

The Planetary System: Active Documents and a Web3.0 for Math.

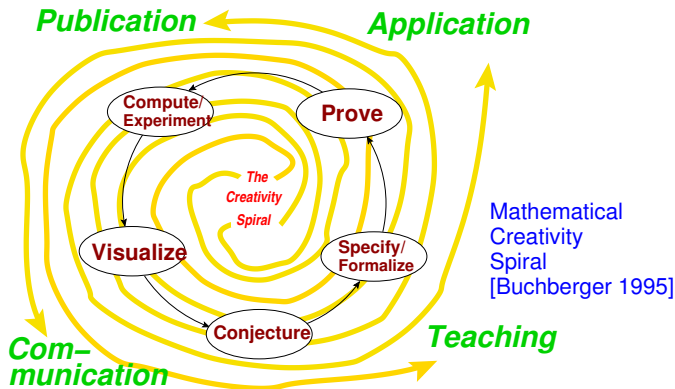
Michael Kohlhasse

<http://kwarc.info/kohlhasse>
Center for Advanced Systems Engineering
Jacobs University Bremen, Germany

May 30. 2012, NIST

Introduction

The way we do math will change dramatically



- Every step will be supported by mathematical software systems
- Towards an infrastructure for web-based mathematics!

eMath 3.0: The Time is Ripe

- Background:
 - Web 2.0 is the term used for the “social Web” (tagging, blogs, wikis, facebook, ...)
 - The “Semantic Web” is a version of the Web, where humans & machines cooperate
 - Web 3.0 is the term used for the “social Semantic Web”.
- We will apply these to eMath (⚠ regular Math may or may not change)

eMath 3.0: The Time is Ripe

- Background:
 - Web 2.0 is the term used for the “social Web” (tagging, blogs, wikis, facebook, ...)
 - The “Semantic Web” is a version of the Web, where humans & machines cooperate
 - Web 3.0 is the term used for the “social Semantic Web”.
- We will apply these to eMath (⚠ regular Math may or may not change)
- Recent Improvements of the Math-on-the Web Environment:
 - MathML3 is out
 - MathML enabled in WebKit (in Safari 5.1, on the road for Chrome, Konqueror)
 - MathML is in HTML5 (without Namespaces though)
 - MathJax has reached 1.0 (Display MathML by JavaScript in many Browsers)





eMath 3.0: The Time is Ripe

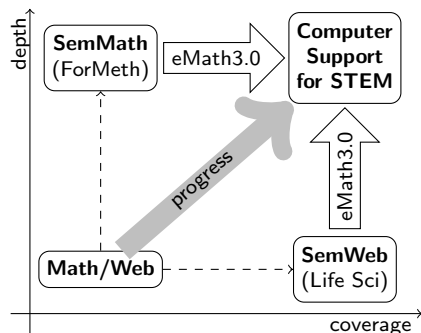
- **Background:**
 - Web 2.0 is the term used for the “social Web” (tagging, blogs, wikis, facebook, ...)
 - The “Semantic Web” is a version of the Web, where humans & machines cooperate
 - Web 3.0 is the term used for the “social Semantic Web”.
- We will apply these to eMath (⚠ regular Math may or may not change)
- **Recent Improvements of the Math-on-the Web Environment:**
 - MathML3 is out
 - MathML enabled in WebKit (in Safari 5.1, on the road for Chrome, Konqueror)
 - MathML is in HTML5 (without Namespaces though)
 - MathJax has reached 1.0 (Display MathML by JavaScript in many Browsers)
- **Recent Improvements of the Semantic Web:**
 - RDF can be embedded into XML via RDFa (linked data export)
 - RDF querying via SPARQL (modulo OWL Ontologies) (semantic search)
 - OMDoc as a mathematical Ontology format (modularity, documentation, full Math)

eMath 3.0: The Time is Ripe

- **Background:**
 - Web 2.0 is the term used for the “social Web” (tagging, blogs, wikis, facebook, ...)
 - The “Semantic Web” is a version of the Web, where humans & machines cooperate
 - Web 3.0 is the term used for the “social Semantic Web”.
- We will apply these to eMath (⚠ regular Math may or may not change)
- **Recent Improvements of the Math-on-the Web Environment:**
 - MathML3 is out
 - MathML enabled in WebKit (in Safari 5.1, on the road for Chrome, Konqueror)
 - MathML is in HTML5 (without Namespaces though)
 - MathJax has reached 1.0 (Display MathML by JavaScript in many Browsers)
- **Recent Improvements of the Semantic Web:**
 - RDF can be embedded into XML via RDFa (linked data export)
 - RDF querying via SPARQL (modulo OWL Ontologies) (semantic search)
 - OMDoc as a mathematical Ontology format (modularity, documentation, full Math)
- **Overview over the talk:**
 - MathML3 brings more semantics (strict content Math, elementary Math)
 - integrating MathML/L^AT_EX into the Web 2.0
 - A L^AT_EX-based Semantic Web for Mathematics

Contributions from KWARC@Jacobs@Bremen

-  STEM Knowledge: more like a Digital Library than the Open WWW 
(reviewed publication \rightsquigarrow less junk, little duplication, partly inaccessible)
- Combination of SemMath and SemWeb
- Expertise in Semantics of STEM Docs
- Expressive Analysis Target Format (OMDoc)
- Software Stack for Semantic Processing
- eSTEM3.0 System Planetary (Active Docs)
- Invasive authoring (Office/L^AT_EX)
- Semantic Analysis for L^AT_EX-based Corpora (arXiv, ZBL, PlanetMath...)
-  We use Math as a test tube for STEM (Science, Tech, Eng, & Math) 

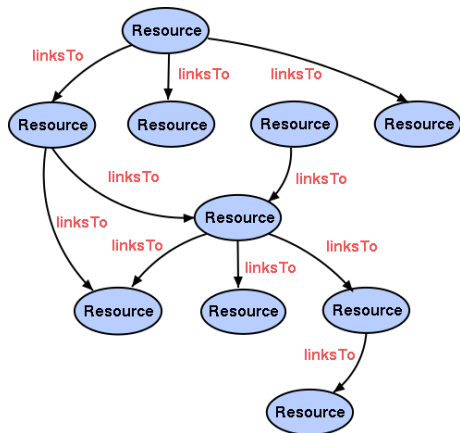


Foundations

The Semantic Web

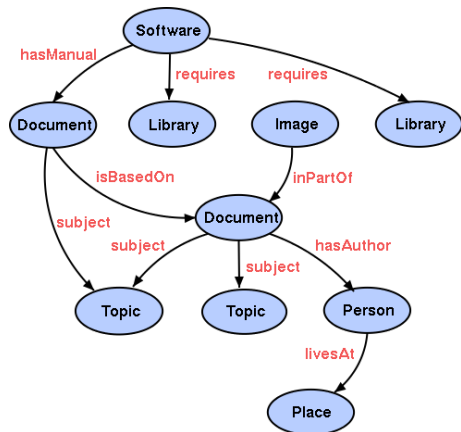
The Current Web

- **Resources:** identified by URI's, untyped
- **Links:** href, src, ... limited, non-descriptive
- **User:** Exciting world - semantics of the resource, however, gleaned from content
- **Machine:** Very little information available - significance of the links only evident from the context around the anchor.



The Semantic Web

- **Resources:** Globally Identified by URI's or Locally scoped (Blank), Extensible, Relational
- **Links:** Identified by URI's, Extensible, Relational
- **User:** Even more exciting world, richer user experience
- **Machine:** More processable information is available (Data Web)
- **Computers and people:** Work, learn and exchange knowledge effectively



What is the Information a User sees?

WWW2002

The eleventh international world wide web conference

Sheraton waikiki hotel

Honolulu, hawaii, USA

7-11 may 2002

1 location 5 days learn interact

Registered participants coming from

australia, canada, chile denmark, france, germany, ghana, hong kong, india, ireland, italy, japan, malta, new zealand, the netherlands, norway, singapore, switzerland, the united kingdom, the united states, vietnam, zaire

On the 7th May Honolulu will provide the backdrop of the eleventh international world wide web conference. This prestigious event ?

Speakers confirmed

Tim Berners-Lee: Tim is the well known inventor of the Web, ?

Ian Foster: Ian is the pioneer of the Grid, the next generation internet ?

What the machine sees

wwwE''E

T[1]↓[1]⊆[1]∖u(∠)∖u[∇∖+u)∖-↓⊇∇↓[⊇]∩⊇[1]∪∖{[∇[∖]∪]

S([∇+u∖∖⊇+))|||)(∖u[↓

H∖∖↓∩↓∩⊆(∖⊇+))⊆USA

ℵ_∞∞↓+†E''E

R[}]∫∪[∇[√+∇u)∫)√+∖u f[∖)∖}{∇∖

+∩∪∇+↓)+⊆]∖+[-+⊆](∠)↓[∩∖↓+∇||⊆{∇+∪[⊆}]∇↓+∖+⊆}(∖+⊆(∖)}||∖)}⊆)∖[∩)+⊆

∇[↓+∖[⊆]u+↓+⊆+√+∖+⊆+↓u+⊆+∖[⊇‡]+↓+∖[⊆u(∩∖[∪(∇↓+∖[f⊆+∖∇⊇+†⊆

∫∖}+√∇[⊆∫⊇]u‡[∇↓+∖[⊆u(∩∖[∪[∩||∖)}[↓⊆u(∩∖[∪[∪+u] f⊆⊆]∪∖+↓⊆+†+∇[

O∖u(∩[u(M+†H∖∖↓∩↓∩⊇)↓√∇[⊆]∩u(∩[-+||[∇[√{u(∩[↓]⊆[1]∖u(

∖∖u[∇∖+u)∖-↓⊇∇↓[⊇]∩⊇[1]∪∖{[∇[∖]∪]↯T) f√∇[∪)}∩∫[1]∖u⊥

S√]-+||∇ f[∖{)∇↓[

T↓[∩∇∖∇ f↓[1]-T)∫∪(∩[↓‡||∖∖⊇)∖⊆[∖u∇∖{u(∩W[[-⊆⊥

I+∖ f∪[∇-I+∖)∪(√)∖[∇[∖{u(∩ g∇)[⊆u(∩[∖]§u}[∖∖∇+u)∖∖∖u[∇∖]u⊥

Solution: XML markup with “meaningful” Tags

[illegible]

What the machine sees of the XML

[illegible]

Need to add “Semantics”

- External agreement on meaning of annotations E.g., Dublin Core
 - Agree on the meaning of a set of annotation tags
 - Problems with this approach: Inflexible, Limited number of things can be expressed

Need to add “Semantics”

- External agreement on meaning of annotations E.g., Dublin Core
 - Agree on the meaning of a set of annotation tags
 - Problems with this approach: Inflexible, Limited number of things can be expressed
- Use Ontologies to specify meaning of annotations
 - Ontologies provide a vocabulary of terms
 - New terms can be formed by combining existing ones
 - Meaning (semantics) of such terms is formally specified
 - Can also specify relationships between terms in multiple ontologies

Need to add “Semantics”

- External agreement on meaning of annotations E.g., Dublin Core
 - Agree on the meaning of a set of annotation tags
 - Problems with this approach: Inflexible, Limited number of things can be expressed
- Use Ontologies to specify meaning of annotations
 - Ontologies provide a vocabulary of terms
 - New terms can be formed by combining existing ones
 - Meaning (semantics) of such terms is formally specified
 - Can also specify relationships between terms in multiple ontologies
- Inference with annotations and ontologies (get out more than you put in!)
 - Standardize annotations in RDF [KC04] or RDFa [BAHS] and ontologies on OWL [w3c09]
 - Harvest RDF and RDFa in to a triplestore or OWL
 -

MathML: Presentation and Content of Mathematical Formulae

Representation of Formulae as Expression Trees

- Mathematical Expressions are build up as expression trees
 - of layout schemata in Presentation-MathML
 - of functional subexpressions in Content-MathML
- Example: $\frac{3}{x+2}$

```
<mfrac>
  <mn>3</mn>
  <mfenced>
    <mi>x</mi>
    <mo>+</mo>
    <mn>2</mn>
  </mfenced>
</mfrac>
```

```
<apply>
  <divide/>
  <cn>3</cn>
  <apply>
    <plus/>
    <ci>x</ci>
    <cn>2</cn>
  </apply>
</apply>
```

Layout Schemata and the MathML Box model

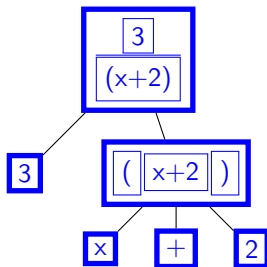


Diagram illustrating the MathML Box model for the expression $3 \cdot (x + 2)^3$ using MathML tags. The expression is represented as a tree of tags. The root tag is `<mfrac>...</mfrac>`. The numerator is `<mn>3</mn>`. The denominator is `<mfenced>...</mfenced>`. The denominator is further divided into three tags: `<mi>x</mi>`, `<mo>+</mo>`, and `<mn>2</mn>`.

Content Mathml: Expression Trees in Prefix Notation

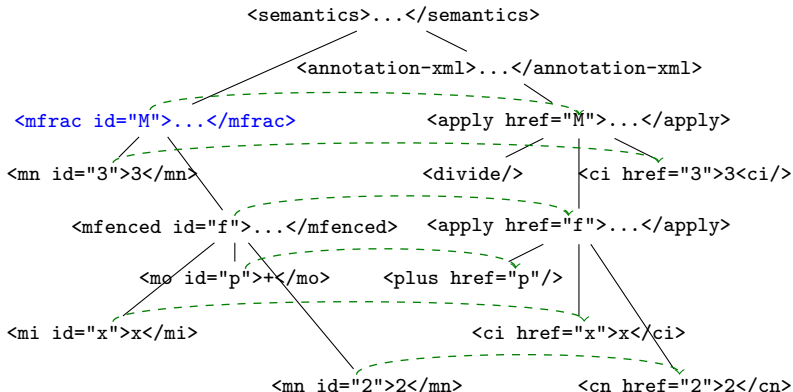
- Prefix Notation saves parentheses (so does postfix, BTW)

$(x - y)/2$	$x - (y/2)$
<pre><apply> <divide/> <apply> <minus/> <ci>x</ci> <ci>y</ci> </apply> <cn>2</cn> </apply></pre>	<pre><apply> <minus/> <ci>x</ci> <apply> <divide/> <ci>y</ci> <cn>2</cn> </apply> </apply></pre>

- Function Application:** `<apply>function arg1 ... argn </apply>`
- Operators and Functions:** ~ 100 empty elements `<sin/>`, `<plus/>`, `<eq/>`, `<compose/>`,...
- Token elements:** `ci`, `cn` (identifiers and numbers)
- Extra Operators:** `<csymbol definitionURL="...">...</csymbol>`

Parallel Markup e.g. in MathML

- Combine the presentation and content markup in one tree and cross-reference



- use e.g. for semantic copy and paste.
(click on presentation, follow link and copy content)

Mixing Presentation and Content MathML

```
<semantics>
  <mrow>
    <mrow><mo>(</mo><mi>a</mi> <mo>+</mo> <mi>b</mi><mo>)</mo></mrow>
    <mo>&InvisibleTimes;</mo>
    <mrow><mo>(</mo><mi>c</mi> <mo>+</mo> <mi>d</mi><mo>)</mo></mrow>
  </mrow>
  <annotation-xml encoding="MathML-Content">
    <apply><times/>
      <apply><plus/><ci>a</ci> <ci>b</ci></apply>
      <apply><plus/><ci>c</ci> <ci>d</ci></apply>
    </apply>
  </annotation-xml>
  <annotation-xml encoding="openmath">
    <OMA><OMS cd="arith1" name="times"/>
      <OMA><OMS cd="arith1" name="plus"/><OMV name="a"/><OMV name="b"/></OMA>
      <OMA><OMS cd="arith1" name="plus"/><OMV name="c"/><OMV name="d"/></OMA>
    </OMA>
  </annotation-xml>
</semantics>
```

Converting the arXiv

The arXMLiv Project: arXiv to semantic XML

- **Idea:** Develop a large corpus of knowledge in OMDoc/PhysML
 - to get around the chicken-and-egg problem of MKM
 - corpus-linguistic methods for semantics recovery (linguists interested)
- **Definition 1 (The Cornell Preprint arXiv)** (<http://www.arxiv.org>)
Open access to ca. 700K e-prints in Physics, Mathematics, Computer Science and Quantitative Biology.
- **Definition 2 (The arXMLiv Project)** (<http://arxmliv.kwarc.info>)
 - use Bruce Miller's \LaTeX XML to transform to XHTML+MathML
 - extend to \LaTeX XML daemon (RESTful web service) (<http://latexml.mathweb.org>)
 - we have an automated, distributed build system (ca. 2 CPU-years)
 - create ca. 12K \LaTeX XML binding files (8 Jacobs students help)
 - use MathWebSearch to index XML version (realistic search corpus)
- More semantic information will enable more added-value services, e.g.
 - filter hits by model assumptions (expanding, stationary, or contracting universe)
 - use linguistic techniques to add the necessary semantics

Why reimplement the T_EX parser? I

- **Problem:** The T_EX parser can change the tokenizer while at runtime (`\catcode`)
- **Example 3 (Obfuscated T_EX)** David Carlisle posted the following, when someone claimed that word counting is simple in T_EX/L^AT_EX

```
\let~\catcode~'76~'A13~'F1~'j00~'P2jdefA71F~'7113jdefPALLF
PA''FwPA;;FPAZZFLaLPA//71F71iPAHHFLPAzzFenPASSFthP;A$$$FevP
A@@FfPARR717273F737271P;ADDFRgniPAWW71FPATTFvePA**FstRsamP
AGGFRruoPAqq71.72.F717271PAY7172F727171PA??Fi*LmPA&&71jfi
Fjfi71PAVVfjbigskipRPWGAUU71727374 75,76Fjpar71727375Djifx
:76jelset&U76jfiPLAKK7172F7117271PAXX71FVLn0SeL71SLRyadR@oL
RrhC?yLRurtKFeLPFovPgaTLtReRomL;PABB71 72,73:Fjif.73.jelset
B73:jfiXF71PU71 72,73:PWs;AMM71F71diPAJJFRdriPAQQFRsreLPAI
I71Fo71dPA!!FRgiePBt'el@ lTLqdrYmu.Q.,Ke;vz vzLqip.Q.,tz;
;Lql.IrsZ.eap,qn.i. i.eLlMaesLdRcna,;!;h htLqm.MRasZ.ilnk,%
s$;z zLqs'.ansZ.Ymi,/sx ;LYegseZRyal,@i;@ TLRlogdLrDsW,@;G
LcYlaDLbJsW,SWXJW ree @rzchLhzsW;;WERcesInW qt.'oL.Rtrul;e
doTsW,Wk;Rri@stW aHAHHFndZPpqar.tridgeLinZpe.LtYer.W,:jbye
```

When formatted by TeX, this leads to the full lyrics of “The twelve days of christmas”. When formatted by L^AT_EXML, it gives

Why reimplement the T_EX parser? II

```
<song>
  <verse>
    <line>On the first day of Christmas my true love gave to me</line>
    <line>a partridge in a pear tree.</line>
  </verse>
  <verse>
    <line>On the second day of Christmas my true love gave to me</line>
    <line>two turtle doves</line>
    <line>and a partridge in a pear tree.</line>
  </verse>
  <verse>
    <line>On the third day of Christmas my true love gave to me</line>
    <line>three french hens</line>
    <line>two turtle doves</line>
    <line>and a partridge in a pear tree.</line>
  </verse>
  <verse>
    <line>On the fourth day of Christmas my true love gave to me</line>
    <line>four calling birds</line>
    <line>three french hens</line>
    <line>two turtle doves</line>
    <line>and a partridge in a pear tree.</line>
  </verse>
  ...
```

Why reimplement the $\text{T}_{\text{E}}\text{X}$ parser? III

- But the real reason is: that we can take advantage of the semantics in the $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$.
- $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{XML}$ does not need to expand macros, we can tell it about XML equivalents.
- **Example 4 (Recovering the Semantics of Proofs)**

Add the following magic incantation to `amsthm.sty.ltxml` ($\text{L}_{\text{A}}\text{T}_{\text{E}}\text{XML}$ binding)

```
DefEnvironment('{proof}', "<xhtml:div class='proof'>#body</xhtml:div>");
```

The `arXMLiv` approach: Try to cover most packages and classes in the arXiv
(Jacobs undergrads' intro to research)

Future Plans for arXMLiv

- **State:** \LaTeX -to-XHTML+MathML Format Conversion works (65% success)
- **Over the summer:** Bump up success rate to 75%, daily downloads, web site, instrumentation, . . .
- **Soon:** Integrate user-level quality control (integrate JS feedback into html)
- **starting Fall:** Extend post-processing by linguistic methods for semantic analysis
 - build semantics blackboard/database for linguistic information (rdf triples)
 - extend build system for arbitrary XML2BB processes
 - invite the linguists over (they leave semantics results in BB)
 - harvest the semantics BB to get OMDoc representations

Current and Possible Applications

- the arxmliv build system <http://arxmliv.kwarc.info>
- the transformation web service <http://tex2xml.kwarc.info>
- L^AT_EXML daemon to avoid perl and L^AT_EX startup times (Deyan Ginev)
 - keep L^AT_EXML alive as a daemon that can process multiple files/fragments (patch memory leaks)
 - a L^AT_EXML client just passes files/fragments along ($\frac{10}{s}$ to $\frac{100}{s}$)
- embedding/editing L^AT_EX in web pages <http://tex2xml.kwarc.info/test>
- a MathML version of the arXiv allows vision-impaired readers to understand the texts
- generalization search (need to know sentence structure for detecting universal variables)
- semantic search by academic discipline or theory assumption (need discourse structure)
- development of scientific vocabularies (over the past 18 years; drink from the source)

Planetary: An Integrated Platform for eMath3.0

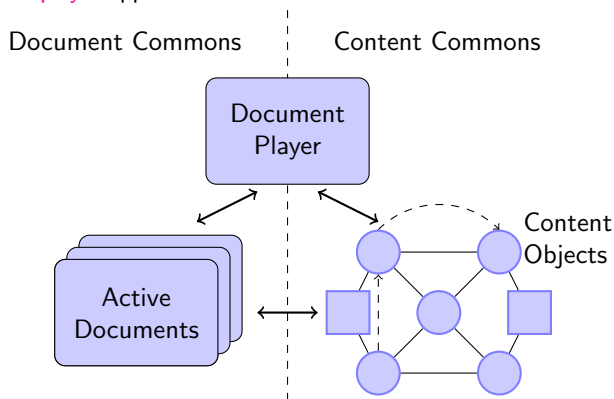
Planetary: A Social Semantic eScience System

The PLANETARY System

- The PLANETARY system is a Web 3.0 system for semantically annotated document collections in Science, Technology, Engineering and Mathematics (STEM).
- **Web 3.0** stands for extension of the Social Web with Semantic Web/Linked Open Data technologies.
- documents published in the PLANETARY system become flexible, adaptive interfaces to a content commons of domain objects, context, and their relations.
- PLANETARY is based on the Active Documents Paradigm (see next)
- **Example 5 (Example installments)**
 - arxivdemo.mathweb.org (presentation/structural Level: arXiv)
 - panta.kwarc.info (semantic level: PantaRhei course system)
 - logicatlas.omdoc.org (fully formal level: Logic Representations)
 - planetbox.kwarc.info (Technology Sandbox)
- The PLANETARY system is finalist in the Elsevier Executable Papers Challenge.

The Active Documents Paradigm

- **Definition 6** The **active documents paradigm (ADP)** consists of
 - *semantically annotated documents* together with
 - background ontologies (which we call the **content commons**),
 - *semantic services* that use this information
 - a **document player** application that embeds services to make documents executable.

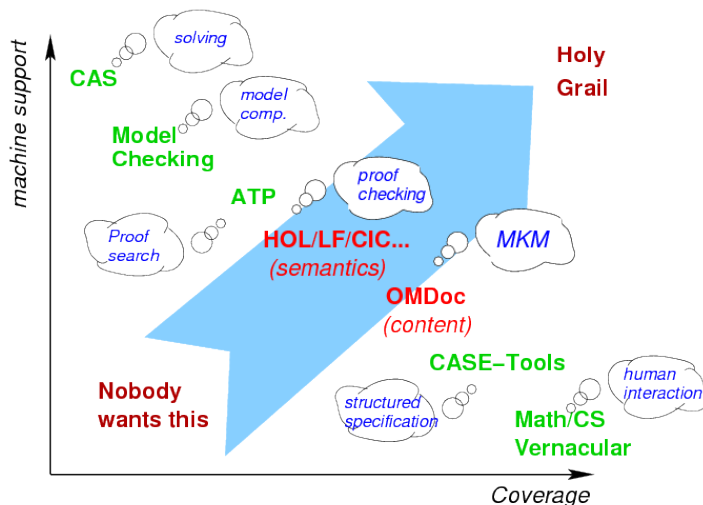


- **Example 7** Services can be program (fragment) execution, computation, visualization, navigation, information aggregation and information retrieval

OMDoc in a Nutshell (three levels of modeling)

<p>Formula level: OpenMath/C-MathML</p> <ul style="list-style-type: none"> • Objects as logical formulae • semantics by ref. to theory level 	<pre> <OMA> <OMS cd="arith1" name="plus"/> <OMS cd="nat" name="zero"/> <OMV name="N"/> </OMA> </pre>
<p>Statement level:</p> <ul style="list-style-type: none"> • Definition, Theorem, Proof, Ex. • semantics explicit forms and refs. 	<pre> <defn for="plus" type="rec"> <CMP>rec. eq. for plus</CMP> <FMP>$X + 0 = X$</FMP> <FMP>$X + s(Y) = s(X + Y)$</FMP> </defn> </pre>
<p>Theory level: Development Graph</p> <ul style="list-style-type: none"> • inheritance via symbol-mapping • theory-inclusion by proof-obligations • local (one-step) vs. global links 	<p>The diagram illustrates the Theory level Development Graph with four main components: Nat-List, List, Nat, and Param. Nat-List and List are at the top, while Nat and Param are at the bottom. Nat-List contains 'cons, nil' and '0, s, Nat, <'. List contains 'cons, nil' and 'Elem, <'. Nat contains '0, s, Nat, <'. Param contains 'Elem, <'. Solid green arrows labeled 'Actualization' point from List to Nat-List and from Param to Nat. Dashed blue arrows labeled 'imports' point from List to Nat-List and from Param to Nat. A yellow box labeled 'Proof Obligations' is connected to Nat and Param by a solid green arrow labeled 'theory-inclusion'. A dashed line also connects the 'Proof Obligations' box to the 'Actualization' arrow between List and Nat-List.</p>

Situating OMDoc: Math Knowledge Management



sT_EX: A Semantic Variant of L^AT_EX

$\text{T}_{\text{E}}\text{X}/\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ as MKM Format: The Notation/Context Problem

- idiosyncratic notations that are introduced, extended, discarded on the fly

$$\lambda X_{\alpha}.X =_{\alpha} \lambda Y_{\alpha}.Y \triangleq \mathbf{I}^{\alpha}$$

meaning of α depends on context: **object type** vs. **mnemonic** vs. **type label**.

- even “standard notations” depend on the context, e.g. binomial coefficients: $\binom{n}{k}$, ${}_nC^k$, C_k^n , and C_n^k all mean the same thing: $\frac{n!}{k!(n-k)!}$ (**cultural context**)
- Notation scoping follows complex rules (**notations must be introduced**)
 - “We will write $\wp(S)$ for the set of subsets of S ” (**for the rest of the doc**)
 - “We use the notation of [BrHa86], with the exception...” (**by reference**)
 - “Let S be a set and $f: S \rightarrow S \dots$ ” (**scope local in definition**)
 - “where w is the...” (**scope local in preceding formula**)
 - Book on group theory in Bourbaki series uses notation [Bou: Algebra]

Observation: **Notation scoping** is different from the one offered by $\text{T}_{\text{E}}\text{X}/\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$

T_EX/L^AT_EX as MKM Format: The Reconstruction Problem

- Mathematical communication relies on the inferential capability of the reader.
- semantically relevant arguments are left out (or ambiguous) to save notational overload
(reader must disambiguate or fill in details.)

$$\log_2(x) \text{ vs. } \log(x) \qquad \mathbb{[A]}^{\mathcal{M}}_{\varphi} \text{ vs. } \mathbb{[A]}$$

- condensed notation: $f(x+1) \pm 2\pi = g(x-1) \mp 2i$ (stands for 2 equations)
- ad hoc extensions: $\#(A \cup B) \leq \#A + \#B$ (exceptions for ∞)
- overt ambiguity: $\sin x/y$ vs. $\frac{\sin x}{y}$ vs. $\sin \frac{x}{y}$ vs. $-1 \leq \sin x/\pi \leq 1$
- size of the gaps varies with the intended readership and the space constraints.
- can be so substantial, that only a few specialists in the field can understand

The $\mathfrak{sT}_E\text{X}$ approach

- The reconstruction and the notation/context problem have to be solved to turn or translate $\text{T}_E\text{X}/\text{L}_E\text{T}_E\text{X}$ into a MKM format
- **Problem:** This is impossible in the general case (AI-hard)
- **Idea:** Enable the author to make structure explicit and disambiguate meanings
 - use the T_EX macro mechanism for this (well established)
 - the author knows the semantics best (at least she understands)
 - the burden is alleviated by manageability savings (MKM on $\text{T}_E\text{X}/\text{L}_E\text{T}_E\text{X}$)
- **Definition 8 ($\mathfrak{sT}_E\text{X}$ Approach)** Semantic pre-loading of $\text{T}_E\text{X}/\text{L}_E\text{T}_E\text{X}$ documents.
 - Introduce semantic macros: e.g. $\backslash\text{union}\{a,b,c\} \rightsquigarrow a \cup b \cup c$
 - Mark up discourse structure: (largely invisible)
e.g. $\backslash\text{begin}\{\text{proof}\}[\text{id}=\text{Wiles},\text{for}=\text{Fermat}]\dots\backslash\text{end}\{\text{sproof}\}$
 - Generate PDF and XML from that (via $\text{L}_E\text{T}_E\text{XML}$ [Miller])

$\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ Modules help with the Notation/Context Problem

- **Note:** *the context of notations coincides with the context of the concepts they denote*
- **Idea:** Use the theory structure for notational contexts
 - The scoping rules of $\mathcal{T}\mathcal{E}\mathcal{X}/\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ follow a hierarchical model:
 - a $\mathcal{T}\mathcal{E}\mathcal{X}$ macro is either globally defined or defined exactly inside the group induced by the $\mathcal{T}\mathcal{E}\mathcal{X}/\mathcal{L}\mathcal{A}\mathcal{T}\mathcal{E}\mathcal{X}$ curly braces hierarchy.
- **Solution:** provide explicit grouping for scope with inheritance.
 - new $\mathcal{S}\mathcal{T}\mathcal{E}\mathcal{X}$ environment module,
 - new macro definition `\symdef`, scoped in module
 - specify the inheritance of `\symdef`-macros in module explicitly
 - `\symdef`-macros are undefined unless in home module or inherited.

STEX Modules: Example

```
\begin{module}[id=pairs]\symdef{pair}[2]{\langle#1,#2\rangle} ... \end{module}

\begin{module}[id=sets]
  \symdef{member}[2]{#1\in #2} % set membership
  \symdef{mmember}[2]{#1\in #2} ... % aggregated set membership
\end{module}

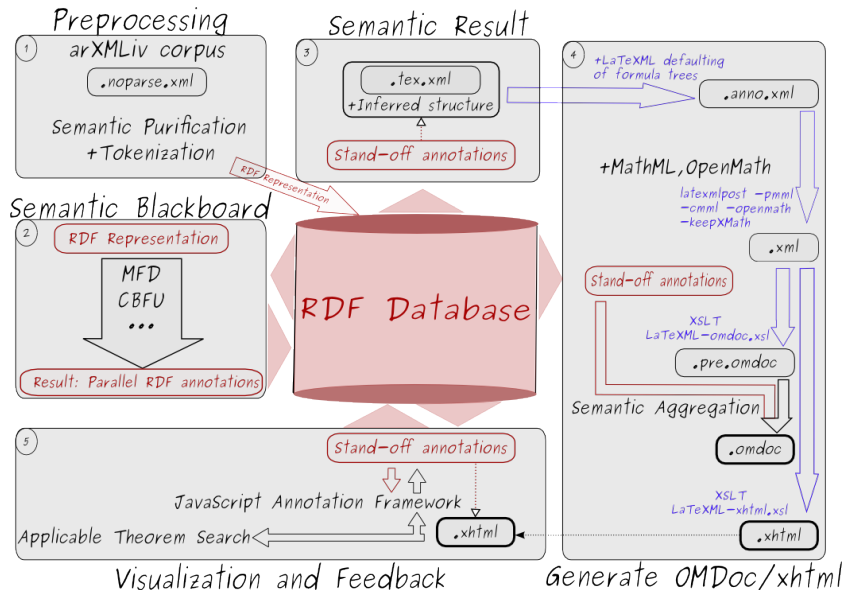
\begin{module}[id=setoid]
  \importmodule{pairs}
  \importmodule{sets}
  \symdef{sset}{\mathcal{S}} % the base set
  \symdef{sopa}{\circ} % the operation symbol
  \symdef{sop}[2]{(#1\sopa #2)} % the operation applied
  \begin{definition}[id=setoid.def]
    A structure  $\langle \text{pair} \setminus \text{sset} \setminus \text{sopa} \rangle$  is called a \defi{setoid}, if  $\langle \text{sset} \rangle$  is closed under  $\langle \text{sopa} \rangle$ , i.e. if  $\langle \text{member} \setminus \text{sop} \setminus \{a\} \setminus \{b\} \rangle \setminus \text{sset}$  for all  $\langle \text{mmember} \setminus \{a,b\} \rangle \setminus \text{sset}$ .
  \end{definition}
\end{module}

\begin{module}[id=semigroup]
  \importmodule{setoid}
  \begin{definition}[id=monoid.def]
    A \trefer{setoid}{setoid}  $\langle \text{pair} \setminus \text{sset} \setminus \text{sopa} \rangle$  is called a \defi{monoid}, if  $\langle \text{sopa} \rangle$  is associative on  $\langle \text{sset} \rangle$ , i.e. if  $\langle \text{sop} \setminus \{a\} \setminus \{ \text{sop} \setminus \{b\} \setminus \{c\} \} \rangle = \langle \text{sop} \setminus \{ \text{sop} \setminus \{a\} \setminus \{b\} \} \setminus \{c\} \rangle$  for all  $\langle \text{mmember} \setminus \{a,b,c\} \rangle \setminus \text{sset}$ .
  \end{definition}
\end{module}
```

The Result of the Example

- **Empirically:** Explicit module structure
 - is a little overhead (can be automated/supported by IDE [JK10])
 - more semantic/portable (but I might be brainwashed)
- In our case study: 320 slides, 160 modules, depth ~ 25

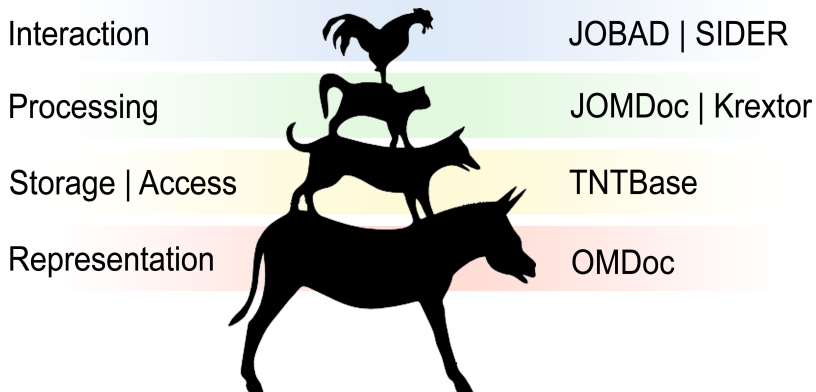
LaMaPUn: Semantic Analysis for Docs with Math (\LaTeX)



Realizing Planetary

Realizing PLANETARY: The KWARC stack

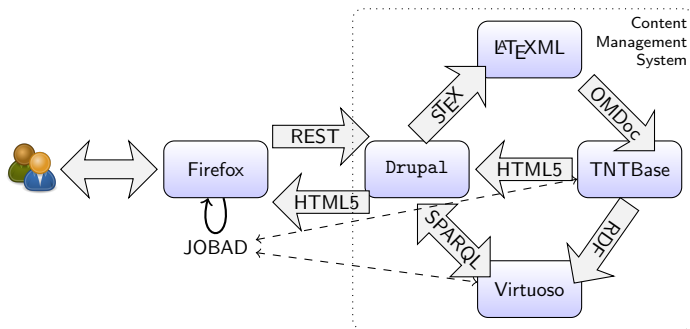
We have already developed the necessary tools/systems over the last decade



PLANETARY is the ideal test bed to integrate them.

Assembling PLANETARY: System Architecture

- PLANETARY functionality can be achieved by integrating existing components.



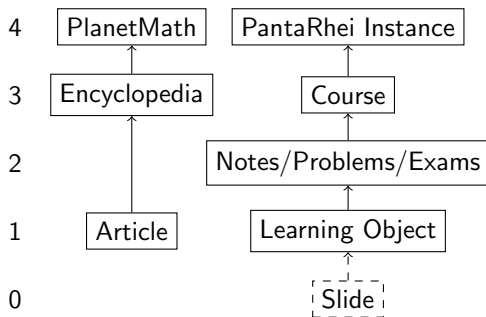
- Drupal for discussions, user management, caching,
- TNTBase for versioned XML storage, OMDoc presentation
- JOBAD integrates semantic services into documents
- Virtuoso is a triple store for semantic relations
- L^AT_EX_ML transforms L^AT_EX/S_TE_X to XHTML+MathML+RDFa

Organization of Content/Narrative Structure

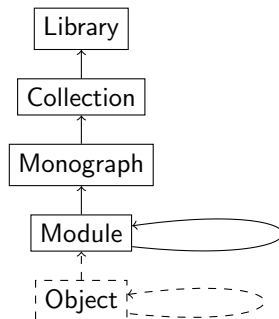
Layers of Documents/Content

- Content and narrative structures come at different conceptual layers

Level **Active Documents**



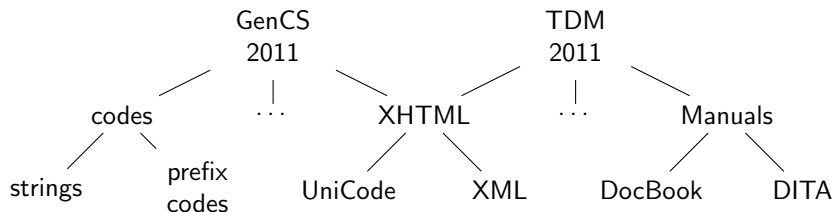
Content Commons



- Different layers support different functionality

Monographs as Module Graphs foster Reuse

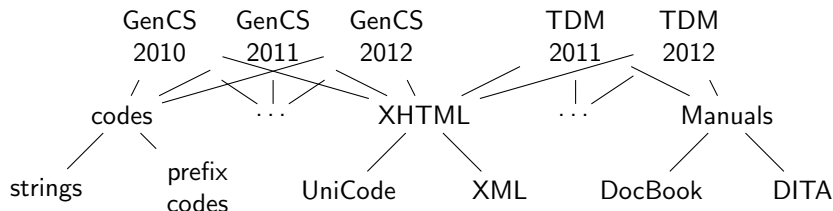
- **Idea:** Modules can be reused in more than one monograph
- **Note:** Similar to, but more general (nesting) than DITA concepts and DITA maps. (but no conditional processing (yet))
- **Example 9** For instance a module on HTML/XML in the courses “General Computer Science” and “Text and Digital Media”.



Observation: These graphs can get quite large: Our corpus has 3300 nodes with 130 roots.

Monographs as Module Graphs foster Reuse

- **Idea:** Modules can be reused in more than one monograph
- **Note:** Similar to, but more general (nesting) than DITA concepts and DITA maps. (but no conditional processing (yet))
- **Example 10** For instance a module on HTML/XML in the courses “General Computer Science” and “Text and Digital Media”.

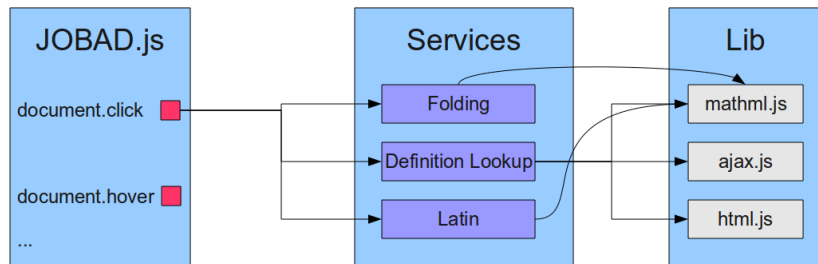


Courses given in different years share most of their content (but not all)

- **Observation:** These graphs can get quite large: Our corpus has 3300 nodes with 130 roots.

JOBAD: Embedding Semantic Services into Web Docs I

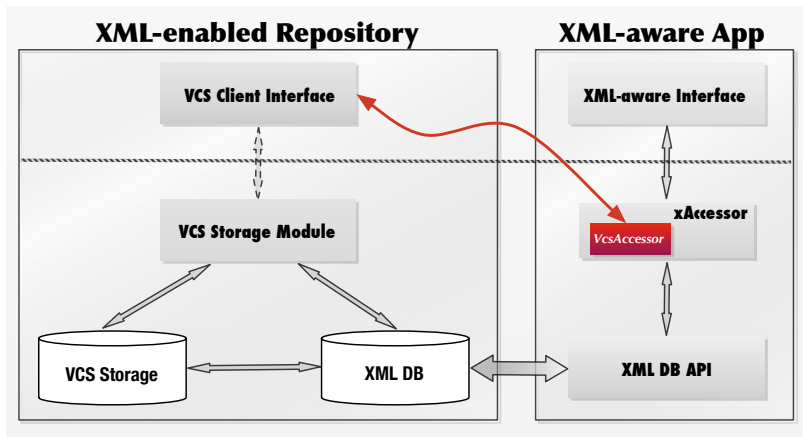
- JavaScript API for (J)OMDoc Based Active Documents
- runs inside client browser (Firefox currently)
- provides client-only or server-based features (extensible framework)
based on semantic annotations in XHTML+MathML+RDFa documents
- Project home page: <https://jomdoc.omdoc.org/wiki/JOBAD>



TNTBase: Versioned Storage for XML

- The TNTBase system is a versioned storage system for XML documents. It combines the functionality and interfaces of Subversion with those of an XML database.

Versioned XML Database



OMDoc in a Nutshell (three levels of modeling)

<p>Formula level: OpenMath/C-MathML</p> <ul style="list-style-type: none"> • Objects as logical formulae • semantics by ref. to theory level 	<pre> <OMA> <OMS cd="arith1" name="plus"/> <OMS cd="nat" name="zero"/> <OMV name="N"/> </OMA> </pre>
<p>Statement level:</p> <ul style="list-style-type: none"> • Definition, Theorem, Proof, Ex. • semantics explicit forms and refs. 	<pre> <defn for="plus" type="rec"> <CMP>rec. eq. for plus</CMP> <FMP>$X + 0 = X$</FMP> <FMP>$X + s(Y) = s(X + Y)$</FMP> </defn> </pre>
<p>Theory level: Development Graph</p> <ul style="list-style-type: none"> • inheritance via symbol-mapping • theory-inclusion by proof-obligations • local (one-step) vs. global links 	<p>The diagram illustrates the Theory level Development Graph with four modules: Nat-List, List, Nat, and Param. Nat-List and List are at the top, while Nat and Param are at the bottom. Nat-List contains 'cons, nil' and '0, s, Nat, <'. List contains 'cons, nil' and 'Elem, <'. Nat contains '0, s, Nat, <'. Param contains 'Elem, <'. Solid green arrows labeled 'Actualization' point from List to Nat-List and from Param to Nat. Dashed blue arrows labeled 'imports' point from List to Nat-List and from Param to Nat. A yellow box labeled 'Proof Obligations' is connected to Nat and Param by a solid green arrow labeled 'theory-inclusion'.</p>

\LaTeX ML: Converting \TeX / \LaTeX Documents to XML

- **Definition 11** \LaTeX ML converts \LaTeX documents to XHTML+MathML
 - re-implement the \TeX parser in perl. (do not expand semantic macros)
 - needs \LaTeX ML bindings for all \LaTeX packages and classes (specify the XML for the emitter)

Case Study: Converting the arXiv into XHTML+MathML
(70% coverage of 550 k documents)

$\text{\texttt{sTeX}}$, a Semantic Variant of $\text{\texttt{TeX/LaTeX}}$

- **Problem:** Need content markup formats for semantic services, but Mathematicians write $\text{\texttt{LaTeX}}$
- **Idea:** Enable the author to make structure explicit and disambiguate meanings
 - use the $\text{\texttt{TeX}}$ macro mechanism for this (well established)
 - the author knows the semantics best (at least she understands)
 - the burden is is alleviated by manageability savings (MKM on $\text{\texttt{TeX/LaTeX}}$)
- **Definition 12 ($\text{\texttt{sTeX}}$ Approach)** Semantic pre-loading of $\text{\texttt{TeX/LaTeX}}$ documents.
 - Introduce semantic macros: e.g. $\text{\texttt{\union\{a,b,c\}}} \rightsquigarrow a \cup b \cup c$
 - Mark up discourse structure: (largely invisible)
e.g. $\text{\texttt{\begin{sproof}[id=Wiles,for=Fermat]... \end{sproof}}}$
 - Generate PDF and OMDoc from that (via $\text{\texttt{LaTeXML}}$ [Mil])

<http://trac.kwarc.info/sTeX/>

Levels of Service in Planetary

PLANETARY at the Presentation/Structural Level

- PLANETARY can make use objects and relations at various levels,
- Example 13 (arXivdemo: Document Structure and Presentational Math)**

The screenshot displays the PLANETARY web interface for an arXiv article. The browser address bar shows the URL `http://arxivdemo.mathweb.org/article/998/nuc1-th.0011027`. The page header includes navigation links: **Dashboard Questions Activity Inbox smirea Articles Books** and a **Sign Out** button. The article title is **Found in: Articles** `nuc1-th.0011027 2011-01-11 16:18:17`. Below the title, there are links for **Edit Details** and **Delete Article**, and a table showing the author `cdavid` and release date `January 11`. On the right side, there are sections for **Make Your Own Articles!** (with links for **Quick-Start Guide** and **Upload a New Article**) and **Who's Online (1)** (showing `smirea` at `2:11PM`). The main content area displays the article title **Low energy scattering and photoproduction of η -mesons on deuterons.** and the section **§ 2. AGS formalism**. The text describes the AGS transition operator U_{11} and the elastic scattering amplitude $f(\mathbf{p}_1, \mathbf{p}_1; z)$. A context menu is open over the text, showing options like **Context Menu (3 Icons)** and **Discussion Thread**. A **FoldingBar** is visible on the left side of the text. An **InfoBar** is visible on the right side. A **Discussion Thread** is also visible on the right side. The page includes mathematical equations and a discussion thread with a comment: "Guest: very long formula... Guest: What is M1 here? I have no idea... maybe file a bug report?". The page footer shows a **0:10** timer.

User Services at the Semantic Level in PLANETARY

Definition Lookup

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

Definition Lookup Results

DEFINITION:

Cartesian product :
 $A \times B := \{ \langle a, b \rangle \mid a \in A \wedge b \in B \}$, call
 $\langle a, b \rangle$ pair .

Semantic Folding

$$s = s_i + v_i \Delta t + \frac{1}{2} a_i (\Delta t)^2$$

Fold
Semantic Fold

$$s = s_i + s_v + s_a$$

contribution from acceleration

Unit Conversion

City A is 9144ft from city B and 5164ft from city C .

Look-up Definition
convert to miles
convert to meters
convert to feet
convert to inches
convert to yards

convert the units

City A is 3048m from city B and 5164ft from city C .

Prerequisites Navigation

Prerequisites Graph for ./slides/dmath/en/sets-operations.tex - Planetary SandBox

http://planetbox.kwarc.info/art... phase 6 date

Dashboard Questions Activity Index kohlhase Articles Books

Prerequisites Graph for ./slides/dmath/en/sets-operations.tex

THIS

sets-relations

sets-introduction

math-talk-definitions

dates

math-talk

highschool

sequences

relation1

setname2

sts

setname1

math-talk, highschool, sequences

Mathematics uses a very effective technique for dealing with conceptual complexity. It usually starts out with discussing simple, basic objects and their properties. These simple objects can be combined to more complex, compound ones. Then it uses a definition to give a compound object a new name, so that it can be used like a basic one. In particular, the newly defined object can be used to form compound objects, leading to more and more complex objects that can be described succinctly. In this way mathematics incrementally extends its vocabulary by add layers and layers of definitions onto very simple and basic beginnings. We will now discuss four definition schemata that will occur over and over in this course.

Find: la Next Previous Highlight all Match case

PantaRhei: Semantic Course Knowledge Exploration

- PantaRhei is a semantic course knowledge exploration system based on the PLANETARY system.

The screenshot shows a web browser window with the address `http://localhost/math/index.php?p=/book/3/1531#`. The page title is "Math/ Forum - MineField". The navigation bar includes "Dashboard", "Questions", "Activity", "Inbox", "tak3r", "Articles", "Books", and "Meheh". The main content area displays a course structure for "General Computer Science" by M. Kohlhasse. The structure includes sections like "1.1 Preface", "1.2 Getting Started with 'General Computer Science'", "1.3 Elementary Discrete Math", and "1.4 Computing with Functions over Inductively Defined Sets". A sidebar on the right lists the course contents, with "1.3.5 Relations and Functions" highlighted. A question box is overlaid on the content, asking "What is this discrete math?" and "It is ...".

Math/ Forum - MineField

File Edit View History Bookmarks Tools Help

Math/ Forum

`http://localhost/math/index.php?p=/book/3/1531#`

Dashboard Questions Activity Inbox tak3r Articles Books Meheh Sign Out

Back Article|Up Article|Next Article

Author M. Kohlhasse

General Computer Science

1.1 Preface

1.2 Getting Started with "General Computer Science"

1.3 Elementary Discrete Math

1.3.1 Mathematical Foundations: Natural Numbers

1.3.2 Talking (and writing) about Mathematics

On our way to understanding functions

Statement 1.3.1

We need to understand sets first.

1.3.3 Talking (and writing) about Mathematics

1.3.4 Naive Set Theory

1.3.5 Relations and Functions

1.4 Computing with Functions over Inductively Defined Sets

1.4.1 Standard ML: Functions as First-Class Objects

1.4.2 Inductively Defined Sets and Computation

Statement 1.4.1

Now, we have seen that "inductively defined sets" are a basis for computation, we will turn to the

Question: What is this discrete math?

It is ...

Cancel Submit

This is the menu!

Book

- 1 General Computer Science Notes
 - 1.1 Preface
 - 1.2 Getting Started with General Computer Science
 - 1.3 Elementary Discrete Math
 - 1.3.1 Mathematical Foundations: Natural Numbers
 - 1.3.2 Talking (and writing) about Mathematics
 - 1.3.3 Talking (and writing) about Mathematics
 - 1.3.4 Naive Set Theory
 - 1.3.5 Relations and Functions**
 - 1.4 Computing with Functions over Inductively Defined Sets
 - 1.4.1 Standard ML: Functions as First-Class Objects
 - 1.4.2 Inductively Defined Sets and Computation
 - 1.4.3 Inductively Defined Sets in SML
 - 1.4.4 A Theory of SML: Abstract Data Types and Term Languages
 - 1.4.5 More SML: Recursion in the Real World
 - 1.4.6 Even more SML: Exceptions and State in SML
 - 1.5 Encoding Programs as Strings
 - 1.6 Boolean Algebra

User Services at the Formal Level in PLANETARY

- Formal Representations Adapted to Distinct User Settings
(Customized via the Dashboard Widget on the Right)

AlgebraTest

- unfold ☒

view OppositeMagmaCommut:MagmaCommut→MagmaCommut

- mag→OppositeMagma ; MagmaCommut/mag
- commut→forallI ($\lambda x.i.(\text{forall } (\lambda x1.i. (x1 \text{ mag}/* x == x \text{ mag}/* x1)))$
 $(\lambda x.i.(\text{forallI } (\lambda x1.i. (x1 \text{ mag}/* x == x \text{ mag}/* x1)) (\lambda y.i.(\text{forall2E}$
 $(\lambda x1.i.(\lambda x2.i. (x1 \text{ mag}/* x2 == x2 \text{ mag}/* x1))) \text{commut } y x)))$

show definitions ☐

show reconstructed types ☒

show implicit arguments ☒

show redundant brackets (high value
= more brackets)

AlgebraTest

- unfold ☒

view OppositeMagmaCommut:MagmaCommut→MagmaCommut

- mag→OppositeMagma ; MagmaCommut/mag
- commut→forallI ($\lambda x.\text{forallI } (\lambda y.\text{forall2E } \text{commut } y x)$)

show definitions ☐

show reconstructed types ☐

show implicit arguments ☐

show redundant brackets (high value
= more brackets)

Accessing Encyclopedias via Ontologies

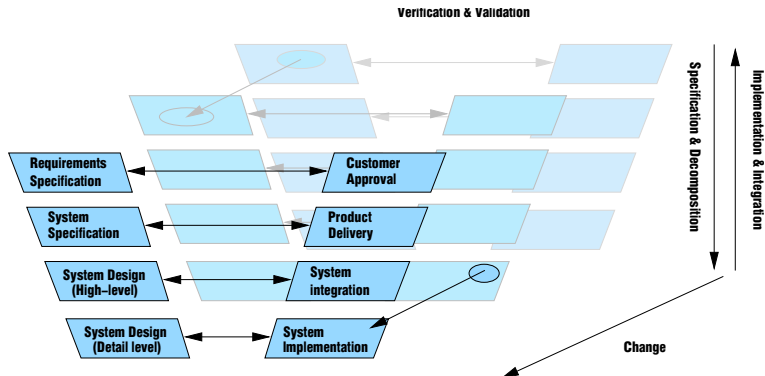
- **Idea:** add classification metadata to articles, harvest as RDF into triplestore, compute access methods via SPARQL queries and SKOS ontology.
- **Example 14 (MSC View in PlanetMath)** use the Math Subject Classification

Discussions Activity Sign In Articles	
top	label
00-xx	General
01-xx	History and biography [See also the classification number -03 in the other sections]
03-xx	Mathematical logic and foundations
subconcept	label
03-00	General reference works [handbooks, dictionaries, bibliographies, etc.]
03-01	Instructional exposition [textbooks, tutorial papers, etc.]
03-02	Research exposition [monographs, survey articles]
03-03	Historical [must also be assigned at least one classification number from Section 01]
article	
	PraeclarumTheorema
	PeircesLaw
	Ampheck
03-04	Explicit machine computation and programs [not the theory of computation]

Ontology-Based Management of Change; A Killer Application for Semantic Techniques

Application: Formal Software Development

- **Idea:** Understand, markup, & version development documents
- **Example 15** For instance in the V Model



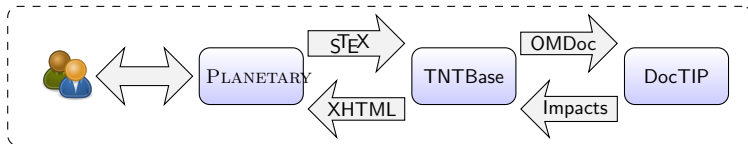
Problem: We need to understand **hybrid documents** (text, math, UML, code)

Management of Change in Planetary

Management of Change in PLANETARY

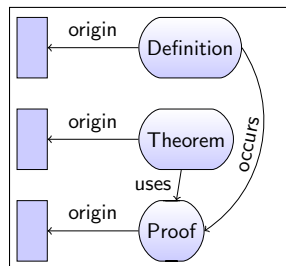
- **Observation**: In an eScience3.0 System, the content is constantly changing.
- **Problem**: How do we maintain **consistency** and **coherence**
- **Idea**: Integrate functionality for **Management of Change**.
 - Make use of the semantic relations already in place in PLANETARY.
 - If A depends on B , then a change in B impacts A .
 - Extend PLANETARY by the DocTIP system from OMoC.

(Joint project with DFKI Bremen).
- Prototypical Integration in PLANETARY available [ADD⁺11]

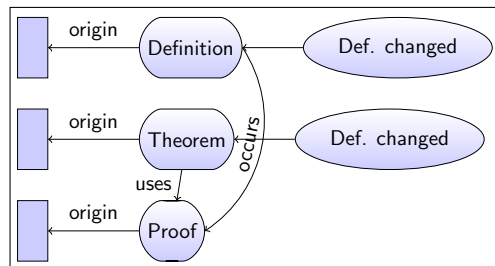


Change Impact Analysis in DocTIP

- **Idea:** If A depends on B , then a change in B impacts A .
- **Definition 16** **Change Impact Analysis (CIA)** is a process for computing potentially impacted fragments in a document collection \mathcal{C} from a change description and semantic relations in \mathcal{C} .
- In DocTIP, CIA is computed by graph rewriting rules on the document ontology.
- **Example 17** CIA propagation rules for OMDoc



(a) Initial Syntax and Semantics



(b) Propagated Impacts after Definition Change

MoC in PLANETARY I

- Extend the commit dialog with CIA

Dashboard Questions Activity Inbox cdavid Articles Books Manage Impacts

Edit Article

Type of Article Encyclopedia Article ▼

CommitMessage reworded, no meaning change

Name balanced-binary-trees

☒ Perform Impact Analysis

RelURL balanced-binary-trees

Commit

Body

```
\begin{module}[id=balanced-binary-trees]
\importmodule{\KWARCS\slides{graphs-trees/en/trees}}{trees}
\importmodule{\KWARCS\slides{graphs-trees/en/graph-depth}}{graph-depth}

\begin{frame}
\frametitle{Balanced Binary Trees}
\begin{itemize}
\item
\begin{definition}[id=binary-tree.def,title=Binary Tree]
A is \termref{cd=trees,name=tree}{tree} is called
\{ \twindefault{binary}{binary}{tree} \}, iff all its
\termref{cd=graphs-intro,name=node}{nodes} have
\termref{cd=graphs-intro,name=out-degree}{out-degree} 2 or 0.
\end{definition}
\item
\begin{definition}[id=balanced-binary-tree.def]
A \termref{name=binary-tree}{binary tree} is called \{ \atwindefault{balanced}{balanced}
{binary}{tree} \} 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 \atwindefault{fully
balanced}{fully}{balanced}{tree}, iff the
```



Kohlhase: Planetary: Web3.0 for Math

64

NIST, May 2012



MoC in PLANETARY II

- The Impact Resolution Dialog

The screenshot shows the PlanetMath web interface. At the top is a navigation bar with links: Dashboard, Questions, Activity, Inbox, cdavid, Articles, Books, Manage Impacts 2, and Sign Out. Below this is a status bar showing 'Commit info: r35 at Sun, 13 Mar 11 12:00:27 -0400' and a message 'reworded, no meaning change' with 'Diff' and 'Quit' buttons. The main content area has two tabs: 'bbs-size' and 'balanced-binary-trees'. Below the tabs are buttons for 'Edit this article' and 'Impact analysis OK'. The article title is 'Size Lemma for Balanced Trees'. The text of the lemma is: 'Lemma 3.1.9: Let $G = \langle V, E \rangle$ be a balanced binary tree of depth $n > 0$, then the set $V_i := \{v \in V \mid \text{dp}(v) = i\}$ of nodes at depth i has cardinality 2^i .' The proof follows: 'Proof : via induction over the depth i . 1. We have to consider two cases 1.1. $i = 0$ 1.1.1. then $V_i = \{v_r\}$, where v_r is the root, so $\#(V_0) = \#(\{v_r\}) = 1 = 2^0$. 1.2 $i > 0$ 1.2.1. then V_{i-1} contains 2^{i-1} vertexes (IH). 1.2.2. By the definition of a binary tree, each $v \in V_{i-1}$ is a leaf or has two children that are at depth i .' A blue box highlights the 'Accept Change' button, which is accompanied by a green checkmark icon.

Searching for Mathematical Formulae

Introduction & Motivation

Why we need a search engine for Mathematics

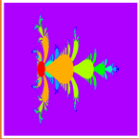
- We have come to rely on the World Wide Web for almost all of our information needs.
- The Internet is only the transport layer. We see a succession of techniques
 - Early Web navigation via explicitly represented hyperlinks
 - Mature Web finding relevant content via search engines
(bag of words techniques for textual content)
 - Semantic Web Representation of meaning and inferring content that is not explicitly represented
- For scientific content, we are still in the “Early Web” phase
 - Need a “Semantic Web for Science” (talk about OMDoc some other time)
 - Today: provide techniques for the “Mature Web”
 - Concretely: a search engine for math. formulae
(a prominent non-textual part of science)

Mathematics Resources on the Web

WOLFRAMRESEARCH**functions.wolfram.com**OTHER WOLFRAM SITES►

Search SiteGoFormula SearchSearch Tips

FUNCTION CATEGORIESVISUALIZATIONSNOTATIONSGENERAL IDENTITIESABOUT THIS SITEContributeEmail CommentsSign the Guestbook



Exp
Exponential function



Mathematica Notation: $\text{Exp}[z]$

Traditional Notation: $\exp(z) = e^z$


VIEW RELATED INFORMATION IN

- The Mathematica Book
- MathWorld

DOWNLOAD FORMULAS FOR THIS FUNCTION

-  Mathematica Notebook
-  PDF File

DOWNLOAD SOURCE FOR VISUALIZATIONS

-  Mathematica Notebook

Elementary Functions ► Exp[z] ► Theorems ▼

► Show All Below

Fourier transformation and Parseval relation (1 formula)

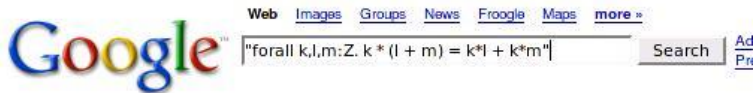
$$\hat{f}(y) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{iyx} dx \Leftrightarrow f(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(y) e^{-iyx} dy,$$

$$\int_{-\infty}^{\infty} f_1(t) f_2(x-t) dt = \int_{-\infty}^{\infty} \hat{f}_1(y) \hat{f}_2(y) e^{-iyx} dy.$$

More Mathematics on the Web

- The Connexions project (<http://cnx.org>)
- Wolfram Inc. (<http://functions.wolfram.com>)
- Eric Weisstein's MathWorld (<http://mathworld.wolfram.com>)
- Digital Library of Mathematical Functions (<http://dlmf.nist.gov>)
- Cornell ePrint arXiv (<http://www.arxiv.org>)
- Zentralblatt Math (<http://www.zentralblatt-math.org>)
- ...
- **Question:** How will we find content that is relevant to our needs
- **Idea:** try Google (like we always do)
- **Scenario:** Try finding the distributivity property for \mathbb{Z}
$$(\forall k, l, m \in \mathbb{Z}. k \cdot (l + m) = (k \cdot l) + (k \cdot m))$$

Searching for Distributivity



Web

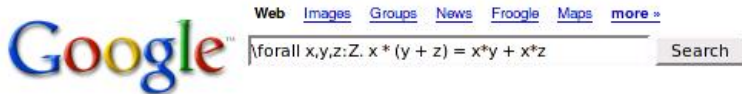
Tip: Try removing quotes from your search to get more results.

Your search - **"forall k,l,m:Z. k * (l + m) = k*l + k*m"** - did not match any documents.

Suggestions:

- ◆ Make sure all words are spelled correctly.
- ◆ Try different keywords.
- ◆ Try more general keywords.

Searching for Distributivity



Web

Untitled Document

... theorem distributive_Ztimes_Zplus: distributive Z Ztimes Zplus. change with ($\lambda \text{forall } x,y,z:Z. x * (y + z) = x*y + x*z$). intros.elim x. ...

[matita.cs.unibo.it/library/Ztimes.ma](http://matita.cs.unibo.it/library/Z%20times.ma) - 21k - [Cached](#) - [Similar pages](#)

Searching for Distributivity



Web

Mathematica - Setting up equations

Try ***Reduce*** rather than ***Solve*** and use ***ForAll*** to put a condition on x, y, and z. In[1]:=

Reduce[ForAll[{x, y, z}, 5*x + 6*y + 7*z == a*x + b*y + c*z], ...

www.codecomments.com/archive382-2006-4-904844.html - 18k - Supplemental Result -

[Cached](#) - [Similar pages](#)

[PDF] arXiv:nlin.SI/0309017 v1 4 Sep 2003

File Format: PDF/Adobe Acrobat - [View as HTML](#)

7.2 Appendix B. Elliptic constants related to $gl(N, \mathbb{C})$ 1 for all $s \leq j$. (4.14). The first condition means that the traces (4.13) of the Lax operator ...

www.citebase.org/cgi-bin/fulltext?format=application/pdf&identifier=oai:arXiv.org:nlin/0309017 -

Supplemental Result - [Similar pages](#)

\documentclass{article} \usepackage{axiom} \usepackage{amssymb ...

i+1) bz := (bz - 2**i)::NNI else bz := bz + 2**i z.bz := z.bz + c z x * y == z ... b,i-1)] be := reduce("...", m)

c = 1 => be c::Ex * be coerce(x): Ex == tl ...

wiki.axiom-developer.org/axiom-test-1/src/algebra/CliffordSpad/src - 20k - Supplemental Result -

[Cached](#) - [Similar pages](#)

Of course Google cannot work out of the box

- Formulae are not words:
 - $a, b, c, k, l, m, x, y,$ and z are (bound) variables. (do not behave like words/symbols)
 - where are the word boundaries for “bag-of-words” methods?
- Idea: Need a special treatment for formulae (translate into “special words”) Indeed this is done ([MY03, MM06, LM06, MG11]) . . . and works surprisingly well (using Lucene as an indexing engine)
- Idea: Use database techniques (extract metadata and index it) Indeed this is done for the Coq/HELM corpus ([Asperti&al'04])
- Our Idea: Use Automated Reasoning Techniques (free term indexing from theorem prover jails)

A running example: The Power of a Signal

- An engineer wants to compute the power of a given signal $s(t)$
- She remembers that it involves integrating the square of s .
- **Problem:** But how to compute the necessary integrals
- **Idea:** call up MathWebSearch with $\int_?^? s^2(t)dt$.
- MathWebSearch finds a document about Parseval's Theorem and $\frac{1}{T} \int_0^T s^2(t)dt = \sum_{k=-\infty}^{\infty} |c_k|^2$ where c_k are the Fourier coefficients of $s(t)$.

Some other Problems (Why do we need more?)

- **Substitution Instances**: search for $x^2 + y^2 = z^2$, find $3^2 + 4^2 = 5^2$
- **Homonymy**: $\binom{n}{k}$, ${}_nC^k$, C_k^n , and C_n^k all mean the same thing (binomial coeff.)
- **Solution**: use content-based representations (MathML, OpenMath)
- **Mathematical Equivalence**: e.g. $\int f(x)dx$ means the same as $\int f(y)dy$ (α -equivalence)
- **Solution**: build equivalence (e.g. α or ACI) into the search engine (or normalize first [Normann'06])
- **Subterms**: Retrieve formulae by specifying some sub-formulae
- **Solution**: record locations of all sub-formulae as well

Term Indexing

Term-Indexing

- **Motivation:** Automated theorem proving (efficient systems)
- **Problem:** Decreasing inference rate (basic operations linear in # of formulae)
- **Idea:** Make use of structural equality between terms (term indexing)
database systems (Algorithms: select, meet, join)

- **Data:** PERSON(hans, manager, 32)

- **Query:** "find all 40-year old persons"

automated theorem proving

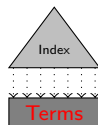
- **Data:** $P(f(x, g(a, b)))$

- **Queries:** "find all literals that are unifiable with $P(f(c, y))$ "

An (additional) index data structure can make the retrieval logarithmic



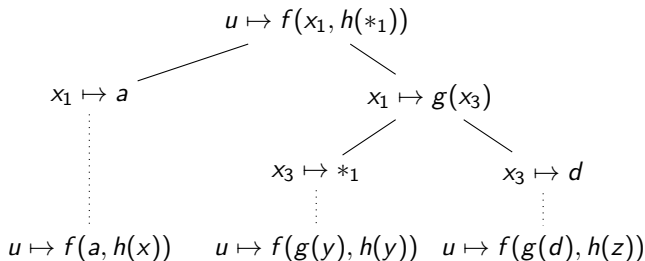
(Algorithm: Unification)



Tree-based Indexing: use structural similarity of terms

Discrimination Tree	Abstraction Tree
<p>Terms as linear chains, shared in trees</p>	<p>Trees of substitution instances</p>

Substitution Tree [Graf '94]



- Variant of abstraction trees that indexes **Substitutions** (Nodes labeled with Substitutions)
- includes **Variable renaming** ($*_i \hat{=} i^{\text{th}}$ variable)
- less redundant than abstraction trees
- allows $n : m$ indexing

Unification-based Search

Unification-Based Querying for Math. Formulae

- **Theory:** Substitution Tree Indexing is a perfect filter for
 - **Variants:** $\{\sigma \mid \sigma \in \mathbf{GEN}(\tau, \rho) \wedge \mathbf{supp}(\sigma) \cap \mathbf{V}^* = \emptyset\}$
 - **Instances:** $\{\sigma \mid \sigma \in \mathbf{UNIF}(\tau, \rho) \wedge \mathbf{supp}(\sigma) \cap \mathbf{V}^* = \emptyset\}$
 - **Generalization:** $\{\sigma \mid \forall x_i \in \mathbf{supp}(\tau). \tau \rho \sigma(x_i) = \rho(x_i)\}$
 - **Unification:** $\{\sigma \mid \forall x_i \in \mathbf{supp}(\tau). \tau \rho \sigma(x_i) = \rho \sigma(x_i) \sigma \text{ mgu}\}$
- **Idea:** Use all of them for querying Formulae mathematical Formulae
 - **Variants:** To find formulae of a given structure
 - **Instances:** To find formulae of a partially remembered structure
 - **Generalization:** To find applicable Theorems
 - **Unification:** A mixture of all three.

MathWebSearch: Search Math. Formulae on the Web

- Idea 1: Crawl the Web for math. formulae (in OpenMath or CMathML)
- Idea 2: Math. formulae can be represented as first order terms (see below)
- Idea 3: Index them in a substitution tree index (for efficient retrieval)
- Problem: Find a query language that is intuitive to learn
- Idea 4: Reuse the XML syntax of OpenMath and CMathML, add variables

Indexing Math Formulae as First-Order Terms?

- Mathematical Expression: $\int_0^\infty s^2(t)dt$.
- Content MathML: Formulae built up by function application, and binding from constants and variables.

```
<math>
  <apply><defint/>
    <apply><interval/><cn>0</cn><infinity/></apply>
    <bind>
      <lambda/>
        <bvar><ci>t</ci></bvar>
        <apply><power/>
          <apply><ci>s</ci><ci>t</ci></apply>
          <cn>2</cn>
        </apply>
      </bind>
    </apply>
  </math>
```

Idea: Extend Substitution Tree Indexing with bound variables and α -renaming

- **Technically:** Use deBruijn Indexes
(bvars as name-less pointers interact well with substitution)

Instantiation Queries

- **Application:** Find partially remembered formulae
- **Example 18** An engineer might face the problem remembering the energy of a given signal $f(x)$
 - **Problem:** hmmm, have to square it and integrate
 - **Query Term:** $\int_{\boxed{\text{min}}}^{\boxed{\text{max}}} \boxed{f}(x)^2 dx$ (\boxed{i} are search variables)
 - **One Hit: Parseval's Theorem** $\frac{1}{T} \int_{-T}^{T} s^2(t) dt = \sum_{k=-\infty}^{\infty} \|c_k\|^2$ (nice, I can compute it)
- This works out of the box (has been working in MathWebSearch for some time)
- **Another Application: Underspecified Conjectures/Theorem Proving**
 - during theory exploration we often have some freedom
 - express that using metavariables in conjectures
 - instantiate the conjecture metavariables as the proof as the proof dictatesapplied e.g. in Alan Bundy's "middle-out reasoning" in proof planing

Generalization Queries

- **Application:** Find (possibly) applicable theorems
- **Example 19** A researcher wants to estimate $\int_{\mathbb{R}^2} |\sin(t) \cos(t)| dt$ from above
 - **Problem:** Find inequation such that $\int_{\mathbb{R}^2} |\sin(t) \cos(t)| dt$ matches left hand side.
 - e.g. Hölder's Inequality: (i are universal variables)

$$\int_{\boxed{D}} |\boxed{f}(x) \boxed{g}(x)| dx \leq \left(\int_{\boxed{D}} |\boxed{f}(x)|^p dx \right)^{\frac{1}{p}} \left(\int_{\boxed{D}} |\boxed{g}(x)|^q dx \right)^{\frac{1}{q}}$$

- **Solution:** Take the instance

$$\int_{\mathbb{R}^2} |\sin(x) \cos(x)| dx \leq \left(\int_{\mathbb{R}^2} |\sin(x)|^p dx \right)^{\frac{1}{p}} \left(\int_{\mathbb{R}^2} |\cos(x)|^q dx \right)^{\frac{1}{q}}$$

Problem: Where do the index formulae come from in particular the universal variables (we'll come back to that later)

Unification Queries

- **Application:** Find applicable theorems for underspecified formulae
- **Example 20** estimate $g^2 \cos(\boxed{x}) + b \sin(\sqrt{y})$
 - this unifies with $\boxed{a} \cos(\boxed{t}) + \boxed{b} \sin(\boxed{t}) \leq ?$
 - result: $g^2 \cos(\sqrt{y}) + b \sin(\sqrt{y}) \leq \sigma(?)$, where σ is the mgu

Problem: Users find it difficult to state the exact unification query

- **Solution:** (from Databases again)
 - express the query in SELECT FROM WHERE form.
 - e.g. **SELECT** instance $\left(\int_{\boxed{D}} \frac{1}{\sqrt{2\pi}} \exp \{ \boxed{B} \} \right)$ **WHERE** $B = \text{variation}(x^2 + jy^2)$
- MathWebSearch preprocessor compiles subqueries into one unification query for efficiency.

Where do the universal variables come from

- **Problem:** we need to have e.g. Hölder's Inequality in the index:

$$\int_D |\boxed{f}(x)\boxed{g}(x)| dx \leq \left(\int_D |\boxed{f}(x)|^p dx \right)^{\frac{1}{p}} \left(\int_D |\boxed{g}(x)|^q dx \right)^{\frac{1}{q}}$$

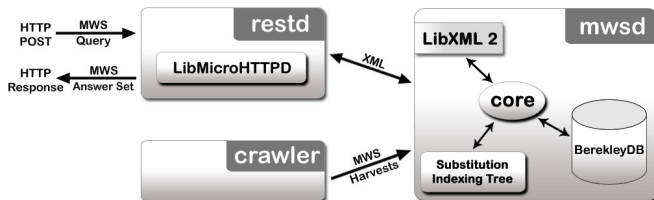
- How do we know what symbols are “universal” (to be instantiated?)
- what is their scope (when are different occurrences of f different?)
- we have no sources with explicit quantifiers, but ([Wikipedia])

*Let (D, Σ, μ) be a measure space and let $1 \leq p, q \leq \infty$ with $1/p + 1/q = 1$.
Then, for all measurable real- or complex-valued functions f and g on D , ...*

- **Solution:** Use techniques from computational linguistics and integrate them into the indexing pipeline. (we have started a bit on the arXiv)

The MathWebSearch System

System Architecture

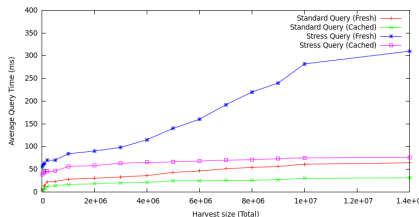


- crawlers for MathML, OpenMath, and OAI repositories. (convert your's?)
- multiple search servers based substitution tree indexing (formula search)
- a RESTful server that acts as a front-end for multiple search servers.
- various front ends tailored to specific applications (search appliances)
 - a Google-like web front end for human users (search.mathweb.org)
 - a \LaTeX -based front-end for the arXiv (<http://arxivdemo.mathweb.org>)
 - special integrations for theorem prover libraries (MizarWiki, TPTP)

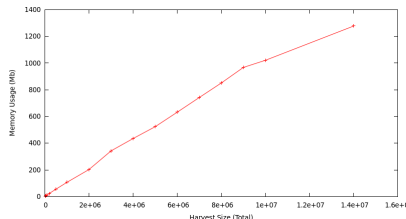
Index statistics

- **Experiment:** Indexing the arXiv (700k documents, $\sim 10^8$ non-trivial formulae)
- **Results:** indexing up to 15 M formulae on a standard laptop

Query Times



Memory Footprint



- query time is constant (~ 50 ms) (as expected; goes by depth \times symbols)
- memory footprint seems linear ($\sim 100 \frac{B}{\text{formula}}$) (expected more duplicates)
- So we need ca 15 GB RAM for indexing the whole arXiv.
- Can index all published Math ($\hat{=}$ $5 \times$ arXiv) on a large server. (ZBL $\hat{=}$ 3M art.)

Instead of a Demo: Searching for Signal Power

Math WebSearch

A SEMANTIC SEARCH ENGINE

Search for:

XML Query

String

int($\lambda x.e^n r$)

QMath:en

$$\int e^n r dx$$

Variables			
Variable	Generic	Any#	Function
r	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
n	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
x	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Arithmetic ...			
Transcendental functions ...			
Calculus			
$\partial_x x$	$\partial^n x$	$\partial_{x,y}(xy)$	
$\int x dx$	$\int_a^b x dx$	$[a, b]$	(a, b)
$\lim_{x \rightarrow x_0} x$	∞	(a, b)	$[a, b)$
∇f	$\nabla^2 v_f$	$\text{curl} v_f$	$\text{div} v_f$
Sets ...			
Logic and relations ...			
Functions			

Search

[Examples](#) | [Help](#) | [API](#) | [About](#) | [Contact](#)

Instead of a Demo: Search Results

[Other integrals \(5 formulas\)](#) (Source)

Other integrals (5 formulas)

Matched term:

$$\int \frac{e^{3z/4}}{(-2+e^{3z/4})\sqrt{-2+e^{3z/4}+e^{3z/2}+5e^{3z/4}-2}} dz = \frac{2}{3} \left(\log(-2+e^{3z/4}) - \log(4\sqrt{-2+e^{3z/4}+e^{3z/2}+5e^{3z/4}-2}) \right)$$

Rank: 100%

[XML Source](#)

Used substitution:

$$n \rightarrow 3z4^{-1}$$

$$r \rightarrow \left(\left((-2) + e^{3z4^{-1}} \right) \left((-2) + e^{3z4^{-1}} + e^{3z2^{-1}} \right)^{1/2} \right)^{-1}$$

$$x \rightarrow z$$

Instead of a Demo: L^AT_EX-based Search on the arXiv

Questions Activity Sign In Books Articles MWS Engine BETA

`\lim_{\var{x}\rightarrow 0}\var{y}`

$\lim_{x \rightarrow 0} y$

```
<m:apply>
  <m:apply>
    <m:csymbol
cd="ambiguous">subscript</m:csymbol>
    <m:limit/>
    <m:apply>
      <m:ci>→</m:ci>
```

Search

Examples - LaTeX queries

Generic subscript search

Specific subscript search

Specific integral search

Physical constant search

All limits approaching zero

Text in math search

1 2 next

$$\chi(t, t_w) = \lim_{h_0 \rightarrow 0} \frac{m[h](t)}{h_0}.$$

Generalized off-equilibrium fluctuation-dissipation relations in random Ising systems

Author: Federico Ricci-Tersenghi <ricci@chimera.roma1.infn.it>

$$\lim_{\mu, \mu_0 \rightarrow 0} I_1^1(\mu, \mu_0, \phi - \phi_0) = \frac{aF_0}{4(c+1)},$$

Behavior of the reflection function of a plane-parallel medium for directions of incidence and reflection tending to horizontal directions

Author: Daphne Stam <d.m.stam@sron.nl>

$$\lim_{\mu, \mu_0 \rightarrow 0} I_1^1(\mu, \mu_0, \phi - \phi_0)$$

Behavior of the reflection function of a plane-parallel medium for directions of incidence and reflection tending to horizontal directions

Instead of a Demo: Applicable Theorem Search in Mizar

definition

```
let k, n be Ordinal;  
pred k divides n means :Def3: :: MTEST1:def 3  
ex a being Ordinal st n = k *^ a;
```

reflexivity

proof

```
let n be Ordinal; :: thesis:  
thus ex a being Ordinal st n = n *^ a ;
```

ATP Proof not found

status: Timeout
Suggest hints, Unification query,

Suggested hints

t73_card_2, t39_ordinal2,

Try SPASS, Export problem to SystemOnTPTP

```
:: thesis:
```

```
end;
```

```
end;
```

But Math consists of more than Formulae

- **Idea:** Text and Formulae co-constrain each other in search would like to search for formulae as well as words
- **Problem:** turnkey text search engines exist, but are incompatible with MathWebSearch
- **Solution:** MaTeSearch combines text and formula search **results**.
 - ① compute text/formula search individually, rank them
 - ② compute ordered intersection of both and display.
- **Problem:** Ranking for formula search is still in its infancy

Conclusions & Future Work

- The time is ripe for eMath 3.0
 - We have the necessary building blocks (integration needed)
 - we have first case studies for eMath 3.0 (more to come)
 - **General Metaphor**: Math Document as an interface to Math Knowledge

Conclusions & Future Work

- The time is ripe for eMath 3.0
 - We have the necessary building blocks (integration needed)
 - we have first case studies for eMath 3.0 (more to come)
 - **General Metaphor**: Math Document as an interface to Math Knowledge
- Active Documents
 - as interfaces to a structured knowledge base (background ontology)
 - via **semantic annotation of documents** (referencing it)

Conclusions & Future Work

- The time is ripe for eMath 3.0
 - We have the necessary building blocks (integration needed)
 - we have first case studies for eMath 3.0 (more to come)
 - **General Metaphor**: Math Document as an interface to Math Knowledge
- Active Documents
 - as interfaces to a structured knowledge base (background ontology)
 - via **semantic annotation of documents** (referencing it)
- The Planetary System as an integrated interaction platform.
 - semantic services for added value
 - semantification support for documents

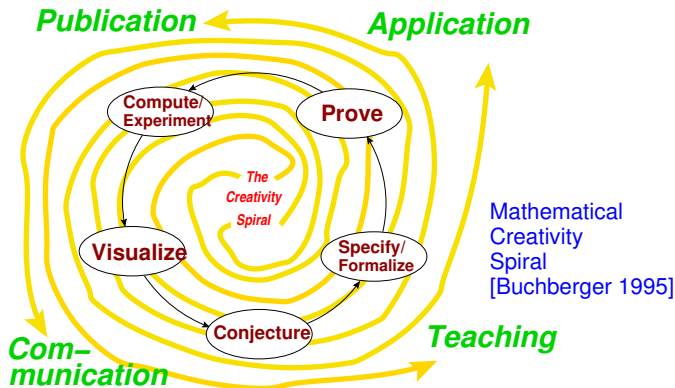
Conclusions & Future Work

- The time is ripe for eMath 3.0
 - We have the necessary building blocks (integration needed)
 - we have first case studies for eMath 3.0 (more to come)
 - **General Metaphor**: Math Document as an interface to Math Knowledge
- Active Documents
 - as interfaces to a structured knowledge base (background ontology)
 - via **semantic annotation of documents** (referencing it)
- The Planetary System as an integrated interaction platform.
 - semantic services for added value
 - semantification support for documents
- Management of Change (to keep mutable collections consistent)
 - Change Impact Analysis based on the existing semantic relations
 - Flexible impact resolution workflows in Planetary

Conclusions & Future Work

- The time is ripe for eMath 3.0
 - We have the necessary building blocks (integration needed)
 - we have first case studies for eMath 3.0 (more to come)
 - **General Metaphor**: Math Document as an interface to Math Knowledge
- Active Documents
 - as interfaces to a structured knowledge base (background ontology)
 - via **semantic annotation of documents** (referencing it)
- The Planetary System as an integrated interaction platform.
 - semantic services for added value
 - semantification support for documents
- Management of Change (to keep mutable collections consistent)
 - Change Impact Analysis based on the existing semantic relations
 - Flexible impact resolution workflows in Planetary
- **Research Interest: Apply this to semi-formal STEM Documents**

The way we do math will change dramatically



- Every step will be supported by mathematical software systems
- Towards an infrastructure for web-based mathematics!



Serge Autexier, Catalin David, Dominik Dietrich, Michael Kohlhase, and Vyacheslav Zholudev.

Workflows for the management of change in science, technologies, engineering and mathematics.

In Davenport et al. [DFRU11], pages 164–179.



Mark Birbeck, Ben Adida, Ivan Herman, and Manu Sporny.

RDFa 1.1 primer.

W3C Working Draft, World Wide Web Consortium (W3C).



James Davenport, William Farmer, Florian Rabe, and Josef Urban, editors. *Intelligent Computer Mathematics*, number 6824 in LNAI. Springer Verlag, 2011.



Constantin Jucovschi and Michael Kohlhase.

sTeXIDE: An integrated development environment for sTeX collections.

In Serge Autexier, Jacques Calmet, David Delahaye, Patrick D. F. Ion, Laurence Rideau, Renaud Rioboo, and Alan P. Sexton, editors, *Intelligent Computer Mathematics*, number 6167 in LNAI. Springer Verlag, 2010.



Graham Klyne and Jeremy J. Carroll.

Resource Description Framework (RDF): Concepts and abstract syntax.

W3C recommendation, World Wide Web Consortium (W3C), 2004.



Paul Libbrecht and Erica Melis.

Methods for Access and Retrieval of Mathematical Content in ActiveMath.
In N. Takayama and A. Iglesias, editors, *Proceedings of ICMS-2006*, number 4151 in LNAI, pages 331–342. Springer Verlag, 2006.
<http://www.activemath.org/publications/Libbrecht-Melis-Access-and-Retrieval-ActiveMath-ICMS-2006.pdf>.



Jozef Misutka and Leo Galambos.

System description: Egomath2 as a tool for mathematical searching on wikipedia.org.
In Davenport et al. [DFRU11], pages 307–309.



Bruce Miller.

LaTeXML: A \LaTeX to XML converter.
Web Manual at <http://dlmf.nist.gov/LaTeXML/>.
seen September 2011.



Rajesh Munavalli and Robert Miner.

Mathfind: a math-aware search engine.
In *SIGIR '06: Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 735–735, New York, NY, USA, 2006. ACM Press.



Bruce R. Miller and Abdou Youssef.

Technical aspects of the digital library of mathematical functions.

Annals of Mathematics and Artificial Intelligence, 38(1-3):121–136, 2003.



OWL 2 web ontology language: Document overview.

W3C recommendation, World Wide Web Consortium (W3C), October 2009.