

X-Technologien

XML and Friends

9. Juli 2001

Jörn Clausen

`joern@TechFak.Uni-Bielefeld.DE`

Übersicht

XML	Extensible Markup Language
SAX	Simple API for XML
DOM	Document Object Model
XSL	Extensible Stylesheet Language
XPath	XML Path Language
XSLT	XSL Transformations
XSL FO	XSL Formatting Objects
SGML	Standard Generalized Markup Language
DSSSL	Document Style Semantics and Specification Language

Markup – Wofür?

- ursprüngliche Domäne: Textdokumente
- Auszeichnung von *Bedeutung* statt *Formatierung*
- Stylesheets beschreiben Formatierung
- Anwendungen: DocBook, CALS, HTML
- heute: strukturierte Beschreibung von Daten aller Art
- Bioinformatik: BioXML (GAME), MAML, GEML, Pise, ...

XML – Extensible Markup Language

- ist ...
 - *keine* Markupsprache
 - *kein* Allheilmittel, trotz des Hypes
 - SGML—
- ist aber auch ...
 - einfach zu erlernen
 - weithin akzeptiert
 - durch viele Implementierungen etabliert

ein Beispiel

```
<?xml version="1.0"?>

<presentation>
  <slide>
    <title use="internal">XML & Friends for Dummies</title>
    <title use="external">Doing B2B with XML & Friends</title>
    <ilist>
      <item>XML is not a markup language
        (unlike HTML)</item>
      <item>XML is a meta language for
        defining markup languages</item>
      <item>XML instances can be <emph>well
        formed</emph> or even
        <emph>validating</emph></item>
    </ilist>
  </slide>
</presentation>
```

Aufbau von XML: Elemente

- öffnendes und schließendes *tag*:

```
<item>XML is not a ...</item>
```

- keine Minimierungsregeln

- leeres Inhaltsmodell:

```
<hr /> statt <hr></hr> statt <hr>
```

- Schachtelung muß „passen“: *well-formed*

Aufbau von XML: Attribute

- Zusatzinformationen zu Elementen

```
<title use="internal">... for Dummies</title>
```

- Attribute im öffnenden tag
- Wertebereich: Aufzählungstyp, Zahlen, Freitext
- schlechte Typisierung
- Design-Frage: Wann Element, wann Attribut?

```
<date y="2001" m="7" d="9" />
```

vs.

```
<date><y>2001</y><m>7</m><d>9</d></date>
```

Aufbau von XML: Entitäten

- Makros und Sonderzeichen
- aus HTML bekannt: `ä` `ê` `©`
- in XML vordefiniert: `&` `<` `>` `'` `"`
- XML verwendet Unicode

Document Type Definition

- Grammatik kann durch DTD beschrieben werden
- `<!DOCTYPE presentation SYSTEM "presentation.dtd">`
- Parser kann XML-Instanz gegen DTD *validieren*
- Definition von
 - Elementen: Inhaltsmodell
 - Attributen: Wertebereich, Vorbelegung
 - Entitäten: character entities, parameter entities
- DTD sehr eingeschränkt
- XML Schema, RELAX NG, ca. 10 weitere Vorschläge

DTD zum Beispiel

```
<!ENTITY % text          "(#PCDATA|emph|bold|ital)*">

<!ELEMENT  presentation  (slide)+>
<!ELEMENT  slide         (title*, ilist)>
<!ELEMENT  title         %text;>
<!ATTLIST  title
           use            (internal|external)  #IMPLIED>

<!ELEMENT  ilist         (item)+>
<!ELEMENT  item          %text;>

<!ELEMENT  emph          %text;>
<!ELEMENT  bold          %text;>
<!ELEMENT  ital          %text;>
```

Name Spaces

- Kombination von verschiedenen XML-Sprachen
- z.B. allgemeine, wiederverwendbare Definition von Tabellen:

```
<pr:presentation xmlns:pr="http://www.slides.org/presentation-1.0">
  <pr:slide>
    <pr:title>increasing use of XML</pr:title>
    <cals:table xmlns:cals="http://www.army.mil/SPEC-CALS">
      <cals:tr>
        <cals:th>year</cals:th>
        <cals:th>applications</cals:th>
      </cals:tr>
      <cals:tr>
        <cals:td>1998</cals:td>
        <cals:td>15</cals:td>
      </cals:tr>
    </cals:table>
  </pr:slide>
</pr:presentation>
```

Verarbeitung von XML

- zwei Verfahren:
 - event-basiert
 - Baumstruktur
- Speicherplatz vs. freier Zugriff
- Geschwindigkeit vs. startup-Zeit
- stream-Fähigkeit

SAX – The Simple API for XML

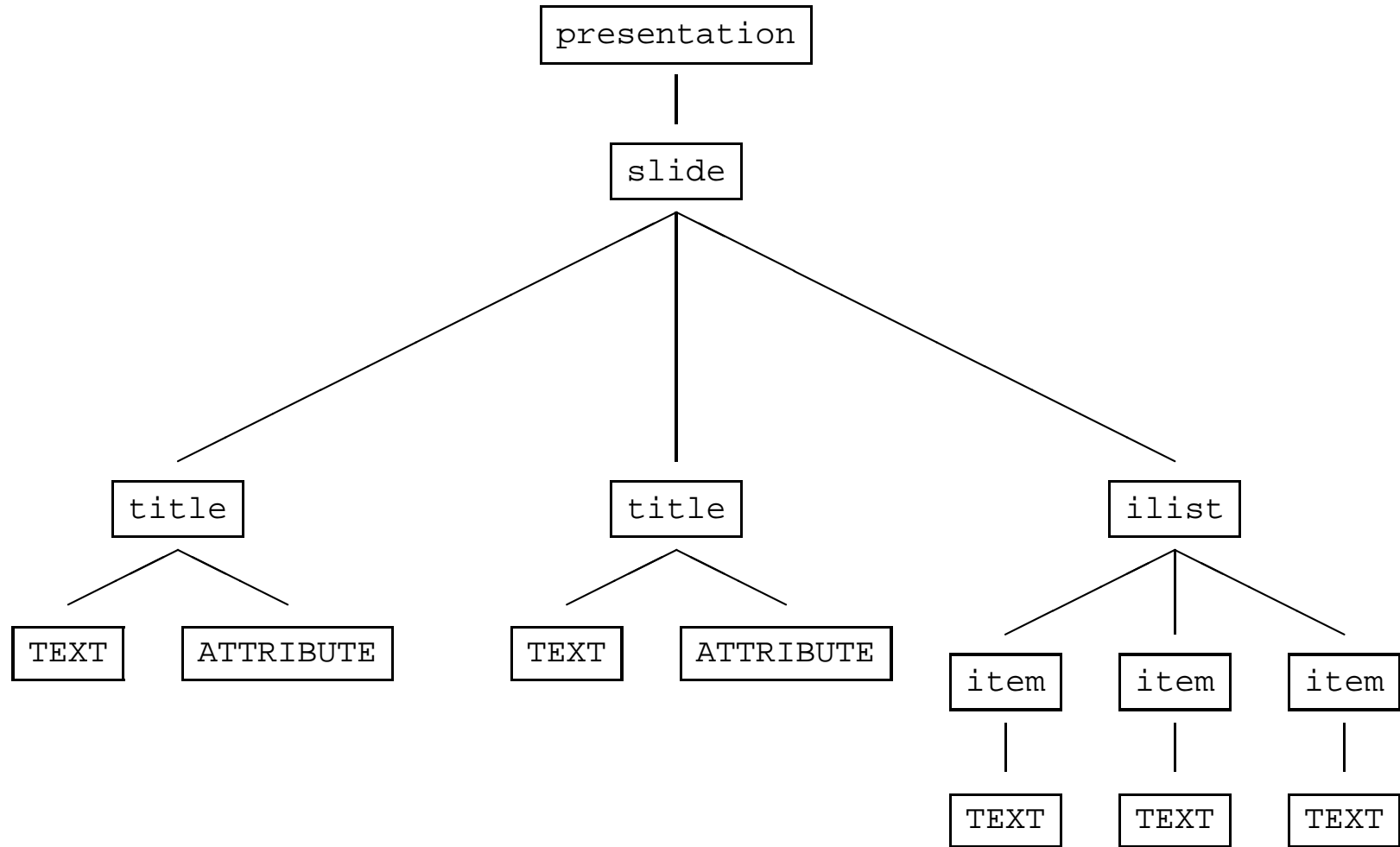
- events: u.a.
 - startDocument, endDocument
 - startElement, endElement
 - characters
- z.B. einfache Übersetzung nach \LaTeX :

<code><slide></code>	<code>→</code>	<code>\begin{slide}</code>
<code></slide></code>	<code>→</code>	<code>\end{slide}</code>
<code><item></code>	<code>→</code>	<code>\item</code>
<code></item></code>	<code>→</code>	
<code><bold></code>	<code>→</code>	<code>\textbf{</code>
<code></bold></code>	<code>→</code>	<code>}</code>

DOM – Document Object Model

- Dokument als Baum
- Elemente, Attribute, Text usw. sind die Knoten
- Klassen: Node, NodeList, Element, Text
- Methoden: `childNodes`, `nextSibling`, `parentNode`, `getElementsByTagName`

das Beispiel als Baum



XSL – Extensible Stylesheet Language

- XML beschreibt Syntax
- Formatierung? Semantik?
- XSL besteht aus mehreren Standards:
 - Zugriff auf bestimmte Knoten (ähnlich DOM)
 - Transformation von XML nach XML
 - XML-Sprache für Layout-Anweisungen

XPath

- enthält
 - Pfadbeschreibung, Bedingungen
 - Funktionen
- *Achsen*: child, parent, sibling, attribute
- Verwendung:

<code>slide</code>	slide-Kind des aktuellen Knotens
<code>slide/title</code>	title-Element unter <code>slide</code>
<code>slide/title/@use</code>	use-Attribut
<code>//title[@use="internal"]</code>	bestimmte title-Elemente

Beispiele für XPath

```
juser@hobel> xpath slides.xml '/presentation/slide/title'  
<title use="internal">XML & Friends for Dummies</title>  
<title use="external">Doing B2B with XML & Friends</title>
```

```
juser@hobel> xpath slides.xml '/presentation/slide'  
<slide>  
  <title use="internal">XML & Friends for Dummies</title>  
  <title use="external">Doing B2B with XML & Friends</title>
```

```
juser@hobel> xpath slides.xml '//title/@use'  
use="internal"  
use="external"
```

```
juser@hobel> xpath slides.xml '//title[@use="internal"]'  
<title use="internal">XML & Friends for Dummies</title>
```

```
juser@hobel> xpath slides.xml '//item[position()=2]'  
<item>XML is a meta language for  
  defining markup languages</item>
```

XSLT – XSL Transformations

- Umformung von XML nach
 - XML: Extraktion von Daten, Vorstufe der Formatierung
 - HTML: Visualisierung
- *style sheet* definiert *templates*
- Verarbeitung durch XSLT-Prozessor
- diverse Implementierungen

das Beispiel in HTML

```
<xsl:template match="/presentation">
  <html>
    <body>
      <xsl:apply-templates/>
    </body>
  </html>
</xsl:template>
```

```
<xsl:template match="slide">
  <xsl:apply-templates select="title[@use='external']"/>
  <xsl:apply-templates select="ilist"/>
</xsl:template>
```

```
<xsl:template match="title">
  <h1>
    <xsl:value-of select="."/>
  </h1>
</xsl:template>
```

XSL Formatting Objects

- XML-Sprache zur Beschreibung von Formatierung
- sehr komplex, aber universell
- FO wird durch XSLT-style-sheet erzeugt
- FO wird durch FO-Engine in Ausgabeformat umgewandelt
- FOP (Apache) erzeugt PDF, FrameMaker, PCL, ASCII, Preview

XML-Ausgabe von blast in HTML

blast of dm|gj|7290022|gb|AAF45489.1|CG13377 gene product against ./vert

Sequence	Score (bits)	E Value
emb RNDBHYDEH Sprague-Dawley D-beta-hydroxybutyrate dehydrogenase mRNA, completecds.	88	4.50265e-17
emb MM17BDTII M.musculus mRNA for 17-beta-hydroxysteroid dehydrogenase type II	77	7.94792e-14
emb MM17BDTII M.musculus mRNA for 17-beta-hydroxysteroid dehydrogenase type II	77	7.94792e-14
emb AF030513 Mus musculus cis-retinol androgen dehydrogenase 1 mRNA, completecds.	59	2.23983e-08
emb RN17BDI1 R.norvegicus mRNA for 17-beta-hydroxysteroid dehydrogenase type 1	48	5.16365e-05
emb RNBDHSDH R.norvegicus mRNA for 17-beta-hydroxysteroid dehydrogenase type I	48	5.16365e-05
emb RN33501 Rattus norvegicus retinol dehydrogenase type III mRNA, completecds.	47	0.000115034
emb MM17BDDEH M.musculus mRNA for 17-beta-hydroxysteroid dehydrogenase type 1	45	0.000256269
emb RN18762 Rattus norvegicus liver microsomal retinol dehydrogenase type Ia mRNA, complete cds.	45	0.000437129
emb AF100930 Oncorhynchus mykiss carbonil reductase/20beta-hydroxysteroiddehydrogenase B mRNA, complete cds.	34	0.590798
emb AF100931 Oncorhynchus mykiss carbonil reductase/20beta-hydroxysteroiddehydrogenase A mRNA, complete cds.	34	0.590798

emb|RNDBHYDEH Sprague-Dawley D-beta-hydroxybutyrate dehydrogenase mRNA, completecds.

Length = 1420
Score = 88 bits (216), Expect = 4.50265e-17
Identities = 63/258 (24%), Positives = 122/258 (47%), Gaps = 3/258 (1%)
SANADSHPSRWLITSADITLGLQLCTHLANKGYRVFAG--MKEAQDSLPAKLLGQMKIREYSEEP IAGTIIPMLDVTREDVLRREATVIIGANLNADERGIARVINTSGSVFRGQVESQNVQQHEMLRTN ILGTL
++ AD+ + VL+T D+ G L HL +KG+ VFAG +KE D+ +RE + + + ++L+V + + +A + + L E+G+ ++N +G G+VE +++ + + N+ GT+R K+F+ LR +GR++ + + G R + + + V+ ++ L
TSQADAASGKAVLVTGCDSSGFGFLAKHLHKGFLVFAAGLLKEQGD-----GVREL-DSLKSDLRIT IQLNVONSEVEKAVETVRSGLKDFEKGMBLWNAGISTFGVEFTSMETYKEVREVNILGTV

emb|MM17BDTII M.musculus mRNA for 17-beta-hydroxysteroid dehydrogenase type II

Length = 1342