# Descriptive Metadata Requirements for Long-term Archival of Digital Product Models

Joshua Lubell
National Institute of Standards and Technology
Gaithersburg, Maryland, USA
1-301-975-3563

lubell@nist.gov

Ben Kassel
Naval Surface Warfare Center, Carderock Division
West Bethesda, Maryland, USA
1-301-227-1142

ben.kassel@navy.mil

Sudarsan Rachuri
National Institute of Standards and Technology
Gaithersburg, Maryland, USA
1-301-975-4264

sudarsan@nist.gov

## ABSTRACT
Digital product engineering information comes from a variety of applications ranging from computer-aided design, engineering, and manufacturing software to product data and lifecycle management systems. As the volume of digital engineering data increases, so does the need to archive this information for long-term use. The challenges of archiving this information are evident in the U.S. Navy's Torpedo Weapon Retriever (TWR) data set. In order to build an effective archival information system from the TWR data, an ingest process must include the generation of descriptive metadata to aid in future access of archived TWR information.

Descriptive metadata is essential for supporting long-term access requirements for ship product model data. Through determining descriptive metadata requirements and by attempting to successfully generate TWR descriptive metadata during ingest, we can (1) gain insight into the feasibility of generalizing ship data ingest and customizing tools developed by the digital library community for engineering-specific applications, (2) help to coordinate ongoing long-term archiving efforts in the engineering design community with the work of digital library researchers and archivists, and (3) establish requirements for long-term archiving of engineering data produced by computer-aided design/engineering and product data/lifecycle management software applications.

## Keywords
Descriptive metadata, descriptive information, product lifecycle management, long-term archiving, digital preservation, torpedo weapon retriever, engineering.

## 1. INTRODUCTION
Digital product engineering information comes from a variety of applications ranging from computer-aided design, engineering, and manufacturing software to product data and lifecycle

management systems. As the volume of digital engineering data increases, so does the need to archive this information for long-term use. Challenges to long-term archival include:

- The variety of engineering data types and complexity of the relationships between the information units comprising these data types.

- Data accuracy requirements where, unlike in other domains, a small anomaly in an engineering design can have great economic and social consequences throughout a product's lifecycle.

- Requirements that the digital models and systems built today be extensible and reusable by subsequent generations of technologists, even though a digital product model may have a longer lifespan than the data formats, application software, and computing platforms used to create the model.

- The need for digital product models to be semantically rich enough to address long-term socio-technical concerns such as forensics (accident and incident investigation) and environmental issues (carbon footprint, disposal).

A series of workshops held at the National Institute of Standards and Technology (NIST) [10][12] and the University of Bath [1] over the last two years resulted in the following recommendations for advancing long-term preservation and reuse of digital engineering information:

- Develop domain-specific preservation criteria and metrics.

- Build a registry representing and classifying engineering and scientific digital objects.

- Determine how best to capture workflows of business and manufacturing processes, and develop software tools to help automate the process capture.

- Collect and preserve archiving case studies. Case studies can provide lessons learned by pointing out examples of poorly organized archived data.

- Collaborate with other groups concerned with long-term access to engineering designs.

- Develop new representation methods for both product and process information.

- Define domain-specific extensions to the Open Archival Information System (OAIS) reference model [2].

- Anticipate future access requirements for managing archived digital objects.

This paper focuses primarily on the last three bullets by presenting an approach to addressing archiving challenges evident in the U.S. Navy's Torpedo Weapon Retriever (TWR) data set. The unstructured collection includes product model databases and thousands of files in a variety of formats encompassing 3D product models, computer-aided design (CAD) data, publications, and drawings. The amount of data and proliferation of formats make it hard to manage and search the collection, let alone ensure the collection's long-term accessibility. The volume of data, diversity of information sources, and multitude of formats make the TWR a good case study for long-term archival of ship information in the general case, as well as for other engineering domains.

In order to build an effective archival information system from the TWR data, an ingest process must address the complexity and diversity of the information sources, as well as requirements for access and reuse. Ingest must therefore include the generation of descriptive metadata to aid in future access of TWR Archival Information Packages (AIPs), content units with associated metadata preserved within an OAIS. Descriptive metadata, referred to as *Descriptive Information* in the Open Archival Information System (OAIS) reference model, is essential for supporting long-term access requirements for ship product model data. OAIS defines descriptive metadata to be "the set of information, consisting primarily of Package Descriptions, which is provided to Data Management to support the finding, ordering, and retrieving of OAIS information holdings by Consumers." Other sources [9][13][14] give somewhat different definitions, but these definitions all agree that descriptive metadata should support searching and discovery of content.

**Table 1. Items affecting descriptive metadata requirements.**

| Influence | Ship Examples |
|---|---|
| Product model data categories | Modeling and simulation, characteristics, systems, geometry |
| Product model subsystems | Molded forms, hull structure, equipment, piping, HVAC (heating, ventilating, and air conditioning), foundations |
| Information access use cases | New ship design, engineering analysis, cost estimation, decommissioning, historical review, bid for construction, logistical support, operation, deployment |

Table 1 shows some factors influencing descriptive metadata requirements, providing ship domain-specific examples of each factor.

Through determining descriptive metadata requirements and by attempting to successfully generate TWR descriptive metadata during ingest, we can (1) gain insight into the feasibility of generalizing ship data ingest and customizing tools developed by the digital library community for engineering-specific applications, (2) help to coordinate ongoing long-term archiving efforts in the engineering design community with the work of digital library researchers and archivists, and (3) establish requirements for long-term archiving of engineering data produced by computer-aided design/engineering and product data/lifecycle management software applications.

## 2. ADDITIONAL DESCRIPTIVE METADATA ISSUES

In this section we discuss in depth some additional factors beyond those listed in Table 1 determining descriptive product model metadata. We first present a classification of types engineering data archive access we call the "3Rs" [11], a generalization of the last row in Table 1. We then discuss the multitude, complexity, and evolution of digital formats for product model data, using the TWR and shipbuilding domains as illustrative examples.

### 2.1 Engineering Informatics and the "3Rs"

The ability to replicate the behavior of the artifact or the experiment in the validation of science and engineering knowledge is crucial. This requires that the information be available in the best form for retrieval and reuse. The need to know a designer's intent becomes important in the context of redesign and reuse of existing parts. Another important aspect of engineering archiving is the ability to store the digital objects at different levels of granularity and abstractions as required by the design decision-making tasks. Without such an ability to compose different digital objects for archiving it would not be possible to maintain the ability to encode reuse or rationale-based access needs.

We therefore consider end-user needs from the point of view of reference, reuse, and rationale – the "3Rs" – to better understand the level of granularity and abstractions required in the definition of digital objects. By "end user" we mean what the OAIS reference model refers to as the *designated community*.

The 3Rs – reference, reuse, and rationale – define a taxonomy of designated community access scenarios. By reference we mean the ability to read the digital object and produce the digital object for proper reproduction in a given display medium (computer display, paper, etc.). We use the term reuse to mean the ability to refer to and modify the digital object in an appropriate system environment (software and hardware). The rationale is the highest level of access in which the end user should be able to refer, reuse, and explain the decisions about the content of the digital object.

The primary driver for the 3Rs is the special retrieval needs for each of these scenarios. For example a collection intended primarily for reference may need to be organized differently than one intended for reuse, where not only the geometric aspects of

the product are sought but also additional information regarding manufacturing, part performance, assembly, and other aspects. In a similar vein, rationale information may have to be packaged differently in that it may include requirements information along with other performance data on the part or the assembly. Given the range of uses and perspectives of the end users, their needs will have a large impact on the process of archiving and retrieval.

## 2.2 Variety of Formats

The number of formats available to represent product model data is seemingly limitless. The major classes are proprietary native, open native, standards-based neutral, and ad-hoc neutral. The selection of the product model data format is dependent upon several variables. These include the type of data defined in the product model, the frequency and length of the data exchange program, the maturity of the data definition specification, and the availability of translators.

The reality is that an archive must capture all of the data required to completely define the product, and in some instances, processes. Therefore the archive must accommodate the full spectrum of data formats. This can cause redundancies, further complicating matters. The repository is an indexed collection of data – nothing less, nothing more. The presence of data is not an indicator of the quality of data. The burden of determining both the utility and quality of the data is the responsibility of the consumer. What the repository *can* do is make this process easier by providing good metadata.

A native format is one created by an application. Typically it is binary, proprietary, prone to change, and its specification is not available to the general public. By maintaining tight control over the definition of the native data format, the software provider can optimize the data to its own requirements. One of the advantages of a native format is the ability to share the data with other users of the same system with a relatively high confidence that the complete model can be exchanged without any loss of data. Generally this is true if all users have the same or compatible software version, if none have customized the software, and if all users have identical or compatible libraries. To further complicate matters, just because the application can successfully read and display the native data does not mean the user can interpret the model. This is because in some cases the context of the model may be defined using procedures peculiar to the project. Assuming users can obtain the hardware, software, documentation, and can figure out how the system works, then the native format provides the most reliable and accurate representation. In the near term, this is not much of a problem, but in the long term it can be. In spite of these barriers to long-term access, archiving the native data is done as a matter of course because it has such a small impact on resources and because it is universally accepted as a good system management practice [6].

The Standard for the Exchange of Product Model Data (ISO 10303) – informally known as STEP (the STandard for the Exchange of Product model data) [7][20] – provides an open and stable means for long-term retention of product information. STEP application protocols (APs) specify information models for a specific engineering domain. STEP physical files (informally known as Part 21 files or STEP files) use an ASCII format defined in ISO 10303-21 [3]. A STEP processor can be any software application capable of interpreting and/or generating STEP physical files, for example a CAD tool capable of importing and exporting STEP files, or a visualization tool that can import STEP data. The objects represented and exchanged using STEP, as well as the associations between these objects, are defined in schemas written in EXPRESS (ISO 10303-11) [4], an information modeling language combining ideas from the entity-attribute-relationship family of modeling languages with concepts from object-oriented modeling.

Although new STEP APs continue to be developed today, EXPRESS and the Part 21 format were developed in the 1980s when few, if any, alternatives existed that had the combined representational capabilities needed for STEP's ambitious scope. Although EXPRESS is a powerful language, it is relatively unknown to most programmers. The Part 21 syntax, although effective for the task at hand, lacks extensibility, can be hard for humans to read, and – perhaps most limiting – is computer-interpretable by a relatively small number of software development toolkits that support STEP.

Developers and users of STEP have long realized that, in order to STEP APs to achieve maximum use, they need to be specified in more modern and Internet-friendly languages and formats [15]. In 3.2, we suggest such a mapping in order to facilitate the development of descriptive metadata for long-term archival.

## 2.3 Format Evolution

The modern warship may arguably be the most complex product known to man. It is large, has a huge number of parts, and may be in service for decades. Over the past several thousand years, ship designs and records have been maintained on paper. It is only recently that ship design, construction, and life cycle support data have been developed and maintained digitally. This has led to several problems. On one hand, digitization allows designers, operators, and logisticians to capitalize on the ship's compartmentalization to perform their tasks at a new level of efficiency. On the other hand, the use of multiple disparate and often incompatible systems has led to a new level of redundancy, conflict, and in some cases a reduction in the availability of information.

A ship data archivist must be prepared to handle design data in different formats, with widely varying definitions and levels of detail. For example, during early stage design, the theoretical molded surfaces of the ship can be captured in a system that is incompatible with the tool used to loft plates[1] and arrange structure. In order to transfer the mold data to the tool used for structural design, an Initial Graphics Exchange Specification (IGES) [19] file may be generated. During this phase the structure is modeled at a relatively low level of detail, referred to as a scantling model. Several interactions may be required between the scantling model and a structural finite element (FE) model. The FE model is used to generate a data file that feeds an FE solver. The solver generates yet another file containing the results. The analyst then modifies the scantling model, which may trigger changes to other models. Simultaneously, these changes may

---

[1] "To loft plates" is a marine term denoting the layout of metal plates for a ship. It comes from the days when sails laid out in a large open room. Later on, the mold loft was the large open room where plate templates were laid out.

require modifications be made to several upstream processes, including hullform design. During the next iteration of this design phase, it is almost certain there will be another change to the hullform, necessitating revisions in all downstream processes.

It is critically important that all of this knowledge be captured in a configuration controlled environment. In this small example, it can be seen that dozens of different files of varying formats will be generated. As time progresses, the level of dynamics in a ship design decreases, but the level of detail and complexity increases. This is further exacerbated by the release of new software versions that frequently are not backward compatible. This problem is not only limited to the closed proprietary native systems. It may possibly apply to open standards-based neutral systems as well if a previous release of a standard is not upwardly compatible to a later version. As technology matures, commercial translators that can read an older data file, even one that complies with an international standard may not be available.

We mentioned in 2.2 that the definition of the product model data and knowledge itself can be found in widely varying sources ranging from native word processing files to requirements databases, to CAD systems, to standards-based neutral representations such as STEP. But, although open non-proprietary standardized formats are desirable for long-term data retention, even information represented using standards-based methods is subject to format evolution. As an example, consider a TWR data long-term archival scenario maximizing the use of STEP to represent detail design product model data.

Although STEP APs such as AP218 (ship structures) and AP227 (plant spatial configuration) define information models well-tailored to the ship domain, these APs have not yet been implemented in commercial off-the-shelf software. On the other hand, other APs such as AP203 (configuration controlled 3D design) and AP214 (core data for automotive mechanical design processes) are supported by today's CAD software applications. As time progresses, so do both the STEP standard as well as the translators provided by the CAD vendors. This means, as shown in Figure 1, that the STEP representation of the data may change over the long term [5].

Initially, the data is created in native format, and the neutral file format is selected as a function of the quality of the available translators. If the desired translators are not available, a compromise will have to be made in order to allow the data to be accessible to the applications used during a specific design phase. The evolution from this point forward could be as follows.

First, geometry is exchanged using any means possible, but most probably using AP214 or AP203. The non-graphical data can be extracted separately and saved in a project-specific Extensible Markup Language (XML) [22] format. The XML often contains only product properties and perhaps minimal product structure. This approach is sufficient to enable minimal exchange of geometric data, graphics, and basic properties. It also has the greatest potential for minimizing the dependency on the product model software supplier.
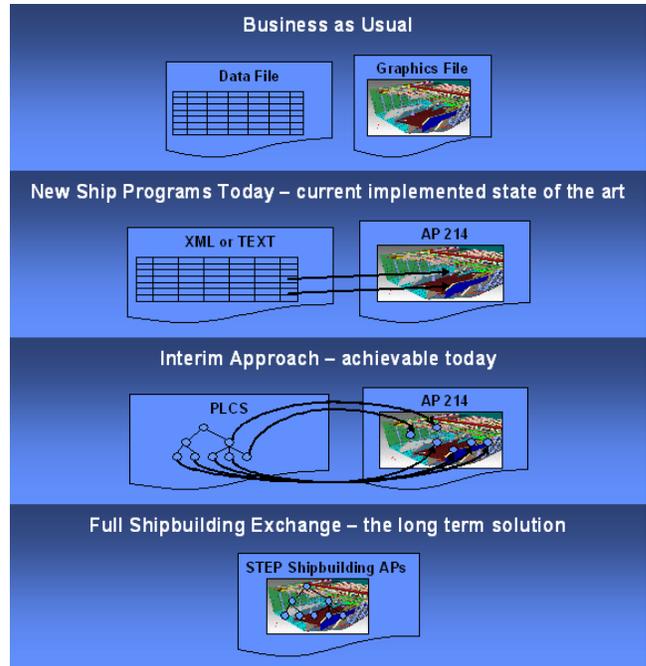


**Figure 1. Evolution of representation formats for ship product models.**

As the archival activity progresses, it is hoped that the product model software supplier will implement AP239 (Product Life Cycle Support) [16], an AP already gaining traction in the United States and United Kingdom defense communities [18], to define product structure, the relationships between objects, and reference data libraries to define an extensible set of properties.

Assuming future vendor support, the ideal long-term solution would be to employ implementations of AP214 for general purpose geometry and AP239 for configuration management and product structure, but to define application-specific product model data using APs with a more specific scope, such as AP218 or AP227.

The lesson learned from this scenario is that, even with a policy favoring STEP for representing product information, the detail design data may manifest itself in many different formats, each conforming to a different AP, or in some cases not conforming to STEP at all. Over the life of the product, format choices will evolve as a function of the product model complexity, the quality of the translators, and the evolution of the standard itself. Therefore, the goal of a single standard to represent product model data may not be realistic, and we should be prepared to encounter multiple formats over the lifecycle of a product.

## 3. SUGGESTED TECHNICAL APPROACH

We suggest a two-pronged approach to addressing the issues mentioned in Table 1 and discussed in Section 2. The first aspect of this approach is to follow the OAIS reference model description of the ingest function, attempting to maximize use of existing tools developed by the digital preservation community. The second aspect is to leverage the robustness and permanence of STEP while, at the same time, applying a transformation from EXPRESS and Part 21 to the Web Ontology Language (OWL) [21] to facilitate the creation of descriptive metadata.

## 3.1 Role of Descriptive Metadata in the Context of OAIS Ingest

Our approach to TWR data archival uses the OAIS reference model. Our goal is to develop a TWR Submission Information Package template, i.e., a "Ship SIP," and to generate accompanying descriptive metadata to aid in future access of TWR Archival Information Packages (AIPs). The generated descriptive information is essential for supporting the various engineering designated community access scenarios classified using the "3Rs." We expect lessons learned and insights gained through attempting to generate descriptive information from the TWR product model data to be a major contribution of our research.

In OAIS, a SIP must contain not only content information, but also associated preservation description information (PDI) to aid in the preservation of the content and packaging information, i.e., metadata delimiting and identifying the content information and PDI. Standards and tools have been developed by digital library researchers to aid in encoding PDI and packaging information. These technologies are generic, i.e., they were not designed specifically with engineering applications in mind. However, they are extensible. Thus our plan is to leverage these existing technologies when feasible to create the Ship SIP, tailoring them as needed to properly address domain-specific requirements.

Figure 2 illustrates the workflow of our approach. In the lower left corner, TWR content is augmented with PDI and rules or packaging information as needed for ingest. The augmented TWR content then serves as input to two activities, each shown as a rounded rectangle. The first is the generation of descriptive information needed for future access of the archived TWR content using methods we will develop. The second is creation of a Ship SIP, using existing third party tools when feasible. The Ship SIP and accompanying descriptive information are then ingested into an archive, again using third party applications as appropriate.
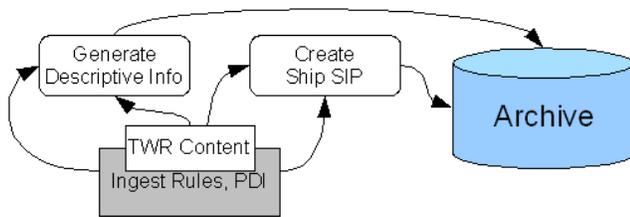


**Figure 2. Ingest workflow.**

## 3.2 Use of Semantic Technology

In [8], a transformation of STEP information models and Part 21 data is presented to represent both geometry and non-geometric data (such as function, behavior, requirements, weight, spring or damper rates, flows, and ground clearance height) in a manner more conducive to generating descriptive metadata. The transformation, called "OntoSTEP," produces semantically enriched product models as output. With the semantic enrichment of STEP, reasoning and inference mechanisms can be applied. These reasoning mechanisms allow us to check the validity of the models, to check the consistency of the instances and to infer new knowledge. This transformation of STEP is developed using OWL-DL, a sublanguage of OWL that is both computationally complete and decidable. OWL-DL has a formal foundation

(description logics) that allows for automated reasoning mechanisms. OWL-DL enables integration with software tools by providing an XML serialization of the ontology. A plug-in to a CAD application could transform STEP geometry into OWL and then allow the insertion, reading and editing of the non-geometric information of the designed product. Geometry and non-geometric information would then be represented in a unique consistent model.

The Part 21 files (EXPRESS instances) are classified using a reasoner according to an OWL translation of the AP214 EXPRESS schema. The outputs of this classification are OWL individuals corresponding to the Part 21 source files. Figure 3 illustrates this scenario using as an example a STEP representation of a gear object. OntoSTEP can enable the 3Rs as follows:

**Reference** – Users can retrieve semantically rich information for 2D and 3D visualization. The visualization information can be annotated with appropriate metadata using OntoSTEP. The language of annotation can be defined using OntoSTEP.

**Reuse** – Users can retrieve geometry and non-geometric information in a semantically rich fashion using OntoSTEP. The information can be reused effectively as it has formal semantics. The annotations can be effectively executed to carry out the instructions embedded in the annotations. The users can also retrieve information based on object properties and can select and modify shapes based on metadata.

**Rationale** – Since a project model created and annotated using OntoSTEP has formal semantics and allows automatic inference, it will be more capable of answering "why" questions and reasoning about design decisions than would an EXPRESS/Part 21 product model.

## 4. CONCLUSIONS

Descriptive metadata is key to successful long-term preservation and reuse of digital product models. Using the 3Rs as a framework and the Navy's TWR data repository as a guiding example, we enumerated several product model data issues that descriptive metadata must address. We then suggested an approach using semantic technology to create the descriptive metadata within the context of the OAIS reference model.

The next step in this research is to define a descriptive metadata schema for TWR and other product model data addressing the requirements discussed. The schema could then be populated using the OntoSTEP approach highlighted in 3.2. As we mentioned in 3.1, defining this schema and developing algorithms for instantiating it will be a major milestone and research contribution.

Although we focus on the ship product model domain in this paper, we expect our results to be relevant to other engineering domains involving complex products with long life cycles. The aerospace industry [17] and – more recently – the automotive industry [23] are adopting a rigorous and standards-based approach to long-term archiving of their product model data. If our research can benefit their respective efforts, then we will have made significant progress toward achieving the objectives stated at the end of Section 1.
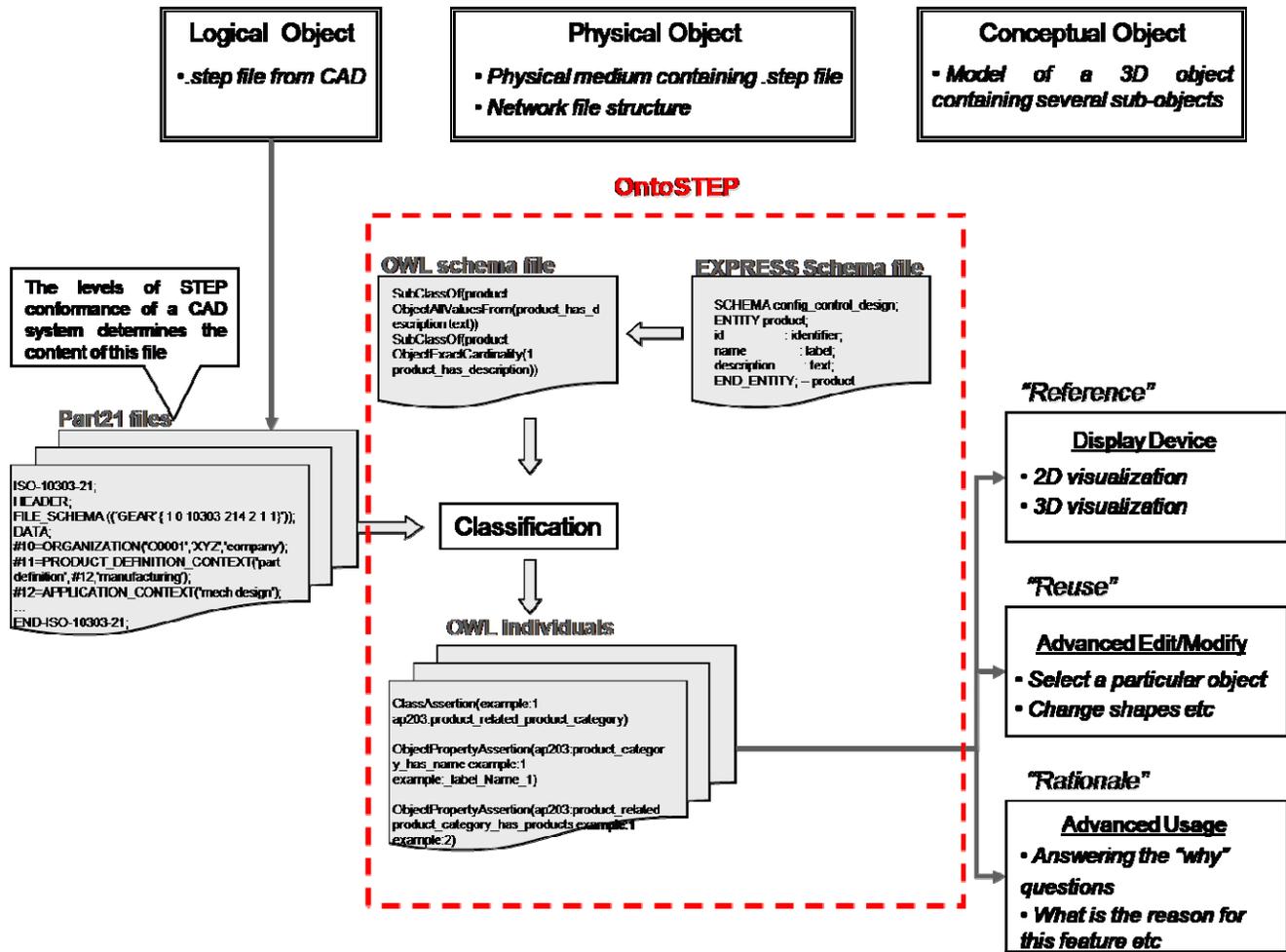
**Figure 3. Product Ontology and 3Rs.**

# 5. REFERENCES

[1] Ball. Ding. *Proceedings of the Atlantic Workshop on Long Term Knowledge Retention*. University of Bath. 13 April 2007. http://www.ukoln.ac.uk/events/ltkr-2007/proceedings/

[2] Consultative Committee for Space Data Systems. CCSDS 650.0-B-1: Reference model for an open archival information system (OAIS). Blue Book. Issue 1. ISO 14721:2003. 2005. http://public.ccsds.org/publications/archive/650x0b1.pdf

[3] ISO 10303-21:2002. Industrial automation systems and integration - Product data representation and exchange - Part 21: Implementation methods: Clear text encoding of the exchange structure.

[4] ISO 10303-11:2004: Industrial automation systems and integration -- Product data representation and exchange -- Part 11: Description methods: The EXPRESS language reference manual.

[5] Kassel. Briggs. *An Alternate Approach to the Exchange of Ship Product Model Data*. In Proceedings of the 2007 SNAME Maritime Technology Conference and Ship Production Symposium (Fort Lauderdale, Florida, November 14-16, 2007). http://www.sname.org/AM2007

[6] Kassel. David. *Long Term Retention of Product Model Data*. In Proceedings of the 2006 SNAME Maritime Technology Conference and Ship Production Symposium (Fort Lauderdale, Florida, October 10-13, 2006). http://www.sname.org/AM2006

[7] Kemmerer (Ed.). *STEP: The grand experience*. NIST Special Publication 939. National Institute of Standards and Technology. 1999. http://www.mel.nist.gov/publications/view_pub.cgi?pub_id=821224

[8] Krima. Barbau. Fiorentini. Rachuri. Sriram. *OntoSTEP: Incorporating Semantics into STEP*. NIST. In preparation.

[9] Library of Congress. *Introduction to Metadata Elements*. http://www.loc.gov/standards/metadata.html

[10] Lubell. Mani. Subrahmanian. Rachuri. *Long Term Sustainment Workshop Report*. NIST Interagency/Internal Report (NISTIR) – 7496. March 2008. http://www.mel.nist.gov/publications/view_pub.cgi?pub_id=824648

[11] Lubell. Rachuri. Mani. Subrahmanian. *Sustaining Engineering Informatics: Towards Methods and Metrics in Digital Curation*. International Journal of Digital Curation. Vol 3. No 2 (2008).
http://www.ijdc.net/index.php/ijdc/article/view/87

[12] Lubell. Rachuri. Subrahmanian. Regli. *Long Term Knowledge Retention Workshop Summary*. NIST Interagency/Internal Report (NISTIR) – 7386. January 2007. http://www.mel.nist.gov/publications/view_pub.cgi?pub_id= 822647

[13] National Information Standards Organization. *Understanding Metadata*. NISO Press (2004).
http://www.niso.org/publications/press/UnderstandingMetad ata.pdf

[14] Oxford Digital Library. *Metadata in the Oxford Digital Library*. http://www.odl.ox.ac.uk/metadata.htm

[15] Peak. Lubell. Srinivasan. Waterbury. *STEP, XML and UML: Complementary Technologies*. Journal of Computing and Information Science in Engineering. Vol. 4. No. 4 (2004).

[16] Product Life Cycle Support (PLCS). http://www.plcs-resources.org/

[17] ProSTEP iViP Association. Long Term Archiving – LOTAR. http://www.prostep.org/lotar

[18] Rachuri. Foufou. Kemmerer. *Analysis of Standards for Lifecycle Management of Systems for US Army – a preliminary investigation*. NIST Interagency/Internal Report (NISTIR) – 7339. August 2006. http://www.mel.nist.gov/publications/view_pub.cgi?pub_id= 822627

[19] U.S. Product Data Association. Initial Graphics Exchange Specification IGES 5.3. Formerly ANS US PRO/IPO-100-1996. http://www.uspro.org/documents/IGES5-3_forDownload.pdf

[20] Wikipedia. ISO 10303.
http://en.wikipedia.org/wiki/ISO_10303

[21] World Wide Web Consortium. OWL Web Ontology Language Overview. W3C Recommendation 10 February 2004. http://www.w3.org/TR/owl-features/

[22] World Wide Web Consortium. Extensible Markup Language (XML) 1.0 (Fifth Edition). W3C Recommendation 26 November 2008. http://www.w3.org/TR/xml/

[23] Yunas.        Long-Term Archiving: Defining System Functionality. AIAG Engineering & Product Development newsletter. Automotive Industry Action Group. February 2009.