

A Framework for Semantic Publishing of Modular Content Objects

Catalin David, Deyan Ginev, Michael Kohlhase, Bogdan Matican, Stefan Mirea

Computer Science, Jacobs University, Germany; <http://kwarc.info/>

Abstract. We present the Active Documents approach to semantic publishing (semantically annotated documents associated with a content commons that holds the background ontologies) and the *Planetary* system (as an active document player).

In this paper we explore the interaction of content object reuse and context sensitivity in the presentation process that transforms content modules to active documents. We propose a “separate compilation and dynamic linking” regime that makes semantic publishing of highly structured content representations into active documents tractable and show how this is realized in the *Planetary* system.

1 Introduction

Semantic publication can range from merely equipping published documents with RDFa annotations, expressing metadata or inter-paper links, to frameworks that support the provisioning of user-adapted documents from content representations and instrumenting them with interactions based on the semantic information embedded in the content forms. We want to propose an entry to the latter category in this paper. Our 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 (see Figure 5). Finally a document player application can embed these services to make documents executable. 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. In this paper we present the ADP with a focus on the *Planetary* system as the document player (see Figure 1)

The *Planetary* system (see [Koh+11; Dav+10; Plab] for an introduction) 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. The system achieves

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

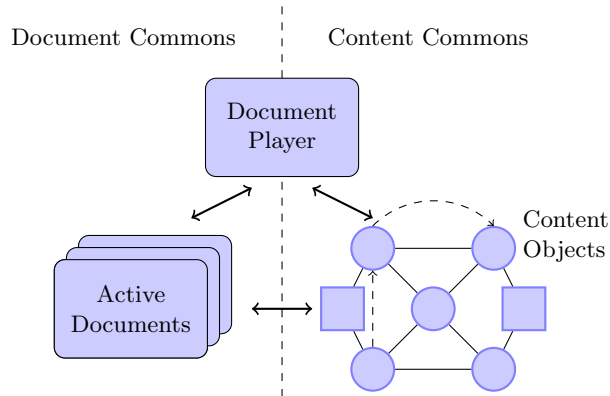


Fig. 1. The Active Documents Paradigm

this by providing embedded user assistance through an extended set of user interactions with documents based on an extensible set of client- and server side services that draw on explicit (and thus machine-understandable) representations in the content commons.

However, the flexibility and power designed into the active documents paradigm comes at a (distribution) cost: Every page that is shown to the user has to be assembled for the user in a non-trivial compilation process (which we call the **presentation** process) that takes user preferences and context into account. On the other hand, if the content is organized modularly, it can be reused across contexts. This presents a completely new set of trade-offs for publishing. One of them is that an investment in modular and semantic representational markup enhances reusability and thus may even lower the overall cost of authoring. We will explore another such trade-off in this paper: optimizing the distribution costs for modular content by “separate compilation”.

In the next section we will look at the organization of the content presented to the user. This will constitute the conceptual backdrop against which we can discuss the issues involved in separate compilation and how we have solved them in the Planetary system.

2 Organization of Content/Narrative Structure

The Planetary system is intended as a *semantic publishing framework*, i.e. as a system providing the baseline capabilities needed for multiple specialized instantiations. We have shown the initial feasibility of the concept in a variety of publicly available case studies² ranging from pre-semantic archives of scientific literature [Arx], over a community-driven mathematical encyclopedia [Plac]

² Note that all of these are research systems under constant development, so your mileage may vary.

and the course system *PantaRhei* [Koh+], to a community portal of formal logics [Plaa]. As a consequence of this, we employ the general, modular knowledge structure depicted in Figure 2.

2.1 Levels of Content/Documents

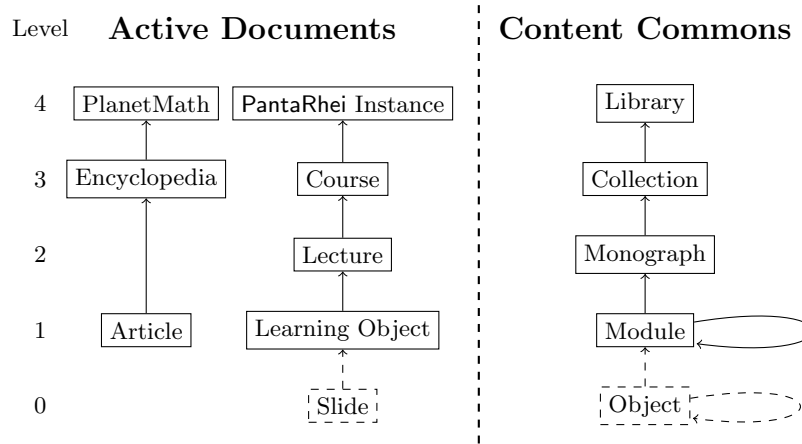


Fig. 2. Levels of Organisation of Content

The lowest level consists of atomic “modules”³, i.e. content objects that correspond to small (active) documents dedicated to a single topic. For a course management system these might be learning objects (either as single modules or module trees), for an encyclopedia these would be the individual articles introducing a topic. Note that technically, we allow modules to contain (denoted by the arrows) other modules, so that larger discourse structures could be formed. For example, sections can be realized as modules referencing other modules of subsections, etc. The next level up is the level of “monographs”, written works on a single subject that have a complete, self-contained narrative structure, usually by a single author or group of authors who feel responsible for the whole monograph. As a content object, a monograph is usually built up from modules, e.g. as a “module tree” that corresponds to sectioning structure of traditional books, but often also includes front and backmatter such as a preface, acknowledgements (both special kinds of modules), table of contents, lists of tables and figures, an index and references (generated from content annotations). Figure 3 shows course notes in the *PantaRhei* system, while other documents at the mono-

³ The level of objects below modules consists of individual statements (e.g. definitions, model assumptions, theorems, and proofs), semantic phrase-level markup, and formulae. Even though it carries much of the semantic relations, it does not play a great role for the document-level phenomena we want to discuss here in this paper.

graph level are articles in a journal, or books in a certain topical section of a library.

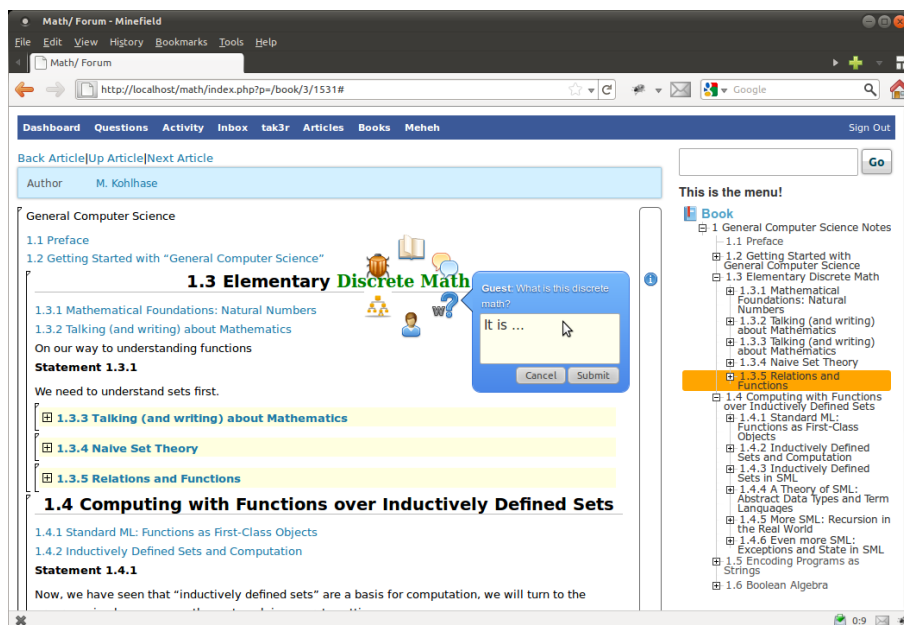


Fig. 3. A Monograph (Course Notes) in the Planetary system

Multiple monographs can be combined into collections, adding special modules for editorial comments, etc. Concrete collections in the document realm are encyclopedias, academic journals, conference proceedings, or courses in a course management system. Finally, the library level collects and grants access to collections, concrete, modern-day examples are digital libraries, installed course management systems, etc. In practice, a library provides a base URI that establishes the web existence of the particular installation. In the Semantic Web world, the library is the authority that makes its resources addressable by URLs.

2.2 Content Objects and their Presentations in Active Documents

To understand the differences between content objects and the documents generated from them in the presentation process, let us consider the example in Figure 4. Even though internally the content objects in Planetary are represented in OMDoc [Koh06], we will use the surface language $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ ⁴ for the example, since this is what the author will write and maintain. $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ is a variant of $\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$

⁴ We speak of an OMDoc surface language for any language that is optimized for human authoring, but that can be converted to OMDoc automatically.

that allows to add semantic annotations in the source. It can be transformed into OMDoc via the L^AT_EXML daemon [GSK11] for management in Planetary; see [Koh08] for details. We are using an example from a mathematical document⁵ since content/presentation matters are most conspicuous there. In our experience, S^TE_X achieves a good balance (at least for authors experienced with L^AT_EX) conciseness and readability for mathematical documents. In particular, since S^TE_X documents such as the one in Figure 4 can be transformed to PDF via the classical pdf_latex for prototyping and proofreading. The semantic editing process can further be simplified by *semantic document development environments* like S^TE_XIDE [JK10], which provides edit-support services like semantic syntax highlighting, command completion/retrieval, and module graph management.

<pre> \begin{module}[id=binary-trees] \importmodule[\KWARCslides{graphs-trees/en/trees}]{trees} \importmodule[\KWARCslides{graphs-trees/en/graph-depth}]{graph-depth} ... \begin{definition}[id=binary-tree.def,title=Binary Tree] A \definiendum[binary-tree]{binary tree} is a \termref[cd=trees,name=tree]{tree} where all \termref[cd=graphs-intro,name=node]{nodes} have \termref[cd=graphs-intro,name=out-degree]{out-degree} 2 or 0. \end{definition} ... \begin{definition}[id=bbt.def] A \termref[name=binary-tree]{binary tree} G is called \definiendumalt[bbt]{balanced} iff the \termref[cd=graph-depth,name=vertex-depth]{depth} of all \termref[cd=trees,name=leaf]{leaves} differs by at most by 1, and \definiendum[fullbbt]{fully balanced}, iff the \termref[cd=graph-depth,name=vertex-depth]{depth} difference is 0. \end{definition} ... \end{module} </pre>
<p>Definition 3.1.7: (<i>Binary Tree</i>) A binary tree is a tree where all nodes have out-degree 2 or 0.</p> <p>Definition 3.1.8: A binary tree G is called balanced iff the depth of all leaves differs by at most by 1, and fully balanced, iff the depth difference is 0.</p>

Fig. 4. Content and Presentation of an Object in S^TE_X

The upper half of Figure 4 shows the content representation of a module on binary trees, and its presentation in Planetary is in the lower box. The first aspect that meets the eye is that the presentation process⁶ adds the textual marker “**Definition 3.1.7**” which is not present in the content representation `\begin{definition}[id=binary-tree.def,title=Binary Tree]`. Note that there are (at least) four issues at hand here pertaining to the presentation of the text marker:

⁵ Actually from a second-semester course on Computer Science [Koh] hosted in **PantaRhei**— an instance of the Planetary system that is optimized for active course notes and discussions.

⁶ We disregard the presentation of formulae in content representation like OpenMath or content MathML into presentation MathML in this paper and refer the reader to [KMR08] for details.

1. The marker “**Definition**” is context-sensitive: The presentation of a Spanish text would have generated “**Definición**”.
2. The number “**3.1.7**” is content-sensitive in a totally different way: it is determined by the document structure, here it is a consequence of being the seventh definition in the first section in chapter 3.
3. The “house style” of a journal might use a different font family for the whole textual marker, for the text of the definition, or add an end marker for a distinctive layout. For instance in mathematical publications, theorems are usually set in italics and proofs use a box on the right of the last line as an end marker.
4. Finally, the whole text marker may be left out altogether in some situations, where a less formal presentation is called for.

Note that all these considerations have to be taken into account when referencing objects like these definitions. More so, these dimensions combine into a unique multi-dimensional point, which identifies the exact presentation of a document fragment. A content reference `\sref{binary-tree.def}` might be presented as “**Def. 3.1.7**”, in the same context as above (again subject to language, house style, etc). Note that here the style (e.g. the keyword) and generated contextual locators (e.g. the number) of the referenced object determines the actual label of the reference⁷. We follow the context dimensions specified in [KK08, Chapter 3], but note that many of the phenomena involve a separate, publishing context dimension (e.g. “house style”).

Another phenomenon related to referencing is induced by the term reference `\termref[cd=graphs,name=vertex]{node}`, which identifies the phrase “node” as a technical term and links it to its defining occurrence by the symbol name (here `vertex`) and the module name (also called content dictionary; here `graphs`). The specified module must be accessible in the current module via the `\importmodule` relation and must contain a definition that contains a definiendum with symbol name `vertex`. The content module in Figure 4 specifies a module/content dictionary with name `balanced-binary-trees`, whose first definition supplies a definiendum with name `balanced-tree` via the `\twindf` macro, which is referenced in the second definition. Note that in the presentation process where term references are displayed e.g. as hyperlinks to the definition the name-based semantic links have to be converted into regular URI references. For this presentational conversion to hyperlinks one utilizes not only the module tree structure (i.e. visibility relationship) but also the library context that provides the base of the URI.

Finally, note that some content objects contribute to the context of other objects higher up in the content hierarchy in Figure 2. A good example for this are the definienda discussed above. In \TeX , they trigger index entries that populate the backmatter of monographs that include the respective module. Section titles populate the frontmatter in a similar way. Concretely, we have a top-level index stub in the backmatter, which “builds” itself from the context. In a sense, the index is an abstract concept with volatile presentation, generated

⁷ a rather peculiar notion of context when viewed from a content-only perspective

from the module tree with the help of the content commons, which answers what objects should be indexed.

3 Separate Compilation

We have seen above that the various contexts (conceptual/document/language) have a significant effect on the presentation. But observe that if all the context-dependent parts of the presentation can be generated (albeit laboriously), the content representations are context-independent and can be reused in different contexts. This makes the content representations very portable. Consider for instance the definitions in our example above. They have been reused not only in eight instances of the “General Computer Science II” course [Koh] in the years 2004-2011 (each time with different numbers due to additions or deletions of preceding material), but also in different courses, e.g. as a recap in a more advanced CS course (without definition marker). But these are not the only contexts: the *Planetary* system can generate “guided tours” (self-contained explanations adapted to the user’s prerequisite knowledge) for any concept in a document. Clearly, we cannot reasonably pre-compute all necessary presentation variants.

Computationally, the described situation is analogous to (and in fact conceptually influenced by) the situation in software design, where large programs are broken up into reusable source modules. As source modules are re-used in many programs, it is important that compilers support a regime of “separate compilation and linking” to make software development tractable: if one of many software modules used in a program changes, only that one module has to be re-compiled and the whole program re-linked. The first factor that enables this is the observation that for compilation of a module only the (relatively stable) signatures⁸ of modules it depends on are needed, not the (relatively change-prone) module implementations. The second factor is that source modules can be compiled into a form, where references to functions imported from other modules are left symbolic and can later be replaced by concrete static references by the linker. We will call such forms of modules **contextable**, since they are contextualized by the linker in the way described.

In the *Planetary* architecture semantic publishing consists of the transformation of content structures encoded in \LaTeX to active documents encoded in XHTML+MathML+RDFa (see Section 3.2 for details). To foster reuse, and make the process tractable, we want to assemble active documents from reusable content modules much in the same way as assembling an executable program from source modules. To make the separate compilation analogy fertile for semantic publishing it is useful to look at the role of context in the separate compilation regime: source modules are compiled into a context-independent form, which is then contextualized by linking compiled modules together into a consistent configuration for a concrete program. In the next two sections we

⁸ Signatures contain the names of functions/procedures, possibly their types, but not their implementations.

examine how the two factors identified as crucial for the separate compilation regime can be obtained in the context of semantic publishing.

3.1 Contextable Presentations

Just as in programming, *separate* compilation of content modules into active documents is impossible without contextable structures in the presentation. It is an original contribution of our work to introduce them in the document setting. Concretely, we make use of the XML styling architecture and computes context-independent presentations that can be contextualized later. For instance, the XHTML header for the first definition in Figure 4 has the following form.

```
<div id="binary-tree.def" class="omdoc-definition" >
  <span class="omdoc-statement-header" >
    <span class="omdoc-definition-number" />7</span>
    <span class="omdoc-statement-title" >Binary Tree</span>
  </span>
  ...
```

We can then add (house) style information via CSS:

```
span.omdoc-statement-header {font-weight:bold}
span.omdoc-statement-title:before {content:"("}
span.omdoc-statement-title:after {content:")"}
span.omdoc-definition-number:after {content:" :"}
span.omdoc-definition-number:before {content:" Definición "}
```

Note that the keywords are not represented explicitly in the XHTML presentation, but added by content declarations in the CSS. This allows to overwrite the default ones via cascaded language-specific CSS bindings, e.g. using

```
span.omdoc-definition-number:before {content:" Definición "}
```

Note furthermore, that the presentation process only adds preliminary statement numbers in the XHTML presentation (here the number 7, since the definition is the seventh statement in the module). In the **Planetary** system, these numbers are dynamically overwritten by values computed from the context; in our example “3.1.7”. The case for references is similar; for the table of contents shown in Figure 3 the presentation generates

```
<div class="omdoc-expandableref" >
  <span class="omdoc-ref-number" >4</span>
  <span class="omdoc-reftitle" >
    <a href="../computing-dmath.omdoc" class="expandable" >
      Computing with Functions over Inductively Defined Sets
    </a>
  </span>
</div>
```

in the table of contents on the right and in the text. The CSS class `omdoc-expandableref` triggers the **Planetary** interaction that expands the references in place to get the expanding ToC and the main document that can be folded/unfolded via the Mathematica-style folding bars on the extreme left.

3.2 Supporting the Logistics of Separate Compilation: Dynamic Linking

The role of the module signatures (think C header files) is taken by \LaTeX module signatures, i.e. auxiliary files generated from \LaTeX content modules that excerpt the information about references, modules and their dependencies; see [Koh08] for details. This information is used to establish a mapping between the content commons and the document commons (see Figure 1) that can be queried for the semantic interaction services embedded into the active documents.

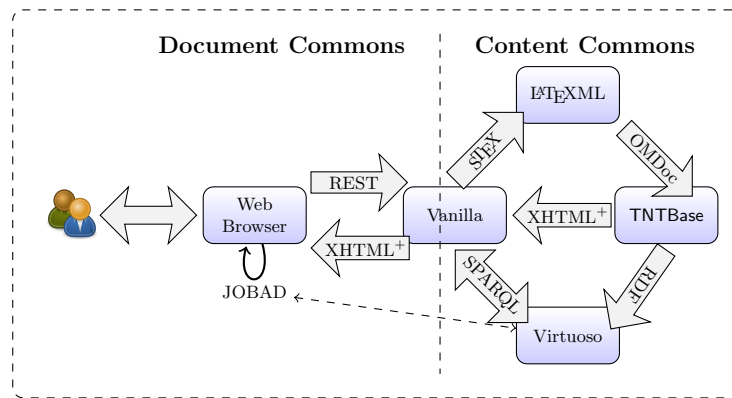


Fig. 5. The Planetary System Ecosystem

Actually, to understand the compilation and linking phases in *Planetary*, consider the system architecture in Figure 5. There, the document commons and content commons are layouted in a generic web browser and encapsulated versioned XML knowledge store *TNTBase* [ZK09] respectively. The *Planetary* system acts as an intermediary between these two:

1. \LaTeX is converted to OMDoc via the \LaTeX XML daemon [GSK11], this is then converted to XHTML+MathML+RDFa. Both transformations are highly dependent on notation information in the content commons, so they are under the control of the *TNTBase* system, which stores the content commons.
2. *Planetary* caches the contextable presentations that are generated by the *TNTBase* system. New presentations are requested from *TNTBase* whenever *a*) the content module in *TNTBase* has changed, and *b*) a user requests a to view the module.
3. *Planetary* hosts a triple store (*Virtuoso*) of structural metadata from the content commons that can be used for semantic services and document-level features, such as different views based on various selection criteria for an encyclopedia.
4. Finally, *Planetary* hosts structural information about the knowledge items at the different levels in Figure 2, used by the linker. In the examples in Fig-

ure 3 and Section 2.2, the numbering is linked into the contextable modules whenever a page is viewed, based on this information. Recall we need this *dynamic (i.e. view-time) linking* as modules are re-used in different document contexts.

4 Conclusion

In this paper we have explored the conceptual and practical decoupling and interaction of content and presentation in the active documents paradigm of semantic publishing. Our main focus rested on the interaction of content object reuse and context sensitivity of the presentation process. To make semantic publishing of highly structured content representations into active documents tractable we have developed a “separate compilation and dynamic linking” regime for transforming highly structured content representations into active documents. The concrete realization in the Planetary system hinges on the development of contextable pre-presentations that are contextualized at document load time.

While the basic architecture has been realized in the Planetary system, there is still a lot to explore in the active documents paradigm and its SCDL implementation. One crucial aspect is that while SCDL makes building active documents tractable, it also leads to the well-known “late binding problems” (aka “DLL Hell”), if modules change without adaptation of the dependent ones. We are currently working on an integration of an ontology-based management of change process [AM10] into the Planetary system (see [Aut+11]). This tries to alleviate late binding problems by analyzing the impacts of a change via the dependency relation induced by the semantic structure of the content commons and supports authors in adapting their work. To complement this, we are currently developing a notion of “versioned references” that support the practice of creating and cultivating “islands of consistency” in the presence of change (see [KK11]). We hope that together, these measures can lead to semantic content management workflows that alleviate the side-effects of the semantic publishing workflow described in this paper.

References

- [AM10] Serge Autexier and Normen Müller. “Semantics-based Change Impact Analysis for Heterogeneous Collections of Documents”. In: *Proceedings of the 10th ACM symposium on Document engineering*. Ed. by Michael Gormish and Rolf Ingold. DocEng ’10. Manchester, United Kingdom: ACM, 2010, pp. 97–106. ISBN: 978-1-4503-0231-9. DOI: <http://doi.acm.org/10.1145/1860559.1860580>. URL: <http://doi.acm.org/10.1145/1860559.1860580>.
- [Arx] *arXMLiv Build System*. URL: <http://arxivdemo.mathweb.org> (visited on 09/27/2010).

- [Aut+11] Serge Autexier et al. “Workflows for the Management of Change in Science, Technologies, Engineering and Mathematics”. submitted. 2011.
- [Cic] *Intelligent Computer Mathematics*. submitted. 2011.
- [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>.
- [GSK11] Deyan Ginev, Heinrich Stamerjohanns, and Michael Kohlhase. “The L^AT_EX XML Daemon: A L^AT_EX Entrance to the Semantic Web”. submitted. 2011. URL: <https://kwarc.eecs.iu-bremen.de/repos/arXMLiv/doc/cicm-systems11/paper.pdf>.
- [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. ISBN: 3642141277. arXiv:1005.5489v1 [cs.OH].
- [KK08] Andrea Kohlhase and Michael Kohlhase. “Semantic Knowledge Management for Education”. In: *Proceedings of the IEEE; Special Issue on Educational Technology* 96.6 (June 2008), pp. 970–989. URL: <http://kwarc.info/kohlhase/papers/semkm4ed.pdf>.
- [KK11] Andrea Kohlhase and Michael Kohlhase. “Maintaining Islands of Consistency via Versioned Links”. submitted. 2011. URL: <http://kwarc.info/kohlhase/submit/mkm11-verlinks.pdf>.
- [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>.
- [Koh] *General Computer Science: GenCS I/II Lecture Notes*. 2011. URL: <http://gencs.kwarc.info/book/1> (visited on 03/03/2001).
- [Koh+] Michael Kohlhase et al. *Planet GenCS*. URL: <http://gencs.kwarc.info> (visited on 09/22/2010).
- [Koh06] Michael Kohlhase. *OMDOC – An open markup format for mathematical documents [Version 1.2]*. LNAI 4180. Springer Verlag, Aug. 2006. URL: <http://omdoc.org/pubs/omdoc1.2.pdf>.
- [Koh08] Michael Kohlhase. “Using L^AT_EX 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>.
- [Plaa] *Logic Atlas and Integrator*. URL: <http://logicatlas.omdoc.org> (visited on 09/22/2010).
- [Plab] *Planetary Developer Forum*. URL: <http://trac.mathweb.org/planetary/> (visited on 01/20/2011).
- [Plac] *PlanetMath.org – Math for the people, by the people*. URL: <http://planetmath.org> (visited on 01/06/2011).
- [ZK09] Vyacheslav Zholudev and Michael Kohlhase. “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/BalisageVol3.Zholudev01.