

# **Publishing on the Web, today and tomorrow**

**Michel Goossens, CERN**

JACoW99

- current problems with HTML;
- HTML 4, a step forward;
- the extensible markup language (XML) effort in general;
- a closer look at the XML and XSL syntax;
- math on the Web: images and Java today, MathML tomorrow;
- conclusion.

## Why HTML has problems

- *Invalid HTML*: applications produce invalid HTML or introduce vendor-specific extensions; browsers do not reject invalid HTML.
- *Broken links*: problem of deleted or moved Web pages (use URL).
- *Fixed grammar*: HTML is a *fixed* SGML DTD (language syntax); no *standard* way to extend, adapt, renew the language.
- *Limited support for meta-data*: what about keywords, search-engines, re-use of the same source file.
- *Absence of structural tags*: HTML tags control mainly appearance not structure; navigation difficult.
- *Data exchange difficulties*: HTML aimed at presentation (Latin1 codes); difficult for data mining, internationalisation, exchange.
- *Absence of modern features*: Web changes rapidly; HTML must adapt quickly to new technology; XHTML and modules.

## HTML, the way ahead

- *Extend HTML*: HTML 4 (4.01) recommendations (Dec. 1997, Nov. 1999), HTML 4. is end of the road (May 1998).
- *Extensible Markup Language (XML)*: XML and its tools offer better application-specificity and data organisation.
- *XHTML*: Modular XML application of HTML: version 1.1 (April 2000), and version 2.0, replacing Images and Anchors with XLink (July 2000). XHTMLBasic for thin clients (April 2000).
- *Cascading Style Sheets*: CSS1 and CSS2 recommendations: decouple presentation and content.
- *Document Object Model (DOM)*: programmatic access to HTML (XML) elements as object data, exposing properties and methods.

**Problem:** Almost no browsers support XML directly at present.

## Extensible Markup Languages

- Since mid 1996 W3C SGML (now renamed XML) Working Group had been developing a “system for defining, validating, and sharing document formats on the Web”.
- XML 1.0 became a W3C Recommendation on 10 February 1998.
- It is the first in a series of standards.
- Others are XSL, XPointer, Xlink, MathML.
- Much other work is going on:  
Chemical Markup Language (CML), Object Data Model (DOM), Music Markup Language (MusicML), Scalable Vector Graphics (SVG), Resource Description Framework (RDF), Weather Observation Markup Format (OMF), Many Ecommerce initiatives, etc.

## What is XML?

- light-weight subset of SGML (ISO/IEC 8879:1986);
- meta-language to described logical structure of document;
- well integrated with the Internet;
- internationalisation built in from the start since basic languages is Unicode (UTF-16 and UTF-8);
- easy to learn, modify and implement;
- logical relations between elements described in *Document Type Definition* (DTD);
- documents can be *correct* (checked w.r.t. DTD);
- documents can be *well-formed* (no DTD needed, nesting correct).

## Structure of an XML document

- simplest complete document (*well-formed*);

```
<coolxml>XML is a cool idea!</coolxml>
```

- the same document in its *valid* form:

```
<?xml version="1.0" standalone="yes"?> <!-- XML PI -->
```

```
<!DOCTYPE coolxml [
```

```
  <!ELEMENT coolxml (#PCDATA)>
```

```
<coolxml>XML is a cool idea!</coolxml>
```

It has three parts:

1. XML processing instruction (version, encoding, stand-alone);
2. document type declaration (internal and external subsets);
3. document instance.

## The Document Type Definition (DTD)

- defines grammar of our *little* language;
- uses its own syntax (*XML-data*, Xschema, DCD propose XML format);
- *elements*: visual structural markup components;

```
<!ELEMENT pome      (prambule, corps)>
```

```
<!ELEMENT prambule (titre, recueil?, auteur, date?)>
```

```
<!ELEMENT titre     (#PCDATA)>
```

- *attributes*: properties of elements;

```
<!ATTLIST line      newpara (yes|no)  "no"
```

```
                indent  CDATA      "0mm"
```

```
                text    CDATA      #REQUIRED
```

```
                color   NMTOKEN   #FIXED "red">
```

- *entities*: internal and external foreign material.

## Internal entities

- defined in DTD, have no associated storage object, always parsed;
- abbreviations (*general entity*).

```
<!ENTITY MML "Mathematical Markup Language">
```

- special characters, accents or symbols (*general character entity*).

```
<!ENTITY gt CDATA "&#62;">
```

```
<!ENTITY cyrya CDATA "&#x044f;">
```

Predefined entities: lt (<), gt (>), amp (&), apos ('), and quot (").

- definition of variables in DTD (*parameter entity*).

```
<!ENTITY % list "UL | OL | DIR | MENU">
```

- definition must precede reference: &symb; or %list;;
- references can be used inside definitions:

```
<!ENTITY XMLS "&MML; and other extensible languages">
```



## External entities

- those that are not *internal*;
- reference data external to current document;
- data from external document can be parsed or declared NDATA (unparsed, e.g., bitmap image or binary file);
- include external file with "SYSTEM" keyword and followed by URL (*Universal Resource Identifier*).

```
<!ENTITY article SYSTEM "/usr/goossens/articles/xmlart.xml">
```

Contents of XML source file included (after parsing) with  
&article;

- include external file with "PUBLIC" keyword and followed by public identifier literal, and system literal (URL):

```
<!ENTITY % html4-strict PUBLIC "-//W3C//DTD HTML 4.0//EN"  
"http://www.w3.org/TR/REC-html40/strict.dtd">
```

The XML application first tries to build URL using public name -//W3C//DTD HTML 4.0//EN, e.g., via the catalog file proposed by OASIS (Organization for the Advancement of Structured Information Standards), otherwise uses explicit URL at end.

- non-parsable data (GIF or JPEG images, binary files) need the definition of a notation to handle data type

```
<!NOTATION GIF SYSTEM
```

```
    "c:\Program Files\Internet Explorer\Ie4.dll" >
```

A gif image is then declared with:

```
<!ENTITY xmlfig1 SYSTEM
```

```
    "http://www.myserver.edu/book-files/figures/xmlfig1"
```

```
    NDATA GIF >
```

and later included, for instance, with:

```
<IMG url="&xmlfig1;">
```

## Doing it with style: XSL

- Historically two approaches : CSS (“HTML-world”) *versus* DSSSL (“computer professionals”);
- CSS: from simple to more complex (CSS1, CSS2), mainly for HTML-based web applications;
- DSSSL-o: subset of DSSSL, allows for transformations between document types, can generate toc’s, indexes, page headers, . . .
- XSL (*Extensible Style Language*) WG set up on January 23rd 1998.
- Uses its own language (not XML), so that it can be used inside attributes.
- Contains XPath (addressing points in XML documents, common with XPointer), XSLT (transformation), and XLSF (formatting objects).
- XPath and XSLT became W3C recommendations on 16 November 1999.

## A small example file

### A simple source file

```
<!DOCTYPE invitation SYSTEM "invitation.dtd">
<invitation>
<front>
<to>Anna, Bernard, Didier, Johanna</to>
<date>Next Friday Evening at 8 pm</date>
<where>The Web Cafe</where>
<why>My first XML baby</why>
</front>
<body>
<par>I would like to invite you all to celebrate
the birth of <emph>Invitation</emph>, my
first XML document child.</par>
<par>Please do your best to come and join me next Friday
evening. And, do not forget to bring your friends.</par>
<par>I <emph>really</emph> look forward to see you soon!</par>
</body>
<back>
<signature>Michel</signature>
</back>
</invitation>
```

### ... and its DTD

```
<!-- invitation DTD -->
<!ELEMENT invitation (front, body, back) >
<!ELEMENT front      (to, date, where, why?) >
<!ELEMENT date       (#PCDATA) >
<!ELEMENT to         (#PCDATA) >
<!ELEMENT where      (#PCDATA) >
<!ELEMENT why        (#PCDATA) >
<!ELEMENT body       (par+) >
<!ELEMENT par        (#PCDATA|emph)* >
<!ELEMENT emph       (#PCDATA) >
<!ELEMENT back       (signature) >
<!ELEMENT signature  (#PCDATA) >
```

## Doing it with XSL

```

<?xml version='1.0'?>
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
                xmlns:fo="http://www.w3.org/1999/XSL/Format"
                default-space="">
    <xsl:variable name="PageMarginTop">75pt</xsl:variable>
    <xsl:variable name="PageMarginBottom">125pt</xsl:variable>
    <xsl:variable name="PageMarginLeft">80pt</xsl:variable>
    <xsl:variable name="PageMarginRight">150pt</xsl:variable>
    <xsl:variable name="BodySize">12pt</xsl:variable>

    <xsl:template match='/'>
        <fo:root xmlns:fo="http://www.w3.org/XSL/Format/1.0">
            <fo:layout-master-set>
                <fo:simple-page-master
                    page-master-name="allpages"
                    margin-top="{PageMarginTop}"
                    margin-bottom="{PageMarginBottom}"
                    margin-left="{PageMarginLeft}"
                    margin-right="{PageMarginRight}"
                    <fo:region-body margin-bottom="100pt"/>
                    <fo:region-after extent="25pt"/>
                </fo:simple-page-master>
            </fo:layout-master-set>
            <fo:page-sequence>
                <fo:sequence-specification>
                    <fo:sequence-specifier-repeating
                        page-master-first="allpages"
                        page-master-repeating="allpages"/>
                </fo:sequence-specification>
                <fo:flow font-family="serif">
                    <xsl:apply-templates/>
                </fo:flow>
            </fo:page-sequence>
        </fo:root>
    </xsl:template>

    <xsl:template match="invitation/front">
        <fo:block font-family="sans-serif" font-size="24pt"
            text-align-last="centered" font-weight="bold"
            space-after.optimum="24pt">
            <xsl:text>INVITATION</xsl:text>
        </fo:block>
        ...
    <xsl:template match="invitation/back">
        <fo:block space-before.optimum="12pt"
            font-weight="bold" text-align-last="end">
            <xsl:text>From: </xsl:text>
            <xsl:value-of select="signature"/>
        </fo:block>
    </xsl:template>
</xsl:stylesheet>

```

## Result with FOP and PassiveTeX

XML  $\xrightarrow{\text{XSL}}$  FO  $\xrightarrow{\text{FOP}}$  PDF

### INVITATION

*To:* Anna, Bernard, Didier, Johanna  
*When:* Next Friday Evening at 8 pm  
*Venue:* The Web Cafe  
*Occasion:* My first XML baby

I would like to invite you all to celebrate the birth of *Invitation*, my first XML document child.

Please do your best to come and join me next Friday evening. And, do not forget to bring your friends.

I *really* look forward to see you soon!

**From: Michel**

XML  $\xrightarrow{\text{XSL}}$  FO  $\xrightarrow{\text{PassiveTeX}}$  PDF

### INVITATION

*To:* Anna, Bernard, Didier, Johanna  
*When:* Next Friday Evening at 8 pm  
*Venue:* The Web Cafe  
*Occasion:* My first XML baby

I would like to invite you all to celebrate the birth of *Invitation*, my first XML document child.

Please do your best to come and join me next Friday evening. And, do not forget to bring your friends.

I *really* look forward to see you soon!

**From: Michel**

CERN

## MathML: An XML vocabulary for math

*MathML will make the Web even better for educational, scientific and technical materials. It also has the potential to make mathematics accessible to those with visual disabilities. It will allow mathematical content to be reused and exchanged with technical computing systems for further manipulation.*

W3C Director Tim Berners-Lee

7 April, 1998

(release of MathML as a W3C Recommendation).

MathML contains elements of three categories: presentation, content, and interface elements.

## Aims of MathML

- Encode mathematical material suitable for teaching and scientific communication at all levels;
- encode both mathematical notation and mathematical meaning;
- facilitate conversion to and from other math formats, both presentational and semantic. Output formats should include: graphical displays, speech synthesizers, computer algebra, other math layout languages ( $\text{T}_{\text{E}}\text{X}$ ), plain text displays (VT100 emulators), print media, braille. Such conversions may entail loss of information.
- allow information intended for specific renderers and applications;
- support efficient browsing for lengthy expressions;
- provide for extensibility;
- be well suited to template and other math editing techniques;
- be human legible, and simple for software to generate and process.



## Presentation and content markup

- **Presentation markup**
  - describe mathematical notation structure;
  - 28 elements, more than 50 attributes;
  - generally, each presentation element corresponds to a single kind of “layout schema” (two-dimensional notational device, such as row, superscript, underscript);
  - variants are expressed via attributes.
- **Content markup**
  - directly describe mathematical objects;
  - about 75 elements and about a dozen attributes;
  - Tokens, basic content, arithmetic, algebra, logic, relations, calculus, set theory, sequences, series, trigonometry, statistics, linear algebra, semantic mappings.

## A first MathML Example

$$x^2 + 4x + 4 = 0$$

### Presentation Markup

```
<mrow><!-- horizontal alignment -->
  <mrow>
    <msup>      <!-- superscript-->
      <mi>x</mi><!-- identifier -->
      <mn>2</mn><!-- number -->
    </msup>
    <mo>+</mo> <!-- operator -->
    <mrow>
      <mn>4</mn>
      <mo>&InvisibleTimes;</mo>
      <mi>x</mi>
    </mrow>
    <mo>+</mo>
    <mn>4</mn>
  </mrow>
  <mo>=</mo>
  <mn>0</mn>
</mrow>
```

### Content Markup

```
<reln>
  <eq/>
  <apply>
    <plus/>
    <apply>
      <power/>
      <ci>x</ci>
      <cn>2</cn>
    </apply>
    <apply>
      <times/>
      <cn>4</cn>
      <ci>x</ci>
    </apply>
    <cn>4</cn>
  </apply>
  <cn>0</cn>
</reln>
```

## Another MathML example

$$A = \begin{bmatrix} x & y \\ z & w \end{bmatrix}$$

### Presentation Markup

```

<mrow>
  <mi>A</mi>
  <mo>=</mo>
  <mfenced open="[" close="]">
    <mtable><!-- table or matrix -->
      <mtr> <!-- table row      -->
        <mtd><mi>x</mi></mtd><-- table -->
        <mtd><mi>y</mi></mtd><-- entry -->
      </mtr>
      <mtr>
        <mtd><mi>z</mi></mtd>
        <mtd><mi>w</mi></mtd>
      </mtr>
    </mtable>
  </mfenced>
</mrow>

```

### Content Markup

```

<reln>
  <eq/>
  <ci>A</ci>
  <matrix>
    <matrixrow>
      <ci>x</ci>
      <ci>y</ci>
    </matrixrow>
    <matrixrow>
      <ci>z</ci>
      <ci>w</ci>
    </matrixrow>
  </matrix>
</reln>

```

## Document strategies for the Web

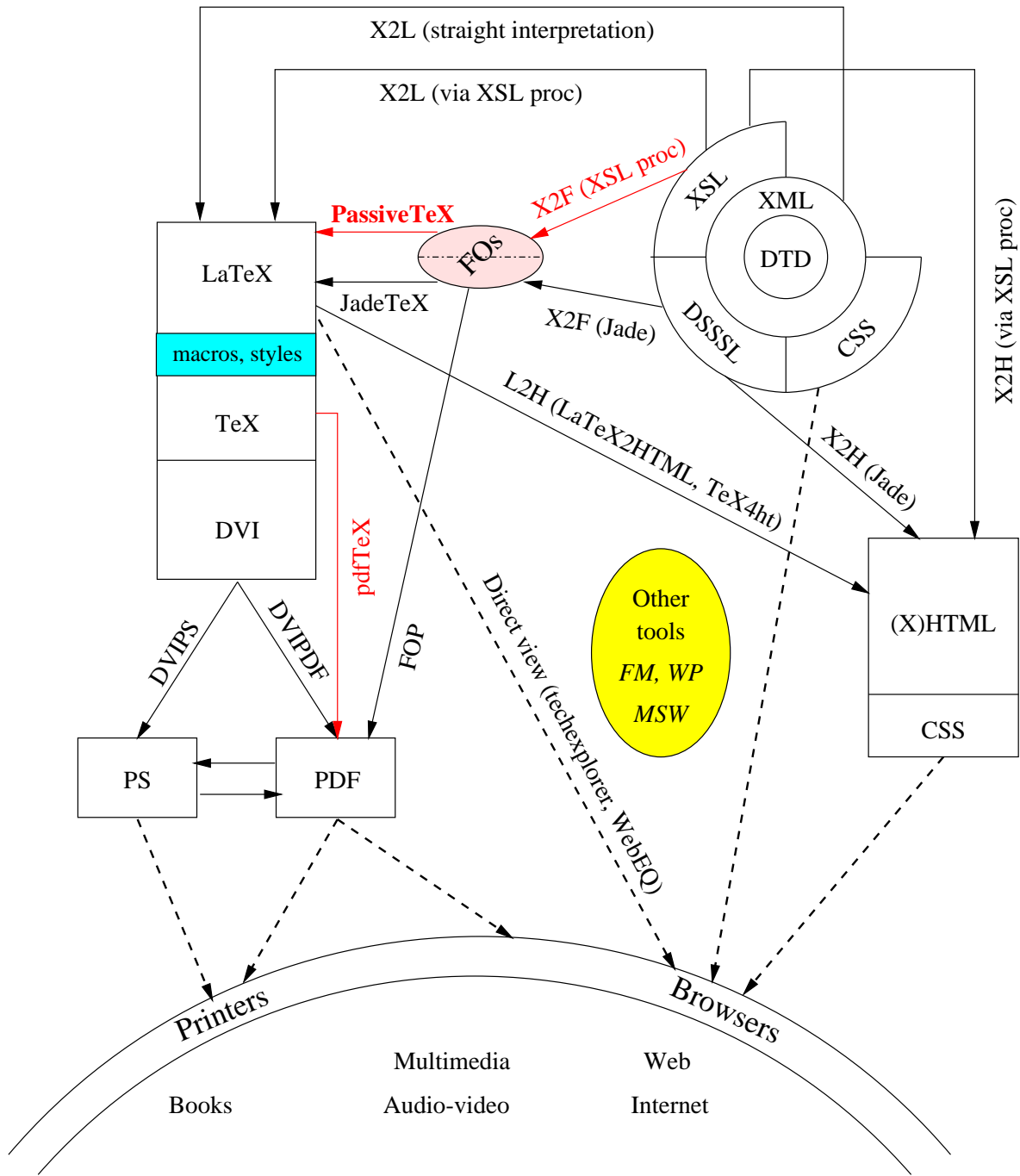
- Millions of documents exist in electronic form. *Reuse* of the large investment in this knowledge base is an important point.
- *New* documents can exploit the advantages of the latest technologies and adapt readily to the target audience (e-commerce, scientific, dictionaries, etc.).
- *Several presentations* of the same source are often a must, and *ad hoc* and different techniques have to be used to optimally exploit the display possibilities of the medium (browser, print, audio-video).
- They exist a *variety* of application domains. Scientific articles, financial data, or ancient Greek poetry need different tools and target different audiences.

## XML as central source repository

- The developers of XML have taken into account the lessons learned in the last decade using SGML and HTML.
- By construction XML is (should be) an ideal tool for dealing with most kinds of data and (multi-lingual) source documents (based on Unicode, seamless integration with modern languages, such as Java, perl, and python).
- Has gotten support from all corners of the Internet world: *Open Source* people as well as commercial players.
- Many free (and not-so-free) tools are available for all conceivable operating systems and purposes. Soon HTML-only browsers will only be a (bad) memory, and most Internet tools will support XML natively.

# XML documents and the Internet

## My (idealized) vision for the *not-too-far* future



CERN

## Batch translation: LaTeX2HTML

- Written in Perl;
- uses  $\text{\LaTeX}$ , dvips, ghostscript, netpbm utilities for the conversion of EPS images or math formulae into GIF or PNG;
- different hierarchical levels of source document put into different HTML files (under user control);
- table of contents, navigational aids;
- lot of packages supported.

## Output example LaTeX2HTML

### Math examples

$$\phi(\lambda) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \exp(u \ln u + \lambda u) du \quad \text{for } c \geq 0 \quad (1)$$

$$\lambda = \frac{\epsilon - \bar{\epsilon}}{\xi} - \gamma' - \beta^2 - \ln \frac{\xi}{E_{\max}} \quad (2)$$

$$\gamma = 0.577215 \dots \quad (\text{Euler's constant}) \quad (3)$$

$$\gamma' = 0.422784 \dots = 1 - \gamma \quad (4)$$

$$\epsilon, \bar{\epsilon} = \text{actual/average energy loss} \quad (5)$$

Since (6) or (7d) should hold for arbitrary  $\delta \mathbf{c}$ -vectors, it is clear that  $\mathcal{N}(A) = \mathcal{R}(B)$  and that when  $y = \mathcal{F}(x)$  one has...

...the [Pythagorians](#) knew infinitely many solutions in integers to  $a^2 + b^2 = c^2$ . That no non-trivial integer solutions exist for  $a^N + b^N = c^N$  with integers  $N > 2$  has long been suspected ([Fermat, c.1637](#)). Only during the current decade has this been proved ([Wiles, 1995](#)).

$$V \pi^{sr} = \left\langle \sum_i M_i \mathbf{V}_i \mathbf{V}_i + \sum_i \sum_{j>i} \mathbf{R}_{ij} \mathbf{F}_{ij} \right\rangle \quad (6)$$

$$= \left\langle \sum_i M_i \mathbf{V}_i \mathbf{V}_i + \sum_i \sum_{j>i} \sum_{\alpha} \sum_{\beta} \mathbf{r}_{i\alpha j \beta} \mathbf{f}_{i\alpha j \beta} - \sum_i \sum_{\alpha} \mathbf{p}_{i\alpha} \mathbf{f}_{i\alpha} \right\rangle$$

CERN



## LaTeX2HTML – advantages

### Quality

- High-level typography and performant navigational aids;
- cross-references intact (figures, tables, sections, equations);
- equations and figures (GIF) very readable;
- large series of user options to control translation process (“global” or “fractioned” translation for formulae);
- consistency in the visual presentation is an important design feature; (nevertheless some imprecise placement in the elements of the sub-formulae remains—also with TeX4ht and tth).

## LaTeX2HTML – advantages (cont.)

### Reliability and lisibility

- Runs without human intervention on thousands of systems worldwide (of course Perl must be installed also!);
- has no problem to handle very large documents (hundreds of pages) in a reliable way;
- LaTeX2HTML generates “standard” HTML files (3.2, 4.0);
- the “ALT” textual representation of the content of images permits an acceptable presentation of the information on non-graphic browsers, such as Lynx.

## LaTeX2HTML – drawbacks

### Execution time

- Can take several minutes for a few tens of pages;
- specially formulae and images take a lot of time;
- impossible to use for “on the fly” translation.

### Diskspace

- Can easily generate hundreds of HTML files and GIF images;
- can double or triple the diskspace needed to store a document;
- maintenance and management non-trivial (export, regeneration).

## LaTeX2HTML – Drawbacks (cont.)

### Time to load the Web pages

- Large number of images to be transferred makes loading Web pages take a long time.

### Installation

- perl, ghostscript, netpbm utilities, . . . must be installed (not so easy for a non-specialist);
- the complete L<sup>A</sup>T<sub>E</sub>X source document must be loaded into memory as a whole, so that a lot of system resources are needed during the translation process (can lead to problems on small systems, such as DOS).

## Batch translation *bis*: TeX4ht

- TeX4ht uses  $\text{T}_{\text{E}}\text{X}$  itself to transform  $\text{T}_{\text{E}}\text{X}$  commands into HTML tags using a dedicated  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$  package;
- does not need an external program, such as perl;
- uses the utilities of the ghostscript system to translate EPS files and mathematics formulae with special symbols;
- some additions in the  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$  sources might be necessary (to give some directives to the program).

## Fragment generated by TeX4ht

3 Vavilov theory - Microsoft Internet Explorer

File Edit View Favorites Tools Help

[\[next\]](#) [\[prev\]](#) [\[prev-tail\]](#) [\[tail\]](#) [\[up\]](#)

### 3 Vavilov theory

Vavilov<sup>5</sup> derived a more accurate straggling distribution by introducing the kinematic limit on the maximum transferable energy in a single collision, rather than using  $E_{\max} = \infty$ . Now we can write<sup>2</sup>:

$$\phi_v(\lambda_v, \kappa, \beta^2) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \phi(s) e^{\lambda_v s} ds \quad c \geq 0$$

$$f(\epsilon, \delta s) = \frac{1}{\xi} \phi_v(\lambda_v, \kappa, \beta^2) \quad \text{where} \quad \begin{aligned} \phi(s) &= \exp[\kappa(1 + \beta^2 \gamma)] \exp[\psi(s)], \\ \psi(s) &= s \ln \kappa + (s + \beta^2 \kappa) [\ln(s/\kappa) + E_1(s/\kappa)] - \kappa e^{-s/\kappa}, \text{ and} \end{aligned}$$

$$E_1(z) = \int_{\infty}^z t^{-1} e^{-t} dt \quad (\text{the exponential integral})$$

$$\lambda_v = \kappa \left[ \frac{\epsilon - \bar{\epsilon}}{\xi} - \gamma' - \beta^2 \right]$$

The Vavilov parameters are simply related to the Landau parameter by  $\lambda_z = \lambda_v / \kappa - \ln \kappa$ . It can be shown that as  $\kappa \rightarrow 0$ , the distribution of the variable  $\lambda_z$  approaches that of Landau. For  $\kappa \leq 0.01$  the two distributions are already practically identical. Contrary to what many textbooks report, the Vavilov distribution *does not* approximate the Landau distribution for small  $\kappa$ , but rather the distribution of  $\lambda_z$  defined above tends to the distribution of the true  $\lambda$  from the Landau density function. Thus the routine GVAVIV samples the variable  $\lambda_z$  rather than  $\lambda_v$ . For  $\kappa \geq 10$  the Vavilov distribution tends to a Gaussian distribution (see next section).

My Computer

## TeX4ht – Advantages and drawbacks

- Mostly the same problems and successes as LaTeX2HTML;
- TeX4ht uses L<sup>A</sup>T<sub>E</sub>X with a few packages (for instance TeX4ht.sty) and some font files, making the installation a lot easier (especially on non-Unix systems, which do not have Perl and the other tools pre-installed);
- still needs dvips and ghostscript;
- compared to LaTeX2HTML the quality can be less good (e.g., placement of the limits on the integrals);
- may need modification in the source files.

## On the fly conversion: tth

- tth is written in C and uses no external programs, thus making the program extremely portable;
- formulae translated into HTML;
- uses the Symbol font available on X-Window, Mac ou PC;
- by default, images are not converted, but they are made available through external links (it is also possible to ask that GIF Images be generated, but that makes the procedure a lot more complicated).



## Fragment generated by tth

Simulation of Energy Loss Straggling - Microsoft Internet Explorer

File Edit View Favorites Tools Help

### 3 Vavilov theory

Vavilov[5] derived a more accurate straggling distribution by introducing the kinematic limit on the maximum transferable energy in a single collision, rather than using  $E_{\max} = \infty$ . Now we can write [2]:

$$f(\varepsilon, \delta s) = \frac{1}{\xi} \phi_v(\lambda_v, \kappa, \beta^2)$$

where

$$\phi_v(\lambda_v, \kappa, \beta^2) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \phi(s) e^{\lambda s} ds \quad c \geq 0$$

$$\phi(s) = \exp[\kappa(1 + \beta^2 \gamma)] \exp[\psi(s)],$$

$$\psi(s) = s \ln \kappa + (s + \beta^2 \kappa) [\ln(s/\kappa) + E_1(s/\kappa)] - \kappa e^{-s/\kappa},$$

and

$$E_1(z) = \int_{\infty}^z t^{-1} e^{-t} dt \quad (\text{the exponential integral})$$

My Computer

**tth – advantages****Performance acceptable**

- Very fast (one pass over the data);
- source documents can be stored in a directory in compressed form (e.g., `.tar.gz`);
- the translation (`gtar zxvf`, then `tth`) can be performed very fast on the fly without any delay.

**Easy maintenance**

- Modifications must only be introduced in the sources which do not need to be retranslated.

**tth – drawbacks****Medium quality**

- Documents with only text or simple math come out not too bad;
- complete mess for complex formulae;
- (almost) no cross-references.

**Modifications in the working environment**

- The user's working visualisation environment must be adapted (X-window files, program setup), otherwise the documents will be unreadable.

## “Online” L<sup>A</sup>T<sub>E</sub>X: techexplorer

- Extension module (*plug-in*) pour Netscape (NS) and MS Internet Explorer (MSIE);
- (free) version (*Introductory Edition*) is available for W95/NT, SGI IRIS 6.2, Sun Solaris 2.5 and IBM AIX 4.1;
- techexplorer reads directly the L<sup>A</sup>T<sub>E</sub>X source document;
- does not (yet) handle all L<sup>A</sup>T<sub>E</sub>X commands (uninterpreted commands are displayed *as-is* in red in the browser);
- techexplorer offers a whole set of hypertext and multimedia extensions;
- attractive tool for preparing documents for and on the Web;
- initial support for MathML exists.

## Fragment displayed with techexplorer

**Vavilov theory**

Vavilov[bib-VAVI] derived a more accurate straggling distribution by introducing the kinematic limit on the maximum transferable energy in a single collision, rather than using  $E_{\max} = \infty$ . Now we can write[bib-SCH1]:

$$f(\epsilon, \delta s) = \frac{1}{\xi} \phi_v(\lambda_v, \kappa, \beta^2)$$

where

$$\phi_v(\lambda_v, \kappa, \beta^2) = \frac{1}{2\pi} \int_{c-i\infty}^{c+i\infty} \phi(s) e^{\lambda_v s} ds \quad c \geq 0$$

$$\phi(s) = \exp[\kappa(1 + \beta^2 \gamma)] \exp[\psi(s)],$$

$$\psi(s) = s \ln \kappa + (s + \beta^2 \kappa) [\ln(s/\kappa) + E_1(s/\kappa)] - \kappa e^{-s/\kappa},$$

and

$$E_1(z) = \int_0^z t^{-1} e^{-t} dt \quad (\text{the exponential integral})$$

$$\lambda_v = \kappa \left[ \frac{\epsilon - \bar{\epsilon}}{\xi} - \gamma' - \beta^2 \right]$$

The Vavilov parameters are simply related to the Landau parameter by  $\lambda_L = \lambda_v / \kappa - \ln \kappa$ . It can be shown that as  $\kappa \rightarrow 0$ , the distribution of the variable  $\lambda_L$  approaches that of Landau. For  $\kappa \leq 0.01$  the two distributions are already practically identical. Contrary to what many textbooks report, the Vavilov distribution *does not* approximate the Landau distribution for small  $\kappa$ , but rather the distribution of  $\lambda_L$  defined above tends to the distribution of the true  $\lambda$  from the Landau density function. Thus the routine GVAVIV samples the variable  $\lambda_L$  rather than  $\lambda_v$ . For  $\kappa \geq 10$  the Vavilov distribution tends to a Gaussian distribution (see next section).

## Equation browsing with Java and WebEQ

- Interprets WebTeX (subset of  $\text{\LaTeX}$ );
- formulae input in  $\text{\TeX}$ -syntax between  $\$. . . \$$  or  $\backslash[ . . . \backslash]$ ;
- special utility (the *Sizer*) translates them into Java applets that can be visualized with WebEQ;
- WebEQ has its own editor, which has support for direct MathML input;
- the WebEQ editor is an ideal teaching tool for learning about MathML since it allows you to translate (not too complicated) formulae from  $\text{\LaTeX}$  into MathML.

## Fragment generated by WebEQ

$$\phi_\nu(\lambda_\nu, \kappa, \beta^2) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \phi(s) e^{\lambda s} ds \quad c \geq 0$$

$$\phi(s) = \exp[\kappa(1 + \beta^2 \gamma)] \exp[\psi(s)],$$

$$\psi(s) = s \ln \kappa + (s + \beta^2 \kappa) [\ln(s/\kappa) + E_1(s/\kappa)] - \kappa e^{-s/\kappa},$$

and

$$E_1(z) = \int_z^\infty t^{-1} e^{-t} dt \quad (\text{the exponential integral})$$

$$\lambda_\nu = \kappa \left[ \frac{\epsilon - \bar{\epsilon}}{\xi} - \gamma' - \beta^2 \right]$$

The Vavilov parameters are simply related to the Landau parameter by  $\lambda_L = \lambda_\nu / \kappa - \ln$ . It can be shown that as  $\kappa \rightarrow \mathcal{C}$ , the distribution of the variable  $\lambda_L$  approaches that of Landau. For  $\kappa \leq 0.0$ : the two distributions

## T<sub>E</sub>X brings typography to XML

- Often typographic quality is a *must*.
- The *print* button of most present-day HTML browsers does not in general give high quality printable copy.
- *Special* and separate procedures are used to prepare the printable output for XML documents.
- With SGML and DSSSL James Clark's Jade (nowadays being extended by members of the DSSSL community as open source library under the name of Openjade, see <http://jade-cvs.avionitek.com/>).
- Jade reads DSSSL style sheets and provides several back ends (FOT, T<sub>E</sub>X, \*ML, RTF).
- Sebastian Rahtz' (David Megginson) `jadetex` generates a printable document (PostScript, PDF).



## What is Passive $\text{T}_{\text{E}}\text{X}$ ?

- Passive $\text{T}_{\text{E}}\text{X}$  is a library of  $\text{T}_{\text{E}}\text{X}$  macros which can be used to process an XML document resulting from an XSL transformation to formatting objects.
- Passive $\text{T}_{\text{E}}\text{X}$  provides a rapid development environment for experimenting with XSL FO, using a reliable pre-existing formatter.
- Running Passive $\text{T}_{\text{E}}\text{X}$  with the pdf $\text{T}_{\text{E}}\text{X}$  variant of  $\text{T}_{\text{E}}\text{X}$  generates high-quality PDF files in a single operation.
- Passive $\text{T}_{\text{E}}\text{X}$  shows how  $\text{T}_{\text{E}}\text{X}$  can remain the formatter of choice for XML, while hiding the details of its operation from the user.

Passive $\text{T}_{\text{E}}\text{X}$  is available at

<http://users.ox.ac.uk/~rahtz/passivetex/>.

## Passive $\text{T}_{\text{E}}\text{X}$ : history and components

Passive $\text{T}_{\text{E}}\text{X}$  derives from and builds on:

- `typehtml`, a  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$  package by David Carlisle, used to typeset HTML directly using  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ ;
- `jadetex`, a package by Sebastian Rahtz that implements the output of the Jade DSSSL processor's  $\text{T}_{\text{E}}\text{X}$  backend;
- A UTF-8 handler by David Carlisle, in conjunction with the catalogue of Unicode/ $\text{T}_{\text{E}}\text{X}$  mappings built up for `jadetex`.

The system components include macros to:

- parse XML input (attributes, entities, etc.), parse UTF-8 input, and MathML elements;
- instantiate formatting object elements and attributes;
- map Unicode to  $\text{T}_{\text{E}}\text{X}$  font layouts.

## Advantages and disadvantages of Passive $\text{\TeX}$

- ☺ rapid development;
  - ☺ well-understood, robust, stable, and freely available page formatter;
  - ☺ fonts, graphics inclusion, hyperlinks, etc. come for free;
  - ☺ mature handling of language issues, including hyphenation;
  - ☺ high-quality math rendering ( $\text{\TeX}$ 's *raison d'être*);
  - ☺ pdf $\text{\TeX}$  variant generates very high-quality PDF.
- 
- ☹ constraint to use  $\text{\TeX}$ 's page makeup model, and force XSL FO to fit it;
  - ☹ as  $\text{\LaTeX}$  is already high-level markup, it is too easy to allow things to fall through and take  $\text{\LaTeX}$  defaults;
  - ☹  $\text{\TeX}$  macro writing is obscure and difficult, so that the system is not transparent for most (non- $\text{\TeX}$ ) programmers;
  - ☹  $\text{\TeX}$  is large and monolithic (difficult to embed in other applications);
  - ☹  $\text{\TeX}$  seems much like a sledgehammer to crack a nut . . .

## Passive $\text{T}_{\text{E}}\text{X}$ : how it handles maths

Passive $\text{T}_{\text{E}}\text{X}$  supports MathML directly. An XSL style sheet can pass  $\langle\text{math}\rangle$  and its children through unchanged, as follows:

```
<xsl:template match="math">
  <xsl:apply-templates mode="math"/>
</xsl:template>

<xsl:template mode="math"
  match="*|@*|comment()|processing-instruction()|text()">
  <xsl:copy>
    <xsl:apply-templates mode="math"
      select="*|@*|processing-instruction()|text()"/>
  </xsl:copy>
</xsl:template>
```

A reasonable subset of presentation MathML is recognized, and produces good output. We show an example later.

## XSL, FO and T<sub>E</sub>X, cons and pros

- ☹ The XSL FO page model is inherited from DSSSL, and is unproven for production-quality print formatting.
- ☹ The XSL specification is unfinished and incomplete.
- ☹ One cannot easily tweak T<sub>E</sub>X's behaviour with this system.
- ☹ The table model of XSL is sufficiently far from T<sub>E</sub>X's that it may require a pre-processor.
  
- ☺ Together with FOP, we now have systems to experiment with, and commercial implementations cannot be far behind (can they?).
- ☺ T<sub>E</sub>X is close to being a XSL FO-capable formatter.
- ☺ With the Omega T<sub>E</sub>X variant (using Unicode internally), we have a native Unicode typesetting system ready and waiting.
- ☺ XSL FO does not threaten TeX — it gives it a reason to survive.

## Passive $\text{T}_{\text{E}}\text{X}$ , typesetting a scientific document

- original source is  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ ;
- translated to XML with  $\text{T}_{\text{E}}\text{X}4\text{ht}$  using  $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ -like *ad hoc* DTD and MathML;
- write an XSL style sheet to treat the “textual elements” of the document;
- “pass through” the math components to the back-end; (XSL, and DSSSL up to a point, do not have a sufficient set of FOs to do math correctly, and although the X3C Math WG is talking to the XSL WG, at present the present consensus seems to be that it is best to treat the MathML directly at the end application level);
- interpret MathML directly in Passive $\text{T}_{\text{E}}\text{X}$ .

$$\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X} \xrightarrow{\text{T}_{\text{E}}\text{X}4\text{ht}} \text{XML} \xrightarrow{\text{XSL}} \text{FO and MathML} \xrightarrow{\text{PassiveT}_{\text{E}}\text{X}} \text{PDF}$$

## The generated MathML source

```
<section id="vavref">
<stitle>Vavilov theory</stitle>

<par>Vavilov<cite refid="bib-VAVI"/> derived a more accurate
straggling distribution by introducing the kinematic limit on the
maximum transferable energy in a single collision, rather than using
<inlinemath><math><msub><mi>E</mi><mrow><mtext>max</mtext></mrow></msub>
<mo>=</mo><mi>&infin;</mi></math></inlinemath>.
```

Now we can write<cite refid="bib-SCH1"/>:

```
<eqnarray ><subeqn><math><mi>f</mi> <mfenced open='(' close=')'>
<mi>&epsi;</mi><mo>,</mo><mi>&delta;</mi><mi>s</mi></mfenced>
<mo>=</mo> <mfrac><mrow><mn>1</mn></mrow>
<mrow><mi>&xi;</mi></mrow>
</mfrac><msub><mi>&phi;</mi><mrow><mi>v</mi></mrow></msub>
<mfenced open='(' close=')'>
<msub><mi>&lambda;</mi><mrow><mi>v</mi></mrow></msub><mo>,</mo>
<mi>&kappa;</mi><mo>,</mo><msup><mi>&beta;</mi><mrow><mn>2</mn></mrow>
</msup></mfenced></math></subeqn></eqnarray>
where
<eqnarray><subeqn><math><msub><mi>&phi;</mi><mrow><mi>v</mi></mrow></msub>
```

```

<mfenced open='(' close=')'>
<msub><mi>&lambda;</mi><mrow><mi>v</mi></mrow></msub><mo>,</mo>
<mi>&kappa;</mi><mo>,</mo>
<msup><mi>&beta;</mi><mrow><mn>2</mn></mrow></msup></mfenced>
  <mo>=</mo>
<mfrac><mrow><mn>1</mn></mrow>
  <mrow><mn>2</mn><mi>&pi;</mi><mi>i</mi></mrow>
</mfrac>
<msubsup><mo>&int;</mo>
<mrow><mi>c</mi><mo>+</mo><mi>i</mi><mi>&infin;</mi></mrow>
<mrow><mi>c</mi><mo>-</mo><mi>i</mi><mi>&infin;</mi></mrow></msubsup>
<mi>&phi;</mi><mfenced open='(' close=')'><mi>s</mi></mfenced>
<msup><mi>e</mi><mrow><mi>&lambda;</mi><mi>s</mi></mrow></msup>
<mi>d</mi><mi>s</mi><mspace width='2cm' /><mi>c</mi><mo>&geq;</mo><mn>0</mn>
  </math></subeqn>

<subeqn><math><mi>&phi;</mi><mfenced open='(' close=')'><mi>s</mi></mfenced>
<mo>=</mo><mo>exp</mo><mfenced open='[' close=']'><mi>&kappa;</mi>
<mrow><mo>(</mo><mn>1</mn><mo>+</mo><msup><mi>&beta;</mi>
  <mrow><mn>2</mn></mrow></msup><mi>&gamma;</mi><mo>)</mo></mrow>
</mfenced><mo>exp</mo><mfenced open='[' close=']'><mi>&psi;</mi>
<mfenced open='(' close=')'><mi>s</mi></mfenced></mfenced>
<mo>,</mo> </math></subeqn>

```



# MathML viewed with Amaya

The screenshot shows the Amaya web browser window titled 'a.html'. The address bar contains the URL 'fafs/cern.ch/user/g/goossens/passivetex/latexexa/a.html' and the title bar shows 'Simulation of Energy Loss Straggling'. The main content area displays the following text and mathematical formulas:

**Urbán model**

The method for computing restricted energy losses with  $\delta$ -ray production above given threshold energy in GEANT is a Monte Carlo method that can be used for thin layers. It is fast and it can be used for any thickness of a medium. Approaching the limit of the validity of Landau's theory, the loss distribution approaches smoothly the Landau form as shown in Figure .

Energy loss distribution for a 3 GeV electron in Argon as given by standard GEANT. The width of the layers is given in centimeters.

It is assumed that the atoms have only two energy levels with binding energy  $E_1$  and  $E_2$ . The particle--atom interaction will then be an excitation with energy loss  $E_1$  or  $E_2$ , or an ionisation with an energy loss distributed according to a function  $g(E) \sim 1/E^2$ :

$$g(E) = \frac{(E_{\max} + I)I}{E_{\max} E^2}$$

The macroscopic cross-section for excitations ( $i = 1, 2$ ) is

$$\Sigma_i = C \frac{f_i \ln(2m\beta^2 \sqrt{E_i}) - \beta^2}{E_i \ln(2m\beta^2 \sqrt{I}) - \beta^2} (1 - r)$$

and the macroscopic cross-section for ionisation is

$$\Sigma_3 = C \frac{E_{\max}}{I(E_{\max} + I) \ln\left(\frac{E_{\max} + I}{I}\right)} r$$

$E_{\max}$  is the GEANT cut for  $\delta$ -production, or the maximum energy transfer minus mean ionisation energy, if it is smaller than this cut-off value. The following notation is used:

- $r, C$  parameters of the model
- $E_i$  atomic energy levels
- $I$  mean ionisation energy
- $f_i$  oscillator strengths

The model has the parameters  $f_i, E_i, C$  and  $r$  ( $0 \leq r \leq 1$ ). The oscillator strengths  $f_i$  and the atomic level energies  $E_i$  should satisfy the constraints

MATH \ P \ BODY \ HTML

thus implying

$$\text{mean } \sigma^2 = \frac{\xi^2 (1 - \beta^2/2)}{\kappa}$$

## 5 Urbán model

The method for computing restricted energy losses with  $\delta$ -ray production above given threshold energy in GEANT is a Monte Carlo method that can be used for thin layers. It is fast and it can be used for any thickness of a medium. Approaching the limit of the validity of Landau's theory, the loss distribution approaches smoothly the Landau form as shown in Figure 2.

It is assumed that the atoms have only two energy levels with binding energy  $E_1$  and  $E_2$ . The particle-atom interaction will then be an excitation with energy loss  $E_1$  or  $E_2$ , or an ionisation with an energy loss distributed according to a function  $g(E) \sim 1/E^2$ :

$$g(E) = \frac{(E_{\text{max}} + I)I}{E_{\text{max}}} \frac{1}{E^2} \tag{1}$$

The macroscopic cross-section for excitations ( $i = 1, 2$ ) is

$$\Sigma_i = C \frac{f_i}{E_i} \frac{\ln(2m\beta^2\gamma^2/E_i) - \beta^2}{\ln(2m\beta^2\gamma^2/I) - \beta^2} (1 - r) \tag{2}$$

## General conclusion

- XML allows you to construct your own markup language, optimized for the task at hand and using tags and attribute names in the user's native language and an encoding adapted to the local computer system.
- A lot of tools are freely (and not so freely) available.
- All major Internet players have adopted the standard.
- “Standard” DTDs are already being prepared by various scientific, commercial, etc., communities, and this will improve communication and interoperability
- XSL is a good companion tool for styling and transforming XML sources. It is well integrated with the Web and builds on DSSSL and CSS (formatting objects, SVG,...).
- In “our” area, preparing scientific documents for the Web is a non trivial task, especially if they contain mathematical information that we want to “share”.

- We have looked at a few “non-optimal” solutions: PDF/PostScript (Java, plugins) or GIF images.
- We think/hope that XML, MathML and companion specifications (SVG, XML Schema, etc.) will allow us to handling scientific documents on the Web gracefully.
- T<sub>E</sub>X is a reliable output engine that can with rather minimal work provide typographically excellent output, especially for typesetting scientific XML documents containing a lot of maths;
- With well-known and “standard” DTDs (TEI, DocBook, ISO12083, etc.) XML can be used as a *lingua franca* to transport documents between various editing and document handling systems. At CERN we plan to test portability between XML, L<sup>A</sup>T<sub>E</sub>X, and FrameMaker/SGML; later also Microsoft Word, Wordperfect, etc. when these applications will have better support for XML and XSL.
- An *Esprit* project TIPS (*Tools for Innovative Publishing in Science*) will look at possible solutions.