

XLingPaper's use of T_EX Technologies

H. Andrew Black, Hugh J. Paterson III

Abstract

We discuss the use of T_EX technologies by XLingPaper, an authoring tool for producing academically oriented publications with features required for linguistic publishing. We present the T_EX modules used and the rationale for the history of XLingPaper development.

1 Introduction

Within the publishing industry, there are several notable products for producing complex documents in beautiful formats. T_EX [23] [24] is one of the well known publishing technologies used to meet these needs. Since 2000, XML-based technologies such as XSL-FO¹ or the T_EXML² project [28] have been used to integrate content and compose complex documents such as textbooks and maintenance manuals. Requirements for composing these large, inter-linked documents birthed the development of tools like XMLmind,³ the <oxyen/> XML Editor,⁴ and Xpublisher.⁵ These tools can be used to compose content within predefined XML structures. XLingPaper, as discussed in [7] [8] [9], seeks to provide a constrained environment in which authors of complex works dealing with language descriptions and linguistic analyses can focus on content structure independently from the styling requirements of documents. In this way the underlying design principle of XLingPaper maximizes the SGML design practice of separating content from presentation. With XLingPaper, authors can keep content structure independent from page layout information and thereby provide maximal transferability between publishing styles. The software does this while providing authors a clear structured interface for authoring content.

XLingPaper is designed to reduce friction in the process of writing, composing, and publishing linguistic papers, grammars, and books by removing common time-sinks related to inconsistent formatting (especially citations, references, and numbered elements like examples). A full list of benefits to all parties in the publishing work flow is available [9].

The XLingPaper software has a growing number of users who have successfully typeset complex documents including:

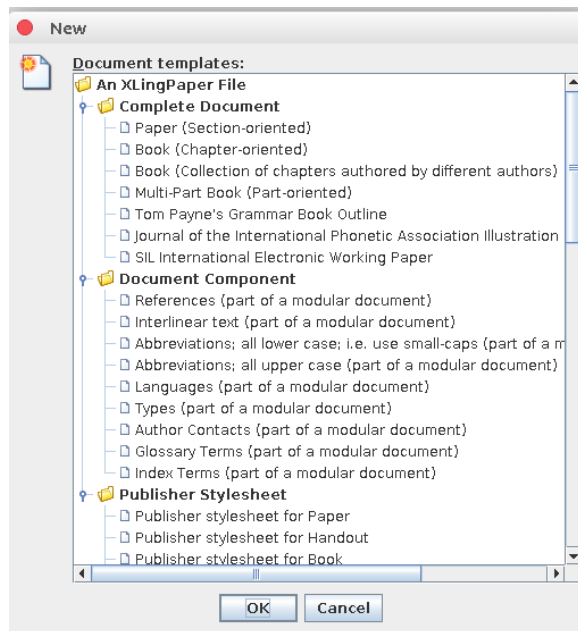


Figure 1: XLingPaper predefined document types via DTD

- master theses [50] [25] [33],
- doctoral dissertations [16] [37],
- textbooks [30],
- linguistic grammars [11],
- books [1] [38],
- journal articles [10], and
- bilingual software documentation [2] [3].

2 What is XLingPaper?

XLingPaper⁶ is a plug-in to the XMLmind XML Editor. XLingPaper benefits from the XMLmind XML Editor's Java-based implementation which allows it to be used on Mac OS X, Windows, and Linux. Via a DTD, XLingPaper defines several document classes (articles, books, chapters, etc., as illustrated in Figure 1), in each case providing document layout sections (paragraphs, examples, endnotes, etc.). Figure 2 illustrates the main screen of the user-interface of XMLmind XML Editor. By using this interface, formatting errors are reduced because users are constrained on where in the document flow they can introduce block and line level document elements. That is, first, authors cannot input page layout instructions directly into the document and second, the introduction of layout sections within the document flow is constrained via the DTD.

¹ w3.org/TR/xsl11

² getfo.org/texml

³ xmlmind.com/xmleditor

⁴ oxygenxml.com/xml_author.html

⁵ xpublisher.com/products

⁶ software.sil.org/xlingpaper

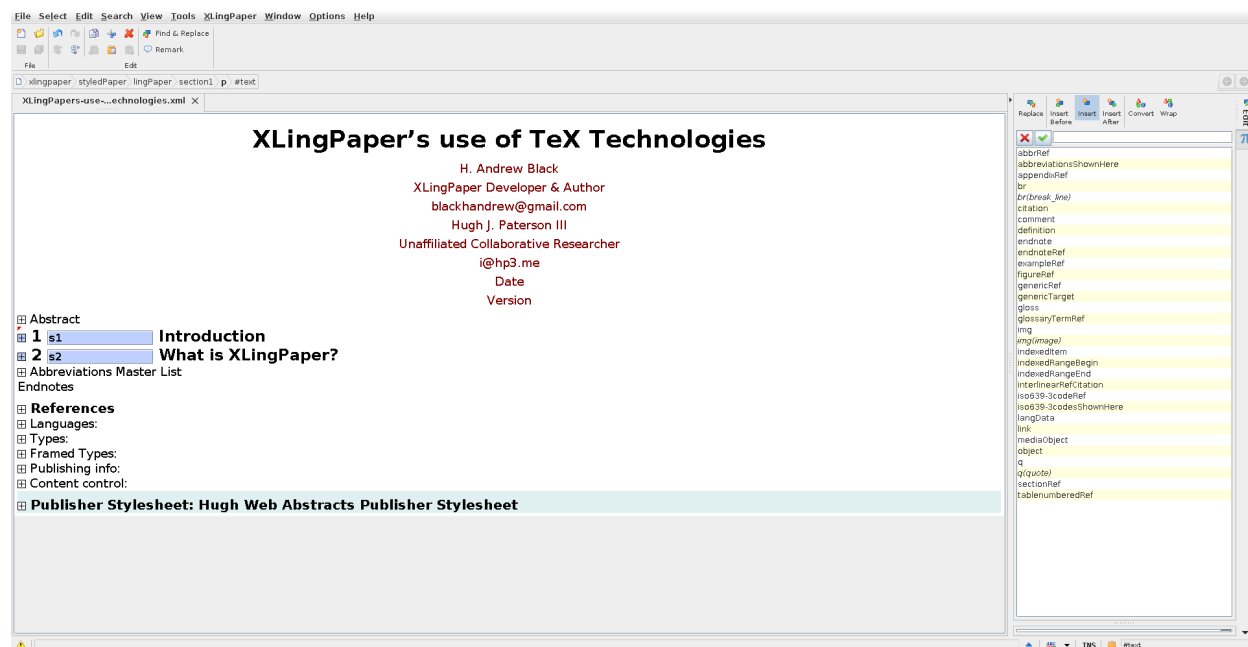


Figure 2: XLPaper’s user interface. Left side: document content editing. Right side: block and line level units available for use at the cursor location.

For many, the PDF format is the quintessential file format for final distribution of publishing outputs. XLPaper supports PDF production; however, as illustrated in Figure 3, XLPaper can also produce documents with at least five outputs, all from the same source document:

- PDF (version 1.5),
- Web pages (HTML 4),
- Microsoft Word (.doc),
- Open Office Writer Document (.odt), and
- ePUB.

XLPaper automatically numbers tables, examples, figures, and sections. It keeps track of internal references to these entities along with citation references, abbreviations, and gloss abbreviations. This keeps numbering and reference links dependable and automated. XLPaper also automatically generates indexes, a table or list of abbreviations used, and a section for references cited (using a custom references implementation).

Unlike most editing programs which are based on either the WYSIWYG paradigm or are unconstrained text editors such as those used to code or produce Markdown, XLPaper (via the XMLmind XML Editor) is a structured editor much more like the block editors we see in tools like MailChimp⁷ or WordPress’s Gutenberg editor,⁸ albeit without

the drag-and-drop features. Rather than visually structuring the document to look the way it is to be formatted, the author “marks up” the items in the document according to their kind. One of the many benefits that using a DTD provides is that there is a “grammar” of what a well-formed linguistic document looks like. This makes moving, replacing, switching, or reordering sections, chapters, tables, figures, and examples less error prone because it prevents users from inadvertently creating ill-formed documents. The following sections of this paper discuss the TeX technologies used by XLPaper.

3 XLPaper and TeX

Linguistic publishing has unique requirements when compared to general publishing. The following sections provide more detail on the linguistic publishing context, design requirements and L^AT_EX packages used by XLPaper.

3.1 TeX and Linguistic document production

TeX has long been embraced by linguists. Peter [35] writes of a personal communication with Don Knuth where Knuth suggests that linguists were some of the earliest adopters outside of mathematicians. Thiele [45] in an interview given in 2007 states that she was typesetting linguistic journals via TeX in 1983—a date prior to the release of Knuth’s book on using the TeX typesetting system [23]. Thiele [44]

⁷ mailchimp.com

⁸ developer.wordpress.org/block-editor

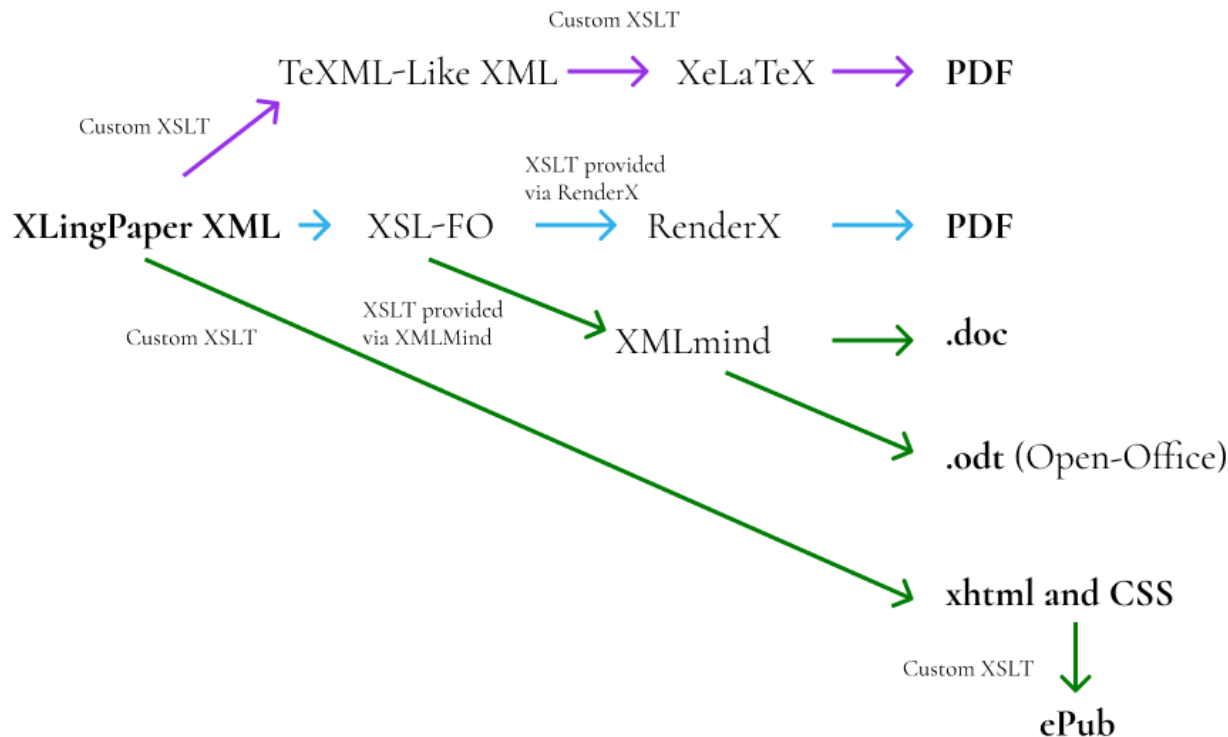


Figure 3: XLPingPaper’s data processing pipeline to multiple formats

gives an early overview of \TeX use in linguistics with mention of significant repositories outside of CTAN. A slightly more recent (2004) update by Peter [35] provides some additional tips and tools for typesetting common information structures in linguistic publishing. The \TeX community has produced many packages which have shaped the visual face of publishing in linguistics, including `tipa`⁹ by Rei [39], which provided access to an excellent typeface for phonetic transcriptions, and `pst-asr`¹⁰ by Frampton [17] for autosegmental representations. Some packages used in linguistic publishing are special purpose but are not exclusive to linguistics. For example, Donnelly [15] describes how to use various packages to draw phonetic pitch traces using \TeX . Peter [35] and Thiele [44] list out and review (through 2004) various packages across several areas of linguistics. Among others, they discuss several packages used to draw syntax trees such as `qtrees`¹¹ and `forest`¹² and specialized packages for presenting examples and interlinear glossed texts such as `expex`.¹³ Their reviews also discuss packages such as

`covington`¹⁴ and `gb4e`¹⁵ whose collections of macros serve a variety of page layout functions targeted at publishing in linguistic topics.

The CTAN repository currently lists fifty-four different \TeX packages for linguistic typesetting,¹⁶ though some of these packages also include capabilities targeted as multi-lingual or multi-script publications or are specific style sheet implementations for publications at linguistic programs at institutions of higher education (there may be more packages which are not tagged but should be). Several of the packages tagged “linguistic” pre-date Unicode [46] but still see significant use. Sometimes it is the case that secondary packages are developed in an attempt to “fix” publishing outputs in different ways to bring Unicode features along with the features of the original package. For example, `tipa` is not Unicode compatible, but the packages `unitipa`¹⁷ and `tipauni`¹⁸ seek to address different implications of not publishing with Unicode while giving access to the beautiful typeface of `tipa`. Understanding the long history of publishing and the interdependency

⁹ ctan.org/pkg/tipa

¹⁰ ctan.org/pkg/pst-asr

¹¹ ctan.org/pkg/qtrees

¹² ctan.org/pkg/forest

¹³ ctan.org/pkg/expex

¹⁴ ctan.org/pkg/covington

¹⁵ ctan.org/pkg/gb4e

¹⁶ ctan.org/topic/linguistic

¹⁷ ctan.org/pkg/unitipa

¹⁸ ctan.org/pkg/tipauni

that packages have (including the order of loading packages) constitutes barriers of adoption to new T_EX users.

We discuss T_EX barriers of adoption for two reasons. First, it exemplifies some of the complexities that X_LingPaper seeks to simplify as it presents authors not just a visual environment for document composition, but also a cohesive output solution. Second, it speaks to the software design process in finding the minimal viable product. That is, *how much (or little) of a software stack is needed to make a usable software product for linguistic publishing?* The T_EX community is divided on this. While the diagrams in linguistic books and journals since the 1980's exemplify many beautiful, sharp, crisp, illustrations created directly in T_EX, many trainers of T_EX tools,¹⁹ but not all,²⁰ have steered authors towards a more generic set of packages which do not include specific diagram creating macros. Rather, they suggest that authors use secondary illustration tools to generate illustrations and then include them as vector PDFs or images. In fact this second method is the document production path that the X_LingPaper philosophy follows. That is, X_LingPaper reduces the complexity of the typesetting task for authors by requiring complex visualizations to be produced via graphical tools. We have found tools like Figma²¹ and Inkscape²² very helpful in the graphic production task. The X_LingPaper product seeks to lower barriers of entry, only produce valid documents, and keep the code base to a minimum.

As mentioned in the discussion of `tipa`, linguistic documents have not always been typeset with Unicode. Unicode was introduced in 1991 and by the early 2000's Unicode along with document and data storage in XML formats were being heralded in academic linguistics as a best practice in order to avoid vendor lock-in, increase interoperability across use cases, and to separate data life-cycles from encoding or software life-cycles [5] [6] [47]. Due to the heavy reliance on Unicode by today's practitioners of language documentation and linguistic work, X_LingPaper specifically uses X_QL^AT_EX and compatible packages to produce PDF outputs. This brings continuity to the text input process for users across their workflows. It also makes importing and using language or phonetically transcribed examples sim-

pler by removing the need to use macros to derive characters.

3.2 Design desiderata for X_LingPaper outputs via T_EX

Three goals have driven the development of X_LingPaper:

- separation of content and style,
- software accessibility (license and size), and
- beautiful multi-format outputs.

Deciding how T_EX technologies fit within the project has been a journey. X_LingPaper development started in 2001 without any use of T_EX technologies. In 2006, X_LingPaper added XSL-FO for PDF production. Prior to 2009, X_LingPaper used RenderX's XEP²³ product to produce PDF documents. As far as we know, there are two cross-platform XSL-FO processors written in Java: RenderX's XEP application and the Apache FOP project.²⁴ Using a Java implementation reduces the size of the required stack because the XMLmind XML Editor requires Java. XSL-FO processors can have various degrees of implementation of the XSL-FO standard. RenderX has some limitations which affect page layout but has more complete coverage than the Apache FOP project which lacks certain required table-oriented capabilities.²⁵ The limitations of RenderX are discussed in Section 6. In 2009 plans were made to add X_QL^AT_EX-based output to X_LingPaper because, while there was a free version of RenderX, the output contained a watermark. By implementing the ability to export to PDF via X_QL^AT_EX, watermarks could be avoided all together. The X_QL^AT_EX method of PDF production is now the default method to produce PDF documents, although the RenderX method is still possible.

Maintaining a separation of content and style in the X_LingPaper environment was a key design requirement. When the X_QL^AT_EX method of PDF production was introduced, X_LingPaper already had a way to format output per a user-created publisher style sheet — allowing great flexibility due to the separation of style and content. Using T_EX technologies meant the developer (Andrew Black) needed to be able to map from an X_LingPaper publisher style sheet to X_QL^AT_EX. It was known that L^AT_EX was the ideal T_EX implementation to target. However, pure L^AT_EX came with predefined output formatting for front matter, chapters, sections, back matter, etc. Pure L^AT_EX, then, would not allow direct control of

¹⁹ Among others, see the Linguistics Dissertation guide for the University of Hawai'i at Mānoa [20], University of Pennsylvania [14], and Language Science Press Guidelines [32].

²⁰ For counter examples see [27] [42] [19] and [36].

²¹ figma.com

²² inkscape.org

²³ renderx.com/tools/xep.html

²⁴ xmlgraphics.apache.org/fop/index.html

²⁵ xmlgraphics.apache.org/fop/compliance.html

formatting of all of these per an X_LingPaper user-defined publisher style sheet. This required overriding these standard features of L^AT_EX with a custom implementation of the T_EX commands needed to control formatting. X_LingPaper takes a custom approach in implementing flexibility here. Appendix A lists the custom commands used.

The programmer of X_LingPaper recently discovered `memoir`²⁶ [48] [49]. As a package, `memoir` accomplishes many of the same tasks and could be considered to replace some of the custom code if it were shown to be easy to implement and that the size of the total X_LingPaper code base would be reduced.

The distributability of the software was also seen as a design requirement. Distributability is understood to have two components: license and accessibility, including size.

From the outset, X_LingPaper was designed to be costless to the end user. It is licensed under the MIT license, and its code is currently available on Github.²⁷ The XMLmind XML Editor had a costless Personal Use License that met this requirement for the vast majority of the target audience of X_LingPaper. The few X_LingPaper users who did not meet the terms of that license most likely would be able to afford to purchase (or have their organization purchase) a professional license of the XMLmind XML Editor. The actual X_LingPaper plug-in has always been free.

The software size of X_LingPaper is a major design influencer. Many of the expected users of X_LingPaper live and work in places around the world where Internet connections are characterized by high costs, low bandwidth capacity, and general unavailability. Therefore, the download required to install X_LingPaper needed to be as small as possible. On Windows the current full X_LingPaper installer is 146MB, and the XMLmind XML Editor installer is 116MB. Both are required. This stands in contrast to the T_EXLive 2010 installer which has a size of about 1.2GB when downloaded and 2.38GB when uncompressed. The size constraint impacts X_LingPaper because its distribution must be independent of larger mainstream T_EX distribution solutions which have a large footprint. This, of course, includes T_EXLive. Therefore the developer identified which L^AT_EX packages and binaries were needed and created a custom installation package which met the required specifications.

X_LingPaper currently uses the following L^AT_EX packages (in alphabetical order):

<code>attachfile2</code>	<code>lineno</code>
<code>booktabs</code>	<code>longtable</code>
<code>calc</code>	<code>lscape</code>
<code>color</code>	<code>mdframed</code>
<code>colortbl</code>	<code>multirow</code>
<code>etoolbox</code>	<code>normalem</code>
<code>fancyhdr</code>	<code>polyglossia</code>
<code>fontspec</code>	<code>setspace</code>
<code>footmisc</code>	<code>tabularx</code>
<code>hyperref</code>	<code>xltxtra</code>

The twenty L^AT_EX packages that are part of the custom X_LingPaper distribution are still rather large (29MB) for someone for whom Internet bandwidth is an expensive and inconsistent commodity.

To reduce bandwidth requirements two assumptions were made which have more or less proven to obtain. The first assumption that the developer made was that the twenty packages and binaries would not need to change over time; in contrast, the second assumption was that X_LingPaper would acquire new features and need bug fixes. These assumptions resulted in an architecture where page layout information expressed in XML is translated via custom T_EX commands to either T_EX directly or to commands understood by L^AT_EX packages distributed with X_LingPaper. This abstraction layer was then executed when the X_LingPaper file was processed. This middle layer has granted X_LingPaper flexibility in adding new code and capabilities while keeping the “heavy” L^AT_EX packages stable. The net result is a “heavy” first install package (116MB), but light-weight upgrade packages (6.21MB). In the thirteen year history of development, there have been a few occasions where upgrades have required the download of new “heavy” packages. One such case was when the ability to use framed units was added. These elements depend on the `mdframed`²⁸ package [12]. The architecture separating stable packages from custom code, however, has generally worked out well and kept update sizes low.

3.3 PDF production

We know of two pathways for converting XML content into PDFs. The first is via XSL-FO, and the second is via T_EXML which converts XML content to T_EX formatted documents for further processing to PDF. Given certain limitations in both XSL-FO and T_EXML, X_LingPaper uses a custom (or third) method. When an author instructs X_LingPaper to

²⁶ ctan.org/pkg/memoir

²⁷ github.com/sillsdev/XLingPap

²⁸ ctan.org/pkg/mdframed

produce PDF output via $\text{Xe}\text{L}\text{A}\text{T}\text{E}\text{X}$, X $\text{L}\text{ing}\text{Paper}$ produces a $\text{T}\text{E}\text{X}\text{ML}$ -like XML file. This is then converted into a $\text{L}\text{A}\text{T}\text{E}\text{X}$ formatted document via a set of XSLT transforms and processed via $\text{X}\text{e}\text{L}\text{A}\text{T}\text{E}\text{X}$ to produce the PDF. Figure 3 contains a diagram of the data handling process.

3.4 $\text{T}\text{E}\text{X}\text{ML}$

$\text{T}\text{E}\text{X}\text{ML}$ was discovered in the process of planning for the transition of the default PDF renderer from RenderX 's XEP to $\text{X}\text{e}\text{L}\text{A}\text{T}\text{E}\text{X}$. Initial analysis conducted in 2009 understood $\text{T}\text{E}\text{X}\text{ML}$ to have two infelicities for use-cases required in linguistic publishing with X $\text{L}\text{ing}\text{Paper}$:

1. $\text{T}\text{E}\text{X}\text{ML}$ has Python as a dependency and the developer did not want to require X $\text{L}\text{ing}\text{Paper}$ users to install a version of Python specifically for $\text{T}\text{E}\text{X}\text{ML}$. This is especially the case since that version of Python might conflict with other installed versions of Python on their operating systems. Moreover, this approach would make the installation package for X $\text{L}\text{ing}\text{Paper}$ much larger due to the inclusion of Python.
2. X $\text{L}\text{ing}\text{Paper}$ users require a high degree of control for white space. The fine grain control of whitespace was not immediately clear how to accomplish with $\text{T}\text{E}\text{X}\text{ML}$.

3.5 Control characters

Even with the use of Unicode in the text of documents, there are some features of typesetting with TEX -based implementations which require the use of control characters. Additionally, XML also has control characters. In TEX these include `[`, `]`, `<`, and `>`. When transforming data between XML and TEX , TEX control characters and commands need to be escaped to ensure proper data processing. This has been implemented via Java since Java was already present in the dependency stack due to the XMLmind XML Editor requiring it. Additionally, some small methods have been written in Java to provide additional access to features via the graphical user interface. Among other things, these include adding rows and/or columns to tables, automatically converting glosses to abbreviation references, and importing references from various XML formats.

3.6 Ling- TEX

One might ask, “Why not add more linguistic related TEX packages to the available stack, or use those instead of creating custom code?” The answer has two simple parts: First, in 2009 the linguistic capabilities of TEX packages were different than they

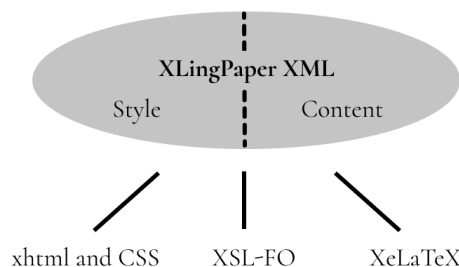


Figure 4: X $\text{L}\text{ing}\text{Paper}$ combines style and content information contained in its custom XML and then exports it into three different formats for further processing.

are today. Second, X $\text{L}\text{ing}\text{Paper}$ is more than a TEX document producer. For example, some authors [2] [3] use X $\text{L}\text{ing}\text{Paper}$ to manage multilingual content on websites.

Besides TEX , X $\text{L}\text{ing}\text{Paper}$ also produces XSL-FO and XHTML/CSS outputs. When new features are considered for inclusion, they must be considered for all output formats.

After excluding $\text{T}\text{E}\text{X}\text{ML}$ as a viable option, and still seeking to create $\text{X}\text{e}\text{L}\text{A}\text{T}\text{E}\text{X}$ -based output, a solution was needed to determine which set of minimal TEX packages would be needed. The Ling- TEX group,²⁹ which also ran the Ling- TEX mailing list from 1995–2018, was discovered.³⁰ Ling- TEX seemed to be the locus of activity in linguistic typesetting via TEX even though other web pages discussing linguistics and TEX also existed, e.g., Essex³¹ and UPenn.³² Today, now that the mailing list is no longer in operation, many of the mailing list participants can be found interacting on the TEX stackexchange.³³

State-of-the-Art for TEX -based linguistic publishing in 2009, as recommended by the Ling- TEX website, suggested using `covington` and `ling-mac`—the list of macros discussed by Thiel in [44]. These macros were used to solve similar use cases, among others, to those already implemented by X $\text{L}\text{ing}\text{Paper}$. Their approaches and outputs, however, had

²⁹ web.archive.org/web/20150702123633/http://heim.ifi.uio.no/~dag/ling-tex

³⁰ ling-tex.ifi.uio.no/narkive.com

³¹ essex.ac.uk/linguistics/external/clmt/latex4ling

³² ling.upenn.edu/advice/latex.html

³³ tex.stackexchange.com

more limitations than what XLingPaper already offered. XLingPaper had the following capabilities for typesetting interlinears:

- no limits on the number of lines within an interlinear grouping;
- no limits on the number of free translation and literal translation lines;
- the ability to include a source reference within the interlinear; and especially
- the ability to tag interlinear items with an ISO 639-3 code for the language used in the interlinear.

At the time the best solution given the state of the \TeX packages available was custom \TeX scripts, although now similar features may be possible via other packages. For example in 2019 Pellard [34] discussed the limiting approaches in various \TeX packages related to interlinear glosses and his solution `typgloss`.³⁴ XLingPaper's examples can be seen in Figures 5–7 which contain output illustrating some of the special capabilities XLingPaper offers.

4 Typesetting tasks XLingPaper users often encounter

Linguistic documents have several formatting needs that other kinds of documents do not. This section discusses some of them.

4.1 Numbered example layouts

Linguistic documents usually have many numbered examples. The prose often refers to examples which are typographically nearby or to previous examples. XLingPaper automatically keeps track of the example identifiers. This is especially important in linguistic publishing because authors, and publishing style sheets, often make use of different kinds of examples, including sub-examples, and table-like design layouts which can contain lists of words along with their glosses (as shown in Figure 6) and interlinear clauses (as shown in Figure 5). Some cases even have headings in portions of the example.

4.2 Interlinear glossed texts

There is a long tradition within linguistics and language study of presenting phrases containing different languages (but the same content) as interlinear texts. Di-Biase-Dyson et al. [13] trace the practice back as far as the 1652 publication of Kircher [22]. More recent publications display significant variation in page layout related to interlinear glossed texts and interlinear examples. Variation exists in three dimensions:

- content elements,
- data-structure of the encoded elements, and
- page layout (visual display of the elements).

A full demonstration of the variation in content and its positioning across common style sheets in linguistics is beyond the scope of this paper. Significant variations include the presence or absence of the following elements:

- index elements such as example numbers or sub-numbers (as shown in Figure 5),
- headings to the interlinear,
- speaker indicator,
- language indicator,
- citation indicator pointing to the larger text from which the example element is taken (see Figure 5 for an example), and
- limits on the number of rows in the original, gloss, translation, and free translation tiers.

Existing \TeX packages approach these content requirements in different ways. As far as we can tell, the following commonly used packages for interlinear glossing all have limitations to some degree. The `expex` package does not offer a content solution for the language code or the citation. The package `langsci-gb4e`,³⁵ a fork of `gb4e`, supports the *Leipzig Glossing Rules*,³⁶ a commonly adopted set of linguistic typesetting conventions. However, while the Leipzig Glossing Rules do call for the language name or identifier to appear on the right hand side of the interlinear glossed text, it does not have a place for the citation. The package `linguex` does not have either language or citation content places built in. With these considerations, it was clear in 2009 that XLingPaper offered more to authors than any single package in the \TeX ecosystem. In order to implement existing XLingPaper features, it meant creating custom \TeX scripts to implement interlinear texts.

There are also some reasons related to data structure for considering XLingPaper over alternatives. Interlinear glossed texts are often stored in one of a few formats: ELAN files,³⁷ FLE_x Text files,³⁸ Standard Format files,³⁹ L^AT_EX files,⁴⁰ custom project-specific XML files, or relational databases such as MySQL, PostgreSQL, or FileMakerPro. Moving content from analysis and markup tools to typesetting tools is an ever present need for linguists. Several

³⁴ github.com/tpellard/typgloss

³⁵ ctan.org/pkg/langsci-gb4e

³⁶ eva.mpg.de/lingua/pdf/Glossing-Rules.pdf

³⁷ archive.mpi.nl/tla/elan

³⁸ software.sil.org/fieldworks

³⁹ software.sil.org/toolbox

⁴⁰ For examples see [41] and [43].

Una frase cuantificadora puede acompañar al sustantivo (véanse [Los Cuantificadores](#) y [Los Números Cardinales](#)). Cuando se presenta esta frase, siempre va delante del núcleo de la frase nominal, como en los ejemplos en (2).

- (2) a. [tcf- NÁa majñuu nákhü iduu iya'
Zila] náā māhjūù" nákù idūū ījā?
LOC entre TOT.cuatro ojo.3SG agua
'De entre los cuatro manantiales'⁵ [Smajiin:6]
- b. [tpl- Gí'doo witsu rakhóó mikhúdú
Tlac] EST.tener.3SG cinco nariz.3SG (EST).picud@
'Tiene cinco esquinas picudas' [FC:5.1]

El cuantificador puede presentarse en construcciones donde no hay sustantivo expreso, como se explica en [Los Cuantificadores](#). Un ejemplo se incluye aquí.

Figure 5: Interlinear example from [29]. Note the example numbers on the left followed by example groups (a) and (b). Each interlinear then also has a language indicator in square brackets. Customization allows for as many rows per group as is required. Finally, on the right the hyperlinked citation to the reference for the source text is indicated.

- (9) a. ko-sis-o [— — —] move forward
b. kɔ-kij-a [— — —] act
c. ko-ʃund-o [— — —] break
d. kɔ-ʃut-a [— — —] become long
e. ko-ʃep-o [— — —] wink
f. kɔ-kɛk-a [— — —] decorate
g. ko-sok-o [— — —] cackle
h. kɔ-mvɔd-a [— — —] suck
i. kɔ-bab-a [— — <] carry

Figure 6: List of words as seen in [37]

tools such as ELAN and FLEEx have well-established workflows for data transfer [40]. FLEEx is often considered the tool of choice for many field linguists, language documenters, and lexicographers. For many linguists entering the field, it is the tool of choice over older tools like Toolbox (which uses standard format files) due to built-in collaborative features and grammar parsers [4]. Interlinear text in FLEEx can be exported into XML and the data used within XLPaper documents. This presents FLEEx users the opportunity to typeset their texts rather easily. Enabled by XML's modular document referencing features, XLPaper documents can reference components. Using the XML document referencing strategy with XML-encoded FLEEx texts allows authors to reflow typesetting outputs easily if they make content changes in their FLEEx environment.

XLPaper does not have a direct ELAN import process. However, we have had reports of linguists using the FLEEx-XLPaper publication pathway to typeset ELAN texts in L^AT_EX documents. One user reports capturing the X_QL^AT_EX document prior to rendering and then copying the relevant T_EX sections to their primary document and adding any required packages required by XLPaper to the header of their primary T_EX document.

Still finally, there is the matter of page layout. The main types of variation in page layout we have seen include the grouping of lines into sets or subsets (see Figure 5 for example), the labeling of sets and subsets, wrapping of interlinear glosses across lines (recall that these may themselves include three or more lines), and the alignment of the various elements of content within the interlinear glosses. We have seen word and morpheme aligned interlinears. XLPaper automatically wraps interlinears which makes the author's job much easier. Figure 7 in FC:1 and FC:2 demonstrate the wrapping of interlinear glossed texts. It does so by formatting each aligned word in an hbox and then having X_QL^AT_EX put them together in a hanging indent paragraph. This is based on the work of Kew & McConnel 1990 [21].

4.3 Gloss abbreviations

Linguists use two types of abbreviations. First, they might use abbreviations for names, titles, or commonly used words. This is much like standard

Rikha²

FC:1

Rikha rígi' najmaꞑ náꞑ yúoo' rꞑ'khꞑ ká',³ rꞑ'khꞑ suan'⁴
 flor.de.calabaza INAN:PROX IMPF.producirse LOC guía.3SG calabaza.especie calabaza.especie

khamí náꞑ yúoo' rꞑ'khꞑ májin'.⁵
 y LOC guía.3SG chilacayote

'La flor de calabaza se da en la guía de la calabaza de Castilla, de la "calabaza espina" y del chilacayote.'

FC:2

Rí rikhoo rꞑ'khꞑ suan', nagí'dyꞑ namídi rí
 SBD:INAN flor.de.calabaza.3SG calabaza.especie IMPF.empezar.3SG.FM ± IMPF.florear SBD:INAN

gꞑn' agóstó.
 luna agosto*

'La flor de la "calabaza espina" empieza a abrir en el mes de agosto.'

FC:3

Mba'jꞑ, muꞑmu' rí'jiꞑꞑ.
 (EST).grande:PL (EST).amarill@ flor.3SG

'Sus flores son grandes y amarillas.'

Figure 7: Wrapped interlinear text as seen in [31].

publishing. The second way that linguists use abbreviations is for indicating the grammatical meaning of pieces of words (morphemes). This second usage is often referred to as ‘glossing’ with the abbreviations referred to as ‘glosses’. One common set of glosses is the *Leipzig Glosses*. Leipzig Glosses, however, are not universally used for several reasons including:

- some authors have established their own tradition within their works which they started prior to the release of the Leipzig Glosses,⁴¹
- the typeset examples are quoted from a database which does not use Leipzig Glosses,
- they are not comprehensive, and
- they are not theoretically sufficient for some linguists.

XLingPaper supports both types. XLingPaper approaches this by providing built-in access to Leipzig Glosses, but also allowing the author to fine-tune a set of abbreviations and their definitions. When producing the output, XLingPaper creates hyperlinks between the abbreviation and its definition. This allows readers to quickly find the meaning of glosses and for the automatic generation of a table or list of abbreviations used.

4.4 Bibliographies

For better or worse XLingPaper has rolled its own bibliography solution. Import options are provided

⁴¹ For examples of the variation and scope of coverage consider the works of Greville Corbett, William Croft, Denis Creissels, and Martin Haspelmath.

Chao, Yuen Ren. 1930. ə sistim əv "toun-letəz" [A system of "tone-letters"]. *Le Maître Phonétique (Troisième Série du Le Maître Phonétique)* 30. 24–27.

赵元任 [Chao, Yuen-Ren]. 1980. 一套标调的字母 (英文). *方言* 1980(2). 81–83.

Chelliah, Shobhana Lakshmi, Willem Joseph de Reuse. 2011. *Handbook of descriptive linguistic fieldwork*. Dordrecht, Netherlands; New York: Springer. doi:10.1007/978-90-481-9026-3

Chen, Yiya & Carlos Gussenhoven. 2015. Shanghai Chinese. *Journal of the International Phonetic Association* 45(3). 321–337. doi:10.1017/S0025100315000043

Cheung, Kwan-hin [張群顯]. 2016. Chao Tone Letters: Original theory Versus Current Practice. In 錢志安, 郭必之 and 鄭嘉彥, *Commemorative Essays for Professor Yuen Ren Chao: Father of Modern Chinese Linguistics* 現代漢語語言學之父 — 趙元任先生紀念論文集, 65–76. 臺北市 [Taipei City]: 文鶴出版有限公司 [Crane Publishing Company].

Figure 8: An XLingPaper bibliography demonstrating mixed Latin and Chinese scripts.

for MODS and EndNote XML formats. This enables users to import from tools like EndNote,⁴² Zotero,⁴³ and JabRef.⁴⁴ XLingPaper uses custom T_EX scripts to output T_EX code for final rendering. It does not rely on B_IB_TE_X or B_IB_LA_TE_X. Figure 8 shows an example of a bibliography created with XLingPaper.

5 Outputs L^AT_EX allow that others do not

While XLingPaper has a large array of linguistically-oriented formatting capabilities across all output formats, there are some that only the X_LL^AT_EX output can produce. This is, of course, due to the formatting power of T_EX and X_LL^AT_EX.

⁴² endnote.com

⁴³ zotero.org

⁴⁴ jabref.org

5.1 Automatically wrapping interlinears

One of the most popular features of XLingPaper is its ability to automatically wrap long interlinear examples and lines in interlinear texts. As seen in Figure 7, wrapping occurs for the glossed text tiers. Not shown, but equally important, is that free translation tiers which are often longer than glossed lines also wrap neatly and coherently.

5.2 Font rendering

X_QL^AT_EX renders fonts extremely well. We show three cases where XSL-FO (via RenderX) and/or XHTML outputs have text rendering issues while X_QL^AT_EX does not.

First, when a line of text contains material rendered in different fonts on the same line, the two fonts may not line up evenly in the vertical direction. See Figure 9. This mismatch is due to the two fonts having different ascender and descender values. In order to overcome this when using XSL-FO, one has to add custom commands to deal with the font that differs from the primary font.

Second, the RenderX way of producing PDF cannot handle stacked diacritics, but the X_QL^AT_EX way does it very well. See Figure 10.

Third, X_QL^AT_EX can even handle special features requiring Graphite⁴⁵ processing. Figure 11 illustrates the special font handling needed for the Awami Nastaliq font. Of the three output renderings, only X_QL^AT_EX renders these correctly. Note that for this particular font, one must use XLingPaper’s X_QL^AT_EX package from T_EXLive 2020.⁴⁶ The Graphite included in the 2010 version of X_QL^AT_EX was not capable of rendering swami well.

5.3 Hyphenation for non-English languages

Since we use the `polyglossia` package, one can write an XLingPaper document in any of the sixty-one non-English languages listed in the `polyglossia` documentation, indicate the language used in the XML attributes and XLingPaper will pass this information to X_QL^AT_EX which will hyphenate according to that language’s hyphenation rules.

5.4 Author contact information

XLingPaper allows one to define a set of contact information for authors. Only the X_QL^AT_EX output is able to format them correctly.

5.5 Vertical fill

For title page material, only the X_QL^AT_EX output allows using vertical fill between items. The other outputs require using `overt`, fixed spacing values.

6 Features other outputs have that the L^AT_EX output does not

X_QL^AT_EX does not allow for custom table cell padding and spacing. Having said that, the developer cannot remember any XLingPaper user ever asking for a way to do this for the X_QL^AT_EX output. It just looks great.

Background color is not available for section titles.

Section 11.17.1.1 “Known limitations of using X_QL^AT_EX” in the XLingPaper user documentation lists known problems.

7 Conclusion

While the XLingPaper approach to composing documents via DTD controlled user interface limitations has great value in and of itself, the fact that it can produce great looking output via X_QL^AT_EX makes the learning curve rewarding. We feel that being able to produce PDF via X_QL^AT_EX has made XLingPaper a fantastic tool for linguists.

A Custom T_EX commands

XLingPaper has a number of custom commands that enable it to handle various tasks in a way that is consistent with our desired outcomes. The following lists some of them in a schematic way:

Command for	Purpose
Table of contents	Store and retrieve page numbers; format the contents.
Lists	Numbered and bulleted lists with control over indents, etc.
Examples	Example number and example content, where the content can be a line, a list of lines, a set of words, a list of a set of words, interlinear, a list of interlinears, etc.

⁴⁵ graphite.sil.org

⁴⁶ software.sil.org/xlingpaper/xelatex-package-from-tex-live-2020/

- (16) a. Farsi: **bozorgan** "leaders" (16) a. Farsi: **bozorgan** "leaders" (16) a. Farsi: **bozorgan** "leaders"
 b. Gilaki: **bozorgan** "leaders" b. Gilaki: **bozorgan** "leaders" b. Gilaki: **bozorgan** "leaders"

Figure 9: Ascender/descender font differences: The RenderX output is on the left; XHTML output is in the middle; the X_qL^AT_EX output is on the right

- (1) a. **Duu** **gúḁ**, **mḁñ** **mlā-gə**. (1) a. **Duu** **gúḁ**, **mḁñ** **mlā-gə**.
 house DEM.3C3 1S.CONTR make.PFV-3C3 house DEM.3C3 1S.CONTR make.PFV-3C3
 ‘That house, it’s I who built it.’ ‘That house, it’s I who built it.’

Figure 10: Stacked diacritics on the third word from the left: The RenderX output is on the left; the X_qL^AT_EX output is on the right.

Indexes	Handle keeping track of XLingPaper’s indexing capability, including page numbers.	Greek	el	ell
Interlinears	Handle lines in an inter-linear text or example, including dealing with an ISO 639-3 code in an interlinear example.	Hebrew	he	heb
Block quotes	Handle special cases needed for block quotes.	Hindi	hi	hin
Table headers	Attempt to calculate a column’s width via its contents.	Hungarian	hu	hun
		Icelandic	is	isl
		Indonesian	id	ind
		Interlingua	ia	ina
		Irish	ga	gle
		Italian	it	ita
		Lao	lo	lao
		Latin	la	lat
		Latvian	lv	lav
		Lithuanian	lt	lit
		Lower Sorbian		dsb
		Malay	ms	msa
		Malayalam	ml	mal
		Marathi	mr	mar
		Nynorsk	nn	nno
		Occitan	oc	oci
		Polish	pl	pol
		Portuges	pt	por
		Romanian	ro	ron
		Russian	ru	rus
		Sanskrit	sa	san
		Scottish	gd	gla
		Serbian	sr	srp
		Slovak	sk	slk
		Slovenian	sl	slv
		Spanish	es	spa
		Swedish	sv	swe
		Syriac		syr
		Tamil	ta	tam
		Telugu	te	tel
		Thai	th	tha
		Turkish	tr	tur
		Turkmen	tk	tuk
		Ukrainian	uk	ukr
		Urdu	ur	urd

B Hyphenation supported languages

Language Name	Two letter code	Three letter code
Albanian	sq	sqi
Amharic	am	amh
Arabic	ar	ara
Asturian		ast
Basque	eu	eus
Bengali	bn	ben
Bulgarian	bg	bul
Catalan	ca	cat
Coptic		cop
Croatian	hr	hrv
Czech	cs	ces
Danish	da	dan
Dutch	nl	nld
English	en	eng
Esperanto	eo	epo
Estonian	et	est
Farsi	fa	fas
Finnish	fi	fin
French	fr	fra
Galician	gl	glg
German	de	deu

- University of Indiana Ph.D. dissertation. pqdtopen.proquest.com/doc/1625743679.html?FMT=ABS.
- [17] Frampton, John. 2006. *Pst-Asr: Tex Macros for Typesetting Autosegmental Representations*. Version:1.1. CTAN. bakoma-tex.com/doc/generic/pst-asr/pst-asr-doc.pdf.
- [18] Frampton, John. 2012. *ExpPex for Linguists: Example Formatting, Glosses, and Reference*. Version: 4.1. mathserver.neu.edu/~ling/tex/expex/base/doc/expex-doc.pdf.
- [19] Freitag, Constantin and Antonio Machicao y Priemer. 2019. *LaTeX-Einführung Für Linguisten*. Berlin, Germany: Humboldt-Universität zu Berlin. doi.org/10.13140/RG.2.2.29299.27682.
- [20] Holton, Gary. 2021. *Writing Your Dissertation with L^AT_EX*. Typescript. Hawai'i. Github.com. gmholton.github.io/files/DissertationWriting.pdf.
- [21] Kew, Jonathan and Stephen McConnel. 1990. *Formatting Interlinear Text*. Occasional Publications in Academic Computing, Number 17. Dallas, Texas: Summer Institute of Linguistics.
- [22] Kircher, Athanasius. 1652. *Œdipus Ægyptiacus, hoc est Vniuersalis Hieroglyphicæ Veterum Doctrinæ temporum iniuria abolitæ Instauratio*. Opus ex omni Orientalium doctrina & sapientia conditum, nec non viginti diuersarum linguarum autoritate stabilitum, Romæ: Ex Typographia Vitalis Mascardi.
- [23] Knuth, Donald Ervin. 1984. *The T_EXbook*. A. Computers & typesetting. Reading, Massachusetts: American Mathematical Society; Addison-Wesley.
- [24] Knuth, Donald Ervin. 1986. *T_EX: The Program*. B. Computers & typesetting. Reading, Massachusetts: Addison-Wesley.
- [25] Lamicela, Andrew Charles. 2020. *Distinguishing Passive from MP2-marked Middle in Koine Greek*. University of North Dakota M.A. thesis. commons.und.edu/theses/3277.
- [26] Lehmann, Christian. 2004. *Interlinear morphemic glossing*. In *Morphologie: Ein internationales Handbuch zur Flexion und Wortbildung Morphology: an international handbook on inflection and word-formation*, edited by Geert E Booij, Christian Lehmann, Joachim Mugdan, and Stavros Skopeteas, 2:1834–57. Handbücher zur Sprach- und Kommunikationswissenschaft Handbooks of Linguistics and communication science 17. Berlin; New York: Walter de Gruyter. doi.org/10.1515/9783110172782.2.20.1834.
- [27] Liter, Adam. 2017. *L^AT_EX Workshop (for Linguists)*. adamliter.org/content/LaTeX/latex-workshop-for-linguists.pdf.
- [28] Lovell, Douglas. 1999. T_EXML: Typesetting XML with T_EX. *TUGboat*. 20 (3): 176–183. tug.org/TUGboat/tb20-3/tb64love.pdf.
- [29] Marlett, Stephen A, (compiler). 2012. La Frase Nominal. In Stephen A. Marlett (ed.) *Los Archivos Lingüísticos Me'phaa*. Instituto Lingüístico de Verano, A.C. [SIL International in Mexico]. mexico.sil.org/publications/i-wpindex/work_papers_-_mephaa_grammar_files.
- [30] Marlett, Stephen A. 2019. *Phonology From the Ground Up: The Basics*. Dallas, Texas: SIL International. sil.org/resources/archives/79207.
- [31] Neri Méndez, Emilia and Stephen A. Marlett. 2011 (Nov). Presentación Analítica del Texto “Flor de Calabaza”. In Stephen A. Marlett (ed.) *Los Archivos Lingüísticos Me'phaa* (versión preliminar). Instituto Lingüístico de Verano, A.C. [SIL International in Mexico]. mexico.sil.org/publications/i-wpindex/work_papers_-_mephaa_grammar_files.
- [32] Nordhoff, Sebastian, and Stefan Müller. 2020. *Language Science Press Guidelines*. Berlin, Germany: Language Science Press. langsci.github.io/guidelines/latexguidelines/LangSci-guidelines.pdf.
- [33] Paterson III, Hugh J. 2021. *Language Archive Records: Interoperability of Referencing Practices and Metadata Models*. University of North Dakota M.A. thesis. commons.und.edu/theses/3937.
- [34] Pellard, Thomas. 2019. Automatic formatting of interlinear glosses with L^AT_EX. *Cipanglossia* cipanglo.hypotheses.org/1221.
- [35] Peter, Steve. 2004. T_EX and Linguistics. *TUGboat* 25 (1): 58–62. tug.org/TUGboat/tb25-1/peter.pdf.
- [36] Machicao y Priemer, Antonio, and Constantin Freitag. 2019. *LaTeX-Einführung Für Linguisten*. Presentation at the MGK Workshop – SFB 1412, Berlin. linguistik.hu-berlin.de/de/staff/amp/latex20sfb/07-141-math2-trees-handout.pdf.
- [37] Rasmussen, Kent. 2018. *A Comparative Tone Analysis of Several Bantu D30 Languages (DR Congo)*. University of Texas Arlington Ph.D. dissertation. hdl.handle.net/10106/27483.

- [38] Rastorgueva, V. S., A. A. Kerimova, A. K. Mamedzade, L. A. Pireiko, and D. I. Edel'man. 2012. *The Gilaki Language*. Edited by Ronald M. Lockwood. Acta Universitatis Upsaliensis; Studia Iranica Upsaliensia 19. Uppsala, Sweden: Acta Universitatis Upsaliensis. urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-182789.
- [39] Rei, Fukui. 1996. TIPA: A System for Processing Phonetic Symbols in LATEX. *TUGboat* 17 (2): 102–14. tug.org/TUGboat/tb17-2/tb51rei.pdf.
- [40] Salfner, Sophie, and Tim Gaved. 2014. Working with ELAN and FLEx Together: An ELAN-FLEx-ELAN Teaching Set. Electronic Manuscript. SOAS, London, England. soas.ac.uk/elar/helpsheets/file122785.pdf.
- [41] Schenner, Mathias, and Sebastian Nordhoff. 2016. Extracting Interlinear Glossed Text from LATEX Documents. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, edited by Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Sara Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, et al., 4044–48. Portorož, Slovenia: European Language Resources Association (ELRA). aclweb.org/anthology/L16-1638.pdf.
- [42] Smith, Zac, Todd Snider, and Mia Wiegand. 2016. *LaTeX and Linguistics - How to Make Your Research Pretty*. Presentation at the Cornell Linguistics Circle, Cornell, New York. conf.ling.cornell.edu/miawiegand/Latex_Slides.pdf.
- [43] So Miyagawa, and Vincent W.J. van Gerven Oei. 2021. Building Web Corpus of Old Nubian with Interlinear Glossing as Digital Cultural Heritage for Modern-Day Nubians. In *The Proceedings of the 11th Conference of Japanese Association for Digital Humanities*, vol. 2021. 144–147. Tokyo: Historiographical Institute, The University of Tokyo. hi.u-tokyo.ac.jp/JADH/2021/Proceedings_JADH2021_rev0905.pdf.
- [44] Thiele, Christina. 1995. T_EX and Linguistics. *TUGboat* 16 (1): 42–44. tug.org/TUGboat/tb16-1/tb46ling.pdf.
- [45] Thiele, Christina. 2007. *Christina Thiele Interview by Dave Walden for the T_EX Users Group*. Transcript. tug.org/interviews/thiele.html.
- [46] Unicode Consortium, ed. 1991. *The Unicode Standard: Worldwide Character Encoding*. Version 1.0. Reading, Massachusetts: Addison-Wesley. unicode.org/versions/Unicode1.0.0.
- [47] Ward, Monica. 2002. Reusable XML Technologies and the Development of Language Learning Materials. *ReCALL* 14 (2): 285–94. doi.org/10.1017/S0958344002000629.
- [48] Wilson, Peter. 2007. The Memoir Class. *TUGboat* 28 (2): 243–46. tug.org/TUGboat/tb28-2/tb89wilson.pdf.
- [49] Wilson, Peter. 2021. *The Memoir Class for Configurable Typesetting: User Guide*. version 3.70. Normandy Park, WA: The Herries Press. texdoc.org/serve/memoir/0.
- [50] Wood, Joyce Kathleen. 2012. *Valence-Increasing Strategies in Urim Syntax*. Graduate Institute of Applied Linguistics M.A. thesis. diu.edu/documents/theses/Wood_Joyce-thesis.pdf.
 - ◇ H. Andrew Black
[blackhandrew \(at\) gmail dot com](mailto:blackhandrew@gmail.com)
 - ◇ Hugh J. Paterson III
[i \(at\) hp3 dot me](mailto:i (at) hp3 dot me)
<http://hp3.me>