

Combining Source, Content, Presentation, Narration, and Relational Representation

Fulya Horozal, Alin Iacob, Constantin Jucovschi,
Michael Kohlhase, Florian Rabe
{f.horozal,a.iacob,c.jucovschi,m.kohlhase,f.rabe}@jacobs-university.de

Computer Science, Jacobs University, Bremen

Abstract. In this paper, we try to bridge the gap between different dimensions/incarnations of mathematical knowledge: MKM representation formats (content), their human-oriented languages (source, presentation), their narrative linearizations (narration), and relational presentations used in the semantic web. The central idea is to transport solutions from software engineering to MKM regarding the parallel interlinked maintenance of the different incarnations. We show how the integration of these incarnations can be utilized to enrich the authoring and viewing processes, and we evaluate our infrastructure on the LATIN Logic Atlas, a modular library of logic formalizations, and a set of computer science lecture notes written in \LaTeX – a modular, semantic variant of \LaTeX .

1 Introduction

The Mathematical Knowledge Management (MKM) community has developed XML-based content representations of mathematical formulae and knowledge that are optimized for machine-to-machine communication. They serve as archiving formats, make mathematical software systems and services interoperable and allow to develop structural services like search, documentation, and navigation that are independent of mathematical foundations and logics.

However, these formats are — by their nature — inappropriate for the communication with humans. Therefore the MKM community uses languages that are less verbose, more mnemonic, and often optimized for a specific domain for authoring. Such human-oriented languages (we call them source languages) are converted — via a complex compilation process — into the content representations for interaction with MKM services, ideally without the user ever seeing them. In addition, we have designed presentation-oriented languages that permit an enriched reading experience compared to the source language.

This situation is similar to software engineering, where programmers write code, run the compiled executables, build HTML-based API documentations, but expect, e.g., the documentation and the results of debugging services in terms of the sources. In software engineering, scalable solutions for this problem have been developed and applied successfully, which we want to transfer to MKM.

The work described here originates from our work on two large collections of mathematical documents: our LATIN logic atlas [KMR09] formalized in the logical framework LF; and our General Computer Science lecture notes written in L^AT_EX. Despite their different flavor, both collections agree in some key aspects: They are large, highly structured, extensively inter-connected, and both authoring and reading call for machine support.

Moreover, they must be frequently converted between representation dimensions optimized for different purposes: a human-friendly input representation (source), a machine-understandable content markup (content), interactive documents for an added-value reading experience (presentation-content parallel markup), a linearized structure for teaching and publication (narration), and a network of linked data items for integration with the semantic web (relational).

In Sect. 2, we first give an overview over these document collections focusing on the challenges they present to knowledge management. Then we design a knowledge representation methodology that integrates these different dimensions in Sect. 3. In Sect. 4 and 5, we show how we leverage this methodology in the authoring and the viewing process both of which benefit from a seamless integration of the knowledge dimensions.

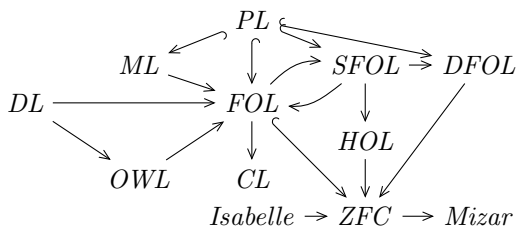
2 Structured Document Collections

2.1 The LATIN Logic Atlas

The LATIN Logic Atlas is a library of formalizations of logics and related formal systems as well as translations between them. It is intended as a reference and documentation platform for logics commonly used in mathematics and computer science. It uses a foundationally unconstrained logical framework based on modular LF and its Twelf implementation [HHP93; PS99; RS09] and focuses on modularity and extensibility.

The knowledge in the Logic Atlas is organized as a graph of LF signatures and signature morphisms between them. The latter are split into inheritance translations (inclusions/imports) and representation theorems, which have to be proved. It contains formalizations of type theories, set theories, and logics. Among them are, for example, propositional (*PL*), first (*FOL*) and higher-order logic (*HOL*), sorted (*SFOL*) and dependent first-order logic (*DFOL*), description logics (*DL*), modal (*ML*) and common logic (*CL*) as illustrated in the diagram below. Single arrows (\rightarrow) in this diagram denote translations between formalizations and hooked arrows (\hookrightarrow) denote imports.

All logics are designed modularly formed from orthogonal theories for individual connectives, quantifiers, and axioms. For example, the classical \wedge connective is only declared once in the whole Logic Atlas, and



the axiom of excluded middle and its consequences reside in a separate signature. We also use individual theories for syntax, proof theory, model theory so that the same syntax can be combined with different interpretations.

As a running example, we introduce a very simple fragment of the formalization of the syntax of propositional logic.

Example 1 (Propositional Logic) The formalization of propositional logic syntax consists of the LF signatures illustrated in Fig. 1. We focus on the structural aspects and omit the details of LF. We four signatures living in two different namespaces. `BASE` declares a symbol for the type of propositions. It is imported into `CONJ` and `IMP` which declare conjunction and implication, respectively, and these are imported `PROP`.

```
%namespace = "http://cds.omdoc.org/logics"
%sig BASE = {
  %% Type of propositions
  o : type.
}

%namespace = "http://cds.omdoc.org/logics/propositional"
%sig CONJ = {%include BASE. ...}
%sig IMP  = {%include BASE. ...}
%sig PROP = {%include CONJ. %include IMP.}
```

Fig. 1. Formalization of Propositional Logic Syntax in Twelf

Overall, the atlas contains over 500 LF signatures, and their highly modular structure yields a large number of inheritance edges. Additionally, the representation theorems include:

- the translation from unsorted to sorted first-order logic (which is almost but not quite an inclusion),
- the translations by relativization of quantifiers from sorted first-order, modal, and description logics to unsorted first-order logic, in most cases including the translation of the model theory,
- the translation from propositional and sorted first-order logic to Andrews-style higher-order logic,
- the negative translation from classical to intuitionistic logic,
- the translation from type theory to set theory that interprets types as sets and terms as elements, including a translation of Isabelle/HOL to ZF set theory,
- the Curry-Howard correspondence between logic, type theory, and category theory.

All translations include translations of the proof theory, which guarantees their proof theoretical soundness due to the type-preservation of signature morphisms.

This leads to a highly interlinked non-linear structure of the atlas. Moreover, it is designed highly collaboratively with strong interdependence between the developers. Therefore, it leads to a number of MKM challenges.

For example, the LF modules are distributed over files and these files over directories. This structure is semantically transparent because all references to modules are made by URIs. The URIs are themselves hierarchical grouping the modules into nested namespaces. It is desirable that these namespaces do not have to correspond to source files or directories. Therefore, the mapping between URIs and URLs is non-trivial and a separate management challenge.

Another problem is that encodings in LF are typically very difficult to read for anybody but the author. In a collaborative setting, it is desirable to interact with the logic graph not only through the LF source syntax but also through browsable, cross-referenced XHMTL+MathML. These should be interactive and for example permit looking up the definition of a symbol or displaying the reconstructed type of a variable. While we have presented such an interface in [GLR09] already, the systematic integration into the authoring process, where the state of the art is a text editor, has so far been lacking.

2.2 Computer Science Lecture Notes

The GenCS corpus consists of the course notes and problems of a two-semester introductory course in Computer Science [Koh] held at Jacobs University by one of the authors in the last eight years. The course notes currently comprise 300 pages with over 500 slides organized in over 800 files; they are accompanied by a database of more than 1000 homework/exam problems. All course materials are authored and maintained in the $\mathfrak{sT}\mathfrak{E}\mathfrak{X}$ format [Koh08], a modular, semantic variant of $\mathfrak{L}\mathfrak{A}\mathfrak{T}\mathfrak{E}\mathfrak{X}$ that shares the information model with OMDoc; see Fig. 2 for an example. In our nomenclature, $\mathfrak{sT}\mathfrak{E}\mathfrak{X}$ is used as a source language that is transformed into OMDoc via the $\mathfrak{L}\mathfrak{A}\mathfrak{T}\mathfrak{E}\mathfrak{X}\mathfrak{M}\mathfrak{L}$ daemon [GSK11]. For debugging and high-quality print the $\mathfrak{sT}\mathfrak{E}\mathfrak{X}$ sources can also be typeset via `pdflatex`, just as ordinary $\mathfrak{L}\mathfrak{A}\mathfrak{T}\mathfrak{E}\mathfrak{X}$ documents. The encoding makes central use of the modularity afforded by the theory graph approach; knowledge units like slides are encoded as “modules” (theories in OMDoc) and are interconnected by theory morphisms (module imports). Modules also introduce concepts via `\definiendum` and semantic macros via `\symdef`, these are inherited via the module import relation.

3 A Multi-Dimensional Knowledge Representation

3.1 Dimensions of Knowledge

In order to address the knowledge management challenges outlined above, we devise a methodology that permits the parallel maintenance of the orthogonal dimensions of the knowledge contained in a collection of mathematical documents. It is based on two key concepts: (i) a hierarchic organization of dimensions and knowledge items in a file-system-like manner inspired by the project view from software engineering, and (ii) the use of MMT URIs [RK10] as a standardized way to interlink both between different knowledge items and between the different dimensions of the same knowledge item.

```

\begin{module}[id=trees]
  \symdef[name=tdepth]{tdepthFN}{\text{dp}}
  \symdef{tdepth}[1]{\prefix\tdepthFN{#1}}
  \begin{definition}[id=tree-depth.def]
    Let  $\text{\defeq{T}}{\text{\tup{V,E}}}$  be tree, then the  $\text{\definiendum [tree-depth]{depth}}$ 
 $\text{\tdepth{v}}$  of a node  $\text{\inset{v}{V}}$  is defined recursively:  $\text{\tdepth{r}}=0$  for
    the root  $\text{r}$  of  $\text{T}$  and  $\text{\tdepth{w}}=1+\text{\tdepth{w}}$  if  $\text{\inset{tup{v,w}}{E}}$ .
  \end{definition}
  ...
\end{module}

\begin{module}[id=binary-trees]
  \importmodule[\KWARCslides{graphs-trees/en/trees}]{trees}
  ...
  \begin{definition}[id=binary-tree.def,title=Binary Tree]
    A  $\text{\definiendum[binary-tree]{binary tree}}$  is a  $\text{\termref[cd=trees,name=tree]{tree}}$ 
    where all  $\text{\termref[cd=graphs-intro,name=node]{nodes}}$ 
    have  $\text{\termref[cd=graphs-intro,name=out-degree]{out-degree}}$  2 or 0.
  \end{definition}
  ...
\end{module}

```

Fig. 2. Semiformalization of two course modules

The MMT URI of a toplevel knowledge item is of the form $g?M$ where g is the namespace and M the module name. Namespaces are URIs of the form $\langle\langle\text{scheme}\rangle\rangle://[\langle\langle\text{userinfo}\rangle\rangle]D_1\dots D_m[:\langle\langle\text{port}\rangle\rangle]/S_1/\dots/S_n$ where the D_i are domain labels and the S_i are path segments. Consequently, $g?M$ is a well-formed URI as well. $\langle\langle\text{userinfo}\rangle\rangle$, and $\langle\langle\text{port}\rangle\rangle$ are optional, and $\langle\langle\text{userinfo}\rangle\rangle$, $\langle\langle\text{scheme}\rangle\rangle$, and $\langle\langle\text{port}\rangle\rangle$ are only permitted so that users can form URIs that double as URLs — MMT URIs differing only in the scheme, userinfo, or port are considered equal.

We arrange a collection of mathematical documents as a folder containing the following subfolders, all of which are optional:

- source** contains the source files of a project. This folder does not have a predefined structure.
- content** contains a semantically marked up representation of the source files in the OMDoc format. Every namespace is stored in one file whose path is determined by its URI. Modules with namespace $D_1. \dots .D_m/S_1/\dots/S_n$ reside in an OMDoc file with path `content/ $D_m/\dots/D_1/S_1/\dots/S_n$.omdoc`. Each module carries an attribute `source="/PATH?colB:lineB-colE:lineE"` giving its physical location as a URL. Here `PATH` is the path to the containing file in the source, and `colB`, `lineB`, `colE`, and `lineE` give the begin/end column/line information.
- presentation** contains the presentation of the source files in the XHTML+MathML format with JOBAD annotations [GLR09]. It has the same file structure as the folder `content`. The files contain XHTML elements whose body has one child for every contained module. Each of these module has the attribute `jobad:href="URI"` giving its MMT URI.
- narration** contains an arbitrary collection of narratively structured documents. These are OMDoc files that contain narrative content such as sectioning and transitions, but no modules. Instead they contain reference elements of the form `<mref target="MMTURI"/>` that refer to MMT modules. It is common but not necessary that these modules are present in the `content` folder.

`relational` contains two files containing an RDF-style relational representation of the content according to the MMT ontology. Both are in XML format with toplevel element `mmtabox` and a number of children. In `individuals.abox`, the children give instances of unary predicates such as `<individual type="IsTheory"uri="MMTURI"source="PATH"/>`. In `relations.abox`, the children give instances of binary predicates such as `<relation subject="MMTURI1"predicate="ImportsFrom"object="MMTURI2"source="PATH"/>`. Usually, the knowledge items occurring in unary predicates or as the subject of a binary predicate are present in the content. However, the object of a binary predicate is often not present, namely when a theory imports a remote theory. In both cases, we use an attribute `source` to indicate the source that induced the entry; this is important for change management when one of the source files was changed.

Example 2 (Continuing Ex. 1)

The directory structure for the signatures from Ex. 1 is given in Fig. 3 using a root folder named `propositional-syntax`. Here we assume that the subfolder `source` contains the Twelf source files `base.elf` which contains the signature `BASE`, `modules.elf` which contains `CONJ` and `IMP`, and `prop.elf` which contains `PROP`.

Based on the MMT URIs of the signatures in the source files, their content representation is given as follows. The signature `BASE` has the MMT URI `http://cds.omdoc.org/logics?BASE`. The other signatures have MMT URIs such as `http://cds.omdoc.org/logics/propositional/syntax?CONJ`. The content representation of the signature `BASE` is given in the OMDoc file `content/org/omdoc/cds/logics.omdoc`. The other content representations reside in the file `content/org/omdoc/cds/logics/propositional/syntax.omdoc`.

The subfolder `presentation` contains the respective XHTML files, `logics.xhtml` and `syntax.xhtml`. All files in Fig. 3 can be downloaded at `https://svn.kwarc.info/repos/twelf/projects/propositional-syntax`.

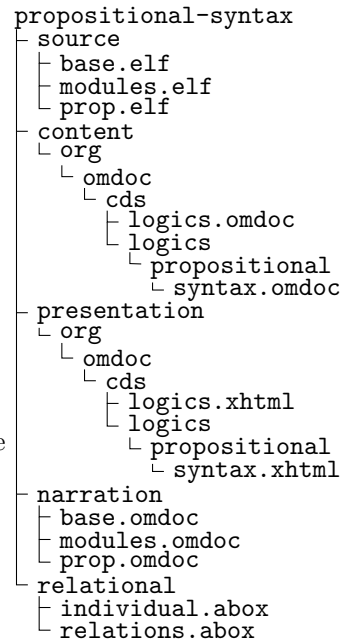


Fig. 3. Files of the Running Example

Our methodology integrates various powerful conceptual distinctions that have been developed in the past. Firstly, our distinction between the source and the content representation corresponds to the distinction between source and binary in software engineering. Moreover, our directory structure is inspired by software projects, such as in Java programming. In particular, the use of URIs

to identify content (binary) items corresponds to identifiers for Java classes. Therefore, existing workflows and implementations from software engineering can be easily adapted, for example in the use of project-based IDEs (see Sect. 4).

Secondly, the distinction between content and presentation has been well studied in the MKM community and standardized in MathML [Aus+03]. In particular, the cross-references from presentation to content correspond to the interlinking of content and presentation in the parallel markup employed in MathML, which we here extend to the level of document collections.

Thirdly, the distinction between content and narrative structure was already recognized in the OMDoc format. The general intuition there is that narrative structures are “presentations” at the discourse level. But in contrast to the formula level, presentations cannot be specified and managed via notation definitions. Instead we add narrative document structure fragments, i.e. document-structured objects that contain references to the content representations and transition texts as lightweight structures to the content commons; see [Mül10] for details and further references.

Finally, the distinction between tree-structured content representation and the relational representation corresponds to the practice of the semantic web where RDF triples are used to represent knowledge as a network of linked data. see [Lan11] for an overview.

3.2 A Mathematical Archive Format

We will now follow the parallelism to software engineering developed in the previous section: We introduce mathematical archives — `mar` files — that correspond to Java archives, i.e., `jar` files [Ora]. We define a mathematical archive to be a zip file that contains the directory structure developed in Sect. 3.1. By packaging all knowledge dimensions in a single archive, we obtain a convenient and lightweight way of distributing multi-dimensional collections of interlinked documents.

To address into the content of a `mar` archive, we also define the following URL scheme: Given a `mar` whole URL is `file:/A` and which contains the source file `source/S`, then the URL `mar:/A/S` resolves to that source file. We define the URL `mar:/A/S?Pos` accordingly if `Pos` is the position of a module given by its line/column as above.

Similarly, to the compilation and building process that is used to create `jar` files, we have implemented a building process for `mar` files. It consists of three stages. The first stage (compilation) depends on the source language and produce one OMDoc file for every source file whose internal structure corresponds to the source file. This is implemented in close connection with dedicated tools for the source language. In particular, we have implemented a translation from LF to OMDoc as part of the Twelf implementation of LF [PS99; RS09]. Moreover, we have implemented a translation from \LaTeX to OMDoc based on the \LaTeX ML daemon [GSK11].

The second stage (building) is generic and produces the remaining knowledge dimensions from the OMDoc representation. In particular, it decomposes

the OMDoc documents into modules and reassembles them according to their namespaces to obtain the content representation. The narrative dimension is obtained from the initial OMDoc representation by replacing all modules with references to the respective content item. We have implemented this as a part of the existing MMT API [KRZ10]. Finally, the API already includes a rendering engine that we use to produce the presentation and the relational representation.

Then the third stage (packaging) collects all folders in a zip archive. For LF, we integrate all three stages into a flexible command line application.

Example 3 (Continuing Ex. 2) The mathematical archive file for the running example can be obtained at <https://svn.kwarc.info/repos/twelf/projects/propositional-syntax.mar>.

3.3 Catalog Services

The use of URIs as knowledge identifiers (rather than URLs) is crucial in order to permit collaborative authoring and convenient distribution of knowledge. However, it requires a catalog that translates an MMT URIs to the physical location, given by a URL, of a resource. Typical URLs are those in a file system, in a mathematical archive, or a remote or local repository. It is trivial to build the catalog if the knowledge is already present in content form where locations are derived from the URI.

But the catalog is already needed during the compilation process: For example, if a theory imports another theory, it refers to it by its MMT URI. Consequently, the compilation tool must already be aware of the URI-to-URL mapping before the content has been produced. However, the compilation tool is typically a dedicated legacy system that natively operates on URLs already and does not even recognize URIs. This is the case for both Twelf and \LaTeX .

Therefore, we have implemented standalone catalog services for these two tools and integrated them with the respective system. In the case of Twelf, the catalog maintains a list of local directories, files, and `mar` archives that it watches. It parses them whenever they change and creates the URI-URL mapping. When Twelf encounters a URI, it asks the catalog via HTTP for the URL. This parser only parses the outer syntax that is necessary to obtain the structure of the source file; it is implemented generically so that it can be easily adapted to other formal declarative languages.

An additional strength of this catalog is that it can also handle ill-formed source representations that commonly arise during authoring. Moreover, we also use the catalog to obtain the line/column locations of the modules inside the source files so that the content-to-source references can be added to the content files.

In the case of \S\TeX , a poor man's catalog services is implemented directly in \TeX : the base URIs of the GenCS knowledge collection (see `\KWARCslides` in Fig 2) is specified by a `\defpath` statement in the document preamble and can be used in the `\importmodule` macros. The `module` environments induce internal \TeX structures that store information about the imports (`\importmodule`)

structure and semantic macros (`\symdef`), therefore these three `\STEX` primitives have to be read whenever a module is imported. To get around difficulties with selective input in `\TeX`, the `\STEX` build process excerpts a `\STEX` signature module `\langle\langle module \rangle\rangle.sms` from any module `\langle\langle module \rangle\rangle.tex`. So `\importmodule\langle\langle module \rangle\rangle.sms` simply reads `\langle\langle module \rangle\rangle.sms`.

4 The Author’s Perspective

The translation from source to a content-like representation has been well-understood. For languages like `\LF`, it takes the form of a parsing and type reconstruction process that transforms external to internal syntax. The translation from internal syntax to an `\OMDoc`-based content representation is conceptually straightforward.

However, it is a hard problem to use the content representation to give the author feedback about the document she is currently editing. For most formal mathematical languages, the state of the art is an emacs mode with syntax highlighting. Only a few systems offer further functionality. For example, the `\Agda` [Nor05] emacs mode can follow cross-references and show reconstructed types of missing terms. The `\Isabelle` [Pau94] `\jEdit` can follow cross-references and show tooltips derived from the static analysis.

A more powerful solution is possible if we always produce all knowledge dimensions using the compilation and building process as described in Sect. 3.2. Then generic services can be implemented easily, each of them based on the most suitable dimension, and we give a few examples in Sect. 4.2.

Note that this is not an efficiency problem: Typically the author only works on a few files that can be compiled constantly. It is even realistic to hold all dimensions in memory. The main problem is an architectural one, which is solved by our multi-dimensional representation. Once this architecture is setup and made available to IDE developers, it is very easy for them to quickly produce powerful generic services.

4.1 Multi-Dimensional Knowledge in an IDE

In previous work, we have already presented an example of a semantic IDE [JK10] based on `\Eclipse`. We can now strengthen it significantly by basing it on our multi-dimensional representation. Inspired by the project metaphor from software engineering, we introduce the notion of a mathematical project in `\Eclipse`.

A mathematical project consists of a folder containing the subfolder from Sect. 3.1. The author works on a set of source files in the `\source` directory. Moreover, the project maintains a `\mathpath` (named in analogy to `\Java`’s `\classpath`) that provides a set of `\mar` archives that the user wishes to include.

The IDE offers the build functionality that runs the compilation and building processes described in Sect. 3.2 to generate the other dimensions from the source dimension. The key requirement here is to gracefully degrade in the presence of errors in the source file. Therefore, we provide an adaptive parser component that consists of three levels:

The **regex level** uses regular expressions to spot important structural properties in the document (e.g. the namespace and signature declarations in the case of LF). This compilation level never fails, and its result is an OMDoc file that contains only the spotted structures and lacks any additional information of the content.

The **CFG parser level** uses a simple context-free grammar to parse the source. It is able to spot more complicated structures such as comments and nested modules and can be implemented very easily within Eclipse. Like the previous level, it produces an approximate OMDoc file, but contrary to the previous level, it may find syntax errors that are then displayed to the user.

The **full parser level** uses the dedicated tool (Twelf or \LaTeX). The resulting OMDoc file includes the full content representation. In particular, in the case of Twelf, it contains all reconstructed types and implicit arguments. However, it may fail in the case of ill-typed input.

The adaptive parser component runs all parser in order, and retains the best OMDoc file any of them returns. This file is then used as the input to produce the remaining content dimensions.

4.2 Added-Value Services

In this section we present several services typically found in software engineering tools which aim at supporting authoring process. We analyze each of these services and show that they can be efficiently implemented by using one or several dimensions of knowledge.

project explorer is a widget giving an integrated view on a project's content by abstracting from the file system location where the sources are defined. It groups objects by their content location, i.e., their MMT URI. To implement this widget, we populate the non-leaf nodes of the tree from the directory structure of the content dimension. The leaf nodes are generated by running simple XPath queries on the OMDoc files.

outline view is a source level widget which visualizes the main structural components. For LF, these include definitions of signatures and namespaces as well as constant declarations within signatures. Double-clicking on any such structural components opens the place in the source code where the component is defined. Alternatively, the corresponding presentation can be opened.

autocompletion assists the user with getting location and context specific suggestions, e.g., listing declarations available in a namespace. Fig. 4a shows an example. Note how the namespace prefix **base** is declared to point to a certain namespace, and the autocompletion suggests only signatures names declared in that namespace. The implementation of this feature requires information about the context where autocompletion is requested, which is obtained from the interlinked source and content dimensions. Moreover, it needs the content dimension to compute all possible completions. In more complicated scenarios, it can also use the relational dimension to compute the possible completions using the relational queries.

hover overlay is a feature that shows in-place meta-data about elements at the position of the mouse cursor such as the full URIs of a symbol, its type or definitions, a comment, or inferred types of variables. Fig. 4b) shows an example. The displayed information is retrieved from the content dimension. It is also possible to display the information using the presentation dimension. **definition/reference search** makes it easy for a user to find where a certain item is defined or used. Although the features require different user interfaces the functionality is very similar, namely, finding relations. Just like in the hover overlay feature, one first finds the right item in the content representation and then use the relation dimension to find the requested item(s). **theory-graph display** provides a graphical visualization of the relations among knowledge items. To implement this feature we apply a filter on the multi-graph from the relations dimension and uses 3rd party software to render it.

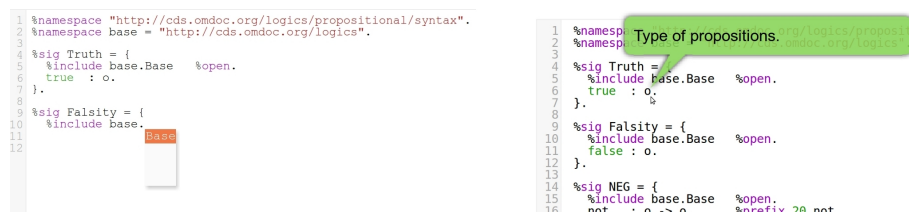


Fig. 4. a) Context aware Auto-Completion b) Metadata information on hover

5 The Reader's Perspective

We have developed the Planetary system (see [Koh+11; Dav+10; Pla] for an introduction) as the reader's complement to our IDE. Planetary is a Web 3.0 system¹ for semantically annotated document collections in Science, Technology, Engineering and Mathematics (STEM). In our approach, *documents published in the Planetary system become flexible, adaptive interfaces to a content commons* of domain objects, context, and their relations.

We call this framework the **Active Documents Paradigm** (ADP), since documents can also actively adapt to user preferences and environment rather than only executing services upon user request. Our

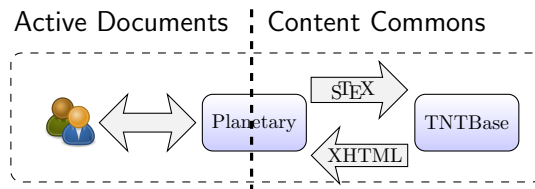


Fig. 5. The Active Documents Architecture

¹ We adopt the nomenclature where Web 3.0 stands for extension of the Social Web with Semantic Web/Linked Open Data technologies.

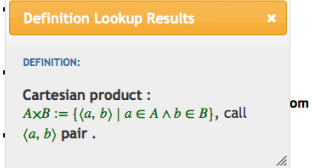
framework is based on *semantically annotated documents* together with semantic background ontologies (which we call the **content commons**). This information can then be used by user-visible, semantic services like program (fragment) execution, computation, visualization, navigation, information aggregation and information retrieval [GLR09].

The Planetary system directly uses all five incarnations/dimensions of mathematical knowledge specified in Sect 3.1. In the **content** incarnation, Planetary uses OMDoc for representing content modules, but authors create and maintain these using \LaTeX in the **source** dimension and the readers interact with the active documents encoded as dynamic XHTML+MathML+RDFa (the **source** incarnation of the material). We use the \LaTeX XML daemon [Mil; GSK11] for the transformation from \LaTeX to OMDoc, this is run on every change to the \LaTeX sources. The basic presentation process [KMR08] process for the OMDoc content modules is provided by the TNTBase system [ZK09].

But Planetary also uses the **narrative** dimension: content modules are used not only for representing mathematical theories, but also for document structures: **narrative modules** consist of a mixture of sectional markup, inclusion references, and narrative texts that provide transitions (narrative glue) between the other objects. Graphs of narrative modules whose edges are the inclusion references constitute the content representations of document fragments, documents, and document collections in the Planetary system, which generates active documents from them. It uses a process of separate compilation and dynamic linking to equip them with document (collection)-level features like content tables, indexes, section numbering and inter-module cross-references; see [Dav+11] for details.

As the active documents use identifiers that are relative to the base URI of the Planetary instance, whereas the content commons uses MMT URIs, the **semantic publishing map** which maintains the correspondence between these is a central, persistent data structure maintained by the Planetary system.

$f \subseteq X \times Y$, is called a **partial function**, iff for all $x \in X$ there is at



This is one instance of the **relational** dimension, another is used in the For instance, the RDFa embedded in the presentation of a formula (and represented in the linking part of the math archive) can be used for definition lookup as shown on the left. Actually the realization of the definition lookup service involves **presentation** (where the service is

embedded) and **content** (from which the definition is fetched to be presented by the service) incarnations as well.

In the future we even want to combine this with the **source** dimension by combining it with a `\symdef-look` service that makes editing easier. This can be thought of as a presentation-

```

\begin{omgroup}[id=sec.contfuncs]{Continuous Functions}
\begin{module}[id=continuous]
\importmodule[./background/functions]{functions}
\importmodule[./background/reals]{reals}
\symdef{continuousfunctions}[2]{\mathcal{C}^0(\#1, \#2)}
\abbrdef{ContRR}[2]{\continuousfunctions.RealNumbers.RealNumbers}
\begin{definition}[for=continuousfunctions]
A function $: \begin{matrix} \text{fun} \\ \text{RealNumbers} \\ \text{cart} \end{matrix} \begin{matrix} \text{A function } f:A \rightarrow B \text{ is a left-total, right-unique} \\ \text{relation in } A \times B \end{matrix}
\end{definition}
\end{module}
\end{omgroup}

```

triggered complement to the editor-based service on the right that looks up `\symdefs` by their definienda.

6 Conclusion

We have presented an infrastructure for creating, storing, managing, and distributing mathematical knowledge in various pragmatic forms. We have identified five aspects that have to be taken into account here: *(i)* human-oriented *source languages* for efficient editing of content, *(ii)* the modular *content representation* that has been the focus of attention in MKM so far, *(iii)* *active presentations* of the content for viewing, navigation, and interaction, *(iv)* *narrative structures* that allow binding the content modules into self-contained documents that can be read linearly, and *(v)* *relational structures* that cross-link all these aspects and permit keeping them in sync.

We have developed and tested this infrastructure on the LATIN logic atlas, a highly modular graph of logics and logic morphisms using Twelf as the source language. Our focus was on the authoring perspective using our semantic IDE. The other experiment that informed the development was the Planetary system, a semantic publishing system that we use for our lecture notes with \LaTeX as a surface language and provides an infrastructure for assembling content modules into document and collection structures.

In the future, we want to combine these systems and perspectives more tightly. For example, we could use Planetary to discuss and review logic formalizations in Twelf, or write papers about the formalizations in \LaTeX . This should not pose any fundamental problems as the surface languages are interoperable by virtue of having the same, very general data model: the OMDoc ontology. By the same token we want to add additional surface languages and presentation targets that allow to include other user groups. High-profile examples include the Mizar Mathematical Language and Isabelle/ISAR.

Finally, there is a sixth aspect that may be added to the math archive infrastructure: discussions. The Planetary system already allows localized discussions on the content modules/presentations, which form an important part of the system content. These would probably be worth saving in math archives.

References

- [Aus+03] R. Ausbrooks et al. *Mathematical Markup Language (MathML) Version 2.0 (second edition)*. Tech. rep. See <http://www.w3.org/TR/MathML2>. World Wide Web Consortium, 2003.
- [Dav+10] Catalin David et al. “eMath 3.0: Building Blocks for a social and semantic Web for online mathematics & ELearning”. In: *1st International Workshop on Mathematics and ICT: Education, Research and Applications*. (Bucharest, Romania, Nov. 3, 2010). Ed. by Ion Mierlus-Mazilu. 2010. URL: <http://kwarc.info/kohlhase/papers/malog10.pdf>.

- [Dav+11] Catalin David et al. “A Framework for Modular Semantic Publishing with Separate Compilation and Dynamic Linking”. 2011. URL: <https://svn.mathweb.org/repos/planetary/doc/sepublica11/paper.pdf>.
- [GLR09] J. Gičeva, C. Lange, and F. Rabe. “Integrating Web Services into Active Mathematical Documents”. In: *Intelligent Computer Mathematics*. Ed. by J. Carette and L. Dixon and C. Sacerdoti Coen and S. Watt. Vol. 5625. Lecture Notes in Computer Science. Springer, 2009, pp. 279–293.
- [GSK11] Deyan Ginev, Heinrich Stamerjohanns, and Michael Kohlhase. “The \LaTeX XML Daemon: A \LaTeX Entrance to the Semantic Web”. submitted. 2011. URL: <https://kwarc.eecs.iu-bremen.de/repos/arXMLiv/doc/cicm-systems11/paper.pdf>.
- [HHP93] R. Harper, F. Honsell, and G. Plotkin. “A framework for defining logics”. In: *Journal of the Association for Computing Machinery* 40.1 (1993), pp. 143–184.
- [JK10] Constantin Jucovschi and Michael Kohlhase. “sTeXIDE: An Integrated Development Environment for sTeX Collections”. In: *Intelligent Computer Mathematics*. Ed. by Serge Autexier et al. LNAI 6167. Springer Verlag, 2010, pp. 336–344. arXiv:1005.5489v1 [cs.OH].
- [KMR08] Michael Kohlhase, Christine Müller, and Florian Rabe. “Notations for Living Mathematical Documents”. In: *Intelligent Computer Mathematics*. 9th International Conference, AISC, 15th Symposium, Calculemus, 7th International Conference MKM (Birmingham, UK, July 28–Aug. 1, 2008). Ed. by Serge Autexier et al. LNAI 5144. Springer Verlag, 2008, pp. 504–519. URL: <http://omdoc.org/pubs/mkm08-notations.pdf>.
- [KMR09] M. Kohlhase, T. Mossakowski, and F. Rabe. *The LATIN Project*. See <https://trac.omdoc.org/LATIN/>. 2009.
- [Koh] *General Computer Science: GenCS I/II Lecture Notes*. <http://gencs.kwarc.info/book/1>. Semantic Course Notes in Panta Rhei. 2011. URL: <http://gencs.kwarc.info/book/1>.
- [Koh08] Michael Kohlhase. “Using \LaTeX as a Semantic Markup Format”. In: *Mathematics in Computer Science* 2.2 (2008), pp. 279–304. URL: <https://svn.kwarc.info/repos/stex/doc/mcs08/stex.pdf>.
- [Koh+11] Michael Kohlhase et al. “The Planetary System: Web 3.0 & Active Documents for STEM”. In: accepted for publication at ICCS 2011 (Finalist at the Executable Papers Challenge). 2011. URL: <https://svn.mathweb.org/repos/planetary/doc/epc11/paper.pdf>.
- [KRZ10] M. Kohlhase, F. Rabe, and V. Zholudev. “Towards MKM in the Large: Modular Representation and Scalable Software Architecture”. In: *Intelligent Computer Mathematics*. Ed. by S. Autexier et al. Vol. 6167. Lecture Notes in Computer Science. Springer, 2010, pp. 370–384.

- [Lan11] Christoph Lange. “Enabling Collaboration on Semiformal Mathematical Knowledge by Semantic Web Integration”. submitted January 31, defended March 13. PhD thesis. Jacobs University Bremen, 2011. URL: <https://svn.kwarc.info/repos/swim/doc/phd/phd.pdf>.
- [Mil] Bruce Miller. *LaTeXML: A L^AT_EX to XML Converter*. URL: <http://dlmf.nist.gov/LaTeXML/> (visited on 03/03/2011).
- [Mül10] Christine Müller. “Adaptation of Mathematical Documents”. PhD thesis. Jacobs University Bremen, 2010. URL: <http://kwarc.info/cmuedler/papers/thesis.pdf>.
- [Nor05] U. Norell. *The Agda Wiki*. <http://wiki.portal.chalmers.se/agda>. 2005.
- [Ora] Oracle. *JDK 6 Java Archive (JAR)*. <http://download.oracle.com/javase/6/docs/technotes/guides/jar>.
- [Pau94] L. Paulson. *Isabelle: A Generic Theorem Prover*. Vol. 828. Lecture Notes in Computer Science. Springer, 1994.
- [Pla] *Planetary Developer Forum*. URL: <http://trac.mathweb.org/planetary/> (visited on 01/20/2011).
- [PS99] F. Pfenning and C. Schürmann. “System Description: Twelf - A Meta-Logical Framework for Deductive Systems”. In: *Lecture Notes in Computer Science* 1632 (1999), pp. 202–206.
- [RK10] F. Rabe and M. Kohlhas. “A Scalable Module System”. To be submitted, see <http://kwarc.info/frabe/Research/mmt.pdf>. 2010.
- [RS09] F. Rabe and C. Schürmann. “A Practical Module System for LF”. In: *Proceedings of the Workshop on Logical Frameworks: Meta-Theory and Practice (LFMTP)*. Ed. by J. Cheney and A. Felty. ACM Press, 2009, pp. 40–48.
- [ZK09] Vyacheslav Zholudev and Michael Kohlhas. “TNTBase: a Versioned Storage for XML”. In: *Proceedings of Balisage: The Markup Conference*. Vol. 3. Balisage Series on Markup Technologies. Mulberry Technologies, Inc., 2009. DOI: 10.4242/BalisageVol13.Zholudev01.