

THE SYNERGY REFERENCE MODEL OF DATA PROVISION AND AGGREGATION

Draft

June 2014

Current Contributors: Martin Doerr¹, Achille Felicetti, Gerald de Jong², Konstantina Konsolaki¹, Barry Norton³, Dominic Oldman³, Maria Theodoridou¹, Thomas Wikman⁴,

¹Institute of Computer Science, FORTH
{[martin, konsolak, maria@ics.forth.gr](mailto:martin,konsolak,maria@ics.forth.gr)}

²Delving B.V.
{[gerald, thomas@delving.eu](mailto:gerald,thomas@delving.eu)}

³British Museum
{[mailto:Bnorton, DOLDMAN@britishmuseum.org](mailto:Bnorton,DOLDMAN@britishmuseum.org)}

⁴The National Archives (Riksarkivet)
{thomas.wikman@riksarkivet.se}

TABLE OF CONTENTS

1	Introduction.....	5
2	Data Provisioning.....	10
2.1	Requirements and assumptions	10
3	User Roles	12
3.1	Primary User Roles.....	13
3.2	Secondary User Roles	15
4	Data Objects	18
4.1	Content Data and Metadata Objects	18
4.2	Schema and logic objects	19
4.3	Control Objects	22
5	Data Provisioning Flow Network	23
6	Analytical Representation of the Data Provisioning Process	25
6.1	Data Provisioning Process	26
6.1.1	Initial Data Delivery.....	28
6.1.1.1	Syntax Normalization	29
6.1.1.2	Mapping Definition.....	32
6.1.1.2.1	Schema Matching	33
6.1.1.2.2	Instance Generation Specification	38
6.1.1.2.3	Terminology Mapping	41
6.1.1.3	Metadata Transformation	43
6.1.1.3.1	Ingest and Storage	46
6.1.2	Update Processing	48
7	Services and S/W components	52
8	References	57

List of Figures

Figure 1: The data provisioning process	9
Figure 2: Working Environment - User Roles	12
Figure 3: Data Objects	18
Figure 4: Data Provisioning Flow Network.....	24
Figure 5: Data Provisioning Process Hierarchy.....	25
Figure 6: The Data Provisioning process	26
Figure 7: Data delivery sub-process.....	28
Figure 8: Syntax Normalization sub-process.....	29
Figure 9: Mapping Definition sub-process	32
Figure 10: Schema Matching sub-process	34
Figure 11: Instance Generation Specification sub-process.....	38
Figure 12 Terminology Mapping sub-process.....	41
Figure 13: Metadata Transformation sub-process.....	43
Figure 14: Ingest and Storage sub-process	46
Figure 15 Update processing sub-process	48
Figure 16 Services and S/W components	53

List of Tables

Table 1: Summary of the Data Provisioning Process.....	27
Table 2: Summary of the Data Delivery sub-process	28
Table 3: Summary of the Syntax Normalization sub-process.....	31
Table 4: Summary of the Mapping Definition sub-process.....	32
Table 5: Summary of the Schema Matching Definition sub-process	37
Table 6: Summary of the Instance Generation Specification sub-process.....	40
Table 7: Summary of the Terminology Mapping sub-process.....	43
Table 8: Summary of the Metadata Transfer sub-process	45
Table 9: Summary of the Ingest and Storage sb-process	47
Table 10: Summary of the Update Processing sub-process.....	51
Table 11: IT objects' Input/Output Documents	56

1 Introduction

This document defines a new reference model for a better practice of data provisioning and aggregation processes, primarily in the cultural heritage sector, but also for e-science. Such processes have become a reality in various, largely disharmonious, forms and in more or less systematic ways in numerous publicly funded projects. The current document defines a consistent set of business processes, user roles, generic software components and open interfaces that will form a harmonious whole. It is based on experience and evaluation of national and international information integration projects. The model is an initiative of the CIDOC CRM Special Interest Group, a Working Group of CIDOC-ICOM, the International Committee for Documentation part of the International Council of Museums. This document is a first draft compiled by a panel of Group members to initiate discussion and further elaboration by the Group and all interested experts and stakeholders in the area. This draft is a skeleton and is in no part complete or elaborated to the intended level of detail and does not represent any authoritative decision or approval of content. The contributors hope that this document is specific and elaborate enough that the reader can appreciate the scope, utility, form and level of specificity of the intended model so that they collaborate in an informed manner. Any and all interested experts are invited to participate and contribute.

The rationale behind this model is the following: Managing heterogeneous cultural heritage data is a complex challenge. Member institutions like galleries, libraries, archives and museums curate different types of collections that, even between similar types of institutions, are documented in different ways using different languages; influenced by different disciplines, objectives and geography, and are encoded using different metadata schemas. However, handling these metadata as a unified whole is vital for progressing new fields of humanities research and discovery, providing more knowledgeable information retrieval and (meta) data exchange, and advancing the field of digital humanities in its various aspects.

The ability to provide users with a uniform interface to access, relate, and combine data while at the same time preserving all the meaning and perspective of the individual data providers, might at first seem like an impossible task. Indeed, the exponential growth of the Web and the extended use of database management systems has brought to the fore the need for the seamless interconnection of large numbers of diverse information sources. In order to provide such uniform access to such heterogeneous and autonomous data sources, complex query and integration mechanisms need to be designed and implemented.

Data aggregation and integration has the potential to create rich resources useful for a range of different purposes, from research and data modeling to education and engagement. There are now significant numbers of projects that aggregate data with these purposes in mind. However, aggregators face two problems.

Firstly, the process of transferring data from source institutions to a central repository can result in a form of data representation stripped of essential information and institutional perspectives. This occurs when mandating target models into which all data sources must fit, regardless of their range, individuality and richness. The generalizations used in these models, designed to facilitate data integration, are too abstract to support the meaningful connections that undoubtedly exist in the data and thus significantly reduce the value of such aggregation initiatives.

The second problem, addressed by this document, relates to the lack of sustainability in the mechanisms and processes through which data is mapped, and the weakness of the partnership between data providers and aggregators inherent in these flawed processes. The mechanisms used for transferring data do not include the full set of necessary processes and tools to create a consistent and good quality outcome and furthermore cannot practically respond to changes in schema and systems on either side of the data provisioning relationship. In order for systems to be sustainable a broader approach is needed that incorporates the experience and knowledge of provider institutions into the infrastructure, in an accessible, beneficial and cost effective manner.

Therefore this document describes a new data provisioning model, the “**Synergy Reference Model**” (specifically the provision of data between providers and aggregators) including associated data mapping components. The intention is to address the design flaws in current models and crucially incorporate, through additional processes and components, the required knowledge and input from providers to create good quality, sustainable aggregations. The funding allocated to humanities aggregation projects over the last two decades has not generated the benefits and progress enjoyed in other sectors who have taken better advantage of digital innovation by using solid and inclusive infrastructures. Unless the value of these infrastructures is clearly demonstrated in the cultural heritage sector, resources for building and developing them will become even more scarce. Unfortunately, numerous systems initiated by projects in both Europe and the United States have failed to understand and identify the relationships and activities necessary to operate collaborative aggregation systems properly and instead have relied on one-sided and centralized approaches using top down modeling and technology led solutions.

The impact of these unresolved issues has been highly detrimental to the landscape of digital humanities, which has thus failed to evolve sufficiently from its fragmented beginnings. Consequently progress has not been made on the essential infrastructure necessary to produce meaningful innovation and fulfill the domains ambitions. Current projects still focus on short term functionality rather than tackle the key issues of sustainability and longevity. Together with an equally debilitating lack of an expressive and real world cultural heritage reference model for data, the humanities has not come close to tackling the challenge of computational reasoning or sophisticated modeling techniques across their vast heterogeneous resources. These resources still remain unexplored and effectively siloed, even when included in aggregation repositories. Indeed such reasoning and modeling activities hold the key to the serious advancement of humanities research; to a degree that would be of serious interest to funding bodies, who are instead becoming increasingly impatient and dissatisfied.

Aggregation systems cannot be viewed as substantially centralized undertakings because aggregators divorced from their providers cannot represent data properly and therefore cannot provide data integration with any real value to any audiences, or, crucially, to the providers themselves. As a result providers are unwilling to commit resources beyond their allocated project funding. If aggregation systems are to have any value they must distribute roles and responsibilities and include the experts who understand the data and know how it should be represented. Data providers curate and understand their information. Aggregators must restrict their involvement to providing homogeneous access to integrated data in such a way that it retains its proper context and its original meaning. Only then can the aggregator provide meaningful services to providers and users. The aggregator provides the information processes (including data improvement and co-referencing) that individual organizations are unable to resource independently, and generate quality services that can be utilized (in return for their data) by the provider. The provider also benefits, if the aggregation is done well and conveys the full meaning of the data, from the ability to use these integrated digital resources to support their own digital strategies. Aggregations based on meaning and context support all organizations small and large and increase the value and relevance of cultural heritage resources for a range of different purposes.

The process of mapping needs support from carefully designed tools and a collaborative knowledge base or “mapping memory” designed to support all organizations with differing levels of resourcing. Together with the CIDOC CRM ontology a new provisioning model is seen as a major step forward in the ability to directly enrich data through collaborative data harmonization, and using the power of multiple data sources to correct, inform and provide greater insight. The challenge is to define a modular architecture that can be developed and optimized by different

developers with minimal inter-dependencies and without hindering integrated UI development for the different user roles involved. The first part of this model is a form of requirements specification, which breaks down in the usual way into a definition of the associated user roles, the primary types of data the system aims to handle, and the complete definition of the processes users of the system carry out to manage the data.

The Synergy Reference Model will be described here in terms of a formal process model. The process model requires the definition of the individual roles, data objects and processes necessary for designing a controlled and managed mapping system. Figure 1 presents a high level view of the data provisioning process. The Provider Institution and the Aggregator Institution agree on the data provisioning and related activities. A “Mapping Manager” is nominated by both parties to oversee the actual data transfer process and forms the third primary role in the Synergy Model. The Provider Institutions own the Provider Records which are transformed to the Aggregator Record Format and are transferred to the Aggregator Institution. The data provisioning process is regarded as an open-ended and on-going task. Throughout the transformation and transfer processes, consistency checks and updates are necessary between all partners and will be supported by the model. The model foresees a series of distinct update processes at all partner sites which trigger each new data transfer.

The details of the Synergy Model will be presented in this document. As is normal in describing processes the descriptions may contain some redundancy. For example, it is common for different roles to be performed by the same people. The processes presented in this document should not be viewed as complete and static but rather are designed to facilitate the growth of other collaborative processes between provider organizations and aggregators. The ability to refer to the same set of stable processes increases the opportunities for organizations to pool and share their resources whether aimed at improving the aggregation service or simply at providing a platform for collaboration outside or connected to the offered aggregation services.

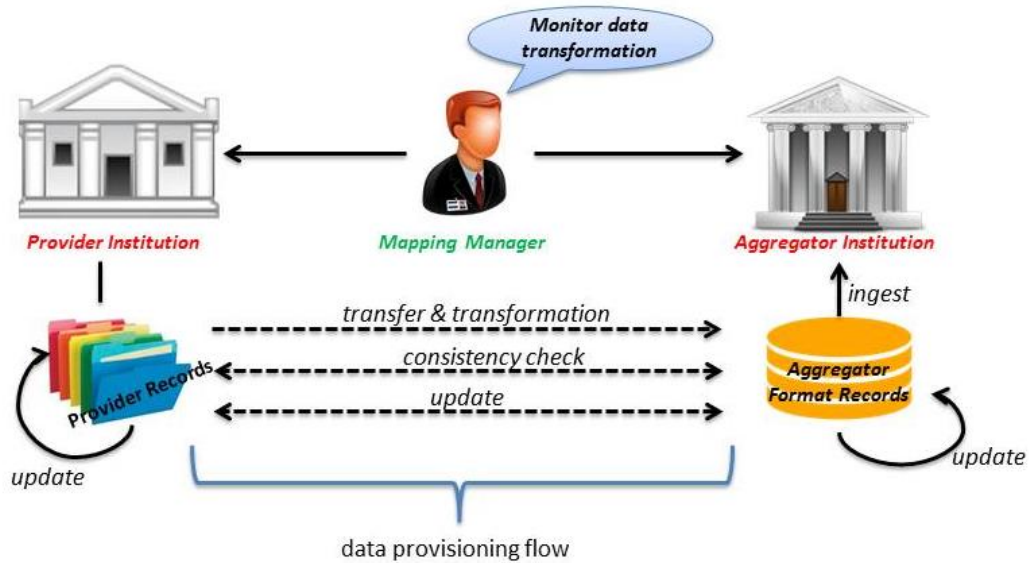


Figure 1: The data provisioning process

The structure of this document is as follows:

In section 2 we present an overview of the data provisioning process to provide a context for the next sections. In the following sections we show analytical model views¹ needed for a detailed understanding of the reference model. In Section 3 we give a detailed description of all the involved roles while in Section 4 we present the data objects that assist the data provisioning process. Section 5 presents the flow of the data objects from the Provider Institution to the Aggregator's Institution. It also provides a high-level overview of the IT objects that either replace manual tasks or assist the user. Section 6 presents the detailed analysis of the data provisioning processes. Finally, section 7 will present in detail the IT objects that assist the mapping process. This section will be extended into a full specification of the interfaces between the various components. All sections are under construction, and only indications of the intended final content.

¹ The modeling of the data mapping components and also of the processes needed for the completion of the mapping is made using Adonis- Business Process Management [1]. Adonis is a freeware tool, useful for the design and documentation of processes.

2 Data Provisioning

The Synergy Reference Model aims at identifying, supporting or managing the processes needed to be executed or maintained when a provider and aggregator agree (see Figure 1):

1. to *transfer data* from the provider to the aggregator,
2. to *transform* their format to the (homogeneous) format of the aggregator,
3. to *curate* the semantic consistency of source and target data and the global referential integrity and
4. to maintain the transferred data *up-to-date* with whatever relevant changes occur in the source and target systems and the employed terminologies.

In the following we present the requirements and assumptions taken into account during the Synergy Model design.

2.1 Requirements and assumptions

The Synergy Reference Model aims to support the management of data between source and target models and the delivery of transformed data at defined times, including updates. This includes a mapping definition, i.e., specification of the parameters for the data transformation process, such that complete sets of data records can automatically be transformed. A provisioning model includes:

1. The ***transfer of data*** (and iterative corrections) until a first consistent state is achieved. This includes transformation of sets of data records submitted to the aggregator, the necessary exception processing of irregular input data between provider and aggregator, ingestion of the transformed records into target system and initial referential integrity processing possibly on both sides.

Note : Referential integrity processing at the aggregator side, out of the context of a particular data submission, is a necessary process but is not part of the model.

2. ***Change detection and update*** processing to provide updated information and ensure semantic consistency in situations where changes have been made in the source or target records, in the provider or aggregator terminologies, in the source or target schemata, in the target instance generation policies and changes in the interpretation of source and target schema recorded in the mapping definition.

It is assumed that the provider has mechanisms that would identify modified records and thus ask for the mapping of these records only. If this is not the case, either the whole mapping is executed, overwriting previous mappings, or smaller units of change are identified and updated. Additionally the target system may also require

some way to clearly identify a modified record. Some aggregation systems may wish to store versions of data for research purposes but in this case the canonical records should always be clear and differentiated. As long as these changes can be identified and accessed then it is a matter for the aggregator to determine their own versioning system.

Only if these processes are sustained can an aggregator provide valid and consistently integrated data over the long term, and thereby deliver the full added value of an aggregation service that would make it attractive for providers and users alike. This sustainability is key in providing benefits for providers and the range of professionals, experts and enthusiasts that would ultimately justify its existence. This report is aware that none of the hundreds of mapping tools and frameworks created in numerous projects has ever systematically addressed this comprehensive scenario.

It is not until data can be analyzed and visualized in an appropriate format and environment that mapping decisions and issues can be made. The system should allow the exploration of data and give some indication of inconsistencies that might exist. The visualization process also provides the ability for both source and target models to be compared. The system may also allow 'test' transformations and provide some of the functionality, described below, to be applied to individual cases to provide some understanding of the requirements necessary to complete a full mapping and transformation.

3 User Roles

The following section describes the key management roles that oversee the process and provide the necessary resources. We distinguish primary and secondary roles associated with the data provisioning process. As depicted in Figure 1, the provider and aggregator organizations consist of performers, the *provider* and the *aggregator employees*, that may be one person or a team. Since the model does not prescribe if the mapping process is managed on account of the provider, the aggregator, or a joint activity, we have named this role as “*provider or aggregator employee*”, that may be one person or a team that is employed by the provider or the aggregator institution or both. Each performer holds different roles during the data provisioning process and as depicted in Figure 2 one performer may have in common one role with another performer. It is obvious that a person can play more than one roles in such a reference model.

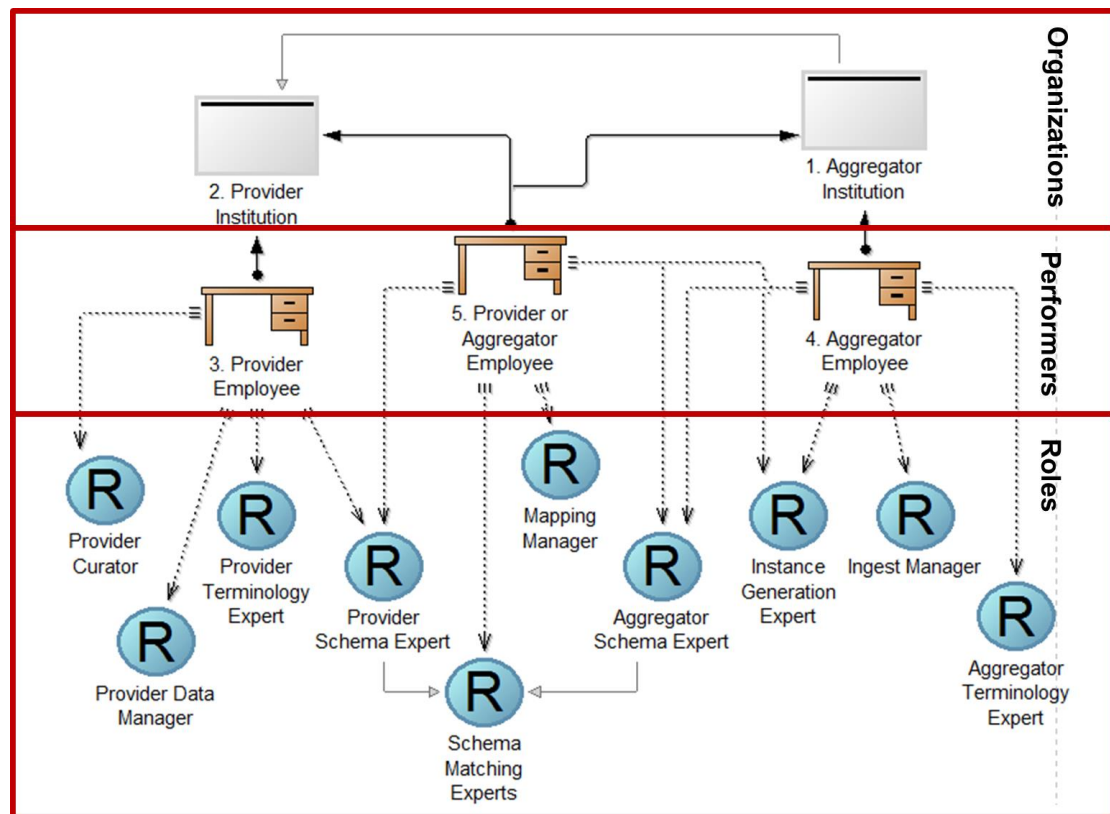


Figure 2: Working Environment - User Roles

3.1 Primary User Roles

As primary user roles we regard the managerially responsible members from the Provider and the Aggregator Institution that agree to perform the data provisioning of the providers' local information systems to the aggregator's integrated access system. These roles are seen in a logical and not a personal way. The model does not describe or exclude that the institution maintaining an aggregation service may also maintain a provider service. In that case, the provider role applies to the respective functions of such an institution. It is not part of the model where the actual data physically reside and how they are replicated or preserved. The model currently does not describe or exclude that an aggregator may hand over data to another aggregator. In that case a chain of interactions should be installed which ensures an information flow between primary providers and primary and secondary aggregators functionally equivalent to the one proposed by this model.

We define the following primary user roles :

Provider Institution

The ***Provider Institution*** maintains one or more Local Information Systems. Following CIDOC CRM v5.0.4 [2]:

“These are either collection management systems or content management systems that constitute institutional memories and are maintained by an institution. They are used for primary data entry, i.e. a relevant part of the information, be it data or metadata, is primary information in digital form that fulfils institutional needs.”

In practice these are systems owned by individual museums, archives, libraries, sites and monument records, academic institutes, private research societies etc., represented by their curators, technologists, documentation personnel and researchers. In other words, “Provider” in the sense of this model is the role of an institution as the authority for the correctness of knowledge represented in the data. Provider Institution systems are also called **source systems** in this text when talking about data transformation and submission. The implementation of provider systems is not part of this model, only certain communication capabilities.

Aggregator Institution

The ***Aggregator Institution*** maintains an Integrated Access System and in this model may also be called simply, ‘target systems’. Following CIDOC CRM v5.0.4 [2] *“These provide a homogeneous access layer to multiple local systems”*. The origin of the information it manages are the Provider Institutions it maintains a business relation

with. It may not alter provided content except for co-reference resolution information, changes of identifier and value representations or schema migrations. It may remove provided content or ask to providers for updates in order to maintain data quality. In case the aggregator institution wishes to add own new content of any form and provenance, the model will treat this part of information as another source system and will regard the respective activities as Provider activities. In other words, “Aggregator” in the sense of this model is the role of an institution of integrating and mediating data without changing meaning.

Aggregator Institution maintains a form of business agreement with providers to send data from local systems to the aggregators’ system, consisting primarily of metadata. The scenario that aggregators may harvest provider information without any formal terms of reference and understanding with the provider, such as the well-known search engine services, is not part of this model. The model will still be of use for such scenarios in a trivial way, but this scenario involves activities that are not the focus of this document and implies aggregators that have no direct knowledge about the meaning of the data they aggregate. This is insufficient for services that seek to harness the rich nature and embedded knowledge of cultural heritage organizations.

Mapping Manager

The ***Mapping Manager*** is the actor responsible for the maintenance of the data transformation process from the provider format to the aggregator format. This role may split into a semantic and a technical part, and may be regarded as an aggregator task, a provider task or a user consortium task. The mapping technology this model aims at, should support scalable management of the data transformation process by the aggregator. Mapping Managers must schedule and negotiate the terms under which data transformation occurs at both ends of the data provisioning system realizing that this may differ between providers.

3.2 Secondary User Roles

In this section we describe the experts whose knowledge or services contribute to the implementation and realization of the data provisioning process. They hold the secondary roles in the Synergy Model and they are employees of either the Provider or the Aggregator Institution.

Provider Curators

The Provider Curators are the employees responsible for curating the content of the source systems. They have the understanding of how their data reflects the real world, or know those who do, and have the means to check the veracity of the representation.

Provider Data Manager

The Provider Data Manager is the employee responsible for managing the relevant IT systems of the provider and for handling the data assets, in contrast to those responsible for creating content. In particular, they are responsible for sending data to the aggregator. The Provider Data Manager may formally be employed by the Provider Institution or the Provider Institution has outsourced this service. For the purpose of this model, it is not necessary to differentiate these situations.

Provider Schema Expert

The curator(s), researcher(s) and/or data manager(s) of the Provider Institution who are responsible for the data creation in their local systems, i.e., the people who know how, following local use and practice, the fields, tables or elements in the schema correspond to the reality described by them. These meanings may have become skewed or misinterpreted over time (the so-called semantic drift). This can be caused by a lack of precision or differentiation in the underlying data models, through misrepresentation in overlaying software and user interfaces or through changing practice. As such these meanings have to be researched and understood before mapping can take place. This is the domain of an increasingly undervalued group of people responsible for the quality and semantic correctness of their data. This group's significance and value would be highlighted by a real, rather than technically artificial, representation of their work. While quality levels will vary from institution to institution, semantic data *harmonization* can contribute to improving the quality of information across all organizations and indeed the whole sector.

Provider Terminology Expert

The Provider Terminology Expert is the curator, maintainer, or other expert on one or more of the terminologies that the provider uses as a reference in the local system. If the terminology is provided by a third party, such as the Getty Research Institute,

there may exist independent external experts conversant with this terminology. If the terminology is local in origin, or even uncontrolled, it is typically the curators or other local data managers (and documentation staff) who, alone, know the meaning of the local terms.

Aggregator Schema Expert

The Aggregator Schema Expert is the expert on the semantics of the schema employed by the aggregator (“integration model”). Some large scale aggregators use a more widely known standard schema, but there is also a growing trend in the linked data world towards lightweight portal or ‘indexing’ aggregation projects that implement narrow custom models. This document refers to aggregation using data modeled according to intelligent frameworks which incorporate the possibility of a wider context and are based on cross disciplinary expert knowledge and description. Typically but not exclusively, this document refers to the CIDOC CRM and extensions of it. The CRM provides a richer but smaller entity model compared to the source models that are mapped to it. An identified issue is that aggregators, despite having expertise in the use of the target schema, have significant gaps in their knowledge about established curatorial practices. This situation can lead to a mismatch of semantics between the provider and the aggregator and therefore the data provision process assumes a greater level of contact with provider representatives than is currently the case. This relationship has a direct effect on increasing the quality of the aggregation.

Aggregator Terminology Expert

The Aggregator Terminology Expert is the curator, maintainer or other expert on one or more of the terminologies that the aggregator uses as reference in the Integrated Access System. Aggregators normally want to avoid engagement with terminology maintenance. They often use more generalized provider independent terminologies rather than (more rarely) take over provider term lists. However, it should be noted that the richness of an aggregation may include the ability to understand why different terminology has been applied in different organizations and that this is also a means by which things can be re-contextualized, to a certain extent, to the real and historical world. Established local terminology should never be replaced by centralized and generalized terminologies, but be correlated to the latter

Ingest Manager

The Ingest manager is responsible for receiving data from the provider and ingesting data into the target system.

Instance Generation Expert

The expert of the aggregator, normally an IT specialist, who is responsible for maintaining the referential integrity of the (meta)data in the Integrated Access

System and who knows how to generate from provider data valid URIs for the Integrated Access System.

Schema Matching Experts

Provider schema experts and an aggregator schema expert collaborate in order to define a schema matching, which is documented in a schema matching definition.

4 Data Objects

We define three categories of information objects that take part in the data provisioning process, (a) the content data and metadata objects, (b) the schema and logic objects and (c) the control objects which are illustrated in Figure 3 and described in detail in the following sections.

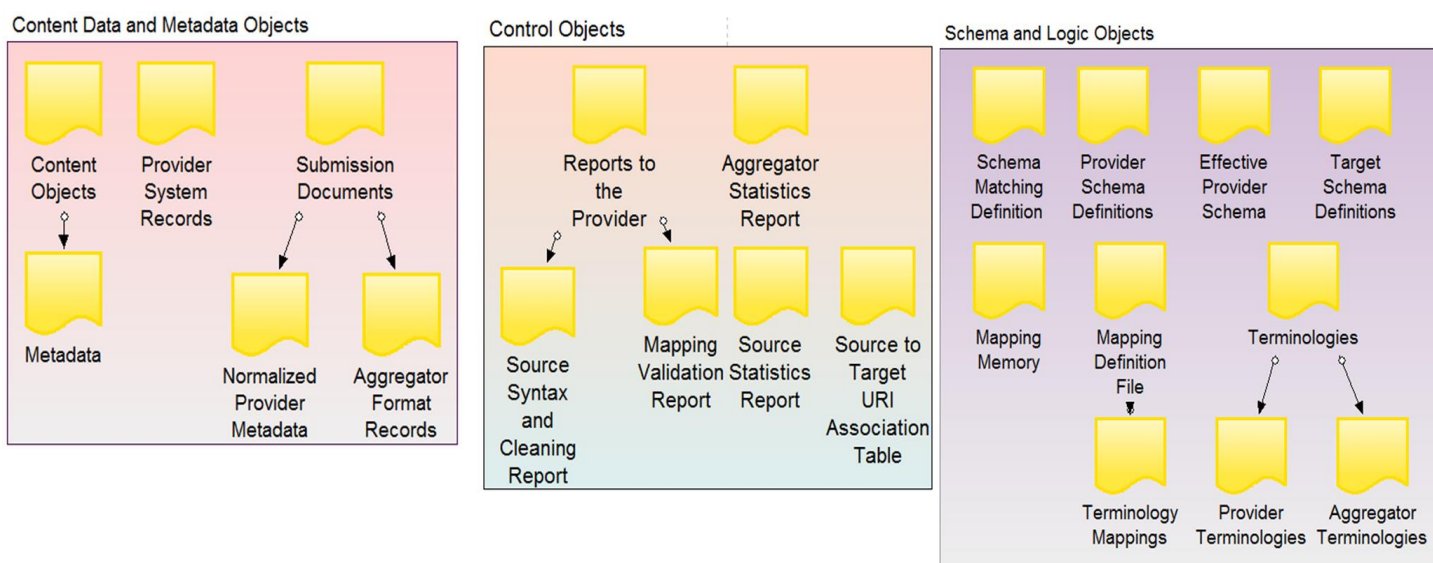


Figure 3: Data Objects

4.1 Content Data and Metadata Objects

Content data and metadata objects consist of all the raw data and metadata. In detail they include:

Content Objects

Individual files or information units with an internal structure that is not described in terms of schema elements of the source or target systems. These are things like images or text documents that are searched with content retrieval indices such as keyword searches, rather than by associative queries. As objects they are typically described by metadata records which are searched by associative queries. What is important in this context is that a content object is not identified by the actual structure it contains but by the way it is treated in the information system (stored either as “blobs” in the database or using references to a file system). Many aggregators do not collect content objects but only link to those resources in the

provider system. If these objects are to be collected then they should be referenced in the mapping and transferred to the target system with the appropriate URL specified in the mapping.

Metadata

Metadata are information units with an internal structure that is described in terms of schema elements of the source or target systems. In our context, these are often data records describing a content object (therefore the term “metadata”). However, bad analogy has also brought the term into use for data describing physical objects. Therefore we define “metadata” here in the same way it is treated in the information system, and not as “data about data”. The metadata records are the common subject of submission to aggregators and therefore of transformation from the source to the target schema.

Submission Documents

The Submission Documents are well-formed documents that are produced either by the syntax normalization process, or by the transformation process.

Provider System Records

The Provider System records are records of the Local Information System. In case they contain fields with local, informal or uncontrolled internal syntax, such as frequently occurring for commented dates, dimensions, or uncertain information, we speak about **Metadata**.

Normalized Provider Metadata

Normalized Provider Metadata is the result of formalizing (“cleaning”) Metadata by extending the provider schema. The result is completely structured data in the effective provider schema, in which each structural element must have a clearly described meaning, i.e., all local methods to subdivide a name or string into meaningful subsections should be expressed by explicit, unambiguous tagging, preferably in XML.

Aggregator Format Records

The Aggregator Format Records are the records in the form to be ingested into the target system.

4.2 Schema and logic objects

Schema and logic objects consist of the schemata, mappings and terminologies of both the source and target systems. In detail they include:

Schema Matching Definition

The Schema Matching Definition contains the mappings of the source schema elements to the target schema paths. This definition must be human and machine readable and is the ultimate communication means on the semantic correctness of the mapping.

Provider Schema Definitions

These include data dictionaries, XML schemata, RDFS/OWL files etc. describing the data structures that are managed and can be searched by associative queries in the source system.

Effective Provider Schema

The Effective Provider Schema is the new source schema definition that comes up, in case local syntax rules exist.

Target Schema Definitions

These include data dictionaries, XML schemata, RDFS/OWL files etc. describing the data structures that are managed and can be searched by associative queries in the target system.

Mapping Memory

A collection of mapping histories of analogous cases collected from the user community.

Mapping Definition

The Mapping Definition comes up by the addition of the instance generation policies to the Schema Matching Definition.

Terminologies

We regard as Terminologies controlled vocabularies of terms that appear as individual data values in the source or target systems and represent *categorical concepts* or “universal”. We do **not** regard reference information about places (gazetteers) and people as terminologies in this document. Matching people and places are regarded as cases of “co-reference resolution” in this document. The term “vocabulary” is **not** used for metadata schemata in this document. Terminologies may be flat lists of words or be described and organized in more elaborate structures as so-called “thesauri” or “knowledge organization systems”, the most popular format now being SKOS. This document distinguishes these structures from an “ontology”, even if the terminology may qualify as such, as long as its use in this context is to provide *data values* and not data structure.

Aggregator Terminologies

The terminologies used by the aggregator as a reference in the Integrated Access System.

Provider Terminologies

The terminologies used by the provider as a reference in the local system.

Terminology Mappings

Terminology Mappings are expressions of exact or inexact (broader/narrower) equivalence between terms from different vocabularies.

In this context, we are primarily interested in the mapping of terms that appear directly or indirectly in mapping conditions of a Schema Matching Definition. In such a mapping condition, a term in the source record is equal to or unequal to a constant, or a narrower term of a constant. This may be expressed in terms of source or target terminology (note that we regard as terminology **only** categorical concepts, not names of particular things, events, places or people).

For instance, take a source schema with a table "Object" and field "object type", to be mapped to the CIDOC CRM. The source schema does not distinguish material from immaterial objects. The target schema we map to however makes the distinction. Then, a source field with value "object type = Vase" may indicate a "Physical Object", and "object type = Image" an "Information Object". As a consequence, the value in the field "object type" determines two alternative interpretations of the table "Object". This is a conditional mapping, in which the mapping of a source schema element depends on the value in some other element of the source record. In general, only categorical terms should affect the schema matching logic.

4.3 Control Objects

Control data objects are the reports and documents that support the data provisioning process and are the products of its different sub-processes. In detail they include:

Reports to the Provider

They include reports useful to the provider in order to monitor the result of the various tasks and to announce all individual inconsistencies in the processed source data which the provider may or should correct.

Source Syntax Report and Cleaning Report

The Source Syntax Report and Cleaning Report is the output report of the syntax normalizer. It contains inconsistencies and errors that occurred during the syntax normalization process.

Source Statistics Report

The Source Statistics Report is the output statistics of the Source Analyzer, used as input to the Source Schema Visualizer. It contains statistic information useful for understanding the provider schema.

Mapping Validation Report

The Mapping Validation Report is the output report of the Metadata Validator Transformer. It contains errors and inconsistencies that occurred during the transformation process.

Aggregator Statistics Report

The Aggregator Statistics Report is the output statistics of the Target Analyzer. It contains statistic information useful for understanding the target schema.

Source to Target URI Association Table

The Source to Target URI Association Table is the output report of the metadata validator transformer tool. It contains source to target URI associations.

Figure 4: Data Provisioning Flow Network

Starting from the Provider Institution, the **Syntax Normalizer** normalizes the provider's records and produces the **Effective Provider Schema** and a report with the errors that occurred during the normalization process.

The next step is performed by the **Schema Matcher** and is completed with the definition of the **Schema Matching Definition**. Different IT objects may assist the user to define the mappings. The **Source Analyzer** provides useful statistics for each field, whereas the **Source Schema Visualizer** and the **Target Schema Visualizer**, respectively, help the user to navigate through all source and targets elements. The **Schema Matcher** is supported with mapping suggestions, provided by the **Mapping Suggester** and produces the Schema Matching Definition. This definition may be viewed with the **Schema Mapping Viewer**.

Subsequently, the **Instance Generation Rule Builder** produces the URI policies and complements the Schema Matching Definition producing the **Mapping Definition**.

The matching may need to interpret provider and aggregator terminologies in order to resolve data dependent mappings. This may be assisted by a **Terminology Mapper**.

Finally, when the mapping definition and the terminology mappings are defined, the **Metadata Validator Transformer** transforms the records to the aggregator format and ingests them to the Aggregator Institution. Since it is usual for schema documentation to be embedded into the structure of the schema the **Target Analyzer** should be able to identify this documentation and expose it to the user when browsing aspects of the target.

6 Analytical Representation of the Data Provisioning Process

In this section we will describe in detail the processes that are involved in data provisioning. The modeling of the data provisioning components and processes is made using Adonis- Business Process Management [1]. The data provisioning process hierarchy is presented in Figure 5.

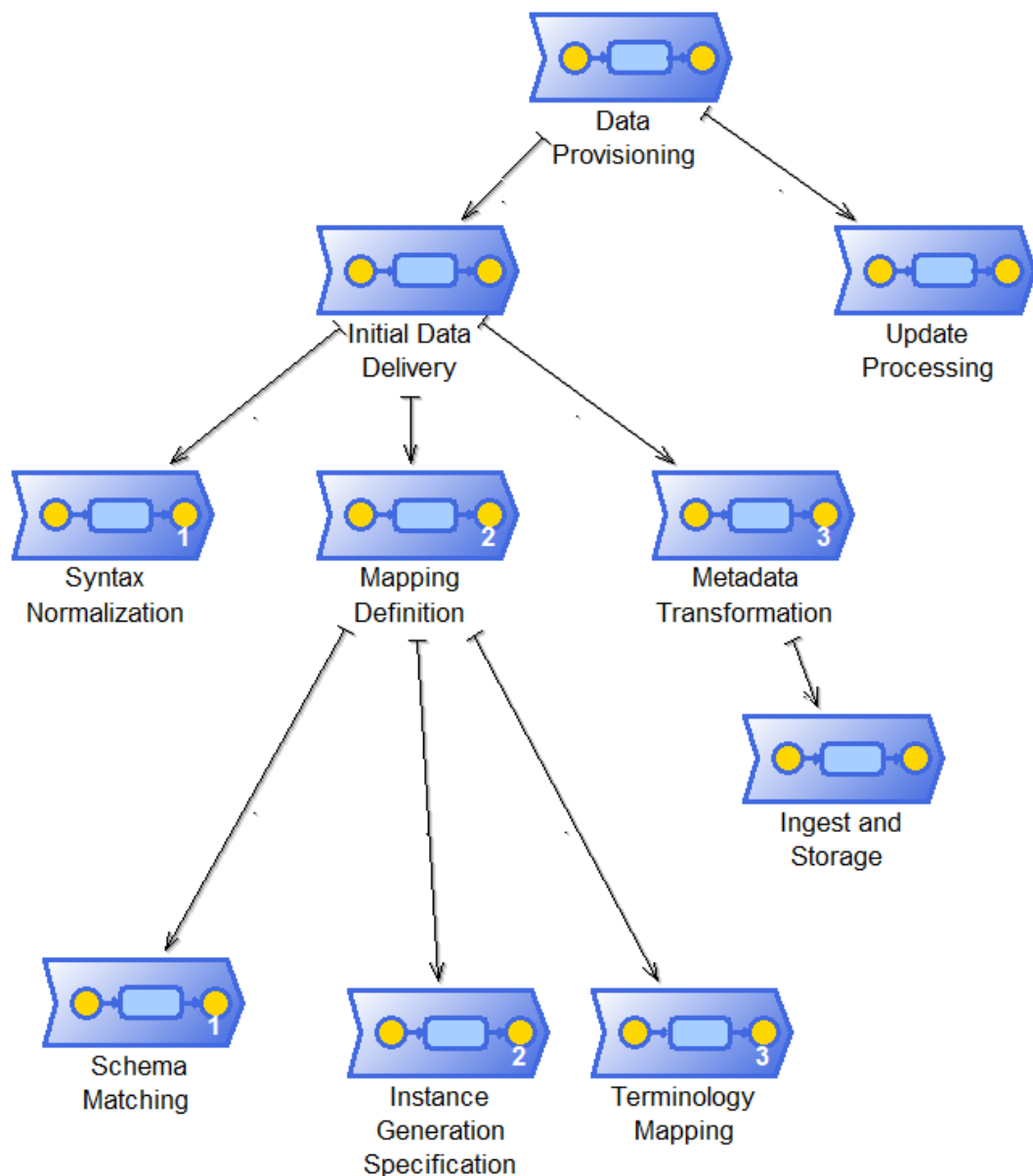


Figure 5: Data Provisioning Process Hierarchy

6.1 Data Provisioning Process

The starting point of the Synergy Model is the Data Provisioning process (Figure 6) which deals with the selection and scheduling of data, including co-reference resolution and updates. A *Mapping Manager* may be responsible for this task that may form part of the agreement with the aggregator that may have been assigned to named representatives on both the provider and the aggregator.

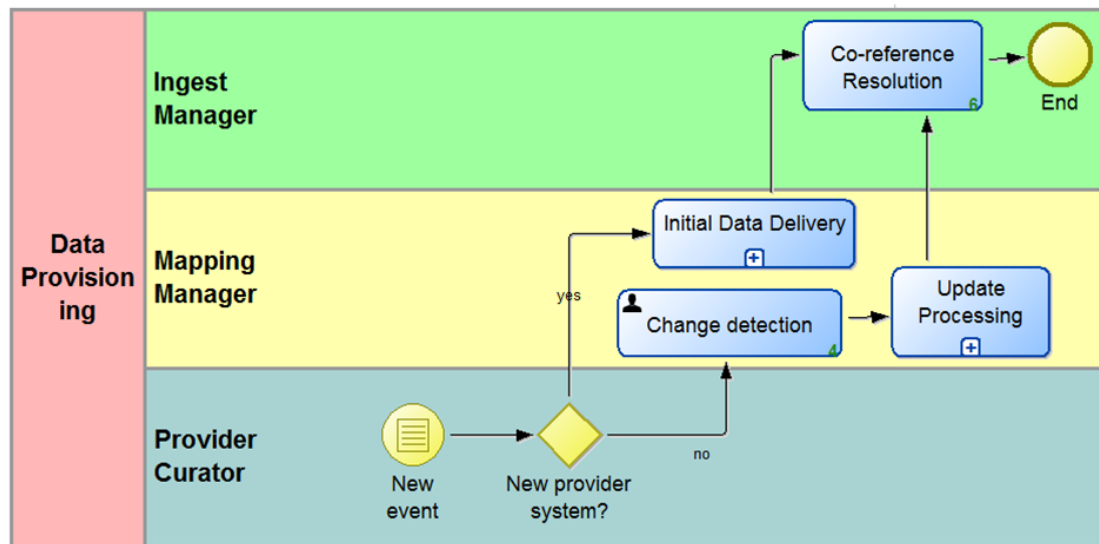


Figure 6: The Data Provisioning process

While it is inevitable that the Synergy Model will be unable to deal with all internal provider processes, it should provide general support through the provision of log files and reports. Reports may be based on queries that include changes of provider system records and other information that indicates the need to resubmit data. Queries would also support provider selection criteria in the source systems. On the target system they may be used to reveal semantic needs in the composition of the aggregation and to derive requests for particular or additional materials from providers.

During the data provisioning process, the Ingest Manager may proceed with referential integrity processing (***co-reference resolution***) i.e. resolving multiple identifiers that denote the same real world thing or object of discourse (co-reference resolution). This is a process in its own right and will be the subject of separate documentation. The main goals are to ensure referential integrity, which is the heart of information integration, and to reduce the number of URIs in use for the same thing. It requires its own dialogue between provider, aggregator and third-party authority managers. Since the aggregator collects more comprehensive knowledge than the providers, it is a natural role of the aggregator. One may regard that the only genuine knowledge of the aggregator, because of the nature of an integrated system, is the co-reference knowledge.

On the other hand, there is a set of characteristic *changes* in the provider – aggregator environment that affect the mapping and may require:

- Re-executing the transformation of records already submitted to the aggregator and updating the transformed records in the target system.
- Resubmission of records from the source system.
- Redefinition of the mapping.

The Mapping Manager must monitor such changes and initiate respective actions. On some occasions, transformation may be affected by a significant change in underlying technology platforms and, while in theory this should not affect the way in which data is mapped to the target schema, such cases may inadvertently prompt changes in the mapping definition and/or require significant re-alignment. This should not be an issue for the system being described.

Table 1 depicts a summary of the tasks presented in Figure 6.

Name	Type ²	Description	Role
Initial Data Delivery	Sub-p	Initial Data Delivery contains 3 sub-processes: Syntax Normalization, Mapping Definition & Metadata Transfer	
Change Detection	T	Monitor the transferred data in order to maintain data up-to-date with whatever relevant changes occur in the source and target systems and the employed terminologies.	Mapping Manager
Update Processing	Sub-p	Restore ability of data transformation and semantic consistency, which comprises changes in the source target records, in the provider or aggregator terminologies, in the source or target schemata, in the target URI policy and in the good practice of interpretation of source and target schema in the mapping definition.	
Co-reference Resolution	T	Referential integrity processing: resolves multiple identifiers that denote the same real world thing or object of discourse. The main goals are to ensure referential integrity, which is the heart of information integration, and to reduce the number of URIs in use for the same thing. It requires its own dialogue between provider, aggregator and third-party authority managers. Since the aggregator collects more comprehensive knowledge than the providers, it is a natural role of the aggregator. One may regard that the only genuine knowledge of the aggregator is the co-reference knowledge.	Ingest Manager

Table 1: Summary of the Data Provisioning Process

² **Sub-p**: A Sub-Process is an activity whose internal details have been modeled in a separate model.

T: A Task is an atomic activity within a process flow. It is used when the work in the process cannot be broken down to a finer level of detail.

6.1.1 Initial Data Delivery

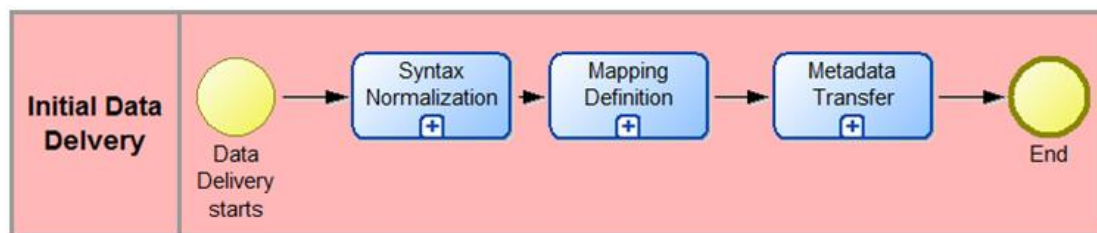


Figure 7: Data delivery sub-process

Initial Data Delivery breaks down into:

- *Syntax Normalization*
- *Mapping Definition*
- *Metadata Transfer*

Table 2 depicts a summary of the tasks presented in Figure 7.

Name	Type	Description
Syntax Normalization	Sub-p	Syntax normalization aims to convert all data structures relevant for the transformation in a standard form since data transformation tools can only deal with a limited set of standard data structures.
Mapping Definition	Sub-p	Mapping definition is the specification of the parameters for the data transformation process, such that complete sets of data records can automatically be transformed, manual exception processing notwithstanding. This includes harmonization between multiple providers.
Metadata Transfer	Sub-p	The actual transfer of data until a first consistent state is achieved. This includes transformation of sets of data records submitted to the aggregator, the necessary exception processing of irregular input data between provider and aggregator, ingestion of the transformed records into target system and initial referential integrity processing possibly on both sides.

Table 2: Summary of the Data Delivery sub-process

6.1.1.1 Syntax Normalization

Syntax normalization aims to convert all data structures relevant for the transformation in a standard form since data transformation tools can only deal with a limited set of standard data structures and thus any non-standard form must be converted to a standard one.

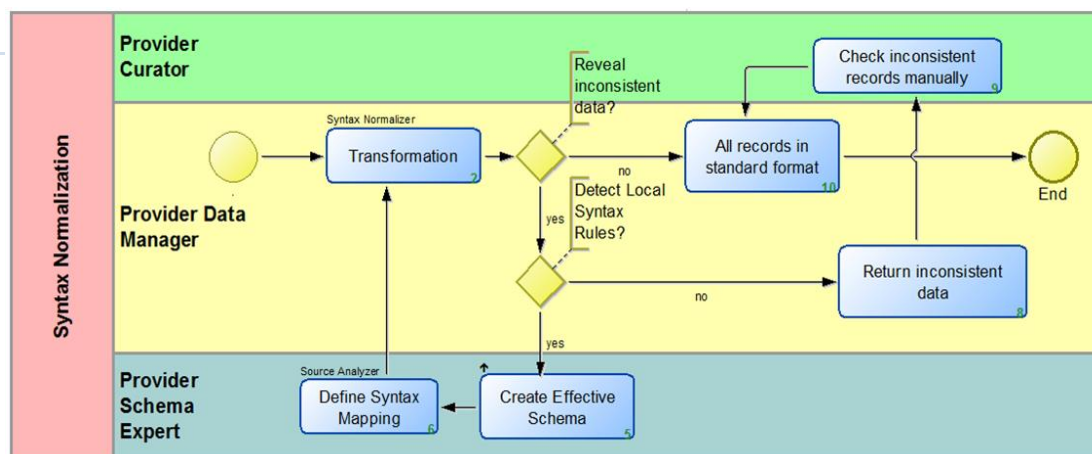


Figure 8: Syntax Normalization sub-process

Although most museum systems now employ structured formats and use relational databases, spreadsheets, XML or even RDF, some systems are still dependent upon unstructured storage such as text documents. Even within structured systems there can be issues related to the use of unstructured text fields and misuse of other fields without validation. Automated data transformation, (i.e. transformation of data from one schema to another without loss of meaning or with controlled loss of meaning by a deterministic algorithm based on a mapping definition), is only possible if the data to be transformed is completely structured. Unstructured storage is out of scope for the mapping system and museums will be encouraged to create structured systems for their data. For data issues within structured systems that require normalization, the following approaches can be used:

1. The institution resolves these issues in the source database or, if exporting to another format, includes syntax normalization as part of the export process.
2. The mapping system provides a library of syntax normalization routines that can be used by the user.
3. The system provides for new syntax normalization routines that can be created by the institution or their agents. These can potentially become part of a central library or resource for collaborative use.
4. The mapping memory suggests an alternative method for mapping the data that does not capture the full semantics.

Local identifiers may have their own syntactic structure, such as inventory numbers, addresses, bibliographic references, date and time etc. It may not be worthwhile to normalize their internal structure to XML prior to the mapping, but instead use specific and additional scripts for instance generation.

The syntax normalization can be done by a technology expert, possibly the same one dealing with instance generation, in collaboration with the source schema expert. Local syntax rules can be so complicated or even deterministic that it is often more effective to use a set of custom filtering routines, resolving one structural feature at a time, and verifying it with the source schema expert. For instance, if italics are used to tag a particular kind of field, it is better to convert first all italics to XML tags.

This process will reveal inconsistencies and alert the provider to issues that may require attention and need to be resolved. There will be a residual of cases so complicated to describe by rules, that manual rewriting is more effective. This is not a bottleneck, or a reason not to proceed. The important thing is to drastically reduce the number of records to be checked and treated manually. Therefore it is equally important to find diagnostic rules for inconsistent cases, as it is to resolve those that can be formally described. If within a system of syntax normalization some inconsistent cases “slip through” undetected, all records may have to be reviewed manually. That would indeed become a bottleneck. Ultimately no mapping tool can mitigate all internal data management issues and organizations wishing to participate in big data initiatives. Organizations would need to address inconsistencies that perhaps have been historically ignored when data was viewed as simply an internal inventory in a closed system.

After syntax normalization, we expect all data structures relevant for the transformation to be in a standard form. Note that in this step NO approximation of the target schema semantics should be attempted. Rather, it must be an exact representation of the data as understood by the provider institution. The CRM is unique in that it represents data as it is understood by the owning organization and does not impose constraints on meaning. There should be no need to manipulate the conceptualization of information to suit the aggregator’s model. The idea is to capture and bring to the fore the semantics as seen and intended by the provider independently of the aggregator.

Table 3 depicts a summary of the tasks presented in Figure 8.

Name	Type ³	Description	Role	IT Object	Output
Transformation	T	Convert data to the structural format described by the effective schema. (syntactic transformation)	Provider Data Manager		
Reveal inconsistent data	EG	Reveal inconsistent data, if any, that need to be mended.			
Detect Local Syntax Rules	EG	Check if local syntax rules exist.			
Create Effective Schema	T	Create a new provider schema definition which contains the formal description of the local syntax rules.	Provider Schema Expert		Effective Provider Schema
Define Syntax Mapping	T	Define the correct syntax mapping according to the new schema.	Provider Schema Expert		
Return inconsistent data	T	Return inconsistent data back to the provider in order to manually check them.	Provider Data Manager		
Check inconsistent records manually	T	Manually review and correct inconsistent data.	Provider Curator		
All records in standard format	T	After syntax normalization, we expect all data structures relevant for the transformation to be in a standard form.	Provider Data Manager		

Table 3: Summary of the Syntax Normalization sub-process

³ **EG**: Exclusive Gateway routes the sequence flow to exactly one of the outgoing branches, when splitting. When merging, it awaits one incoming branch to complete before triggering the outgoing flow.

6.1.1.2 Mapping Definition

Mapping definition is the specification of the parameters for the data transformation process, such that complete sets of data records can automatically be transformed, manual exception processing notwithstanding. This includes harmonization between multiple providers.

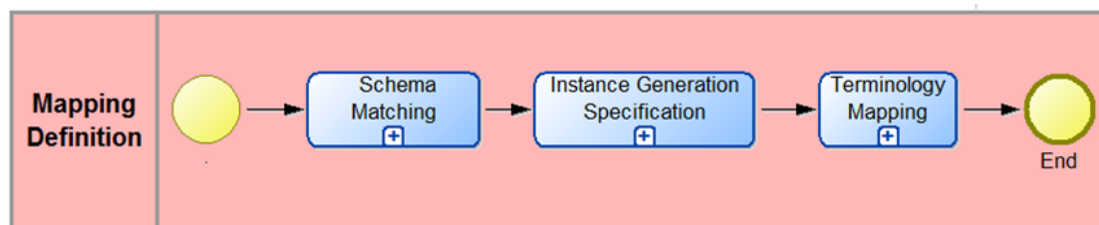


Figure 9: Mapping Definition sub-process

Mapping definition consists of the:

- Schema matching
- Instance generation specification
- Terminology mapping

The Mapping Manager may be responsible for issuing and coordinating these tasks.

Table 4 depicts a summary of the tasks presented in Figure 9.

Name	Type	Description
Schema Matching	Sub-p	Source schema experts and a target schema expert define a schema matching which is documented in a schema matching definition. In order to do so, all source schema elements must be well understood and mapped to target schema paths.
Instance Generation Specification	Sub-p	Define the URI generation policies for each instance of a target schema class referred to in the matching.
Terminology Mapping	Sub-p	Define the terminology mappings between source and target terms.

Table 4: Summary of the Mapping Definition sub-process

6.1.1.2.1 Schema Matching

The *provider schema experts* together with a *target schema expert* (e.g., a CIDOC CRM expert) define a schema matching which is documented in a ***schema matching definition***. This definition must be human and machine readable and is the ultimate definition of the semantic correctness of the mapping. The collaboration between these experts must be well organized and is the bottle-neck of the data provisioning process.

In order to define a schema matching, all source schema elements must be well understood and mapped to equally well understood target schema paths. Both tasks need two Independent tools to visualize source and target schemata and source metadata records. Adequate navigation and viewing techniques must facilitate both overviews and an understanding of the details.

The matching process must lead the user through all source elements in order to make a mapping decision. This may be supported by tools *suggesting mappings* (automated mapping). The automated mapping tools should recalculate their proposals with each new mapping decision. They should make use of “*mapping memories*” of analogous cases collected from the user community. An aggregator may maintain *mapping guidelines* together with provider and user consortia.

As described in section 4.2, the schema matching may need to interpret provider or aggregator terminologies in order to re-solve data dependent mappings (where values might determine the mapping). In the *schema matching definition*, we generally foresee *mapping conditions* that a term in the source record is equal to or unequal to a constant, or a narrower term of a constant. This may be expressed in terms of *source or target terminology*. In order to resolve these specifications at *record transformation time*, partial *terminology mappings* of source and target terminology must exist and may be linked to mapping conditions. The terminology mapping needs to be done only to the degree needed to resolve the conditions of the schema matching. If the provider terminology is hierarchical, the effort can be drastically reduced.

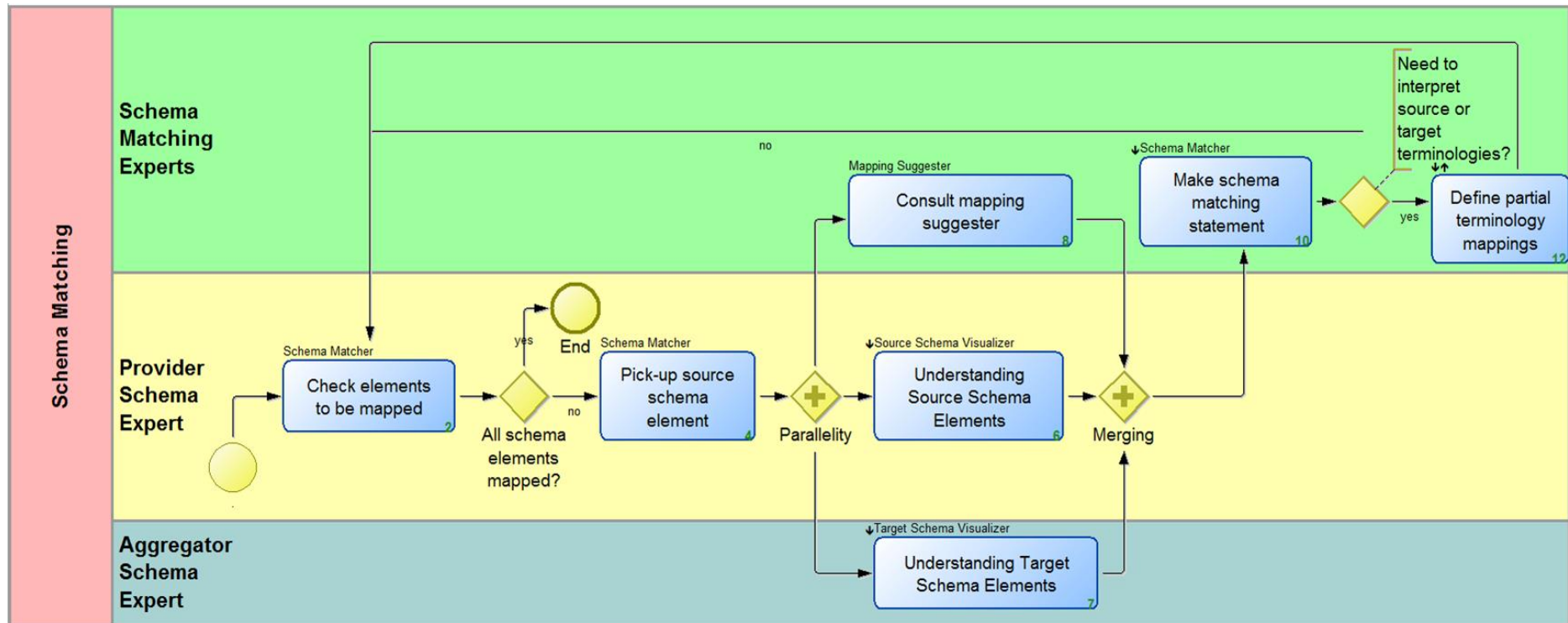


Figure 10: Schema Matching sub-process

Consider a situation in which a provider has a flat authority with different terms that require different mapping. The potential solutions are as follows:

1. Conditions are provided on an individual basis or, where the majority of terms conform to a particular mapping conditions are provided for the exceptions.
2. The terms are grouped into the types that determine their mapping approximating to a single level thesaurus. Conditions can then be applied to those groups.
3. The authority is re-organised into a fully realized thesaurus and conditions applied to branches of the hierarchy.
4. The authority is mapped against the authority/ thesaurus of the aggregator.

As stated above this mapping may only need to be partial in order to resolve particular automated mapping decisions. Once the mapping is complete then additional semi-automated co-referencing can take place using reasoning across all the aggregator's datasets.

All related tools should take into account the need for *incremental mappings* after source or *target schema definition* up-dates, *terminology* updates and *mapping guideline* updates and guide the user through relevant changes. These could be highly sophisticated and granular or relate to modifications at a record level triggering a full record update that would include the update.

Some source data or source schema elements may not allow for matching decisions or create a mapping with more general semantics than might be ideal. In situations where the mapping is highly generalised (and not in accordance with mapping memory) the mapping may be automatically identified as needing improvement. In some situation the mapping cannot be made with an improvement to the source by the provider. It must be possible to define filters for these data that run before and at transformation-time and feed back to the provider

Table 5 depicts a summary of the tasks presented in Figure 10.

Task Name	Type	Description	Role	IT Objects	Input / Output Document
Check elements to be mapped	T	Check schema matching definition in order to examine if all source schema elements are mapped to target schema paths.	Provider Schema Expert	Schema Matcher	
All schema elements mapped	EG	This process continues till all source elements are mapped to a target path.			
Pick-up source schema element	T	Select the source schema path to be mapped.	Provider Schema Expert	Schema Matcher	
Understanding Source Schema Elements	T	All source schema elements must be well understood. It is important to use tools for the visualization of source schema and source metadata	Provider Schema Expert	Source Schema Visualizer	Input: 1. Normalized Provider Metadata 2. Provider Schema Definitions 3. Effective Provider Schema
Understanding Target Schema Elements	T	All target schema elements must be well understood. It is important to use tools for the visualization of target schema.	Aggregator Schema Expert	Target Schema Visualizer	Target Schema Definitions
Consult mapping suggester	T	The mapping decision may be supported by tools suggesting mappings. The automated mapping tools should recalculate their proposals with each new mapping decision by some user. They should make use of “mapping memories” of analogous cases collected from the user community. An aggregator may maintain	Schema Matching Experts	Mapping Suggester	

		mapping guidelines together with provider and user consortia.			
Make schema matching statement	T	Map a source schema element to a target schema path.	Schema Matching Experts	Schema Matcher	Input: 1. Effective Provider Schema 2. Provider Schema Definitions
Need to interpret provider or aggregator terminologies	EG	The matching may need to interpret provider or aggregator terminologies in order to re-solve data dependent mappings. For instance, a field “object type = Vase” may indicate a “Physical Object”,and “object type = Image” an “Information Object”. Such difference will lead to a thorough reinterpretation of most other fields describing such an object.	Provider Curator		
Define partial terminology mappings	T	In the schema matching definition, we generally foresee mapping conditions that a term in the source record is equal to or unequal to a constant, or a narrower term of a constant. This may be expressed in terms of source or target terminology. In order to resolve these specifications at record transformation time, partial terminology mappings of source and target terminology must exist.	Schema Matching Experts		Input: 1. Aggregator Terminologies 2. Provider Terminologies Output: Terminology Mappings

Table 5: Summary of the Schema Matching Definition sub-process

6.1.1.2.2 Instance Generation Specification

After the matching process, an appropriate URI schema must be applied for each target class instance. These URIs are ones defined by a combination of information including the namespace used by the provider, the type of URI (object, terminology, etc.) and the mapping function used. It may also be affected (customized) by particular policies of the provider using language that best fits the data, typically defined by information managers. Some URIs may be based on third party URI definitions and may require a look-up against appropriate online resources. Changes in provider instance generation policies may result in changing the definitions without affecting the schema matching. As a result, the application of URIs may be a combination of both automatic and manual steps.

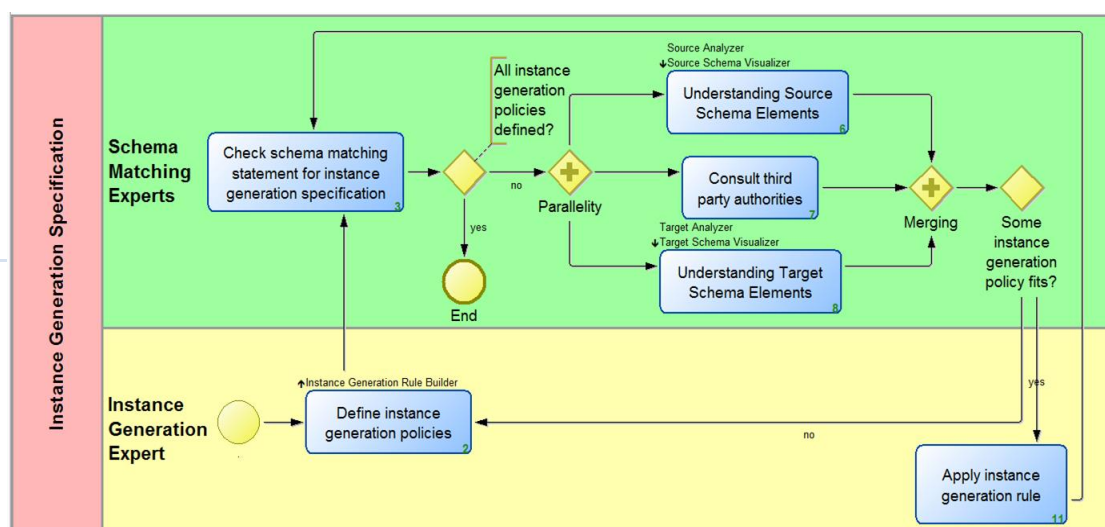


Figure 11: Instance Generation Specification sub-process

Two situations may exist. Firstly, the provider is interested in representing their data using their own URI schema. For example, they may be using the data internally as part of their own data harmonization strategy and may also wish to create their own reusable interface. In this situation they may wish to use their own URI schema. Alternatively they may be happy to use the schema applied by the aggregator who may have a URI policy for the one or the other category of items acceptable for the provider. It is likely that the former will become more common and the provider's data should be represented using the providers URIs. If different aggregators use different instance generation policies then co-reference resolution would be needed between URIs that effectively hold the same information. It is *not the intention* of this model to propose a best practice of URI policies, but to make *the distinctions necessary* to handle consistently and globally resolve the effects of reasonable URI policies.

It is already the case that data providers who publish linked data use the URIs of third parties instead of minting their own URIs. For example, third party URIs are used for

vocabularies and ontology support. URI generation policies that use third party resources (“authority records”) about persons or places at transformation time should ensure that adequate exception handling is built in, in case of lookup failure.

The execution of URI generation rules may also reveal inconsistent or not normalized data of the provider at transformation time, or before. Inconsistent data filters must be foreseen in the generation rules. The source metadata records may be analyzed before transformation time for such cases. Providers must be informed about inconsistent cases, and given the possibility to run an organized, sustainable process to improve the source data. It may be possible to define preliminary workarounds to maintain the submission process, i.e. characteristic data patterns replacing dirty data that can later be recognized and updated at aggregator side without resubmission. Otherwise, inconsistent records are held back until they are updated.

Changes of instance generation policies of the aggregator may result in the need to update the URI generation rules. Changes of naming and identifier policies at the provider side may also make a redefinition of URI generation rules necessary. It must be possible to do that without affecting the schema matching definition file.

Table 6 depicts a summary of the tasks presented in Figure 11.

Name	Type	Description	Role	IT Objects	Input / Output Document
Define URI generation policies	T	The URI generation policies for each instance of a target schema class referred to in the matching must be defined, such as for persons, objects, events, place, and formats of time. The URI generation policies can be introduced in an abstract form as rules or references to code signatures implementing specific rules.	Instance Generation Expert	Instance Generation Rule Builder	Output: Mapping Definition File
Check schema matching statement for instance generation specification	T	Examine each instance of a target schema class to apply a URI generation rule.	Schema Matching Experts		
All instance generation policies defined?	EG	Check each instance of the target schema class, if a URI generation rule is applied.			
Understanding Source	T	All source schema elements must be well	Schema Matching	1. Source	Input: 1.Normalized

The Synergy Reference Model

Schema Elements		understood. It is important to use tools for the visualization of source schema and source metadata records.	Experts	Analyzer 2. Source Schema Visualizer	Provider Metadata 2.Provider Schema Definitions 3.Effective Provider Schema
Consult third party authorities	T	Some URI generation policies may include look-up of on-line resources (“authority records”) about persons or places.	Schema Matching Experts		
Understanding Target Schema Elements	T	All target schema elements must be well understood. It is important to use tools for the visualization of target schema	Schema Matching Experts	1. Target Analyzer 2. Target Schema Visualizer	Input: Target Schema Definitions
Some instance generation policy fits?	EG	Some instance generation policy fits?			
Apply instance generation rule	T	Apply the instance generation rule, in case a instance generation policy fits.	Instance Generation Expert		

Table 6: Summary of the Instance Generation Specification sub-process

6.1.1.2.3 Terminology Mapping

Terminology mapping can be a huge task. Providers may use anything from intuitive lists of uncontrolled terms up to highly structured third party thesauri. However, most of the provider terminology is very specialized and more important as information element when metadata records are displayed than as search term in the target system. As long as the terms are in the same natural language, most terms can just be copied from the source records into the transformed target records. If they are in other languages, aggregators may choose to translate terms, possibly preserving also the provider terms. It may be useful to associate provider terms with broader terms of some standard terminology the aggregator employs as search terms. Since all this can happen even after metadata record transformation, it does not affect the mapping process itself.

In this model, we are only interested in the consistency of the mapping process when the choice of a target class or property depends on a term. For that sake, we can extract from the schema matching definition the terms appearing in mapping conditions. We distinguish two cases:

- **equality/inequality condition:** These terms (constants) must be taken from the provider terminology and no action is needed.
- **broader term condition:**
 1. If the constant term is given in the provider terminology, the narrower term hierarchy for each constant term is used if it exists, otherwise it must be “invented”. The latter is one case of terminology mapping.
 2. If the constant term is given in the aggregator terminology, for each

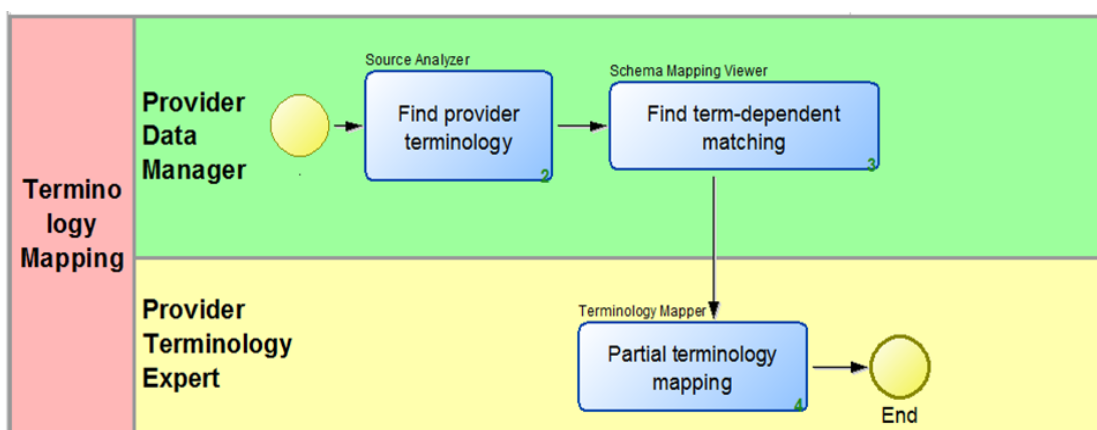


Figure 12 Terminology Mapping sub-process

3. constant term the narrower terms in the source terminology must be identified. This is the second case of terminology mapping.

In any case, the aggregator terminology should have a thesaurus structure, albeit a small vocabulary of high-level terms. Sometimes it may be more effective to merge provider terms with aggregator terms, i.e., replace equivalent terms and insert all other provider terms as narrower terms of aggregator terms.

In this case, the term values used to execute the schema matching conditions of the provider terminology should be replaced by the updated aggregator terminology before transforming the respective records and in the schema matching definition file. This will allow for better controlling the mutual consistency of mappings between different providers. Notwithstanding, the original provider terminology could, possibly should, be added to the source records and be carried over to the target records in separate fields/properties.

The terminology mapping may reveal inconsistent data of the provider, such as spelling errors or unauthorized terms. Inconsistent data filters must be foreseen for terminology. The source metadata records should be analyzed before transformation time for such cases. Providers must be informed about inconsistent data cases, and given the possibility to run an organized, sustainable process to improve the source data.

It may be possible to define preliminary workarounds to maintain the submission process, i.e. characteristic data patterns replacing dirty data that can later be recognized and updated at aggregator side without resubmission. Otherwise, “dirty records” are held back until they are updated.

Mapping identifiers of persons and places to those used at the aggregator system is in general ineffective due to the large number of such identifiers, most of which are only known at local level. It is more effective to update the provider with identifiers (URIs) of persons and places referred to by more than one provider (or third party authority, such as viaf.org). After these steps, metadata records are ready for transformation.

Table 7 depicts a summary of the tasks presented in Figure 12.

Task Name	Type	Description	Role	IT Objects
Find provider Terminology	T	The provider data manager uses the source analyzer tool in order to extract the terms appearing in the source schema.	Provider Data Manager	Source Analyzer
Find term-dependent matching	T	Extract from the schema matching definition the terms appearing in mapping conditions.	Provider Data Manager	Schema Mapping Viewer
Partial terminology	T	Define a partial terminology mapping of source and target	Provider Terminology Mapper	Terminology Mapper

mapping		terminology.	gy Expert	
---------	--	--------------	-----------	--

Table 7: Summary of the Terminology Mapping sub-process

6.1.1.3 Metadata Transformation

Once the mapping definition has been finalized (and all syntax errors resolved) the data and mapping information needs to be submitted, transformed and stored in the aggregators system.

The mapping manager will be informed of the submission to initiate the transformation process, provide final validation and store both the raw data (optional) and transformed data into the target system. The submitted metadata records must be identified with a unique identifier, checksum and modify date-time. The submission management must be able to recognize any change of a source record by the metadata record metadata.

Transformation

The transformation process itself may run completely automatically. However, it is possible that further issues not realized by the provider will materialize. These are issues that the aggregator may be able to resolve on a temporary or permanent basis but in any event may require that records are referred back to the provider for further analysis and resolution. It is possible (given sufficient trusted expertise) that personnel under the aggregator mapping manager can take manually correct issues in an interactive process. The result is a set of valid target records.

Instance Matching

The instance generation algorithm of the automatic data transformation process may employ an initial instance matching process in order

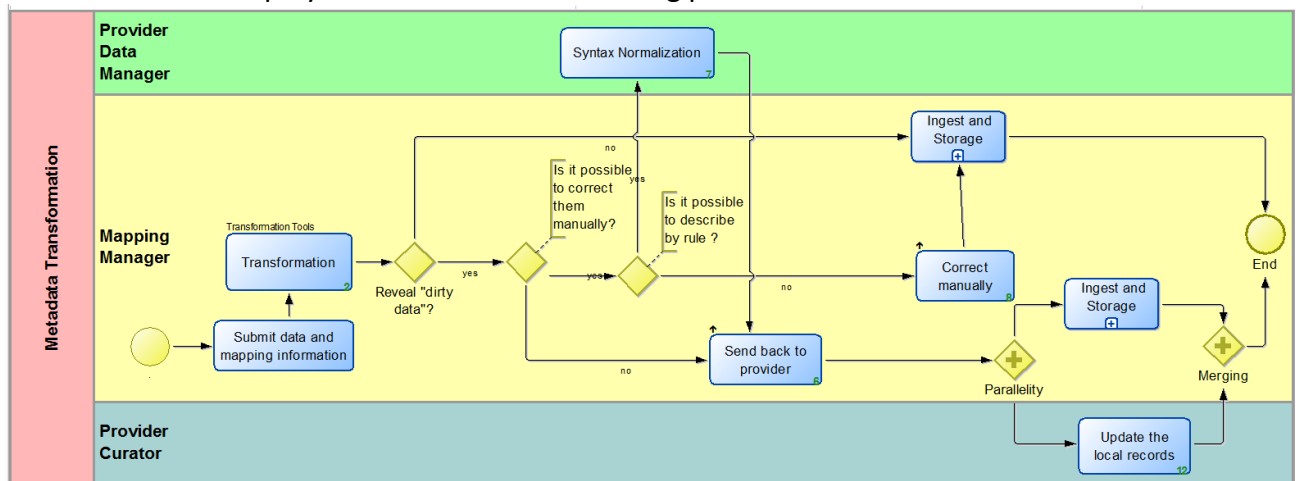


Figure 13: Metadata Transformation sub-process

to reuse existing URIs in the target system. This also holds for third party authority systems with URIs that the aggregator uses as reference (such as viaf.org). In any case, such matching should be reported back to the provider for potential internal

use of such URIs. The provider should not replace local identifiers of items under his authority of knowledge (collection objects, local persons, places, events) with ones matched in the aggregators system. Alternative URIs identified by the aggregator should only be used in addition to established local data for which the provider can verify the referred thing, and aggregation should not be used to homogenize provider data.

Table 8 depicts a summary of the tasks presented in Figure 13.

Name	Type	Description	Role	IT Objects	Input / Output Document
Submit data and mapping information	T	Once the mapping definition has been finalized, the data and mapping information needs to be submitted.	Mapping Manager		
Transformation	T	Transform source metadata to target records, ready to be ingested to the target system. The transformation process itself may run completely automatically.	Mapping Manager	Transformation Tools	
Reveal "dirty data"?	EG	Transformation process may reveal inconsistent data that need to be mended.			
Is it possible to correct them manually?	EG	Is it possible to correct inconsistent data manually?			
Is it possible to describe by rule?	EG	Check if local syntax rules exist.			
Send back to provider	T	Sort out dirty records and send them back to the provider for processing.			
Syntax Normalization	T	Convert data to the structural format described by the rules.	Provider Data Manager		
Correct manually	T	If possible and given sufficient, trusted expertise, experts under the control of the mapping manager may correct some of them manually in an interactive process. The result is a set of valid target records.	Mapping Manager		Output: Aggregator Format Records
Ingest and Storage	Sub-p				
Update the local records	T	The provider should better update his records.	Provider Curator		

Table 8: Summary of the Metadata Transfer sub-process

6.1.1.3.1 Ingest and Storage

Once records are transformed, an automated translation for source terms using a terminology map may follow. The transformed records will then, be ingested into the target system.

An Ingest Manager should also store all source metadata records for the transformed information. This is considered very good practice and supports transparency important for academic projects. The Ingest Manager must preserve a link to the identity and version of the source record it is derived from. Some aggregators additionally provide source data to users of the target system as part of query results (e.g., the German Digital Library). Source

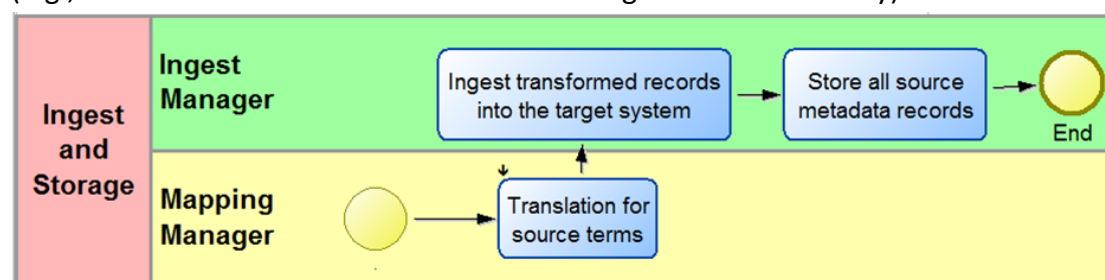


Figure 14: Ingest and Storage sub-process

records may be syntactically normalized for that purpose.

If a new version of a transformed source record is ingested in the target system, the target record representing the previous version must be removed for the purposes of canonical searching. Some aggregators will keep previous versions for historical and academic purposes but these versions should be separate and correctly described to avoid any confusion. However many aggregators will have a deletion policy and it is advisable that providers keep old versions. Otherwise the provider and aggregator may agree terms to manage and preserve old versions as an additional service. In this case providers must make sure that the data is recoverable in the same format that it was submitted in addition to the aggregator’s model.

Table 9 depicts a summary of the tasks presented in Figure 14.

Task Name	Type	Description	Responsible Role	Input / Output Document
Translation for source terms	T	An automated translation for source terms using a terminology map may follow.	Mapping Manager	Input: Terminology Mappings
Ingest transformed records into the target	T	The transformed records will be ingested into the target system.	Ingest Manager	

system				
Store all source metadata records	T	An aggregator should store all source metadata records which are going to be transformed, or which are transformed and have been ingested to the target system. The aggregator must preserve the link to the identity and version of the source record it is derived from. Some aggregators return also fitting source records in query results (e.g., the German Digital Library). Source records may be syntactically normalized for that purpose.	Ingest Manager	

Table 9: Summary of the Ingest and Storage sb-process

6.1.2 Update Processing

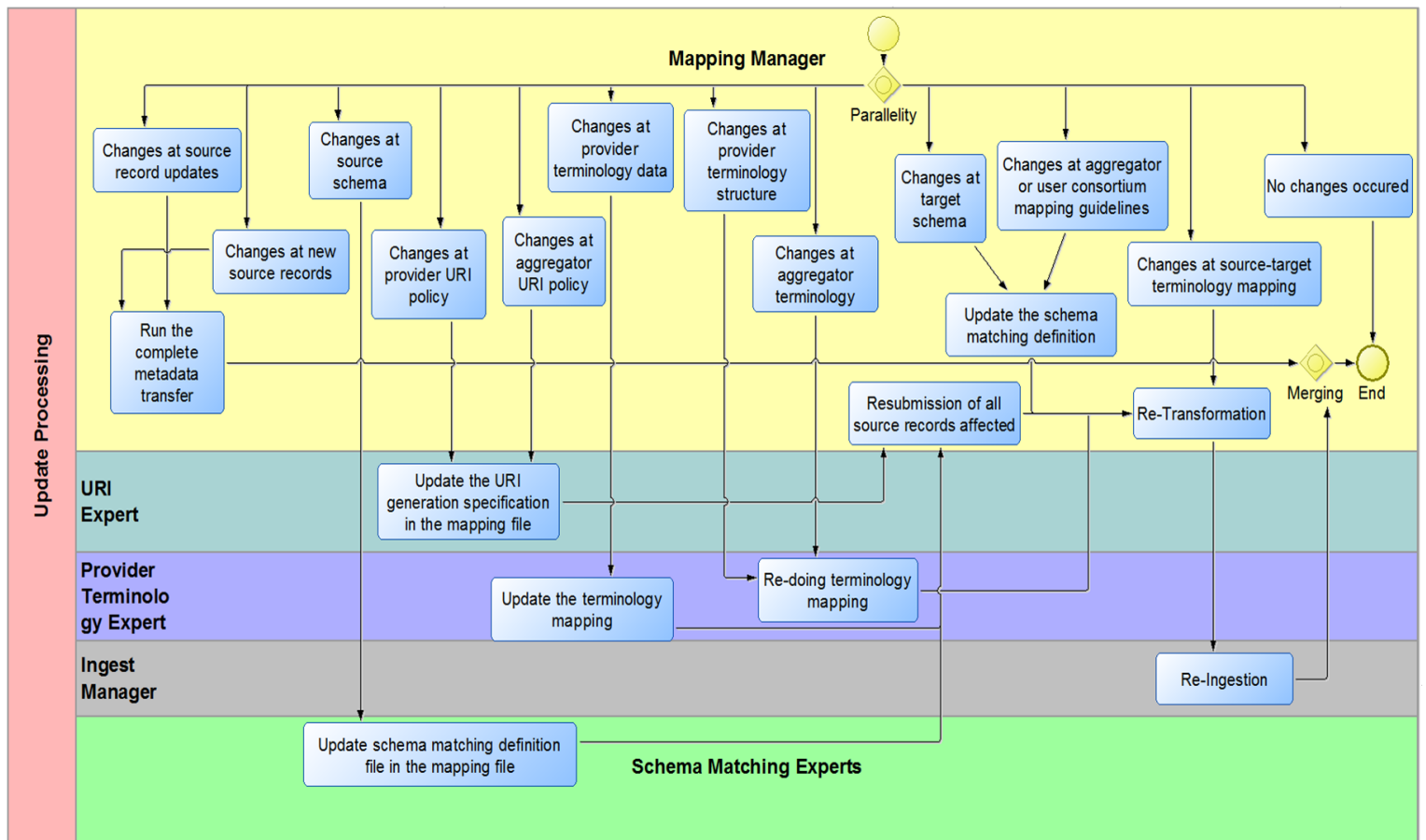


Figure 15 Update processing sub-process

The mapping manager must monitor all changes that may affect the consistency of provider and aggregator data.

Those changes are:

1. New source records
2. Source record updates (new versions)
3. Source schema changes
4. Provider changes identifier policy (for people, objects, events, places, time) and updates his records
5. Provider changes terminology data (terms or authority) and updates his records
6. Provider changes terminology structure (broader term links etc.)
7. Target schema changes

8. Aggregator changes URI policy
9. Aggregator changes terminology (terms or authority)
10. Aggregator or user consortium changes mapping guidelines
11. Source-target terminology mapping changes

The changes of number 1 and 2 require running the complete metadata transfer with the changed or new source records but using the existing mapping definition.

Changes of kind 3 require updating the schema matching definition in the mapping file, resubmission of all source records affected, transformation and ingestion, replacing the target records transformed from the previous version of these source records.

Changes of kind 4 require updating the URI generation specification in the mapping file, resubmission of all source records affected, transformation and ingestion, replacing the target records transformed from the previous version of these source records.

Changes of kind 5 require updating the terminology mapping, resubmission of all source records affected, transformation and ingestion, replacing the target records transformed from the previous version of these source records

The changes of kind 6 and 9 require a rework of the terminology mapping, retransformation and re-ingestion of all (stored) source records already transferred which refer to the respective terms.

Changes of kind 7 and 10 require updating the schema matching definition, retransformation and re-ingestion of all (stored) source records already transferred which refer to the respective terms.

Changes of kind 8 require updating the URI generation specification in the mapping file, retransformation and re-ingestion of all (stored) source records already transferred which refer to the respective terms.

The changes of kind 11 require retransformation and re-ingestion of all (stored) source records already transferred which refer to the respective terms.

Table 10 depicts a summary of the tasks presented in Figure 15.

Name	Type	Description	Role
Changes at source record updates	T	Changes at source record updates (new versions)	Mapping Manager
Changes at provider	T	Provider changes identifier policy	Mapping

URI policy		(for people, objects, events, places, time) and updates his records	Manager
Changes at aggregator URI policy	T	Provider changes identifier policy (for people, objects, events, places, time) and updates his records	Mapping Manager
Changes at target schema	T	Target schema changes	Mapping Manager
Changes at aggregator or user consortium mapping guidelines	T	Aggregator or user consortium changes mapping guidelines	Mapping Manager
Changes at source-target terminology mapping	T	Source-target terminology mapping changes	Mapping Manager
Retransformation	T	Retransformation of the changed terms	Mapping Manager
Re-ingestion of all (stored) source records	T	Re-ingestion of all (stored) source records already transferred which refer to the respective terms.	Ingest Manager
Re-doing terminology mapping	T	Redoing the terminology mapping	Provider Terminology Expert
Changes at aggregator terminology	T	Aggregator changes terminology (terms or authority)	Mapping Manager
Update the schema matching definition	T	Update the schema matching definition	Schema Matching Experts
Run the complete metadata transfer	T	Run the complete metadata transfer with the changed or new source records but using the existing mapping definition.	Mapping Manager
Changes at source schema	T	Changes at source schema	Mapping Manager
Update the URI generation specification in the mapping file	T	Update the URI generation specification in the mapping file.	Instance Generation Expert
Resubmission of all	T	Resubmission of the affected source records	Mapping Manager

source records affected			
Update the terminology mapping	T	Update the URI generation specification in the mapping file,	Provider Terminology Expert
Update schema matching definition file in the mapping file	T	Update the schema matching definition file in the mapping file	Schema Matching Experts
Transformation	T	Transform source records affected to target records.	Mapping Manager
Ingestion	T	Ingestion, replacing the target records transformed from the previous version of these source records.	Ingest Manager
Changes at provider terminology structure	T	Provider changes terminology structure (broader term links etc.)	Mapping Manager
Changes at provider terminology data	T	Provider changes terminology data (terms or authority) and updates his records	Mapping Manager
Changes at new source records	T	Changes at new source records	Mapping Manager
No changes occurred	T	No changes occurred	Mapping Manager

Table 10: Summary of the Update Processing sub-process

7 Services and S/W components

This section describes in detail the software components foreseen by this model. The intention of this model is only to define interoperable interfaces between the components, such that an effective monitoring and workflow control system can interact with the components, and combinations from arbitrary providers in arbitrary technologies can interact to enable the whole process with all its details. It should further enable rich enough variations of particular workflows and extension of functionality. Each component may be implemented with different levels of sophistication, from simple commands to fancy graphic manipulations. Some components may only support limited functionality, e.g., transformation from XML to XML. In such cases, the workflow system should be able to plug in on demand alternative components for other formats, such as E-R to XML, E-R to RDF, XML to RDF etc.

THIS SECTION WILL BE EXTENDED BY EXACT INTERFACE DEFINITIONS.

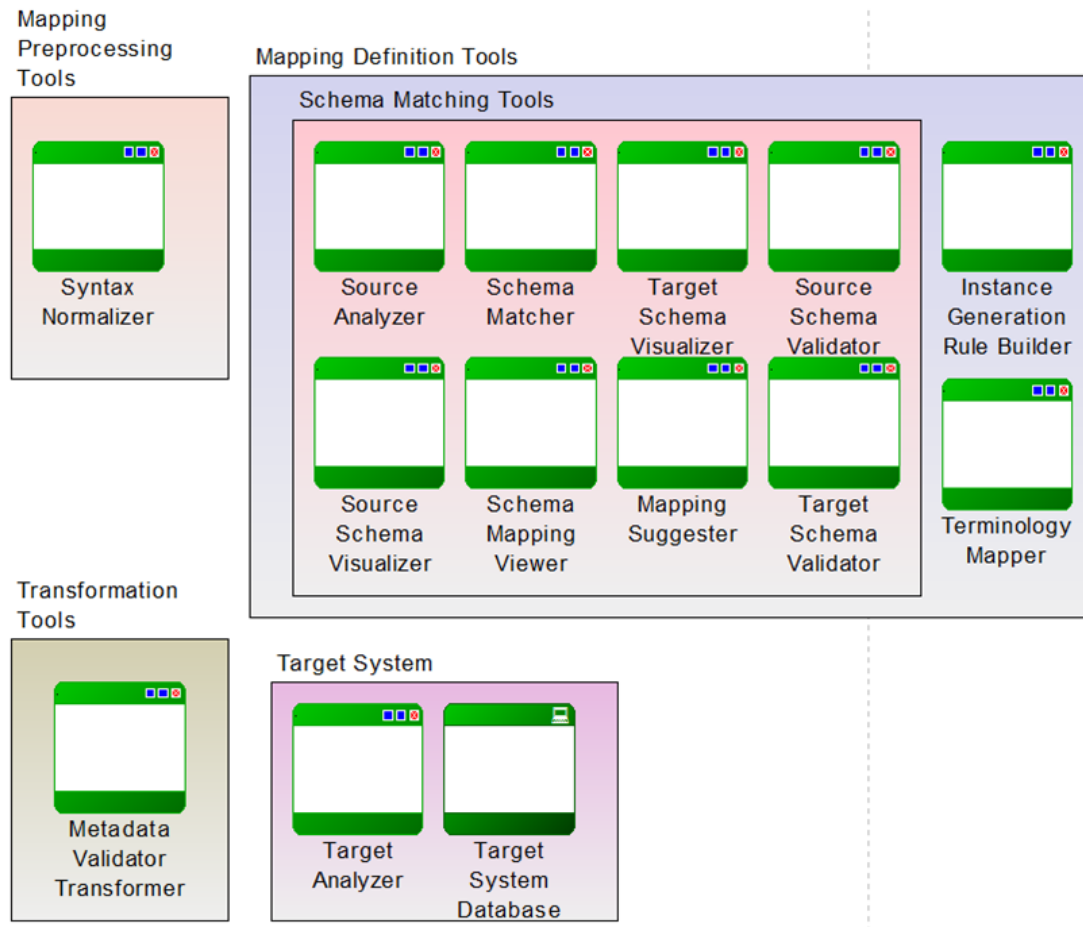


Figure 16 Services and S/W components

Figure 16 illustrates the IT objects that assist the data provisioning process. We do not regard IT objects as self-contained and opposed to user processes, but IT objects are regarded as being part of the user processes replacing or supporting manual work. The following objects are distinguished:

Mapping Definition Tools

A toolset able to transform content and metadata, in various, heterogeneous formats, to a normalized data model such as the CIDOC CRM.

Schema Matching Tools

Tools to support the schema matching process. Takes as input two ontologies and determines the alignment result between entities of the input ontologies.

Source Analyzer

- This component should be developed to deal with different source formats, including relational database tables, XML exports and potentially other formats. In this respect it should be open to the development of additional modules designed to support other formats of source data.

- The system should incorporate and display in context the scope notes, definitions, examples etc. that are provided with the target schema.
- The system should provide visual representations of the schemas within their hierarchical context. The hierarchical view should be expandable and collapsible at each level and overall (expand all)
- The system should show instances of mapped data to a configurable number
- The system should provide statistics that show the number of successful and unsuccessful or empty mappings for each field.
- The user should be able to browse the values in the source and target fields for a configurable number of examples.
- The system should be able to expose random values from a sample of source and mapped fields.
- The system supports the inclusion and ongoing addition (through a user interface) of data transformation functions necessary to 'clean' data necessary for mapping.

Source Schema Visualizer

Tools to visualize source schema definition and source metadata records. Adequate navigation and viewing techniques must allow for overviews and understanding of details.

Schema Matcher

Loops through all source elements in order to make a mapping decision. It may be supported by tools suggesting mappings and should recalculate their proposals.

Mapping Suggester

Tools suggesting mappings. These tools must take into account "mapping memories" of analogous cases collected from the user community.

Target Schema Visualizer

Tools to visualize target schema definition. Adequate navigation and viewing techniques must allow for overviews and understanding of details.

Schema Mapping Viewer

Tools to visualize the schema matching definition file. Adequate navigation and viewing techniques must allow for overviews and understanding of details.

Instance Generation Rule Builder

Defines new instance generation rules for each independent node.

Mapping Preprocessing Tools

Tools used at the preprocessing stage for the normalization of data.

Syntax Normalizer

Converts any data structure that has no standard format to a standard one, like XML, RDF OWL, RDBMS or at least spreadsheets.

Terminology Mapper

Maps source terms to target terms.

Transformation Tools

Tools used to transform data from one standard format to another one.

Metadata Validator Transformer

Performs the transformation process and validates the transformed records before their ingestion at the target system.

Target System

Provide a homogeneous access layer to multiple local systems. The information they manage resides primarily on local systems and (meta) data are sent on a regular base or in a single action by several providers.

Target Analyzer:

- It is usual for schema documentation to be embedded into the structure of the schema. The system should be able to identify this documentation and expose it to the user when browsing aspects of the target. It should also be possible for a user to add additional notes (stored separately) that can be used in the future by the provider or the target schema owner. These notes may be useful for future improvements.
- The schema should be presented in a way that is easy to navigate and should expose relevant associated information when viewing any particular property. For example, it should be easy to isolate and browse information on related or sub-properties. **Note:** The mapping memory might include information about related properties that are often confused and regularly misused. This information could be made available so that the user is alerted to subtle differences in the schema. It also may be possible to isolate parts of the target schema that are identified as relevant to a particular type of mapping. For example, is the user is mapping acquisition information the system might show only those properties that are relevant.
- Hide/expand ISA sub-trees
- Hide/expand path graphs

Target System Database

The database used by the target system.

Table 11 depicts the input and output data objects of each function component.

IT component	Input Document	Output Document
Mapping Suggester	Mapping Memory	
Metadata Validator Transformer	<ol style="list-style-type: none"> 1. Co-Reference Store 2. Normalized Provider Metadata 3. Terminology Mappings 4. Mapping Definition File 	<ol style="list-style-type: none"> 1. Mapping Validation Report 2. Aggregator Format Records
Schema Matcher		Schema Matching Definition
Source Analyzer	<ol style="list-style-type: none"> 1. Normalized Provider Metadata 2. Effective Provider Schema 	<ol style="list-style-type: none"> 1. Provider Field Use Statistics 2. Provider Terminologies 3. Effective Provider Schema
Source Schema Visualizer	<ol style="list-style-type: none"> 1. Effective Provider Schema 2. Provider Field Use Statistics 	
Syntax Normalizer	<ol style="list-style-type: none"> 1. Raw Metadata 2. Target Schema Definitions 	<ol style="list-style-type: none"> 1. Source syntax and Cleaning Report 2. Effective Provider Schema 3. Normalized Provider Metadata
Target Analyzer		<ol style="list-style-type: none"> 1. Authority References Hints 2. Aggregator Statistics Report
Target Schema Visualizer	Target Schema Definitions	
Terminology Mapper	Provider Terminologies Aggregator Terminologies	Terminology Mappings
Instance Generation Rule Builder	<ol style="list-style-type: none"> 1. Schema Matching Definition 2. Provider Field Use Statistics 	Mapping Definition

Table 11: IT objects' Input/Output Documents

8 References

[1] ADONIS-BOC Group, [Online]. Available: <http://www.adonis-community.com/>. [Accessed 2014]

[2] CIDOC-CRM, "The official release of CIDOC-CRM," [Online]. Available: http://www.cidoc-crm.org/docs/cidoc_crm_version_5.0.4.pdf. [Accessed 2014].