

Data Foundation and Terminology

Working Group

Responsible RDA Working Group Co-Chairs:

Gary Berg-Cross – Research Data Alliance Advisory Council, Washington D.C. USA

Raphael Ritz - Max Planck Institute for Plasma Physics, Germany

Peter Wittenburg – Max Planck Institute for Psycholinguistics, Germany

What is the Problem?

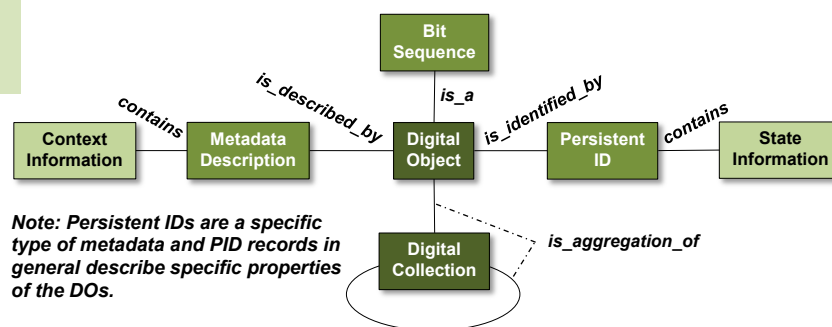
Unlike the domain of computer networks where the TCP/IP and ISO/OSI models serve as a common reference point for everyone, there is no common model for data organisation, which leads to the fragmentation we are currently seeing everywhere in the data domain. Not having a common language between data communities, means that working with data is very inefficient and costly, especially when integrating cross-disciplinary data. As Bob Kahn, one of the Fathers of the Internet, has said, “Before you can harmonise things, you first need to understand what you are talking about.”

When talking about data or designing data systems, we speak different languages and follow different organization principles, which in the end, result in enormous inefficiencies and costs. We urgently need to overcome these barriers to reduce costs

are endless solutions that create enormous hurdles when federating. To give an idea of the scale of the problem, almost every new data project designs yet more new data organisations and management solutions.

We are witnessing increasing awareness of the fact that at a certain level of abstraction, the organisation and management of data is independent of its content. Thus, we need to seriously change the way we are creating and dealing with data to increase efficiency and cost-effectiveness.

For the physical layer of data organisations, there is a clear trend towards convergence to simpler interfaces (from file systems to SWIFT-like interfaces⁷). For the virtual layer information, which includes persistent identifiers, metadata of different types including provenance information, rights information, relations between digital objects, etc., there



This diagram describes the essentials of the basic data model that the DFT group worked out in a simplified way. Agreeing on some basic principles and terms would already make a lot of difference in data practices.

What were the goals?

The goals of this Working Group (WG) were:

- Pushing the discussion in the data community towards an agreed basic core model and some basic principles that will harmonize the data organization solutions.
- Fostering an RDA community culture by agreeing on basic terminology arising from agreed upon reference models.

⁷ <https://wiki.openstack.org/wiki/Swift>

What is the solution?

Based on 21 data models presented by experts coming from different disciplines and about 120 interviews and interactions with different scientists and scientific departments, the DFT WG has defined a number of simple definitions for digital data in a registered⁸ domain based on an agreed conceptualisation.

These definitions include for example:

- **Digital Object** is a sequence of bits that is identified by a persistent identifier and described by metadata.
- **Persistent Identifier** is a long-lasting string that uniquely identifies a Digital Object and that can be persistently resolved to meaningful state information about the identified digital object (such as checksum, multiple access paths, references to contextual information etc.).
- A **Metadata description** contains contextual and provenance information about a Digital Object that is important to find, access and interpret it.
- A **Digital Collection** is an aggregation of digital objects that is identified by a persistent identifier and described by metadata. A Digital Collection is a (complex) Digital Object.

A number of such basic terms have been defined and put into relation with each other in a way that can be seen as spanning a reference model of the core of the data organisations.

What is the impact?

The following benefits will come from wide adoption of a harmonized terminology which will be expanded stepwise:

- Members of the data community from different disciplines can interact more easily with each other and come to a common understanding more rapidly.
- Developers can design data management and processing software systems enabling much easier exchange and integration of data from their colleagues in particular in a cross-disciplinary setting (full data replication for example could be efficiently done if we can agree on basic organization principles for data).
- It will be easier to specify simple and standard APIs to request useful and relevant information related to a specific Digital Object. Software developers would be motivated to integrate APIs from the beginning and thus facilitate data re-use, which currently is almost impossible without using information that is exchanged between people.
- It will bring us a step closer to automating data processing where we can all rely on self-documenting data manipulation processes and thus on reproducible data science.

When can we use this?

The definitions have been discussed at RDA Plenary 4 meeting (Sept 2014) and will become available as a document and on a semantic wiki to invite comments and usage at January 2015. RDA and the group members will take care of proper maintenance of the definitions. For more information see <https://rd-alliance.org/group/data-foundation-and-terminology-wg.html> and http://smw-rda.esc.rzg.mpg.de/index.php/Main_Page

In the next phase of the work, more terms will be defined and interested individuals will have the opportunity to comment via the semantic wiki.

⁸ There will always exist data in private, temporary stores, which will not be made accessible in a standard way.

Data Type Registries Working Group

Responsible RDA Working Group Co-Chairs:

Larry Lannom - Corporation for National Research Initiatives, Virginia USA

Daan Broeder - Max Planck Institute for Psycholinguistics, Netherlands

What is the Problem?

Often researchers receive a file from colleagues, follow a link, or otherwise encounter data created elsewhere that they would like to make use of in their own work. However, they may not know how to work with it, interpret it or visualise its content, being unfamiliar with the specifics of the structure and/or meaning of the data, ranging from individual observations up to complex data sets. Frequently, researchers need to stop here since it requires too much work to look for explanations, tools, and where tools exist, install them.

When sharing data across disciplines, we often get files which we cannot process easily. Dragging such a file on the DTR would immediately yield results and reduce effort.

What was the goal?

The goal of the DTR WG was to allow data producers to record the implicit details of their data in the form of Data Types and to associate those Types, each uniquely identified, with different instances of datasets. Data consumers can then resolve the Type identifiers to Type information for gaining knowledge of the implicit assumptions in the data, finding available services that can be used for this kind of data, and any other useful information that can be used to understand and process the data, without additional support from data producers. DTRs are meant to provide machine-readable information, in addition to presenting human readable information.

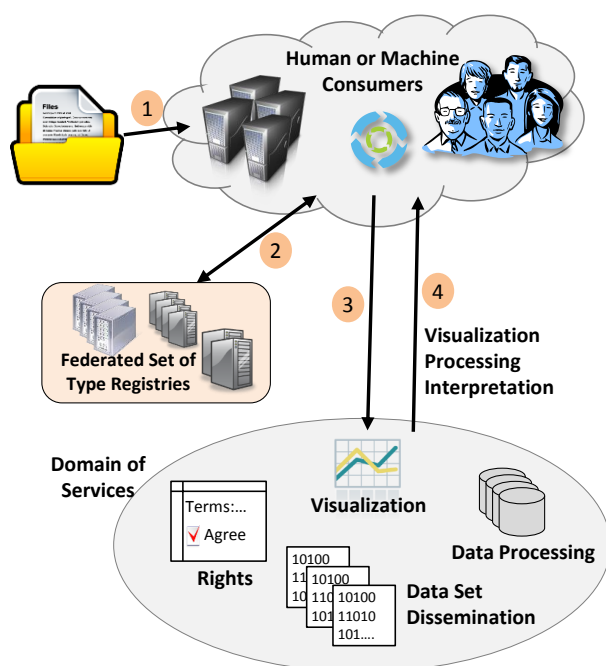
What is the solution?

DTRs offer developers or researchers the ability to add their type definitions in an open registry and, where useful, add references to tools that can operate on them. For example, a user who received an unknown file could query a DTR and receive back a pointer to a visualisation service able to display the data in a useful form. A fully automated system could use a DTR, much like the MIME type system enables the automatic start of a video player in the browser once a video file has been identified. We envision humans taking advantage of Data Types in DTRs through the type definitions that clarify the nuanced and contextual aspects of structured datasets.

Data Types in DTRs can be used to extend or expand existing types, e.g., MIME types, which provide only container-level parsing information. They can additionally describe experimental context, relationships between different portions of data, and so on. Data Types are deliberately intended to be quite open in terms of registration policies.

Two examples may illustrate the benefits of the DTR solution:

1. Researchers dealing with data (e.g. in a cross-disciplinary, cross-border context) find an unknown data type and can immediately process and/or visualize its content by using the DTR service.
2. Machines that want to extract the checksum information of a data object from a PID record to check whether the content is still the same. Without knowing the details of the PID service provider, the machine could ask for CKSM for example, since this is an information type which all PID service providers agreed upon and registered in the DTR.



This diagram indicates how the Data Type Registry (DTR) is working. A user or machine receives an unknown type (1) which can be a file or a term for example. The DTR is contacted and returns information about an available service (2) that will allow the user or machine to continue processing the content (3, 4) such as visualizing an image without asking prior knowledge from the user. This will make cross-disciplinary and cross-border work much more efficient and enable data driven science even to those who are

available here: <http://typeregistry.org/>. We expect software to become available for download around the end of 2014. Please check the information on the DTR WG's web page at <https://www.rd-alliance.org/group/data-type-registries-wg.html> for updates.

This simple model will be the start for designing DTRs, with the intention to extend the specifications according to priorities and usage.

What is the impact?

The potential impact on scientific practices is substantial. Unknown data types as described above can be exploited without any prior knowledge and thus an enormous gain in time and/or in interoperability can be achieved. In a similar way to the MIME types that allow browsers to automatically select visualization software plug-ins when confronted with a certain file type extension, scientific software can make use of the definitions and pointers stored in the DTR to continue processing without the user acquiring knowledge beforehand. DTRs pave the way to automatic processing in our data domain, which is becoming increasingly complex, without putting additional load on the researchers.

Of course, a price needs to be paid in that type creators need to enter the required information into a DTR. We assume that there will be a federation of such DTRs setup to satisfy different needs.

When can we use this?

The first groups are building software to implement such a DTR concept and make the software available. The RDA PID Information Type (PIT) Working Group is already using the first DTR prototype version in its API. The latest version of a DTR prototype is made

PID Information Types Working Group

RDA Working Group Co-Chairs:
Tobias Weigel – DKRZ, Germany

Timothy Dilauro – John Hopkins University, Maryland, United States

What is the Problem?

Numerous systems and providers to register and resolve Persistent Identifiers (PIDs) for Digital Objects and other entities have been designed in the past and are used today. However, almost all

Due to high demand, a variety of trusted PID service providers have been set up already, yet all of the different attributes associated with the registered PIDs make life of a software developer a nightmare. We need to harmonize the major information types and suggest a common API, so that if we request the checksum we simply have to program one piece of software

of them differ in the way they allow researchers to associate additional information, such as for proving identity and integrity with the PID. For application developers this is an unacceptable situation, since for all providers a different Application Programming Interface (API) needs to be developed and maintained. Given that a researcher has found a useful file, but first wants to prove whether it is indeed the same stream of bits after some years, he should be able to request the checksum independent of the provider holding the PID. How should he do this not knowing whether the provider offers this information and if so, how to request it? We can overcome such extreme inefficiencies only if all providers

agree on a common API, register their information types in a common data type registry and agree on some core types, such as the checksum.

What were the goals?

The goals of this WG were:

- Coming to a core set of information types and register (and define) them in a commonly accessible Data Type Registry
- Providing a common API and prototypical implementation to access PID records that employ registered types

What is the solution?

The PIT group accomplished the following:

- Defined and registered a number of core PID information types (such as checksum)
- Developed a model to structure these information types
- Provided an API, including a prototypical server implementation that offers services to request certain types associated with PID records by making use of registered types.

The set of core information types currently provided can help to illustrate cross-discipline usage scenarios. It can also act as an example for a community-driven governance process creating and governing more user-driven types. PID service providers and community experts need to come together regularly and add types to the data type registry to make full use of the possibilities of the results of the PIT group.

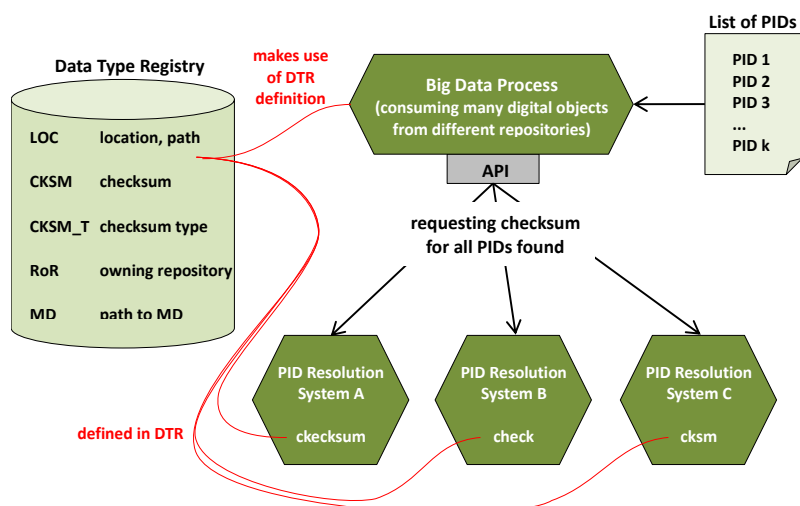
It is now essential to convince PID service providers such as those using the Handle System (DOI, EPIC, etc.) to adopt the API to unify access. In the diagram below, we give an example of the usage and potential of the suggested solution.

What is the impact?

We need to envisage the situation in a few years, when the amount and complexity of data has been increased in all sciences and there is a greater need to rely on automatic processes, as human intervention means loss of efficiency.

In such scenarios, communities can exploit the wealth of the data domain relying on semantic interoperability between all relevant actors for example for Big Data analytics. The above example is just one small usage scenario that would be enabled if the relevant PID service providers accept the results of the PIT WG and harmonize their approach. Application software writing would be reduced dramatically since only one API would be supported and one module would be sufficient for retrieving the checksum, for example, and checking identity and integrity.

Assume that you got a list of PIDs referring to data you want to use in a computation, that these PIDs are being registered at different providers and that you first want to check whether all data objects are still the same. You simply want to provide one module that reads a PID from the list and submits a request to the appropriate resolver to send the checksum. If all actors refer to the same entry in the DTR interoperability is given, i.e. one module would be sufficient to retrieve the checksums independent



The strengthening of PID information types could also move the existing identifier systems and the overall idea of identification into a more central and fundamental position as suggested by DFT's core model of a Digital Object, leading to an enormous increase in efficiency when dealing with data.

When can we use this?

First groups are building software to implement a first prototype based on the defined PIT API. This first prototype works together with the DTR prototype and both are publicly available, but not designed for production use. We expect another update of the prototypes to become available for download at the end of 2014.

Please check the information on the PIT group's web-page at <https://www.rd-alliance.org/group/pid-information-types-wg.html>.

It is now time to convince the PID service providers to adopt the solution.

Practical Policy Working Group

Responsible RDA Working Group Co-Chairs:
Reagan Moore, RENCI, North Carolina, USA
Rainer Stotzka, Karlsruhe Institute of Technology, Germany

What is the Problem?

Repositories' responsibilities of data stewardship and processing require a highly automated, safe and documented process. However, at this time, repositories design and implement these processes in a method that does not support this requirement.

Current practice in managing and processing data collections are determined by manual operations and ad-hoc scripts making verification of the results an almost impossible task. Establishing trust and a reproducible data science requires automatic procedures which are guided by practical policies. Collecting typical policies, evaluating them and providing best practice solutions will help all

With the increasing amount and complexity of data, repositories should not continue to use manual interventions and ad-hoc scripts any longer since they prevent us to establish trust.

All operations or chains of operations that have these capabilities and are enforced on collections of data objects should have "Practical Policies" (PP), which should be stated in simple languages and turned into robust and tested executable code. PPs are at the basis of reproducible science, an important element in the chain of building trust and one of the core elements in repository certification processes.

What were the goals?

The goals of this WG were:

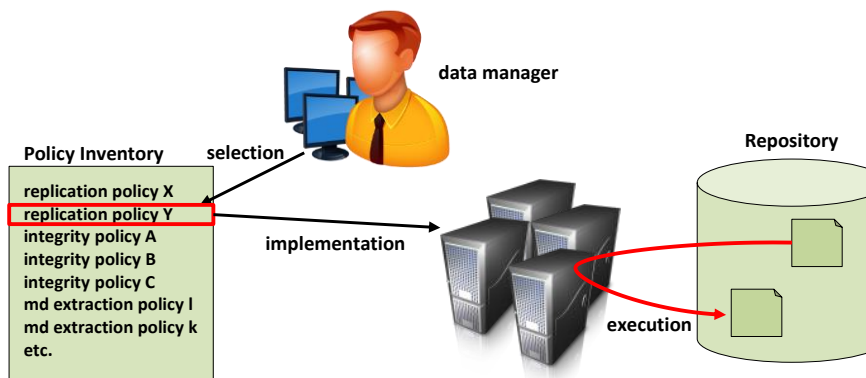
- Defining computer actionable PPs that enforce proper management and stewardship, automating administrative tasks, validating assessment criteria, and automating types of scientific data processing
- Identifying typical application scenarios for practical policies such as replication, preservation, metadata extraction, etc.
- Collecting, registering and comparing existing practical policies
- Enabling sharing, revising, adapting and re-using of such practical policies and thus harmonizing practices, learning from good examples and increasing trust

Since these goals were broad in scope, PP WG focused its efforts on a few application scenarios for the collection and registration process.

What is the solution?

In order to identify the most relevant areas of practice, the PP WG conducted a survey as a first step. The analysis of the survey resulted in 11 highly important policy areas which were tackled first by the WG: 1) contextual metadata extraction, 2) data access control, 3) data backup, 4) data format control, 5) data retention, 6) disposition, 7) integrity (incl. replication), 8) notification, 9) restricted searching, 10) storage cost reports, and 11) use agreements.

Participants and interested experts were asked to describe their policy suggestions in simple semi-formal descriptions. With this information, the WG developed a 50-page document covering the



simple descriptions, the beginning of a conceptual analysis and a list of typical cases such as extract metadata from DICOM, FITS, netCDF or HDF files.

Due to unexpected circumstances, the WG

will continue until Plenary 5 (March 2015). It will focus on further analysing, categorising and describing the offered policies. Currently, volunteers are reviewing the policies and different groups have started to implement some of these policies in environments such as iRODS and GPFS. The goal is to register prototypical policies with suitable metadata so that people can easily find what they are looking for and re-use what they found at abstract, declarative or even at code level. At this point, there is still much work to be done to reach a stage where the policies can be easily used.

What is the impact?

The impact is huge. In the ideal case, data managers or data scientists can simply plug-in useful code into their workflow chains to carry out operations at a qualitatively high level. This will improve the quality of all operations on data collections and thus increase trust and simplify quality assessments. Large data federation initiatives such as EUDAT and DATANET Federation Consortium (US) are very active in this group, since they also expect to share code development/maintenance, thus saving considerable effort by re-using tested software components. Research Infrastructure experts that need to maintain community repositories can simply re-use best practice suggestions, thus avoiding ending up in traps. In particular, when these best practice suggestions for practical policies are combined with proper data organisations, as suggested by the Data Foundation and Terminology Working Group, powerful mechanisms will be in place to simplify the data landscape and make federating data much more cost-effective.

The diagram indicates the final goal of the PP WG. A policy inventory will be made available with best practices examples. Data managers will have the ability to select and implement the procedures most relevant to them.

When can we use this?

The document mentioned above already provides a valuable resource to get inspiration and perhaps make use of suggested policies, thus improving people's own ideas or to even making profit from developed code.

Once evaluated, properly categorised and described, the real step ahead will be registering practical policies in suitable registries, so that data professionals can easily re-use them, if possible even at code level. The group intends to progress to this step by the end of March 2015 for a number of policy areas, making use of the policy registry developed by EUDAT.

For more details on the PP WG, see <https://www.rd-alliance.org/group/practical-policy-wg.html>

Revolutionising Data Practices

Gary Berg-Cross, Keith Jeffery, Rob Pennington, Peter Wittenburg

What is the Problem?

A large survey from mainly RDA Europe and EUDAT (including about 120 interviews and interactions with data professionals from various departments engaged in various research disciplines)

The task of DFIG is to design a flexible and dynamic framework of essential components and services, identifying those that enable efficient, cost-effective and reproducible data science and making these known and available to researchers and data scientists. The goal is to make it possible for scientific users to easily integrate their scientific algorithms into such a data fabric without needing to master the underlying details.

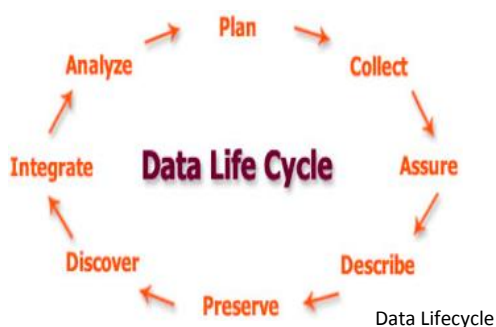
demonstrated that the way we manage and process data is very inefficient and too expensive. In addition, data science generally is not reproducible as some reports have shown which is contrary to good practices and thus not acceptable.

Despite insights from computer science and excellent individual solutions from advanced infrastructure projects we lack a broad and systematic approach to understand the

components, their services and their interfaces that are needed to change our data practices in a way that the deficits will be overcome and to make them available to every researcher. A number of RDA groups are working already on such components, yet doing it in a somewhat isolated way. There is a wide agreement that this needs to be changed urgently.

What are the Goals?

The Data Fabric Interest Group (DFIG) has been setup to address the design of such a framework as a whole, to locate the various activities on the landscape of components, to indicate gaps and to understand how the various groups need to interact to come to an interoperable flexible framework. The intention is thus not to design a relatively fixed architecture of a system that fulfills a particular set of functions, but a flexible framework that can be configured by changing components to meet varying needs, and thus is technology-independent. The framework identifies the minimal set of components required to let any system based on the framework function.



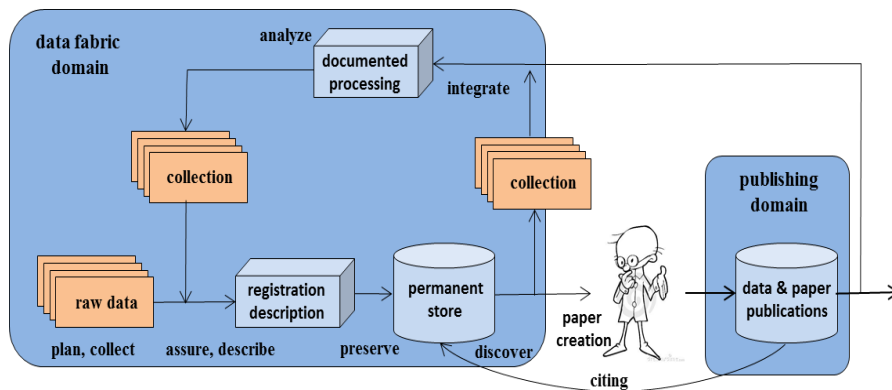
To meet these goals we need to analyse large scale lighthouse infrastructure projects - which are mostly discipline-based developed exemplary solutions - and identify commonalities. DFIG does not start from scratch, but can build on the knowledge already gathered.

DFIG also needs to look at all phases of the lifecycle as schematically indicated by the diagram above.

What is the Solution?

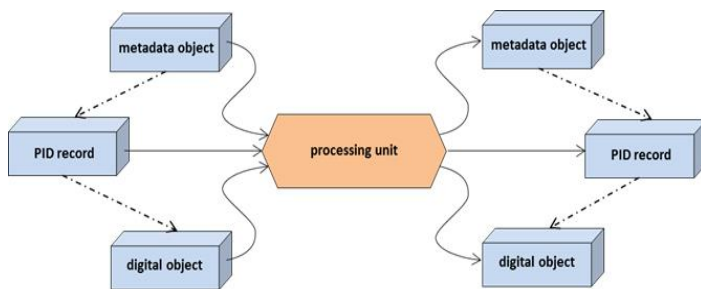
DFIG needs to define a basic and flexible machinery framework that (when implemented as systems) makes data science reproducible, fulfils the G8+O5 recommendations and the need to carry out data management and processing much more efficiently. Recognizing that data intensive science is faced with increasingly large volumes and complexity of data we need to turn to processing which is guided by actionable and documented policies, in which all steps adhere to basic organizational

principles are self-documenting, i.e. provide provenance metadata and are (as much as feasible) autonomic.



The diagram above indicates the data machinery which is being executed in some form in all data intensive scientific work. The relations to the phases in the previous diagram are indicated. Raw data (which can also be long tail data created on a notebook) will

be brought into the accessible domain of data by registering it (assigning Persistent Identifiers), describing it by metadata and depositing it into a permanent and accessible repository which will be distributed. Using metadata scientists will now create new (virtual) collections by making selections which then will be subject to some kind of processing – be it management, curation or analytic. New collections are being created that which again are described, registered and deposited.



If all processing steps follow principles as schematically indicated above where new data and metadata is being generated extending the old objects, we will achieve the kind of self-documentation that is required. To unload the scientist DFIG needs to identify the components that are required to put such machinery in place and that allows researchers to

simply plug-in their scientific algorithms so that they do not need to know about all the details of the machinery. We realize that achieving this, being compliant with the G8+O5 principles (searchable, accessible, interpretable, re-usable) and putting it in place so that everyone can take profit from it is a long road that requires a step-wise approach. But we need to start working on this today and convince software builders to follow these principles. RDA activities need to have this overall picture in mind where the act of publishing papers and data is an integrated phase requiring some explicit steps.

What is the Impact?

The impact of implementing such machinery based on a flexible framework is huge and will revolutionize data intensive science. It can be compared with optimizing the publication and citation machinery as we have seen over the past decades.

When can we use it?

Like with Internet where broad uptake happened about 15 to 20 years after the invention and optimization of the TCP/IP framework, RDA will stepwise optimize the way to deal with data in the various phases. Here the first working and interest groups in RDA take already now important steps and also large lighthouse infrastructure projects facing the inefficiencies daily have designed solutions which need to be analysed and considered carefully. Like with Internet we need to define the basic and essential components now that will allow us adding components and services

RDA Europe Data Practice Analysis

Editors: Peter Wittenburg, Herman Stehouwer*

dependent on insights and technological advancements.

What did we do?

For the RDA Europe Data Practice Analysis Programme we held a large number of interviews with data scientists/practitioners from various communities. We interviewed these people about various aspects of their data environment including data acquisition, data processing, the computational environment, services and tools, and the data related policies being applied.

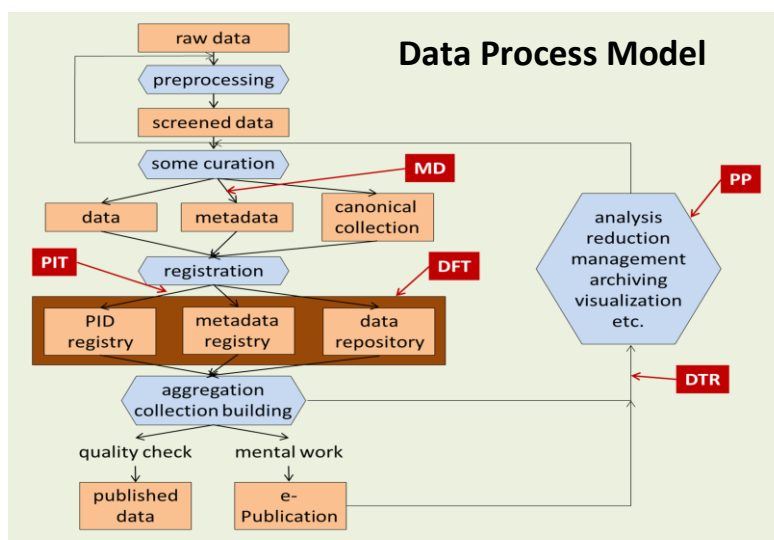
We interviewed 24 communities, and attended more than 70 community meetings. We combined these observations with the interviews and observations made in the EUDAT project, in the Radieschen project, and in the first RDA Europe Science Workshop. Based on these sources of information we came to a large number of observations, which are summarized here in form of the dominant underlying data process model, and 12 key observations.

Key Messages

Support Open Access
Ensure (Meta-) Data Quality
Explicit Structure & Semantics
Change to Documented Methods
Help Increasing Trust
Educate/Train Data Professionals

Data Process Model

The process model in the figure emerges as the dominant underlying process model that most data



scientists/practitioners are implicitly using when processing data. In practice the methods used in the departments deviate slightly from this generic model in various ways, but it summarizes what is being done at an abstract level very well. Furthermore, most often parts of the data processing are implicitly handcrafted with ad-hoc solutions rather than by following an explicit model.

The model helps us to clarify our observations and to identify specific steps as they relate to data, specifically: Data is scientifically meaningful and relevant after the pre-processing step; data is ready

for upload to a repository after the curation step; data is ready for re-use after the registration step; and data is ready for citation after the publishing step. Currently most researchers do not distinguish between these steps explicitly. Explicitly separating these steps of the data process would increase efficiency and decrease cost.

This model shows similarities to existing models of data processing (such as the Kahn/Wilensky, CLARIN, EUDAT, ENVRI, EPOS, and DICE models⁹), and it can be used to place the observations made in the analysis program as well as to talk about a data management system. In the diagram we also placed where the topics of the first RDA Working Groups can be located.

12 Observations

1. **ESFRI** projects and the recent developments within **e-Infrastructure** have had a strong and positive influence on data management practices.
2. **Open Access** is supported everywhere as a basic recommendation. However in practise there are many barriers that still need to be lowered.
3. **Trustworthiness** is a key issue and new methods are urgently required to establish trust in the entire data processing chain.
4. **Legacy Data** is a problem in many communities, however even new data is often badly documented and organised, thus we are creating continuously new legacy data which will cost much effort to integrate them in the accessible data domain. There is 1) a lack of knowledge about principles of proper data organisation; 2) a lack of experts, time and money who could change practices; 3) a lack of off-the-shelf software methods for improved data management and access.
5. **Big Data** is driving many new scientific requirements that dictate the thorough adoption of this paradigm in increasing numbers of departments. However, big data only scales when data management and access methods are used that scale.
6. **Data Management** needs to move towards including the logical layer of information, i.e. metadata, PIDs, rights, relations to other data, etc. At the end the current file-system based methods are too inefficient and costly. A large amount of researchers' time is wasted in finding the right data objects, interpreting them and creating meaningful collections.
7. **Metadata** practise needs to be improved in order to help discovery and reuse (especially after some time). Guidance and ready-to-use packages and software are required to improve the situation.
8. **Lack of Explicitness** is an issue in relation to data, which hinders efficient machine-based processing of data. This lack ranges from non-registered digital objects (i.e. lacking PIDs), data integrity information (such as checksums), collection descriptions, encoding systems, format/syntax, and semantics up to the level of software components. Appropriate registration authorities and mechanisms do exist, but often they are unknown or not used.
9. **Centres** for managing data across communities are a clear trend. Such centres and repositories need to be established to provide a long-term reliable service to all researchers. Creating virtual collections or carrying out distributed processing jobs is still an unsolved issue. Some aspects of distributed authentication and authorization are still not in place at European level and

Establish Trust
Quality and Integrity of data
Availability of high-quality metadata
Sustainable services and PIDs
Clear Responsibilities and Funding

⁹ references

distributed computing, although mentioned increasingly often, is not a well-understood scenario.

10. **Education & Training** is a clear need in order to address the lack of data professionals. This lack hampers changes and progress everywhere.
11. **Lack of Knowledge** and trusted information on services that are being offered (registries, data, storage, curation, analytics, etc.) is an issue. We have a large number of possibilities, but many can't cope with the information flood and have a hard time making selections. A more structured and trusted approach of offering information would have great impact.
12. **RDA** needs to ensure it is a true grass-roots organisation. It needs to provide demonstration cases, and give help and support to research communities.