

HTML to \LaTeX transformation

Frederik R. N. Schlupkothen

Abstract

\LaTeX was created as an authoring language that enables authors to keep typesetting quality standards while preparing their printed matter, assuming that the output context is known from the very beginning of the writing process. As a counter-concept, XML is focused on keeping output flexible, providing mechanisms to manage and control the logical structure of documents. Combining the strengths of both ecosystems has been discussed frequently in the past. This article aims to contribute to this discussion by introducing a mapping from HTML to \LaTeX , the two most widespread document description languages in their respective fields of application.

1 Introduction

The Extensible Markup Language (XML) has become the native markup language for structured information in (distributed) document processing architectures. It provides a core syntax that has been adopted by many document description languages. A broad software ecosystem and further standards have emerged to realize and facilitate the processing of XML-based documents. Among them, the Extensible Stylesheet Language Transformations (XSLT) described in [8] offers a standardized mechanism to translate different XML-based languages into one another. This has made XML the language of choice for cross-media publishing workflows.

However, while the XML processing is fully covered by viable tools, the final document production of printed matter from XML may still be a problem. A potential solution is to integrate the \TeX typesetting engine (described in [10]) in XML-based processing chains. For this task, \TeX XML, an XML representation of \TeX commands, was introduced in [12] and its “production proof” implementation discussed in [14]. As any other XML language, \TeX XML documents can be produced via XSLT. So the last gap to fill in \TeX XML-based workflows is to define XSL stylesheets that realize the transformation between specific XML and \TeX document description languages (see [14]).

Here we introduce a mapping between two document description languages that are well known in their respective fields of application: The HyperText Markup Language (HTML), as the transformation’s source language, is the core language of the World Wide Web and has a history that is closely tied to XML. It offers layout-oriented markup semantics primarily for textual content and is used amongst others

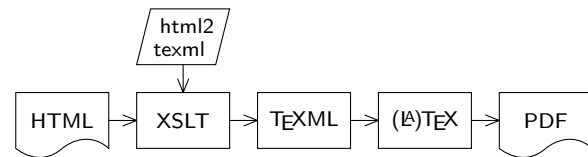


Figure 1: HTML to PDF processing

to describe web pages (see [7]), electronic books (see [4]), and printed matter (see e.g. [9, 13]). \LaTeX , as the transformation’s output, is the abstraction of \TeX ’s typesetting commands to logical markup.

Figure 1 shows the XML-based production workflow that produces the Portable Document Format (PDF) from HTML by the use of the \TeX typesetting engine. The particular processing steps are described via an example of emphasized text:

1. The HTML source document describes the emphasized text “dolor” as: `dolor`
2. The XSLT processor queries the stylesheet “html2texml” for a matching transformation template. The matching template defines the transformation of the HTML’s “em” element to the corresponding \LaTeX command in its \TeX XML representation:

```

<xsl:template match="html:em">
  <tex:cmd name="emph">
    <tex:param>
      <xsl:apply-templates />
    </tex:param>
  </tex:cmd>
</xsl:template>

```

3. The result of the XSLT transformation is the following \TeX XML element:

```

<cmd name="emph">
  <param>dolor</param>
</cmd>

```
4. The \TeX XML processor converts the \TeX XML element to a \LaTeX command: `\emph{dolor}`
5. The \LaTeX processor typesets: *dolor*

The prior example shows that defining the transformation between two languages needs insight into the differences in conceptual syntax and semantic coverage of source and destination language. While HTML is native to the XML syntax, \TeX XML is reproducing the syntax logic inherited from \LaTeX . While \LaTeX is native to printed matter, HTML was initially designed for electronic resources. Section 2 introduces the underlying concepts of the HTML and \LaTeX markup syntax. Section 3 introduces the mapping between HTML and \LaTeX markup semantics. Finally, section 4 shows a complete document with its corresponding representations in HTML and \LaTeX .

2 Markup syntax

The following two subsections give an overview of the very basic HTML and L^AT_EX syntax. They do not introduce the full syntax but focus on the aspects needed within this article. Full descriptions can be found in [3] for XML, the underlying markup language of HTML as described in [6], and in [10, 11] for L^AT_EX.

2.1 Basic syntax of HTML

An HTML document consists of *elements* that are either *empty* or *non-empty*. The boundaries of a non-empty element is marked by a *start-tag* and a *end-tag*. Tag delimiters are the < and > characters. The element *type* is defined in the start-tag by its *name*. The end-tag repeats the element name preceded by a / character. The element's *content* is enclosed between the start- and end-tag. The content consists of *character data* (i.e. text), subordinated *child* elements, or both. An element without content is called empty and is either described by a start-tag that is directly followed by its end-tag or by an *empty-element-tag*. An empty-element-tag has the same form as a start-tag but ends with a / character. The following example shows a non-empty 'p'-element with mixed content:

```

      child element
      ┌───────────┴───────────┐
<p> Lorem ipsum <em>dolor</em> sit amet. </p>
└──────────────────────────┘
start-tag          content          end-tag

```

An element can possess *attributes*. Attributes are noted in the start-tag or empty-element-tag behind the element name. Attributes consist of a *name-value*-pair. The attribute-value is given between two ' or " characters and is assigned to its name by a preceding = character. The following example shows an empty 'img' element with a 'src' attribute of the value 'uri':

```

      attribute
      ┌──────────┴──────────┐
<img src = "uri" />
└──┬──┘ └──┬──┘
  name  value

```

There is exactly one *root* element that includes all the document's content. The tag placement within the document follows the rules of mathematical brackets. The examples below show possible tag placements by means of 'a' and 'b' elements:

sequence	<a>
subordination	<a>
<i>syntax error</i>	<a>

2.2 Basic syntax of L^AT_EX

A L^AT_EX document consists of *commands* that describe either output characters (i.e. characters to typeset), special characters (e.g. the ~ character for a non-breaking space), or *control sequences*. There are two types of control sequences: *control words* and *control symbols*. A control word starts with a \ character followed by its *name* that consists of one or more letters (i.e. lower- or uppercase letters 'a' to 'z') and is terminated by either a space or another non-letter. A control symbol starts with a \ character followed by one non-letter. A command can possess optional and required *parameters* that are set by *arguments*. Optional parameter arguments are noted after the command name between square brackets, and required parameter arguments between curly braces. The following example shows a 'usepackage'-command with an optional parameter set to 'utf8' and a required parameter set to 'inputenc':

```

                                parameters
                                ┌──────────┴──────────┐
\usepackage [utf8] {inputenc}
                                └──┬──┘ └──┬──┘
                                  optional  required

```

Furthermore there are two special types of commands: *environments* and *declarations*. Environments are pairs of 'begin'- and 'end'-commands that enclose the environment's content. The environment name is provided as the first required argument of the corresponding 'begin'- and 'end'-commands. The arguments of the environment are noted as further arguments of the 'begin'-command. Declarations influence the behavior of following commands. The *scope* (i.e. range of effect) of most declarations is limited to its enclosing environment or *group*. The group delimiters are the { and } characters. The placement of group delimiters and environment commands follows the rules of mathematical brackets. The examples below show possible placements by example of a group and an 'x'-environment:

sequence	{ } \begin{x} \end{x}
subordination	{ \begin{x} \end{x} }
<i>syntax error</i>	{ \begin{x} } \end{x}

3 Markup correspondence

The following sections introduce a possible mapping between HTML elements and L^AT_EX commands in the order of the HTML module descriptions in [1]. For an XSLT implementation transforming HTML to T_EXML, the following mapping tables show the resulting L^AT_EX commands for expository purposes.

3.1 Core Modules

The HTML Core Modules assemble the markup that is common to all HTML dialects that are derived from module-based HTML. This core markup for high level structures, basic text, hyperlinks, and lists of HTML documents and its corresponding L^AT_EX commands are described in the following subsections.

3.1.1 Structure Module

The HTML Structure Module defines the high level markup of a document. The `html`-element is the document's root containing the meta-information (`head`) and the actual content (`body`) of a document. L^AT_EX follows a similar separation with its preamble and `document`-environment. Table 1 below shows the corresponding commands.

Table 1: HTML to L^AT_EX structure mapping

HTML	L ^A T _E X
<code><html>{...}</code>	<code>\documentclass{report}{...}</code>
<code><head>{...}</code>	<code>{...}</code>
<code><title>{...}</code>	<code>\title{...}</code>
<code><body>{...}</code>	<code>\begin{document}{...}</code>

3.1.2 Text Module

The HTML Text Module defines the basic text markup to describe heading, block, and inline elements. Most of these elements have equivalent commands in L^AT_EX, but not all. In these cases the ‘ \mapsto ’ symbol indicates the default formatting in HTML where the Presentation Module described in section 3.2.1 might be used for an alternative, not corresponding semantically, mapping.

Headings The HTML Text Module defines six levels of headings (`h1` to `h6`). L^AT_EX offers a specific heading hierarchy that depends on the given document class. Table 2 below shows the corresponding heading commands for the `report` document class.

Table 2: HTML to L^AT_EX heading mapping

HTML	L ^A T _E X
<code><h1>{...}</code>	<code>\chapter{...}</code>
<code><h2>{...}</code>	<code>\section{...}</code>
<code><h3>{...}</code>	<code>\subsection{...}</code>
<code><h4>{...}</code>	<code>\subsubsection{...}</code>
<code><h5>{...}</code>	<code>\paragraph{...}</code>
<code><h6>{...}</code>	<code>\subparagraph{...}</code>

Blocks The HTML Text Module defines elements to mark text groups as paragraphs (`p`), contact information (`address`), quotations (`blockquote`), generic

groups (`div`), and preformatted text (`pre`). Table 3 shows the corresponding L^AT_EX commands (using the L^AT_EX core package `alltt` for preformatted text).

Table 3: HTML to L^AT_EX block mapping

HTML	L ^A T _E X
<code><p>{...}</code>	<code>{...} \par</code>
<code><address>{...}</code>	\mapsto <i>italic</i>
<code><blockquote>{...}</code>	<code>\begin{quote}{...}</code>
<code><div>{...}</code>	<code>{...}</code>
<code><pre>{...}</code>	<code>\begin{alltt}{...}</code>

Inlines The HTML Text Module defines markup for text fragments. This includes abbreviations (`abbr`) and acronyms (`acronym`), citations (`cite`), quotations (`q`), and definitions (`dfn`), program code (`code`), sample output (`samp`), arguments (`var`), and input (`kbd`), regular (`em`) and strong (`strong`) emphases, generic fragments (`span`), and forced line breaks (`br`). Table 4 below shows the corresponding L^AT_EX commands (using the `glossaries` package for abbreviations and acronyms, the `csquotes` package for quotations, and the `listings` package for code).

Table 4: HTML to L^AT_EX inline mapping

HTML	L ^A T _E X
<code><abbr>{...}</code>	<code>\acrshort{...}</code>
<code><acronym>{...}</code>	<code>\ac{...}</code>
<code><cite>{...}</code>	<code>\cite{gen-id}</code> <code>\bibitem{gen-id}{...}</code>
<code><q>{...}</code>	<code>\enquote{...}</code>
<code><dfn>{...}</code>	\mapsto <i>italic</i>
<code><code>{...}</code>	<code>\lstinline ... </code>
<code><samp>{...}</code>	\mapsto <i>teletype</i>
<code><var>{...}</code>	\mapsto <i>italic</i>
<code><kbd>{...}</code>	\mapsto <i>teletype</i>
<code>{...}</code>	<code>\emph{...}</code>
<code>{...}</code>	\mapsto bold
<code>{...}</code>	<code>{...}</code>
<code>
</code>	<code>\newline</code>

3.1.3 Hypertext Module

The HTML Hypertext Module defines markup to describe hyperlinks. They are described by source anchors (`a`) that reference to contents inside or outside of the document via Unified Resource Identifiers (URIs). Referenceable document fragments are marked by common ‘`id`’ attributes that can be applied to all elements. The use of traversable hyperlinks is an adequate solution in the context of electronic documents; its mapping to corresponding

L^AT_EX commands by means of the `hyperref` package is shown in Table 5. However, in the context of printed matter a solution with references by e.g. visual key or page numbers might be more appropriate.

Table 5: HTML to L^AT_EX hypertext mapping

HTML	L ^A T _E X
<code><a href="<uri"><...></code>	<code>\href{<uri>}{<...>}</code>
<code><a href="<id"><...></code> <code>id="<id"></code>	<code>\hyperref[<id>]{<...>}</code> <code>\label{<id>}</code>

3.1.4 List Module

The HTML List Module defines markup to describe ordered (`ol`) and unordered (`ul`) lists as sequences of list items (`li`) and furthermore markup to describe definition lists (`dl`) that are composed of sequences of term (`dt`) and description (`dd`) pairs. Table 6 below shows corresponding L^AT_EX commands.

Table 6: HTML to L^AT_EX list mapping

HTML	L ^A T _E X
<code><...></code>	<code>\begin{enumerate}<...></code>
<code><...></code>	<code>\begin{itemize}<...></code>
<code><...></code>	<code>\item <...></code>
<code><dl><...></code>	<code>\begin{description}<...></code>
<code><dt><...></code>	<code>\item[<...>]</code>
<code><dd><...></code>	<code><...>\</code>

3.2 Text Extension Modules

The HTML Text Extension Modules assemble additional text markup to control text rendering, maintenance, and direction for HTML documents. These and the corresponding L^AT_EX commands are described in the following subsections.

3.2.1 Presentation Module

The HTML Presentation Module defines markup to control the text rendering. It provides elements to render text in/as bold (`b`) and italic (`i`) style, typewriter (`tt`), super- (`sup`) or subscripted (`sub`), larger (`big`) or smaller font (`small`). Additionally the module provides an element to render horizontal rules (`hr`). L^AT_EX offers corresponding commands with the exception of ‘`textsubscript`’ that relies on the `subscript` package. The `relsize` package offers commands to realize relative font sizes (as intended by the ‘`big`’ and ‘`small`’ elements in HTML). Table 7 shows a possible mapping.

3.2.2 Edit Module

The HTML Edit Module defines editing-related markup. It provides elements to mark content as deleted

Table 7: HTML to L^AT_EX presentation mapping

HTML	L ^A T _E X
<code><...></code>	<code>\textbf{<...>}</code>
<code><i><...></code>	<code>\textit{<...>}</code>
<code><tt><...></code>	<code>\texttt{<...>}</code>
<code><sup><...></code>	<code><...></code>
<code><sub><...></code>	<code>\textsubscript{<...>}</code>
<code><big><...></code>	<code>{\larger <...>}</code>
<code><small><...></code>	<code>{\smaller <...>}</code>
<code><hr /></code>	<code>\hrulefill</code>

(`del`) or inserted (`ins`). The `changes` package offers semantically corresponding L^AT_EX commands as shown in Table 8 below. However, if L^AT_EX is used as final output format only, a more stable solution might be to simply output contents of ‘`ins`’-elements, but not those of ‘`del`’-elements.

Table 8: HTML to L^AT_EX edit mapping

HTML	L ^A T _E X
<code><...></code>	<code>\deleted{<...>}</code>
<code><ins><...></code>	<code>\added{<...>}</code>

3.2.3 Bi-directional Text Module

The HTML Bi-directional Text Module defines markup to declare text direction changes. It provides an attribute to control the direction of text (`dir`) that can be applied to all elements including a special element (`bdo`) to override the current text direction. The `bidi` package offers corresponding L^AT_EX commands. Table 9 below shows the corresponding commands for inline text. However, the `bidi` package defines a set of new environments which replace common L^AT_EX commands (e.g. lists and footnotes) which makes the general mapping between elements and commands more complex. Furthermore the combination with other common packages (e.g. `hyperref` or `longtable`) remains problematic. So a more stable solution might be to omit bi-directional text controls during the transformation process and to apply such changes manually in the L^AT_EX document.

Table 9: HTML to L^AT_EX bidi mapping

HTML	L ^A T _E X
<code><bdo dir="ltr"><...></code>	<code>\LR{<...>}</code>
<code><bdo dir="rtl"><...></code>	<code>\RL{<...>}</code>

3.3 Forms Modules

The HTML Forms Modules define markup to describe interactive forms that can define, organize, and receive (textual) input and selections. The `hyperref`

package implements most HTML form elements for L^AT_EX. As with hyperlinks, the use of interactive forms is adequate for electronic documents; their mapping by means of the `hyperref` package is shown in Table 10 below. However, in the context of printed matter, an alternative solution as given e.g. by the `formular` package might be more suitable.

Table 10: HTML to L^AT_EX forms mapping

HTML	L ^A T _E X
<code><form>{...}</code>	<code>\begin{Form}{...}</code>
<code><input /></code>	<code>\TextField{<label>}</code>
<code>type="password"</code>	<code>\TextField[password]{<label>}</code>
<code>type="checkbox"</code>	<code>\CheckBox{<label>}</code>
<code>type="button"</code>	<code>\PushButton{<label>}</code>
<code>type="radio"</code>	<code>\ChoiceMenu[radio]{<label>}{=}</code>
<code>type="submit"</code>	<code>\Submit{<label>}</code>
<code>type="reset"</code>	<code>\Reset{<label>}</code>
<code>type="file"</code>	<code>\TextField[fileselect]{<label>}</code>
<code>type="hidden"</code>	<code>\TextField[hidden]{<label>}</code>
<code>type="image"</code>	<code>\Submit[submitcoordinates]{}</code>
<code><select>{...}</code>	<code>\ChoiceMenu{<label>}{<options>}</code>
<code><option>{...}</code>	<code>{...}</code>
<code><textarea>{...}</code>	<code>\TextField[multiline]{<label>}</code>
<code><button>{...}</code>	<code>\Submit{<...>}</code>
<code>type="button"</code>	<code>\PushButton{<...>}</code>
<code>type="reset"</code>	<code>\Reset{<...>}</code>
<code><fieldset>{...}</code>	<code>{...}</code>
<code><label>{...}</code>	<code>{...}</code>
<code><legend>{...}</code>	<code>{...}</code>
<code><optgroup>{...}</code>	<code>{...}</code>

3.4 Table Modules

The HTML Table Modules define markup to describe tables (`table`) by organizing their data (`td`) and header (`th`) cells in rows (`tr`). These rows can be grouped into table headers (`thead`), footers (`tfoot`), and bodies (`tbody`). Column-based markup is realized by standoff elements (`col` and `colgroup`). A table caption (`caption`) can provide a short description of the table contents.

L^AT_EX table definitions differ in two essential aspects from HTML: (i) The total number of table columns has to be given explicitly to a L^AT_EX table environment. This is not necessary in HTML but calculated continuously by the rendering engine at processing time. (ii) L^AT_EX table cells that span several rows (by means of the `multirow` package) cover the adjacent cells in the following rows; therefore empty cells need to be inserted in the following rows. This is not necessary in HTML but the rendering

```

table = {rowi | rowi = {cellij}}
grid = {sloti | sloti = ⟨x, y⟩}
procedure TABLE(table)
  grid ← {∅}
  for all rowi | i = 1 .. n do
    y ← i
    x ← 1
    for all cellij | j = 1 .. m do
      while grid ∋ ⟨x, y⟩ do
        EMPTY CELL(x, y)
        x ← x + 1
      end while
      for ycell ← 0 .. rowspan(cellij) - 1 do
        for xcell ← 0 .. colspan(cellij) - 1 do
          grid ← grid ∪ ⟨x + xcell, y + ycell⟩
        end for
      end for
      x ← x + colspan(cellij)
    end for
  end procedure

```

Figure 2: HTML table cell positioning algorithm

engine automatically shifts the cells of the following rows according to the reading direction.

Hence for the transformation of HTML tables to L^AT_EX this information (total number of table columns and position of additional empty cells) need to be precalculated. Therefore the transformation process has to include parts of the HTML table processing model described in [2]. This model describes an HTML table as a set of cells that are positioned on a two-dimensional grid of slots. The algorithm shown in Figure 2 calculates the cell positioning and illustrates how the additional empty cells are inserted; hence the total number of table columns is given by the maximum x -coordinate within the final grid. Table 11 shows the mapping of HTML table elements to corresponding L^AT_EX commands by means of the `longtable` package.

3.5 Image Module

The HTML Image Module defines markup to embed external images. The `graphicx` package offers a corresponding L^AT_EX command as shown in Table 12.

3.6 Further Modules

The HTML specification describes further modules that define markup to realize dynamic and interactive document content, mechanisms to control layout, and deprecated markup for backwards compatibility with legacy HTML. Due to the focus of this article on the transfer of the logical structure of HTML documents

Table 11: HTML to L^AT_EX table mapping

HTML	L ^A T _E X
<code><caption><...></code>	<code>\caption{<...>}</code>
<code><table><...></code>	<code>\begin{longtable}{<col>}<...></code>
<code><td><...></code>	<code><...>&</code>
<code><th><...></code>	<code>\bf{<...>&</code>
<code> colspan=""</code>	<code>\multicolumn{}{1}{<...>}</code>
<code> rowspan=""</code>	<code>\multirow{}{*}{<...>}</code>
<code><tr><...></code>	<code><...>\</code>
<code><col /></code>	
<code><colgroup><...></code>	
<code><tbody><...></code>	<code><...></code>
<code><thead><...></code>	<code><...>\endhead</code>
<code><tfoot><...></code>	<code><...>\endfoot</code>

Table 12: HTML to L^AT_EX image mapping

HTML	L ^A T _E X
<code><img src="<uri>" /></code>	<code>\includegraphics{<uri>}</code>

to L^AT_EX, the mapping of these specialized modules is not described in detail. However, in specific use cases the support of these modules might be desired. The following hints might serve as a starting point to implement a transformation of these modules' features to L^AT_EX.

The HTML Applet, Object, Scripting, and Intrinsic Events modules define markup that introduces scripting facilities to manipulate dynamically the document content. At present this is notably realized through the JavaScript programming language, which is partially integrated with L^AT_EX by means of the `insdljs` package.

The HTML Client- and Server-side Image Map modules define markup for interactive and hyperlinked images. This functionality can be potentially realized in L^AT_EX by means of the `TikZ` package.

The HTML Frames and Iframe modules define markup to insert one document into another. This can be realized in L^AT_EX with the `\input` and/or `\include` commands.

The HTML Style Sheet and Style Attribute modules define markup to integrate layout definitions realized through Cascading Style Sheets (CSS). CSS has its own syntax and description logic—its transformation to L^AT_EX is a topic all its own, which has been outlined e.g. in [15].

4 An example

Figure 3 shows an example page taken from a bird guide. On the next page, Figure 4 shows a possible

coding of this page using HTML, and Figure 5 its corresponding representation in L^AT_EX.

1

Gannet

Birds of the open ocean, Gannets breed on small islands off the NW coast of Europe. They move away from land after nesting to winter at sea. The young migrate south as far as W Africa. Gannets feed on fish by plunging from 25m. They nest in large, noisy colonies. The nest is a pile of seaweed. A single egg is incubated for 44 days. The young bird is fed by both parents and flies after 90 days.



Size	Larger than any gull
Adult	White, black wing-tips, yellow nape
Juvenile	Grey, gradually becoming white over 5 years
Bill	Dagger-like
In flight	Cigar-shaped with long, narrow, black-tipped wings
Voice	Usually silent, growling "urr" when nesting
Lookalikes	Skuas, Gulls and Terns

Figure 3: An example document, derived from [5]

5 Conclusion

While HTML is increasingly becoming the common document description language for different output media (web, print, e-books, ...), the problem of creating well-typeset documents from HTML is not yet fully solved within the XML ecosystem. The article at hand has introduced a mapping from HTML elements to corresponding L^AT_EX commands, in order to use the T_EX typesetting engine for this task.

With the multitude of existing L^AT_EX extensions released as packages, almost any HTML description can be ported to L^AT_EX and typeset according to its original logic. Unfortunately, the use of L^AT_EX packages often comes with a catch: while many HTML structures can be used recursively (e.g. nested lists or tables), L^AT_EX packages tend to override existing commands giving them a new meaning (e.g. the new-line command is redefined in table environments to end a row). These context-dependent syntax-changes can make a mapping potentially error-prone for deep document structures.

```

<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.1//EN" "http://www.w3.org/TR/xhtml11/DTD/xhtml11.dtd">
<html xmlns="http://www.w3.org/1999/xhtml">
  <head><title>Gannet</title></head>
  <body>
    <h1>Gannet</h1>
    <p>Birds of the open ocean, Gannets breed on small islands off the <abbr title="northwest"
      >NW</abbr> coast of Europe. They move away from land after nesting to winter at sea. The
      young migrate south as far as <abbr title="west">W</abbr> Africa. Gannets feed on fish by
      plunge-diving from 25<abbr title="meters">m</abbr>. They nest in large, noisy colonies. The
      nest is a pile of seaweed. A single egg is incubated for 44 days. The young bird is fed by
      both parents and flies after 90 days.</p>
    <div></div>
    <table>
      <tr><th>Size</th><td>Larger than any gull</td></tr>
      <tr><th>Adult</th><td>White, black wing-tips, yellow nape</td></tr>
      <tr><th>Juvenile</th><td>Grey, gradually becoming white over 5 years</td></tr>
      <tr><th>Bill</th><td>Dagger-like</td></tr>
      <tr><th>In flight</th><td>Cigar-shaped with long, narrow, black-tipped wings</td></tr>
      <tr><th>Voice</th><td>Usually silent, growling <q>urr</q> when nesting</td></tr>
      <tr><th>Lookalikes</th><td>Skuas, Gulls and Terns</td></tr>
    </table>
  </body>
</html>

```

Figure 4: HTML source, describing the document shown in Figure 3

```

\documentclass{report}
% preamble ...

\begin{document}
\chapter{Gannet}

Birds of the open ocean, Gannets breed on small islands off the \acrshort{NW}coast of Europe. They
move away from land after nesting to winter at sea. The young migrate south as far as \acrshort{W}
Africa. Gannets feed on fish by plunge-diving from 25\acrshort{m}. They nest in large, noisy
colonies. The nest is a pile of seaweed. A single egg is incubated for 44 days. The young bird is
fed by both parents and flies after 90 days.\par

\includegraphics{gannet.jpg}

\begin{longtable}{ll}
\toprule
\bf{}Size      & & Larger than any gull \\
\bf{}Adult    & & White, black wing-tips, yellow nape \\
\bf{}Juvenile & & Grey, gradually becoming white over 5 years \\
\bf{}Bill     & & Dagger-like \\
\bf{}In flight & & Cigar-shaped with long, narrow, black-tipped wings \\
\bf{}Voice    & & Usually silent, growling \enquote{urr} when nesting \\
\bf{}Lookalikes & & Skuas, Gulls and Terns \\
\bottomrule
\end{longtable}

\end{document}

```

Figure 5: L^AT_EX output, exported from HTML shown in Figure 4

The mappings introduced in this article have been developed in the context of an XSLT implementation within a \TeX -based workflow, but do not rely on it and can be implemented through other approaches as well. However, the principle of \TeX XML, to provide a processor that transforms specific \TeX commands from a generic XML representation to the \TeX format, realizes a separation between the task of format transformation and the task of defining appropriate mappings. This facilitates the definition and adaption of markup correspondences as has e.g. been done by extending the HTML mapping with a third party stylesheet that defines the transformation from MathML to \LaTeX .

Acknowledgment

The author would like to thank Prof. Dr. Karl-Heinrich Schmidt for giving valuable advice and Gilles Bülow for integrating this document processing workflow into our daily tasks for testing purposes.

References

- [1] Daniel Austin, Shane McCarron, Subramanian Peruvemba, Masayasu Ishikawa, and Mark Birbeck. XHTML modularization 1.1 — second edition. W3C recommendation, W3C, July 2010. <http://www.w3.org/TR/2010/REC-xhtml-modularization-20100729/>.
- [2] Robin Berjon, Steve Faulkner, Travis Leithead, Erika Doyle Navara, Edward O'Connor, Silvia Pfeiffer, and Ian Hickson. HTML 5. W3C candidate recommendation, W3C, August 2013. <http://www.w3.org/TR/2013/CR-html5-20130806/>.
- [3] Tim Bray, François Yergeau, C. M. Sperberg-McQueen, Jean Paoli, and Eve Maler. Extensible markup language (XML) 1.0 (fourth edition). W3C recommendation, W3C, August 2006. <http://www.w3.org/TR/2006/REC-xml-20060816/>.
- [4] Markus Gylling, William McCoy, Erika J. Etemad, and Matt Garrish. EPUB content documents 3.0. IDPF recommended specification, IDPF, October 2011. <http://www.idpf.org/epub/30/spec/epub30-contentdocs-20111011.html>.
- [5] Renate Henschel, John Bateman, and Judy Delin. Automatic genre-driven layout generation. In *Proceedings of the 6th "Konferenz zur Verarbeitung natürlicher Sprache" (KONVENS) Conference*, Saarbrücken, September 2002.
- [6] Masayasu Ishikawa and Shane McCarron. XHTML 1.1 — module-based XHTML — second edition. W3C recommendation, W3C, November 2010. <http://www.w3.org/TR/2010/REC-xhtml11-20101123/>.
- [7] Ian Jacobs, David Raggett, and Arnaud Le Hors. HTML 4.01 specification. W3C recommendation, W3C, December 1999. <http://www.w3.org/TR/1999/REC-html401-19991224/>.
- [8] Michael Kay. XSL transformations (XSLT) version 2.0. W3C recommendation, W3C, January 2007. <http://www.w3.org/TR/2007/REC-xslt20-20070123/>.
- [9] Sanders Kleinfeld. The case for authoring and producing books in (X)HTML5. In *Proceedings of Balisage: The Markup Conference 2013*, volume 10 of *Balisage Series on Markup Technologies*, Montréal, August 2013.
- [10] Donald Ervin Knuth. *The \TeX book*, volume A of *Computers & Typesetting*. Addison-Wesley, March 1986.
- [11] Leslie Lamport. *\LaTeX : A Document Preparation System*. Addison-Wesley, 2nd edition, November 1994.
- [12] Douglas Lovell. \TeX XML: Typesetting XML with \TeX . *TUGboat*, 20(3):176–183, September 1999. <http://tug.org/TUGboat/tb20-3/tb64love.pdf>.
- [13] Shane McCarron. XHTML-print — second edition. W3C recommendation, W3C, November 2010. <http://www.w3.org/TR/2010/REC-xhtml-print-20101123/>.
- [14] Oleg Parashchenko. \TeX XML: Resurrecting \TeX in the XML world. *TUGboat*, 28(1):5–10, March 2007. <http://tug.org/TUGboat/tb28-1/tb88parashchenko.pdf>.
- [15] S. Sankar, S. Mahalakshmi, and L. Ganesh. An XML model of CSS3 as an \LaTeX - \TeX XML-HTML5 stylesheet language. *TUGboat*, 32(3):281–284, December 2011. <http://tug.org/TUGboat/tb32-3/tb102sankar.pdf>.

◇ Frederik R. N. Schlupkothen
University of Wuppertal
Rainer-Gruenter-Str. 21
D-42119 Wuppertal
Germany
schlupko (at) uni-wuppertal dot de