**

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

888 It is envisaged that there will be "implicit upcasting" between the Integer and the Numeric
889 data types and between the corresponding types of Parameters. This means that wherever in
890 the language it is possible to use a Constant<Numeric>, a Constant<Integer> is allowed as
891 well. Similarly, wherever it is possible to use a Component<Numeric>, a Component<Integer>
892 is allowed as well. Obviously, the opposite is not allowed. In these cases, in the description of
893 the single Operators in the Part 2, the Numeric type is indicated, provided that there are no
894 particular constraints on using Integers.

895 The VTL is strongly typed, in the sense that any Parameter belongs to one of the possible
896 types.

897 The various Operators have specific constraints in terms of number and types of Parameters
898 (see the corresponding sections in the Part 2).

899 Also a VTL Expression is assumed to correspond to a Parameter type, which is the type of its
900 output Parameter. The type of an Expression can be calculated at compile time.

901 An Expression can be input of an Operator, provided that the Parameter type of the (result of
902 the) Expression is compliant with the Operator constraints.

903 The Operators constraints in terms of number and types of Parameters are statically checked
904 (at compile time) so that type errors are not possible at runtime. Moreover, only type-safe
905 upcast conversion for Integers into Numerics is performed.

906 Type errors result in **compile time exceptions** preventing the Transformations from being
907 used (exchanged, executed ...).

908

909 The operations on the Data Sets

910 General rules

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

917 Let us call Data Set Operator a generic VTL Operator which produces a Data Set. If we assume
918 that F is a Data Set Operator, D_r is its result Data Set and D_i ($i=1, \dots, n$) are its input Data Sets, the
919 general form of a Transformation based on F can be written as follows:

$$920 \quad D_r := F(D_1, D_2, \dots, D_n)$$

921 Operator F composes the Data Points of D_i ($i=1, \dots, n$) to obtain the Data Points of D_r .

922 For making this operation, F follows a number of default behaviours described here.

923 In general the Data Sets D_i ($i=1, \dots, n$) and consequently their Data Points may have any number of
924 Identifier, Measure and Attribute Components, nevertheless the VTL Data Set Operators may
925 require specific constraints on the Data Structure Components of their input Data Sets¹².

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

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

930 In particular, as for the number of operands, a **Data Set Operator** is called “**unary**” if it allows
931 only one Data Set as input operand (e.g. minimum, maximum, absolute value ...) and “**n-ary**” if
932 it requires more than one Data Set as input operand (e.g. sum, product, merge ...). The **n-ary**
933 Operators require a preliminary matching between the Data Points of the various input Data
934 Sets.

935 The **Data Sets** may be also usefully categorized with reference to the number of their Measure
936 Components. A Data Set is called “**mono-measure**” if it has just one Measure Component and
937 “**multi-measure**” if it has two or more Measure Components. For the multi-measure Data
938 Sets it may be necessary to specify which measures should be considered in the operation.

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

¹² 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.

944 **The Identifier Components and the Data Points default matching**

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

948 The “n-ary” VTL Data Set Operators compose more than one input Data Sets. A simple
949 example is: $D_r := D_1 + D_2$

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

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

957 $D_1 =$ United States Data

958

<i>Ref.Date</i>	<i>Meas.Name</i>	<i>Meas.ValUe</i>
2013	Population	200
2013	Gross Prod.	800
2014	Population	250
2014	Gross Prod.	1000

959
960
961
962

963 $D_2 =$ European Union Data

964

<i>Ref.Date</i>	<i>Meas.Name</i>	<i>Meas.ValUe</i>
2013	Population	300
2013	Gross Prod.	900
2014	Population	350
2014	Gross Prod.	1000

965
966
967
968

969

970 The desired result of the sum is the following:

971 $D_r =$ United States + European Union

972

<i>Ref.Date</i>	<i>Meas.Name</i>	<i>Meas.ValUe</i>
2013	Population	500
2013	Gross Prod.	1700
2014	Population	600
2014	Gross Prod.	2000

973
974
975
976

977

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

981 The operation above is assumed to happen under a **strict constraint**: the input Data Sets
982 must have the same Identifier Components. The result will also have the same Identifier
983 Components as the operands.

984 Some Data Set operations (including the sum) may be possible also under a more **relaxed**
985 **constraint**, that is if the Identifier Components of one Data Set are a superset of those of the
986 other Data Set.

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

990 $D_1 =$ European Countries

<i>Ref.Date</i>	<i>Country</i>	<i>Population</i>
2012	U.K.	60
2012	Germany	80
2013	U.K.	62
2013	Germany	81

996 $D_2 =$ Europe

<i>Ref.Date</i>	<i>Population</i>
2012	480
2013	500

1000

1001 In order to calculate the percentage of the population of each single country on the total of
1002 Europe, the Transformation will be:

1003
$$D_r := D_1 / D_2 * 100$$

1004 The Data Points will be matched according to the Identifier Components common to D_1 and D_2
1005 (in this case only the Ref.Date), then the operation will take place.

1006 The result Data Set will have the Identifier Components of both the operands:

1007 $D_r =$ European Countries / Europe * 100

<i>Ref.Date</i>	<i>Country</i>	<i>Population</i>
2013	U.K.	12.5
2013	Germany	16.7
2014	U.K.	12.4
2014	Germany	16.2

1013 In the Part 2, dedicated to the description of the library of Operators, it is specified whether
1014 the Operators require the **strict** or the **relaxed** constraint (see the "Constraints" subsections).

1015 More formally, let F be a generic n-ary VTL Data Set Operator, D_r the result Data Set and D_i
1016 $(i=1, \dots, n)$ the input Data Sets, so that: $D_r := F(D_1, D_2, \dots, D_n)$

1017 The “strict” constraint requires that the Identifier Components of the D_i $(i=1, \dots, n)$ are the same.
1018 The result D_r will also have the same Identifier components.

1019 The “relaxed” constraint requires that at least one input Data Set D_k exists such that for each
1020 D_i $(i=1, \dots, n)$ the Identifier Components of D_i are a (possibly improper) subset of those of D_k . The
1021 output Data Set D_r will have the same Identifier Components of D_k .

1022 The n-ary Operator F will produce the Data Points of the result by matching the Data Points of
1023 the operands that share the same values for the common Identifier Components and by
1024 operating on the values of their Measure Components according to its semantics.

1025 Behaviour for Measure Components

1026 As already mentioned, given $D_r := F(D_1, D_2, \dots, D_n)$, the input Data Sets D_i $(i=1, \dots, n)$ may have any
1027 number of Measure Components. Therefore to enforce the desired behaviour it is necessary
1028 to understand which Measures the Operator is applied to. This Section shows the general VTL
1029 assumptions about how Measure Components are handled, while the behaviour of the single
1030 operators is described in the Part 2.

1031 The most simple case is the **application of unary Operators to mono-measure Data Sets**,
1032 which does not generate ambiguity; in fact the Operator is intended to be applied to the only
1033 Measure of the input Data Set. The result Data Set will have the same Measure, whose values
1034 are the result of the operation.

1035 For example, let us assume that D_1 contains the salary of the employees (the only Identifier is
1036 the Employee ID and the only Measure is the Salary):

1037 $D_1 =$ Salary of Employees

1038

<i>Employee ID</i>	<i>Salary</i>
A	1000
B	1200
C	800
D	900

1039
1040
1041
1042

1043
1044 The Transformation $D_r := D_1 * 1.10$ applies to the only Measure (the salary)
1045 and calculates a new value increased by 10%, so the result will be:

1046 $D_r =$ Increased Salary of Employees

1047

<i>Employee ID</i>	<i>Salary</i>
A	1100
B	1320
C	880
D	990

1048
1049
1050
1051
1052

1053 In case of **unary Operators applied to a multi-measure Data Set**, the Operator F is by
1054 default intended to be applied separately to all its Measures, unless differently specified. The
1055 result Data Set will have the same Measures as the operand.

1056 For example, given the import and export by reference date:

1057 $D_1 = \text{Import \& Export}$

<i>Ref.Date</i>	<i>Import</i>	<i>Export</i>
2011	1000	1200
2012	1300	1100
2013	1200	1300

1062 The Transformation $D_r := D_1 * 0.80$ applies to all the Measures (e.g. to
1063 both the Import and the Export) and calculates their 80%:

1064 $D_r = 80\% \text{ of Import \& Export}$

<i>Ref.Date</i>	<i>Import</i>	<i>Export</i>
2011	800	960
2012	1040	880
2013	960	1040

1070 If there is the need to **apply an Operator only to specific Measures**, the membership (#)
1071 operator can be used, which allows referencing specific Components within a Data Set. The
1072 syntax is: *dataset_name#component_name* (for a better description see the corresponding
1073 section in the Part 2).

1074 For example, in the Transformation $D_r := D_1\#\text{Import} * 0.80$

1075 the operation applies only to the Import (and calculates its 80%):

1076 $D_r = 80\% \text{ of the Import, } 100\% \text{ of the Export}$

<i>Ref.Date</i>	<i>Import</i>	<i>Export</i>
2011	800	1200
2012	1040	1100
2013	960	1300

1081 Note that in the example above, the Import is kept and left unchanged. In fact by default all the
1082 Measures are kept in the result, even the ones that are not operated on. If there is the need to
1083 keep only some Measures, the “keep” clause can be used (see the Part 2).

1084

1085 In case of **n-ary Operators**, by default **the operation is applied on the Measures of the**
1086 **input Data Sets having the same names**, unless differently specified. To avoid ambiguities
1087 and possible errors, the input Data Sets are constrained to have the same Measures and the
1088 result will have the same Measures too.

1089 For example, let us assume that D_1 and D_2 contain the births and the deaths of the United
1090 States and the European Union respectively.

1091 $D_1 =$ Births & Deaths of the United States

<i>Ref.Date</i>	<i>Births</i>	<i>Deaths</i>
2011	1000	1200
2012	1300	1100
2013	1200	1300

1096 $D_2 =$ Birth & Deaths of the European Union

<i>Ref.Date</i>	<i>Births</i>	<i>Deaths</i>
2011	1100	1000
2012	1200	900
2013	1050	1100

1101

1102 The Transformation $D_r := D_1 + D_2$ will produce:

1103 $D_r =$ Births & Deaths of United States + European Union

<i>Ref.Date</i>	<i>Births</i>	<i>Deaths</i>
2011	2100	2200
2012	2500	2000
2013	2250	2400

1108

1109 The Births of the first Data Set have been summed with the Births of the second to calculate
1110 the Births of the result (and the same for the Deaths).

1111 If there is the need to **apply an Operator on Measures having different names**, the
1112 “rename” clause can be used to make their names equal (for a complete description of the
1113 clause see the corresponding section in the Part 2).

1114

1115 For example, given these two Data Sets:

1116 D_1 (Residents in the United States)

<i>Ref.Date</i>	<i>Residents</i>
2011	1000
2012	1300
2013	1200

1120

1121

1122

D_2 (Inhabitants of the European Union)

1123

<i>Ref.Date</i>	<i>Inhabitants</i>
2011	1100
2012	1200
2013	1050

1124

1125

1126

1127

1128 A Transformation for calculating the population of United States + European Union is:

1129 $D_r := D_1[\text{rename Residents as Population}] + D_2[\text{rename Inhabitants as Population}]$

1130 The result will be:

1131

D_r (Population of United States + European Union)

1132

<i>Ref.Date</i>	<i>Population</i>
2011	2100
2012	2500
2013	1250

1133

1134

1135

1136

1137 Note that the number and the names of the Measure Components of the input Data Sets are
1138 assumed to match (following their renaming if needed), otherwise the Expression is
1139 considered in error.

1140 In case the Measure Components of the input Data Sets match only partially, the Measure
1141 structure must be properly adapted through the features for structure manipulation (e.g. the
1142 *keep* and the *calc* clauses, see below and in the relevant sections in the Part 2).

1143 If there is the need to **apply an Operator only to specific Measures**, the membership (#)
1144 operator can be used as in the case of unary Operators. Even in this case, by default all the
1145 Measures are kept in the result, even the ones that are not operated on; if there is the need to
1146 keep only some Measures, the “keep” clause can be used (see the Part 2).

1147 Finally, it may be needed to **apply different Operators on different Measures**. This is
1148 possible through the *merge* Operator in combination with the *keep* and *calc* clauses (this
1149 offers a wide variety of possibilities, see the specific sections in the Part 2).

1150 Roughly speaking, *merge* allows the production of a Data Set having the union of the
1151 Components of the input Data Sets (in a similar way to the SQL join), *keep* selects the
1152 Components to keep in the result, *calc* defines specific operations for specific Components.

1153 As a first example, let D_1 and D_2 be two multi-measure Data Sets, both having I as the common
1154 Identifier Component and M_1 and M_2 as Measures. Suppose that we want to calculate D_r
1155 having the Measures M_3 and M_4 , where the former is the sum of the M_1 of the input Data Sets
1156 and the latter is the difference of the M_2 . This can be obtained as:

1157 $D_r :=$
1158 **merge**(D_1 , D_2 , on($D_1\#I = D_2\#I$), return($D_1\#I$ as I ,

```

1159         D1#M1 as M11, D2#M1 as M12, D1#M2 as M21, D2#M2 as M22) )
1160         [calc M11 + M12 as M3, M21 - M22 as M4][keep I, M3, M4]

```

1161 The *merge* operator joins D_1 and D_2 , applying the general key matching behaviour on the
1162 Identifier Component I (the resulting rows). The *return* keyword, which is part of the *merge*
1163 operator (see the Part 2), specifies which columns to return in the result, which will have I as
1164 Identifier Component and four Measure Components, obtained from D_1 and D_2 (two from
1165 each). The *calc* clause calculates the sum and the difference between the right pairs of
1166 measures. Finally, *keep* maintains only the desired Components.

1167 As another example, assume that D_1 and D_2 are two mono-measure Data Sets, both having I as
1168 Identifier Component and M_1 as Measure Component. Suppose that we want to calculate D_r
1169 having two Measures, M_2 obtained as the sum of the M_1 of the input Data Sets and M_3 obtained
1170 as their difference. This can be achieved as:

```

1171         Dr :=
1172         merge(D1, D2, on(D1#I = D2#I),
1173         return(D1#I as I, D1#M1 as M11, D2#M1 as M12) )
1174         [calc M11 + M12 as M2, M11 - M12 as M3][keep I, M2, M3]

```

1175 The *merge* operator joins D_1 with D_2 , the *return* keyword produces a temporary multi-
1176 measure Data Set where M_{11} and M_{12} have been copied from D_1 and D_2 respectively. Those
1177 Measure are in turn summed (into M_2) and subtracted (into M_3). The *keep* maintains only the
1178 desired Components.

1179 Finally, note that **each Operator may be applied on Measures of certain data types**,
1180 corresponding to its semantics. For example *abs* and *round* will require the Measures to be
1181 numeric, while *substr* will require them to be a string. Expressions which violate this
1182 constraint are obviously considered in error.

1183 For example consider the Transformation: $D_r := abs(D_1)$

1184 As already described, this expression is assumed to apply the *abs* Operator (i.e. absolute
1185 value) to all the Measures Components of D_1 . If all these Measures are quantitative the
1186 expression is considered correct, on the contrary, if at least one Measure is of an incompatible
1187 data type, the expression is considered in error. The general description of the VTL data types
1188 is given above while the description of the data types on which each operator can be applied
1189 is given in the Part 2.

1190 Order of execution

1191 VTL allows the application of many Operators in a single expression. For example:

```

1192         Dr := D1 + D2 / (D3 - D4 / D5)

```

1193 When the order of execution of the Operators is not explicitly defined (through the use of
1194 parenthesis), a default order of execution applies.

1195 In the case above, according to the VTL precedence rules, the order will be:

- 1196 I. D_4 / D_5 (default precedence order)
- 1197 II. $D_3 - I$ (explicitly defined order)
- 1198 III. D_2 / II (default precedence order)
- 1199 IV. $D_1 + III$ (default precedence order)

1200 The default order of execution depends on the precedence and associativity order of the VTL
1201 Operators and is described in detail in the Part 2.

1202 **Missing Data**

1203 The awareness of missing data is very important for correct VTL operations, because the
1204 knowledge of the Data Points of the result depends on the knowledge of the Data Points of the
1205 operands. For example, assume $D_r := D_1 + D_2$ and suppose that some Data Points of D_2
1206 are unknown, it follows that the corresponding Data Points of D_r cannot be calculated and
1207 are unknown too.

1208 Missing data can take up two basic forms.

1209 In the first form, **the lack of information is explicitly represented**. This is the case of Data
1210 Points that show a “missing” value for some Measure or Attribute Components, which denotes
1211 the absence of a true value for a Component. The “missing” value is not allowed for the
1212 Identifier Components, in order to ensure that the Data Points are always identifiable.

1213 In the second form, **the lack of information remains implicit**. This is the case of Data Points
1214 that are not present at all in the Data Set. For example, given a Data Set containing the reports
1215 to an international organization relevant to different countries and different dates, and having
1216 as Identifier Components the Country and the Reference Date, this Data Set may lack the Data
1217 Points relevant to some dates (for example the future dates) or some countries (for example
1218 the countries that didn’t send their data) or some combination of dates and countries.

1219 The interpretation of the Data Points that are not present in the Data Set may be different in
1220 different cases. There are situations in which it is not correct to assume that such Data Points
1221 are “unknown”. As a matter of fact, there exist significant cases in which the “known” Data
1222 Points having a prefixed value (e.g. the “zero” value) are intentionally omitted, so that:

- 1223 • It is not possible to conclude that the missing Data Points are unknown;
- 1224 • it may be required to consider the missing Data Points as known and having such a
1225 prefixed value.

1226 The most common case of this kind is the “zero” value for quantitative data. According to a
1227 common practice, in fact, in high volume sparse data (i.e. when most of the Data Points have
1228 the value “zero”), the Data Points equal to “zero” are intentionally omitted, because it would
1229 be highly cumbersome or even unbearable to represent them explicitly. In these cases it may
1230 be correct to assume that the missing Data Points are “known” and have the value “zero”. This
1231 situation will be called hereinafter “implicit zero”.

1232 On the contrary, if the Data Points assuming the value “zero” are explicitly represented, it is
1233 correct to assume that the missing Data Points are “unknown”. This situation is called “explicit
1234 zero”.

1235 For some quantitative Operators, the current version of VTL allows both implicit and explicit
1236 zero operations. In the former case, if a calculation finds missing Data Points for an operand,
1237 the corresponding result is regularly calculated assuming for them the value “zero”. In the
1238 latter case, on the contrary, the result is considered “unknown”.

1239 For the sake of clarity, the VTL introduces distinct operators for the two cases. For example,
1240 the VTL algebraic operators (+, -, *, /) operate in implicit zero mode, while there are other
1241 corresponding operators (++, --, **, //) which perform the same operation in explicit zero
1242 mode.

1243 In practice, considering the case $D_r = F(D_1, D_2)$, if a Data Point P_1 of D_1 does not match with
1244 any Data Point P_2 of D_2 (i.e. there does not exist a P_2 of D_2 having the same value for the
1245 Identifier Components as P_1 of D_1), both the kinds of operators assume a fictitious matching
1246 Data Point P_{2F} , whose Measure Components are assigned the value “zero” by the former kind
1247 of operators ($+$, $-$, $*$, $/$) and the value “unknown” (NULL) by the latter ($++$, $--$, $**$, $//$).

1248 Coming back to the case of the **explicit representation of the “missing” values**, there can
1249 exist more missing values having different meanings. For example, possible meanings are
1250 “non-reported data” (the value should have been reported but it is absent), “nil data” (the data
1251 is negligible or zero), “not applicable data” (data is missing as expected) and so on. At the
1252 moment there is no standardization of the missing values and different organizations may use
1253 different sets of missing values (the goal of standardizing the missing values is out of the
1254 context of this work). Moreover, the needed missing values may change.

1255 A common practice to deal with missing values is to use just one value for the Measure
1256 Components having the generic meaning of “unknown” (the NULL literal) and introducing
1257 dedicated Attribute Components to better qualify the meaning as “non-reported”, “nil”, “not
1258 applicable” and so on.

1259 The VTL supports this practice through the NULL literal and the propagation rules of the
1260 Attribute Components, which are described below.

1261 The general properties of the NULL are the following ones:

- 1262 • **Data type:** NULL is type-less; this means that it is an allowed value for a Component of
1263 any data type (e.g. Numeric, String, Boolean ...)
- 1264 • **Testing.** A specific Operator (**isnull**) allows to test if a value is NULL returning a
1265 Boolean value (TRUE or FALSE).
- 1266 • **Comparisons.** Whenever a NULL value is involved in a comparison ($>$, $<$, $>=$, $<=$, in , not
1267 in , $between$) the result of the comparison is NULL.
- 1268 • **Mathematical operations.** Whenever a NULL value is involved in a mathematical
1269 operation ($+$, $-$, $*$, $/$, ...), the result is NULL.
- 1270 • **String operations.** In operations on Strings, NULL is considered an empty String (“”).
- 1271 • **Boolean operations.** VTL adopts 3VL (three-value logic). Therefore the following
1272 deduction rules are applied:

1273	TRUE	or	NULL	→	TRUE
1274	FALSE	or	NULL	→	NULL
1275	TRUE	and	NULL	→	NULL
1276	FALSE	and	NULL	→	FALSE

- 1277 • **Conditional operations.** The NULL is considered equivalent to FALSE; for example in
1278 the control structures of the type (*if (p) -then -else*), the action specified in *-then* is
1279 executed if the predicate p is TRUE, while the action *-else* is executed if the p is FALSE
1280 or NULL;
- 1281 • **Filter clauses.** The NULL is considered equivalent to FALSE; for example in the filter
1282 clause [*filter p*], the Data Points for which the predicate p is TRUE are selected and
1283 returned in the output, while the Data Points for which p is FALSE or NULL are
1284 discarded.
- 1285 • **Aggregations.** The aggregations (like *sum*, *avg* and so on) return one Data Point in
1286 correspondence to a set of Data Points of the input. In these operations the input Data
1287 Points having a NULL value are in general not considered. In the average, for example,

1288 they are not considered both in the numerator (the sum) and in the denominator (the
1289 count). Specific cases for specific operators are described in the respective sections.
1290 • **Implicit zero.** Arithmetic operators assuming implicit zeros (+, -, *, /) may generate
1291 NULL values for the Identifier Components in particular cases (superset-subset
1292 relation between the set of the involved Identifier Components). Because NULL values
1293 are in general forbidden in the Identifiers, the final outcome of an expression must not
1294 contain Identifiers having NULL values. As a momentary exception needed to allow
1295 some kinds of calculations, Identifiers having NULL values are accepted in the partial
1296 results. To avoid runtime error, possible NULL values of the Identifiers have to be fully
1297 eliminated in the final outcome of the expression (through a selection, or other
1298 operators), so that the operation of “assignment” (:=) does not encounter them.

1299 If a different behaviour is desired for NULL values, it is possible to **override** them. This can be
1300 achieved with the combination of the *calc* and *isnull* operators.

1301 For example, suppose that in a specific case the NULL values of the Measure Component M_1 of
1302 the Data Set D_1 have to be considered equivalent to the number 1, the following
1303 Transformation can be used to multiply the Data Sets D_1 and D_2 , preliminarily converting
1304 NULL values of $D_1 \# M_1$ into the number 1. For detailed explanations of *calc* and *isnull* refer to
1305 the specific sections in the Part 2.

1306
$$D_r := D_1 [\text{calc if}(\text{ISNULL}(M_1) \text{ then } 1 \text{ else } M_1) \text{ as } M_1] * D_2$$

1307 **The Attribute Components**

1308 Given as usual $D_r := F(D_1, D_2, \dots, D_n)$ and considering that the input Data Sets D_i ($i=1, \dots, n$) may
1309 have any number of Attribute Components, there can be the need of calculating the desired
1310 Attribute Components of D_r . This Section describes the general VTL assumptions about how
1311 Attributes are handled (specific cases are dealt with in description of the single operators in
1312 the Part 2).

1313 It should be noted that the Attribute Components of a Data Set are dependent variables of the
1314 corresponding mathematical function, just like the Measures. In fact, the difference between
1315 Attribute and Measure Components lies only in their meaning: it is intended that the
1316 Measures give information about the real world and the Attributes about the Data Set itself
1317 (or some part of it, for example about one of its measures).

1318 The VTL has a different default behaviour for Attributes and for Measures.

1319 As specified above, Measures are kept in the result by default, whereas Attributes are
1320 assigned a characteristic called “**virality**”, which determines if the Attribute is kept in the
1321 result by default or not: a “**viral**” Attribute is kept while a “**non-viral**” Attribute is not kept
1322 (the default behaviour is applied when no explicit indication about the keeping of the
1323 Attribute is provided in the expression).

1324 A second aspect is the “virality” of the Attribute in the result. By default, a viral Attribute is
1325 considered viral also in the result.

1326 A third aspect is the operation performed on an Attribute. By default, **the operations which**
1327 **apply to the Measures are not applied to the Attributes**, so that the operations on the
1328 Attributes need a dedicated specification. If no operations are explicitly defined on an
1329 Attribute, a default calculation algorithm is applied in order to determine the Attribute’s
1330 values in the result.

1331 As already mentioned, when the default behaviour is not desired, a different behaviour can be
 1332 specified by means of the proper use of the *keep*, *calc* and *attrcalc* clauses. In particular,
 1333 through these clauses, it is possible to override the virality (to keep a *non-viral* Attribute or
 1334 not to keep a *viral* one), to alter the virality of the Attributes in the result (from *viral* to *non-*
 1335 *viral* or vice-versa) and to define a specific calculation algorithm for an Attribute (see the
 1336 detailed description of these clauses in the Part 2).¹³

1337 Hence, the **default Attribute propagation rule** behaves as follows:

- 1338 • the non-viral Attributes are not kept in the result and their values are not considered;
- 1339 • the viral Attributes of the operand are kept and are considered viral also in the result;
 1340 in other words, if an operand has a viral Attribute V, the result will have V as viral
 1341 Attribute also;
- 1342 • The Attributes, like the Measures, are combined according to their names, e.g. the
 1343 Attributes having the same names in multiple Operands are combined, while the
 1344 Attributes having different names are considered as different Attributes;
- 1345 • the values of the Attributes which exist and are viral in only one operand are simply
 1346 copied (obviously, in the case of unary Operators this applies always);
- 1347 • the Attributes which exist and are viral in multiple operands (i.e. Attributes having the
 1348 same names) are combined in one Attribute of the result (having the same name also),
 1349 whose values are calculated according to the default calculation algorithm explained
 1350 below;

1351 Extending an example already given for unary Operators, let us assume that D_1 contains the
 1352 salary of the employees of a multinational enterprise (the only Identifier is the Employee ID,
 1353 the only Measure is the Salary, and there are two other Components defined as viral
 1354 Attributes, namely the Currency and the Scale of the Salary):

1355 $D_1 = \text{Salary of Employees}$

<i>Employee ID</i>	<i>Salary</i>	<i>Currency</i>	<i>Scale</i>
A	1000	U.S. \$	Unit
B	1200	€	Unit
C	800	yen	Thousands
D	900	U.K. Pound	Unit

1361
 1362 The Transformation $D_r := D_1 * 1.10$ applies only to the Measure (the salary)
 1363 and calculates a new value increased by 10%, the viral Attributes are kept and left unchanged,
 1364 so the result will be:

¹³ 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

1365

D_r = Increased Salary of Employees

1366

1367

1368

1369

1370

<i>Employee ID</i>	<i>Salary</i>	<i>Currency</i>	<i>Scale</i>
A	1100	U.S. \$	Unit
B	1320	€	Unit
C	880	yen	Thousands
D	990	U.K. Pound	Unit

1371

1372

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.

1373

1374

1375

1376

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:

1377

1378

1379

1380

1381

1382

1383

1384

1385

- 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.

1386

1387

1388

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.

1389

1390

1391

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).

1392

D_1 = Births & Deaths of the United States

1393

1394

1395

1396

<i>Ref.Date</i>	<i>Births</i>	<i>Deaths</i>	<i>Estimate</i>
2011	1000	1200	False
2012	1300	1100	False
2013	1200	1300	True

1397

D_2 = Birth & Deaths of the European Union

1398

1399

<i>Ref.Date</i>	<i>Births</i>	<i>Deaths</i>	<i>Estimate</i>
2011	1100	1000	False
2012	1200	900	True
2013	1050	1100	False

1400

1401 Assuming the weights 1 for “false” and 2 for “true”, the Transformation $D_r := D_1 + D_2$
1402 will produce:

1403 $D_r = \text{Births \& Deaths of United States} + \text{European Union}$

<i>Ref.Date</i>	<i>Births</i>	<i>Deaths</i>	<i>Estimate</i>
2011	2100	2200	False
2012	2500	2000	True
2013	2250	2400	True

1408 Note also that:

- 1409 • if the attribute *Estimate* was non-viral in both the input Data Sets, it would not be kept
1410 in the result
- 1411 • if the attribute *Estimate* was viral only in one Data Set, it would be kept in the result
1412 with the same values as in the viral Data Set

1413 The VTL default Attribute propagation rule (here called A) ensures the following properties
1414 (in respect to the application of a generic VTL operator “§” on the measures):

1415 **Commutative law (1)**

1416 $A(D_1 \S D_2) = A(D_2 \S D_1)$

1417 The application of A produces the same result (in term of Attributes) independently of
1418 the ordering of the operands. For example, $A(D_1 + D_2) = A(D_2 + D_1)$. This may seem
1419 quite intuitive for “sum”, but it is important to point out that it holds for every
1420 operator, also for non-commutative operations like difference, division, logarithm and
1421 so on; for example $A(D_1 / D_2) = A(D_2 / D_1)$

1422 **Associative law (2)**

1423 $A(D_1 \S A(D_2 \S D_3)) = A(A(D_1 \S D_2) \S D_3)$

1424 Within one operator, the result of A (in term of Attributes) is independent of the
1425 sequence of processing.

1426 **Reflexive law (3)**

1427 $A(\S(D_1)) = A(D_1)$

1428 The application of A to an Operator having a single operand gives the same result (in
1429 term of Attributes) that its direct application to the operand (in fact the propagation
1430 rule keeps the viral attributes unchanged).

1431 Having these properties in place, it is always possible to avoid ambiguities and circular
1432 dependencies in the determination of the Attributes’ values of the result. Moreover, it is
1433 sufficient without loss of generality to consider only the case of binary operators (i.e. having
1434 two Data Sets as operands), as more complex cases can be easily inferred by applying the VTL
1435 Attribute propagation rule recursively (following the order of execution of the operations in
1436 the VTL expression).

1437 With regard to this last aspect, the VTL assumes that the **order of execution** of the operations
1438 in an expression is determined by the precedence and associativity rules of the Operators
1439 applied on the Measures, as already explained in the relevant section. The operations on the

1440 Attributes are performed in the same order, independently of the application of the default
1441 Attribute propagation rule or user defined operations.

1442 For example, recalling the example already given:

$$1443 \quad D_r := D_1 + D_2 / (D_3 - D_4 / D_5)$$

1444 The evaluation of the Attributes will follow the order of composition of the Measures:

1445 I. $A(D_4 / D_5)$ (default precedence order)

1446 II. $A(D_3 - I)$ (explicitly defined order)

1447 III. $A(D_2 / II)$ (default precedence order)

1448 IV. $A(D_1 + III)$ (default precedence order)

1449 Storage and retrieval of the Data Sets

1450 The Storage

1451 As mentioned, the general form of Transformation can be written as follows:

$$1452 \quad D_r := F(D_1, D_2, \dots, D_n)$$

1453 In practice, the right-hand side is a mathematical expression like the one described above:

$$1454 \quad D_r := D_1 + D_2 / (D_3 - D_4 / D_5)$$

1455 As already shown, this expression implies the calculation of many Data Sets in different steps:

1456 I. (D_4 / D_5)

1457 II. $(D_3 - I)$

1458 III. (D_2 / II)

1459 IV. $(D_1 + III)$

1460 Calculated Data Sets are assumed to be non-persistent (temporary), as well as D_r , to which is
1461 assigned the final result of the expression (step IV).

1462 A temporary result within the expression can be only input of other operators in the same
1463 expression.

1464 Parameter D_r , which the result of the whole expression is assigned to, can be directly
1465 referenced as operand by other Transformations of the same VTL program (a VTL program is
1466 a set of Transformations, that is a Transformation Scheme, aimed to obtain some meaningful
1467 results for the users, supposed to be executed in the same run).

1468 The *Put* command is used to specify that a result must be persistent. Any step of the
1469 calculation can be made persistent (including all the steps).

1470 The *Put* has two parameters, the first is the (partial) result of the calculation that has to be
1471 made persistent (a non-persistent parameter of *Dataset* type), the second is the reference to
1472 the persistent Data Set, for example:

$$1473 \quad D_r := Put(D_1 + D_2 / (D_3 - D_4 / D_5), "PDS1")$$

1474 means that the overall result of the expression is stored in the persistent Data Set having
1475 name PDS1. The expression:

$$1476 \quad D_r := Put(D_1 + D_2 / Put((D_3 - D_4 / D_5), "PDS1"), "PDS2")$$

1477 Specifies that $(D_3 - D_4 / D_5)$ is stored in *PDS1* and the overall result in *PDS2*.

1478 **The Retrieval**

1479 Considering again the general form of Transformation:

$$1480 \quad D_r := F(D_1, D_2, \dots, D_n)$$

1481 the “n” Data Sets D_i ($i=1, \dots, n$) are the operands of the Expression and their values have to be
1482 retrieved.

1483 The generic D_i may be retrieved either as the temporary result of another Transformation (of
1484 the same VTL program) or from a persistent data source. In the former case D_i is the name of
1485 the left-hand parameter (D_r) of the other Transformation. In the latter, D_i is the reference to a
1486 persistent Data Set (see the following sections).

1487 A specific Operator (Get) ensures powerful features for accessing persistent data (see the
1488 detail in the Part 2). A direct reference to a persistent Data Set is equivalent to the application
1489 of the Get command.

1490 The Operators Get and Put are also called “commands” because they allow the interaction
1491 with the persistent storage.

1492 **The references to persistent Data Sets**

1493 In defining the Transformations, persistent Data Sets can be retrieved or stored by means of
1494 the Get and Put commands respectively.

1495 As described in the VTL IM, the Data Set is considered as an artefact at a logical level,
1496 equivalent to a mathematical function having independent variables (Identifiers) and
1497 dependent variables (Measures and Attributes). A Data Set is a set of Data Points, which are
1498 the occurrences of the function. Each Data Point is an association between a combination of
1499 values of the independent variables and the corresponding values of the dependent variables.

1500 Therefore, the VTL references the conceptual/logical Data Sets and does not reference the
1501 physical objects where the Data Points are stored. The link between the Data Set at a logical
1502 level and the corresponding physical objects is out of the scope of the VTL and left to the
1503 implementations.

1504 Also the versioning of the artefacts of the information model, including the Data Sets, is out of
1505 the scope of the VTL and left to the implementations.

1506 The VTL allows reference through commands (Get and Put) to any persistent Data Set defined
1507 and identified according the VTL IM. For correct operation, knowledge of the Data Structure of
1508 the input Data Sets is essential, in order to check the correctness of the expression and
1509 determine the Data Structure of the result. For this reason, the VTL requires that at
1510 compilation time the Data Structures of the referenced Data Sets are available.

1511 In addition, to simplify some kind of operations, the VTL makes it possible to reference also
1512 Cartesian subsets of the already defined Data Sets (i.e. sub Data Sets specified as Cartesian
1513 products of Value Domain Subsets of some Identifier Components).

1514 This is consistent with the IM, because any subset of the Data Points of a Data Set may be
1515 considered in its turn a Data Set, and with correct VTL operations, because the Data Structure
1516 of a sub Data Set is deducible from the Data Structure of the original Data Set, once that the
1517 specification of the subset is given.

1518 Note however that it is not possible to reference directly a non-Cartesian sub Data Set (i.e. a
1519 sub Data Set that cannot be obtained as a Cartesian product of Value Domain Subsets). As any

1520 other kind of Data Set, however, non-Cartesian subsets can be obtained through an
1521 Expression, as partial or final results.

1522 For example, in case of unit data, given the Data Set “Legal Entity” having as Identifiers of the
1523 Country, the IssuerOrganization, and the LegalEntityNumber, the VTL allows direct reference
1524 to either the whole Data Set or a sub-Data Set obtained specifying some countries, and/or
1525 issuers, and/or numbers. By specifying a single value for each identifier it is possible to
1526 reference even a single Legal Entity (i.e. a single Data Point).

1527 In case of Dimensional Data Sets, assuming that the Country and the Date are the Identifiers, it
1528 is possible to reference the sub Data Sets corresponding to one or some countries, to one or
1529 some dates, and to a combination of them. If the dates are periodical, the sub Data Set
1530 corresponding to one country is a time-series. The sub Data Set corresponding to a certain
1531 date is a cross-section. The sub Data Set corresponding to one country and one date is a single
1532 Data Point. Therefore the VTL allows direct reference to dimensional data, time-series, cross-
1533 sections, and single observations.

1534 In conclusion, a VTL reference to a persistent (sub)Data Set is composed of two parts:

- 1535 • The identification of the Data Set (mandatory)
- 1536 • The specification of a subset of it (optional)

1537 **The Identification of a persistent Data Set**

1538 The identification of the persistent Data Sets to read from (Get) or to store into (Put) follows
1539 the general rules of identification of the persistent artefact (see the corresponding section
1540 above).

1541 Therefore, the Data Set identifier is the **Data Set Name**, which is unique in the environment.
1542 As different environments can use the same Data Set Names for their artefacts, the Data Set
1543 Name can optionally be qualified by a **Namespace**, to avoid name conflicts.

1544 In case the Data Set identifier has a Namespace, a separator character can be chosen (and
1545 configured in the system) among the non-alphanumeric ones. A typical, and recommended,
1546 choice is the slash (“/”) symbol. If the Data Set identifier does not have a Namespace, the same
1547 namespace as the respective Transformation is assumed.

1548 Examples of good references to Data Sets are:

1549 “NAMESPACE/DS_NAME” (explicit Namespace definition)

1550 “DS_NAME” (the Namespace of the Transformation is assumed)

1551 **The specification of a subset of a persistent Data Set**

1552 The VTL allows the retrieval or storage of a subset of a predefined Data Set by filtering the
1553 values of its Identifier Components.

1554 Two basic options are allowed in the grammar of this specification:

- 1555 • A **full notation (query string)**, specifying both the Identifiers and the values to be
1556 filtered (e.g. Date= 2014, Country=USA, Sector=Public ...); in this case the filtering
1557 condition is preceded by the “?” symbol.
- 1558 • A **short notation (ordered concatenation)**, specifying only the values to be filtered
1559 (e.g. 2014.USA.Public); in this case the filtering condition is preceded by the “/”
1560 symbol; the values have to be specified following a predefined order of the Identifiers.

1561 The **query string** is a postfix syntax specifying the filter in case the order of the identifiers is
1562 not defined beforehand or not known.

1563 The filter is specified by concatenating the filtering conditions on the Identifiers, expressed in
1564 any order and separated by "&". If a filtering condition is not specified for an Identifier, the
1565 latter is not constrained and all the available values are taken. For example:

1566 I. DS_NAME?DATE=2014&COUNTRY=USA&SECTOR=PUBLIC

1567 In the example above, **single values** are specified for each filtering condition.

1568 In the same way, it is also possible to specify **multiple values** for some filtering conditions,
1569 separating the values by the "+" keyword (list). For example, to take the years 2013 and 2014
1570 and the countries USA and Canada:

1571 II. DS_NAME?DATE=2013+2014&COUNTRY=USA+CANADA&SECTOR=PUBLIC

1572 Finally, where the Values have an order like the one for the "Date" data type, it is possible to
1573 specify ranges of values for some filtering conditions, separating the first and last values of
1574 the range by the "-" keyword (range). For example, to take all the years from 2008 to 2014:

1575 III. DS_NAME?DATE=2008-2014&COUNTRY=USA+CANADA&SECTOR=PUBLIC

1576 The **ordered concatenation** is a simplified syntax to specify the filter in case the order of the
1577 identifiers is defined beforehand and known.

1578 The filter is specified by concatenating the filtering conditions in the predefined order of the
1579 Identifiers; the filtering conditions do not require the specification of the name of the
1580 Identifier, which can be deduced by their predefined order, therefore only the values are
1581 specified, separated by ".", i.e. a dot. If a value is omitted, the corresponding Identifier is not
1582 constrained and all the available values are taken. For example, (assuming that the order on
1583 the identifiers is 1-Date, 2-Country, 3-Sector):

1584 I. DS_NAME/2014.USA.PUBLIC

1585 This definition in the query string syntax corresponds to:

1586 DS_NAME?DATE=2014&COUNTRY=USA&SECTOR=PUBLIC

1587 II. DS_NAME/.USA.PUBLIC

1588 This definition filters all the available years for the USA and the public sector, and
1589 in the query string syntax corresponds to:

1590 DS_NAME?COUNTRY=USA&SECTOR=PUBLIC

1591 III. DS_NAME/..PUBLIC

1592 This definition filters all the available years and countries for the public sector and
1593 in the query string syntax corresponds to:

1594 DS_NAME?SECTOR=PUBLIC

1595 If needed, the list ("+") and/or range ("-") keywords can be used to specify lists or range of
1596 values respectively. For example:

1597 IV. DS_NAME/2008-2014.USA+CANADA.PUBLIC

1598 This definition in the query string syntax corresponds to:

1599 DS_NAME?DATE=2008-2014&COUNTRY=USA+CANADA&SECTOR=PUBLIC

1600

1601 Conventions for the grammar of the language

1602 General conventions

1603 A VTL program is a set of Transformations executed in the same run, which is defined as a
1604 Transformation Scheme.

1605 Each Transformation consists in a *statement* that is an assignment of the form:

1606 `variable parameter := expression`

1607 “:=” is the assignment operator, meaning that the result of the evaluation of the *expression* in
1608 the right-hand side is assigned to the *variable parameter* in the left-hand side (which is the
1609 output parameter of the assignment).

1610 Examples of assignments are (assuming that *ds_i* (*i=1...n*) are Data Sets):

- 1611 • `ds_1 := ds_2`
- 1612 • `ds_3 := ds_4 + ds_6`

1613 Variable Parameter names

1614 The variable parameters are non-persistent (temporary).

1615 The names of the variable parameters are alphanumeric (starting with an alphabetic
1616 character). Also non alphabetic characters (“_”,“-”) are allowed, but not in the first position.
1617 Parameter names are case-sensitive.

1618 Examples of allowed names for the parameters are: `par1`, `p_1`, `VarPar_ABCD`, `paraMeterXY`.

1619 Reserved keywords

1620 Certain words are reserved **keywords** in the language and cannot be used as parameter
1621 names, they include:

- 1622 - all the names of the operators / clauses
- 1623 - all the symbols used by the language (assignment “:=”, parenthesis “(,)”, “[,]”,
1624 ampersand “&”, hash “#” ...)
- 1625 - true
- 1626 - false
- 1627 - all
- 1628 - imbalance
- 1629 - errorlevel
- 1630 - condition
- 1631 - msg_code
- 1632 - dataset
- 1633 - script

1634 Expressions

1635 The expression is the right-hand side of an assignment and can be built in a number of
1636 alternative options.

1637 In general, an Expression may be the result of the application of an operator to another (sub)-
1638 expression. This may be done recursively, as already shown in the examples in the sections
1639 above. In other words, an Expression can be an operand of an Operator, resulting in another
1640 Expression that in turn may be an operand of an Operator, and so on.

1641 The basic and simplest types of Expressions correspond to the types of Constants and
1642 Parameters, for example an Expression can be:

1643 **A Constant**, that is a literal of any data type. Examples of Constant Expressions are:

1644 String: 'hello world', 'string'
1645 Numeric: 12.34, 0.0, 23.2E+4
1646 Integer: 2, 0, 45
1647 Boolean: true, false
1648 Date: 2012-01-31

1649 **A Constant Set** of any data type. Examples of Constant Sets Expressions are:

1650 String: ('k', '7', '1')
1651 Numeric: (12.34, 0.0, 23.2E+4)

1652 **A Constant List** of any data type. Examples of Constant List Expressions are:

1653 Numeric: [12.34, 0.0, 23.2E+4]
1654 Boolean: [true, false, false]

1655 **A Data Set** of any data type. Examples of Data Set Expressions are:

1656 Reference to temporary Data Sets: ds_1, DatasetA, X, Y
1657 Reference to a persistent Data Set: Namespace/DS_Name

1658 **A Component** of any data type. Examples of Component Data Set Expressions are:

1659 Component of a temporary Data Sets: ds_1#date, X#country
1660 Component of a persistent Data Set: Namespace/DS_Name#sector
1661 In the context of a single Data Set: date, country, sector

1662 **A Value Domain Subset** of any data type. Examples of Value Domain Subset Expressions are:

1663 Reference to a persistent V.D.S.: Namespace/VDS_Name
1664

1665 The other types of Expressions correspond to the ways an expression can be built from the
1666 basic types. For example an Expression can be:

1667 **The application of an Operator to other Expressions** (as explained above). Some examples
1668 are:

1669 Functional style: length(D1), round(D2, 4)
1670 Non functional style: D1+D2, D1 and (D2 or D3),

1671 **The application of a Clause to an Expression.** This is the same as above as for the semantic,
1672 and it is only different for the syntax, because the clauses are operators that use a postfix
1673 style. Some examples are the following ones (the D symbols denote Data Set names and the C
1674 symbols the Component names):

1675 D1[rename C1 as C2]
1676 D2+D4[keep C1, C2, C3]
1677 D3*(D2+D4)[calc C2*C3 as C5]

1678 **Comments**

1679 VTL allows comments within the statements in order to provide textual explanations of the
1680 operations. Whatever is enclosed between `/*` and `*/` shall not be processed by VTL parsers, as
1681 it shall be considered as comment.

1682 For example:

```
1683 /* Set constant for '\pi' */  
1684 numpi := 3.14  
1685 popA := populationDS + 1 /* Assign temp Dataset popA */
```

1686 **Constraints and errors**

1687 VTL supports a number of errors, which can occur in different situations; errors are divided
1688 into three main categories **compile time**, **runtime**, **validation**. Each category is divided in
1689 turn in subcategories, containing the specific errors.

1690 An error is identified by the string “VTL-“ followed by a four digit code CSEE, where:

- 1691 - C identifies the category (0: compile time, 1: runtime, 2: validation)
- 1692 - S identifies the subcategory
- 1693 - EE identifies the specific error in the subcategory

1694 While the three categories (and subcategories for compile errors) are standardized with
1695 codes reported in the remainder of this section, an encoding for specific errors (identified by
1696 the last two digits, EE) is not enforced here and can be independently defined by the adopting
1697 organization.¹⁴

1698 A compile time error prevents an expression from being used (exchanged, executed ...) and
1699 results in an exception reporting the error code (VTL-0XXX) and the wrong expression to the
1700 definer.

1701 In contrast, when a runtime error is raised, it can cause:

- 1702 a) an abnormal termination of the running VTL program, with an exception reporting the
1703 error code (VTL-1XXX) and the wrong expression to the user
- 1704 b) the current expression to be discarded, without generating any exception
- 1705 c) only the violating Data Point to be discarded, without generating any exception.

1706 The choice between these three behaviours should be dependent on the runtime system and
1707 is not part of the language, nor linked to the error codes.

1708 Validation errors are errors resulting from data validation (e.g. *check* operator), which can be
1709 stored in Datasets and used for further elaboration. Indeed, validation errors are not VTL
1710 errors and do not influence the use of the expression or the normal execution of a VTL
1711 program.

1712 **Compile Time errors (VTL-0xxx)**

1713 The VTL grammar specifies the rules to be followed in writing expressions. The VTL language
1714 allows the detection at compile time of the possible violation of the **correct syntax**, the use of

¹⁴ However, notice that in a following version of the language, a standardization is foreseen also for subcategories and specific error codes.

1715 **wrong types** as parameters for the operators or the **violation of any of the static**
1716 **constraints of the operators** (with respect to the rules described in the Part 2).

1717 A VTL compiler has to be able to detect all the syntax errors, help the user understand the
1718 reason and recover. Three subcategories are predetermined (see below). The specific error
1719 can be represented by the adopting organization with any code ranging from 00 to 99
1720 (examples are: unclosed literal string; unexpected symbol, etc.)

1721 Syntax errors (VTL-01xx)

1722 A violation of the VTL syntax with respect to the syntax templates of operators in name of
1723 operators or number of operands.

1724 Examples of syntactically invalid expressions are:

1725 R := C1 + -the second operand is missing

1726 R := C1 exist_in_all C2 - the correct syntax is "exists_in_all".

1727 R := if k1>4 then else K3 + 3 - the "then" operand is missing

1728

1729 Type errors (VTL-02xx)

1730 A violation of of the types of the operands allowed for the operators.

1731 Examples of expressions that are type-invalid are:

1732 R := C1 + '2' - if C1 has a measure component that is not <String>

1733 R := C1 + C2 - if C1 has a MeasureComponent<String> and C2 has a
1734 MeasureComponent<Numeric>

1735 R := C1 / 5 - if C1 has a MeasureComponent<String>.

1736 R:= if (K1 > 3 and k1 < 5) then 0 else "hello" - the "then" and the "else"
1737 operands must be of the same type

1738 Since the language is strongly typed, all type violations can be reported at compile time.

1739

1740 Static constraint violation errors (VTL-03xx)

1741 Every operator may have additional constraints. They are reported in the respective
1742 "Constraints" sections in the Part 2. Some of them are static, in the sense that they can be
1743 checked at compile type.

1744 A constraint violation error is the violation of a static VTL constraint .

1745 Examples of expressions that violate static constraints are:

1746 R := C1 + C2 - if the IdentifierComponents of C1 and C2 are not the same or
1747 are not contained in the ones of the other operator.

1748 R := 3 + 5 - in the plus (+) operator, at least one operand must be a Dataset.

1749

1750 **Runtime errors (VTL-1xxx)**

1751 These are the errors that can be detected only at runtime, typically because they are
1752 generated by the data.

1753 Examples are the classical mathematical constraints on operators arguments (negative or
1754 zero logarithm argument, division by zero, etc.).

1755 Particular types of runtime errors are:

- 1756 • presence of **duplicate** Data Points to be assigned to a Data Set (it is not allowed that
1757 two Data Points in a Data Set have the same values for all the Identifier Components
1758 because the Data Point identification would be impossible)
- 1759 • presence of a **NULL value** in an Identifier Component of a Data Point.

1760 These two errors result in a runtime exception only if the inconsistent Data Points are
1761 assigned (:=) to a Data Set in the left-hand side of a Transformation or are stored in a
1762 persistent Data Set. In other words, if such Data Points are only partial and temporary results
1763 inside the expression on the right-hand side, no runtime exceptions will be raised provided
1764 that the anomalies (duplications or NULLS) are removed before the execution of the
1765 assignment or the Put command.

1766 Examples of expressions generating runtime errors are:

1767 `R := C1 / C2` – where C2 is 0 for any observation

1768 `R := substr(A, 2, 5)` – if A is 1 character long, causing an “out of range”

1769 `R := C1` – if C1 contains NULL values for some IdentifierComponents.

1770 Notice that the assignment causes the runtime error; the fact that C1 contains a NULL value
1771 for an IdentifierComponent is accepted as partial and temporary result in the right-hand side
1772 of the expression.

1773 `R := C1` – if C1 contains duplicates on an IdentifierComponent. Also in this
1774 case, notice that the assignment causes the runtime error; the fact that C1 contains a duplicate
1775 is accepted as partial and temporary result in the right-hand side of the expression.

1776 A VTL runtime environment will be able to detect a wide number of runtime errors. The
1777 specific errors can be divided into subcategories by the adopting organization; moreover, the
1778 specific error can be represented with any code ranging from 00 to 99.

1779

1780 **Validation errors (VTL-2xxx)**

1781 They represent the outcome of a failed user-defined validation. The code can be used for
1782 further elaboration or to report discrepancies.

1783 Error codes can be associated with the single validations with the *check* operator, whose last
1784 parameter is *errorCode*. This is the code to be used for each Data Point having FALSE for its
1785 MeasureComponent.

1786 For example:

1787 `R := check(C1 >= C2, all, 2601)`

1788 Checks if C1 is greater or equal than C2 and, if not the case, stores the code 2601 in the
1789 *errorCode* attribute.

1790

1791

1792

1793

1794

C1		
K1	K2	M1
1	A	1000
2	B	200

1795

1796

1797

1798

1799

1800

C2			
K1	K2	K3	M1
1	A	X	1000
2	B	Y	350
2	B	Z	150

1801 and produces:

1802

1803

1804

1805

1806

1807

R				
K1	K2	K3	CONDITION	ERRORCODE
1	A	X	TRUE	
2	B	Y	FALSE	2601
2	B	Z	TRUE	

1808

1809 A set of VTL validation rules, will be able to detect a wide number of validation errors. The
 1810 specific errors can be divided into subcategories by the adopting organization; moreover, the
 1811 specific error can be represented with any code ranging from 00 to 99.

Governance, other requirements and future work

1812

1813 The SDMX Technical Working Group, as mandated by the SDMX Secretariat, is ensuring the
1814 technical maintenance of the Validation and Transformation Language through a dedicated
1815 VTL task-force. The VTL task-force is open to the participation of experts from other
1816 standardisation communities, such as DDI and GSIM.

1817 As the language is designed to be usable within different standards (SDMX, DDI, GSIM), a
1818 wider body could in future take on the task of synchronising and coordinating any parallel
1819 development. The detailed elements of a wider governance would need to be developed and
1820 shared with the other interested communities (e.g. GSIM, DDI, ESS, ESCB,...). Each community
1821 has its own governance rules and processes, and attention should be given to creating a
1822 system which may ensure a good representation of users' needs together with sound
1823 technical governance.

1824 A number of comments, suggestions and other requirements have been submitted to the VTL
1825 task force in order to enhance the current VTL 1.0 version. The outcome of a preliminary
1826 discussion of these requirements is presented here.

1827 The governance of the extensions

1828 According to the requirements, it is envisaged that the language can be enriched and made
1829 more powerful in future versions according to the evolution of the business needs. For
1830 example, new operators and clauses can be added, and the language syntax can be upgraded.

1831 The VTL governance body will take care of the evolution process, collecting and prioritising
1832 the requirements, planning and designing the improvements, releasing future VTL versions.

1833 The release of new VTL versions is considered as the preferred method of fulfilling the
1834 requirements of the user communities. This way, in fact, the possibility of exchanging
1835 standard validation and transformation rules would be preserved to the maximum extent
1836 possible.

1837 In order to fulfil specific calculation features not yet supported, the VTL provides for a specific
1838 operator (Evaluate) whose purpose is to invoke an external calculation function (routine),
1839 provided that this is compatible with the VTL IM and data types.

1840 The operator "Evaluate" (also "Eval") allows defining and making customized calculations
1841 (also reusing existing routines) without upgrading or extending the language, because the
1842 external calculation function is not considered as an additional operator. The expressions
1843 containing Eval are standard VTL expressions and can be parsed through a standard parser.
1844 For this reason, when it is not possible or convenient to use other VTL operators, Eval is the
1845 recommended method of customizing the language operations.

1846 However, as explained in the section "Extensibility and Customizability" of the "General
1847 Characteristics of VTL" above, calling external functions has some drawbacks in respect to
1848 the use of the proper VTL operators. The transformation rules would be not understandable
1849 unless such external functions are properly documented and shared and could become
1850 dependent on the IT implementation, less abstract and less user oriented. Moreover, the
1851 external functions cannot be parsed (as if they were built through VTL operators) and this
1852 could make the expressions more error-prone. External routines should be used only for

1853 specific needs and in limited cases, whereas widespread and generic needs should be fulfilled
1854 through the operators of the language.

1855 While the “Eval” operator is part of VTL, the invoked external calculation functions are not.
1856 Therefore they are considered as customized parts under the governance, and are
1857 responsibility and charge of the organizations which use it.

1858 Another possible form of customization is the extension of VTL by means of non-standard
1859 operators/clauses. This kind of extension is deprecated, because it would compromise the
1860 possibility of sharing validation rules and using common tools (for example, a standard parser
1861 would consider an expression containing non-standard operators as in error).

1862 Organizations possibly extending VTL through non-standard operators/clauses would
1863 operate on their own total risk and responsibility, also for any possible maintenance activity
1864 deriving from VTL modifications.

1865 Relations with the GSIM Information Model

1866 As explained in the section “VTL Information Model”, VTL 1.0 is inspired by GSIM 1.1 as much
1867 as possible, in order to provide a formal model at business level against which other
1868 information models can be mapped, and to facilitate the implementation of VTL with
1869 standards like SDMX, DDI and possibly others.

1870 GSIM faces many aspects that are out of the VTL scope; the latter uses only those GSIM
1871 artefacts which are strictly related to the representation of validations and transformations.
1872 The referenced GSIM artefacts have been assessed against the requirements for VTL and, in
1873 some cases, adapted or improved as necessary, as explained earlier. No assessment was made
1874 about those GSIM artefacts which are out of the VTL scope.

1875 In respect to GSIM, VTL considers both unit and dimensional data as mathematical functions
1876 having a certain structure in term of independent and dependent variables. This leads to a
1877 simplification, as unit and dimensional data can be managed in the same way, but it also
1878 introduces some slight differences in data representation. The aim of the VTL Task Force is to
1879 propose the adoption of this adjustment for the next GSIM versions.

1880 The VTL IM allows defining the Value Domains (as in GSIM) and their subsets (not explicitly
1881 envisaged in GSIM), needed for validation purposes. In order to be compliant, the GSIM
1882 artefacts are used for modelling the Value Domains and a similar structure is used for
1883 modelling their subsets. Even in this case, the VTL task force will propose the explicit
1884 introduction of the Value Domain Subsets in future GSIM versions.

1885 VTL is based on a model for defining mathematical expressions which is called
1886 "Transformation model". GSIM does not have a Transformation model, which is however
1887 available in the SDMX IM. The VTL IM has been based on the SDMX Transformation model,
1888 with the intention of suggesting its introduction in future GSIM versions.

1889 Some misunderstanding may arise from the fact that GSIM, DDI, SDMX and other standards
1890 also have a Business Process model. The connection between the Transformation model and
1891 the Business Process model has been neither analysed nor modelled in VTL 1.0. One reason is
1892 that the business process models available in GSIM, DDI and SDMX are not yet fully
1893 compatible and univocally mapped.

1894 It is worth nothing that the Transformation and the Business Process models address
1895 different matters. In fact, the former allows defining validation and calculation rules in the
1896 form of mathematical expressions (like in a spreadsheet) while the latter allows defining a
1897 business process, made of tasks to be executed in a certain order. The two models may
1898 coexist and be used together as complementary. For example, a certain task of a business
1899 process (say the validation of a data set) may require the execution of a certain set of
1900 validation rules, expressed through the Transformation model used in VTL. Further progress
1901 in this reconciliation is a task which needs some parallel work in GSIM, SDMX and DDI, and
1902 could be reflected in a future VTL version.

1903 Future directions

1904 **Structural Validation**

1905 We can distinguish two general types of validation according to their goals: “structural
1906 validation” and “content validation”, i.e. validation of the information content. The former can
1907 be defined as the assurance that data observations are compliant with the desired data
1908 structure, the latter that the data give a good representation of the phenomena under
1909 investigation.

1910 As both DDI and SDMX provide for structural metadata which allow structural validation, the
1911 VTL Task Force discussed whether VTL has to support structural validation or not. The
1912 conclusion was affirmative, considering that the use of different kinds of structural metadata
1913 is not homogeneous among organizations and among implementing standards and that it
1914 could be useful to support all kind of validations using the same method.

1915 It has been acknowledged, however, that this makes it possible to express structural
1916 validation rules in two alternative ways: through structural metadata or through VTL rules.
1917 Obviously, different choices by different organizations might compromise the possibility of
1918 exchanging, understanding and applying validation rules defined by others: the two forms of
1919 expressing structural validation rules should be made equivalent, in order to make it possible
1920 to transform one into the other, if needed.

1921 This VTL 1.0 version supports structural validation but does not provide yet for an
1922 equivalence with and easy conversion of structural metadata. This topic is intended to be
1923 covered in future work for a following VTL version.

1924 **Reusable rules**

1925 A main requirement expressed in the VTL public consultation is to allow generic and reusable
1926 rules, in order to apply the same rule in many cases. A typical example is to check that the
1927 values of a certain variable belong to a certain set of values.

1928 In VTL 1.0, such rules have to be written for each case. As structural metadata are typically
1929 reusable, only the structural validation rules defined through structural metadata are
1930 reusable at the moment.

1931 Reusable rules will be supported in a following VTL version, also through the use of “macro”
1932 operators (new operators defined by combining the existing ones).

1933

1934 **Other operators**

1935 In the VTL public consultation, some other kinds of operators have been requested (in
1936 addition to the “macro” operators already mentioned). For example, it was highlighted the
1937 lack of operators to manipulate dates and times, to convert different units of measure and to
1938 deal with time series. Operators of these kinds will be introduced in a following VTL version.

1939 It was also underlined that sometimes it is not easy to understand how to perform some kind
1940 of data manipulation. For example the possibility of converting the codes (from a coding
1941 system to another) is in some way “hidden” in the hierarchy operator. Cases of this kind may
1942 lead to a more explicit documentation or the introduction of more specific operators.

1943 Annex 1 - EBNF

1944 The VTL 1.0 language is also expressed in EBNF (Extended Backus-Naur Form).

1945 EBNF is a standard¹⁵ meta-syntax notation, typically used to describe a Context-Free grammar
1946 and represents an extension to BNF (Backus-Naur Form) syntax. Indeed, any language
1947 described with BNF notation can be also expressed in EBNF (although expressions are
1948 typically lengthier).

1949 Intuitively, an EBNF consists of terminal symbols and non-terminal production rules.
1950 Terminal symbols are the alphanumeric characters (but also punctuation marks, whitespace,
1951 etc.) that are allowed singularly or in a combined fashion. Production rules are the rules
1952 governing how terminal symbols can be combined in order to produce words of the language
1953 (i.e. legal sequences).

1954 More details about EBNF notation can be found on:

1955 http://en.wikipedia.org/wiki/Extended_Backus-Naur_Form

1956 Properties of VTL grammar

1957 VTL can be described in terms of a Context-Free grammar¹⁶, with productions of the form $V \rightarrow$
1958 w , where V is a single non-terminal symbol and w is a string of terminal and non-terminal
1959 symbols.

1960 VTL grammar aims at being unambiguous. An ambiguous Context-Free grammar is such that
1961 there exists a string that can be derived with two different paths of production rules,
1962 technically with two different leftmost derivations.

1963 In theoretical computer science, the problem of understanding if a grammar is ambiguous is
1964 undecidable. In practice, many languages adopt a number of strategies to cope with
1965 ambiguities. This is the approach followed in VTL as well. Examples are: the presence of
1966 *associativity* and *precedence* rules for infix operators (such as addition and subtraction); the
1967 existence of compulsory *else* branch in *if-then-else* operator.

1968 These devices are reasonably good to guarantee the absence of ambiguity in VTL grammar.
1969 Indeed, real parser generators (for instance YACC¹⁷), can effectively exploit them, in particular
1970 using the mentioned associativity and precedence constrains as well as the relative ordering
1971 of the productions in the grammar itself, which solves ambiguity by default.

¹⁵ ISO/IEC 14977

¹⁶ http://en.wikipedia.org/wiki/Context-free_grammar

¹⁷ <http://en.wikipedia.org/wiki/Yacc>