# TeXFolio — a framework to typeset XML documents using TeX

Rishikesan Nair T., Rajagopal C.V.,
Radhakrishnan C.V.

## Abstract

TeXFolio is a web-based framework on the cloud
to generate standards-compliant, hyperlinked, book-
marked PDF output directly from XML sources with
heavy math content, using TeX. TeXFolio is a com-
plete journal production system as well. It can pro-
duce strip-ins which are alternate GIF or SVG images
of MathML content. In addition, DOI look-up, HTML
rendering of XML/MathML, and the whole dataset
generation according to the customer's specification
are also possible. Customer-specific validation tools
can be integrated in this system.

## 1 Introduction

TeXFolio is a web-based complete journal production
system that accepts XML documents as input and
generates a variety of outputs per user directives. At
the moment, TeXFolio accepts documents tagged per
the NLM/JATS or Elsevier Journal Article DTD. The
typesetting engine is TeX, which allows hyperlinked,
bookmarked and standards-compliant PDF outputs
of infinite variants in terms of look and feel. It further
permits TeX authors to directly edit their documents
at the proof stage and a master copy editor can pass
it for publishing with minimal loss of time.

Although the underlying engine of TeXFolio
hasn't deviated from the genre of free/libre software,
the computing paradigms have markedly shifted to
those in vogue to make the system modernized and
competitive in terms of usability, technologies and
performance. In fact, TeXFolio allows users to under-
take any stage of work anywhere in the world, owing
to its absolute compatibility with cloud and mobile
computing. That is further augmented by the usage
of LaTeX3 methodologies and programming to perfect
the production system to an efficient, automated and
accurate one.

## 2 The workflow

### 2.1 XML first

The input can be either XML or LaTeX. XML must be
per NLM/JATS DTDs or the Elsevier Journal Article
DTD. TeXFolio ingests these two types of XML
documents and generates a LaTeX file by applying
corresponding an XSLT stylesheet over the document.
This LaTeX file is used to edit the content and/or to
generate PDF output.

The processing cycle of XML $\rightarrow$ LaTeX $\rightarrow$ edit $\rightarrow$
PDF can be repeated any number of times, as shown
in the schematic diagram provided in Figure 1.

Whenever the LaTeX file is edited, the user can
either generate PDF or go through another cycle of
XML $\rightarrow$ LaTeX $\rightarrow$ PDF to ensure high fidelity between
XML and PDF, thereby making it a truly XML-first
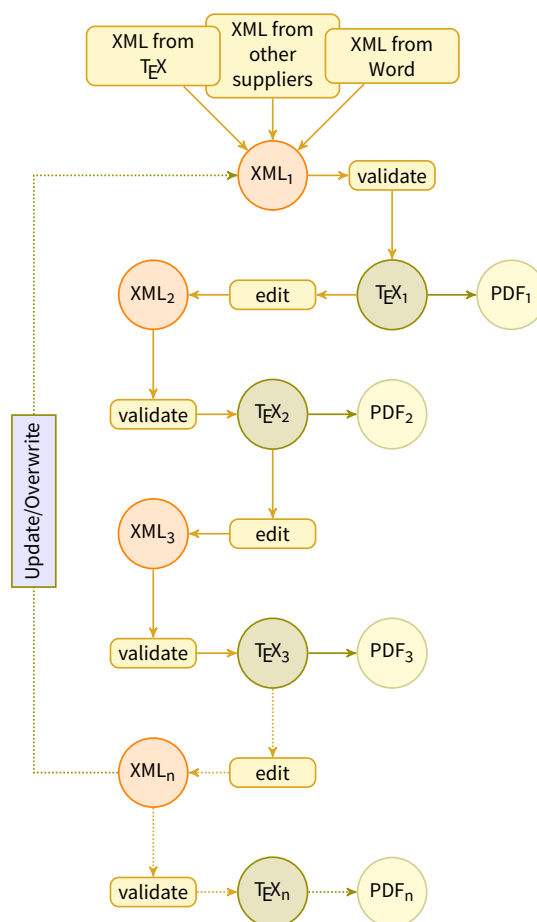workflow. We call it the "XROUND" process.



**Figure 1**: XML-first workflow

### 2.2 LaTeX input

For LaTeX input, we need to restructure the docu-
ment to augment XML generation since the front
matter is very verbose and structured in a granular
way in XML documents. The bibliography must be
provided as a BibTeX database. Barring the front
and back matters, the main body of the document
does not pose any problem during XML generation.
The system can digest all the author's macro defi-
nitions, `newtheorem`-type declarations, or any other
declarations provided by the AMS math packages
and other standard LaTeX packages.

TEX4ht is our preferred engine to translate LaTeX documents to XML format among all the tools and available software around. Since TEX is being used to process the document, it can easily and effectively handle all the macros and user definitions or declarations seamlessly. See Figure 2 for a schematic diagram of the workflow.
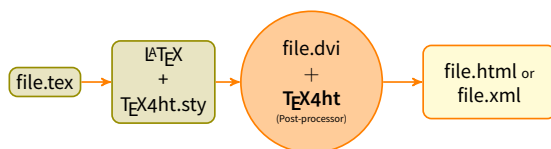


**Figure 2**: LaTeX to XML generation diagram.

## 3 Features

A list of the most commonly used features built into TEXFolio are given here.

### 3.1 General

**Cloud** With the development of TEXFolio, we are able to do text processing in the cloud rather than on a personal computer. This is critically important in the current scenario. Previously, text processing software was installed on each user's desktop and updating the systems with the most recent changes always caused problems. In the text processing world, much software is subject to regular updates and ensuring these updates was a herculean task. With the deployment of TEXFolio, we need to update the software on the server only and users' desktops need not be touched.

**TEX** For the whole process (i.e., generate XML, PDF, strip-ins, etc.), we are using TEX and friends. TEX being the most sophisticated software for high-quality typesetting, TEXFolio ensures high-quality output especially for mathematics, computer science, economics, engineering, linguistics, physics, statistics.

**Low learning curve** With the help of user-friendly and an attractive interface, even a novice user can use it without much learning. Since standard LaTeX commands are used, a normal LaTeX user can quickly start using it.

**Cross-platform** Since it is browser-based, users with different operating systems can access it without any difficulty.

**Browser, sole software** A desktop with any current browser and an Internet connection is all that is required to access TEXFolio. Instead of a desktop machine, if you have a Raspberry Pi, that is more than enough.

**Self-publishing** Supporting self-publishing is another feature of TEXFolio. TEXFolio serves well the "author as publisher" since it accepts TEX documents as input and can easily generate PDF, HTML 5, NLM/JATS or Elsevier Journal Article XML outputs if the sources are marked up per the `elsarticle` or `stm` document classes. These are the deliverables required for a web platform. As mentioned above, the bibliography needs to be provided as a BibTEX database.

### 3.2 Inputs

TEXFolio can currently accept the following input formats. The first is for a LaTeX-first workflow whereas the second and third are for an XML-first workflow.

**LaTeX** In the case of a LaTeX-first workflow, the input source has to be structured according to `neptune.cls`, `elsarticle.cls` or `stm.cls`.

**NLM/JATS XML** The user loads the XML file in one of these DTDs. TEXFolio generates a TEX file from this loaded JATSXML. From now on, the source will be this machine-generated TEX file. Using this, the user can paginate, make changes and/or generate different deliverables from this source.

**Elsevier Journal Article XML** The user loads an Elsevier DTD XML file. TEXFolio generates a TEX file from the input. Just as with NLM/JATS, from now on, the source will be the machine-generated TEX file, and the user can generate different deliverables from this source.

### 3.3 Outputs

Whatever the input (i.e., LaTeX or XML), the user can generate the following output files from the source in a few seconds.

**PDF** The PDF output generated will be according to the standards. You may generate a web version PDF as well as a print version PDF. The Web version will be hyperlinked, bookmarked and according to the PDF/A-1 standards whereas the print ready or fat PDF will be according to PDF/X-1a standards. It can easily be configured by a developer if the user wants a PDF according to another ISO standard.

**HTML5 + MathML** This is another deliverable or output which can be generated from the LaTeX or XML workflow. MathJax is supported.

**XML** The user can always generate a client version XML file in either the LaTeX-first or XML-first workflow. The XML used for generating PDF can have processing instructions embedded if any vital instructions for TEX was lost in the
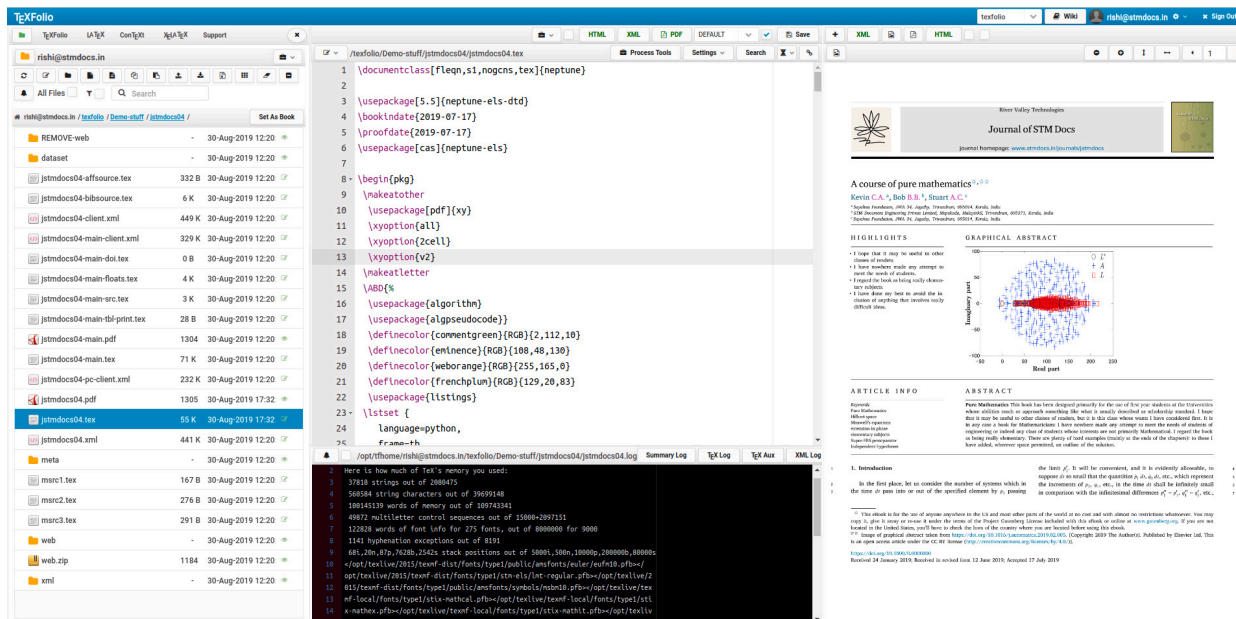
**Figure 3**: Main page of TEXFolio: File manager (left), text editor (center), output viewer (right), log windows (bottom), tools (top).

translation to XML. Before delivering to the client, these processing instructions can also be removed, and any other necessary changes made. Some publishers even require the XML document to be on a single physical line. Currently, it can ship either NLM/JATS XML + MathML or Elsevier Journal Article XML + MathML.

**e-pub** It is not one of the standard outputs at present, but is easily configurable.

**MathML** Math in the above XML documents will be tagged as MathML, SVG, GIF/PNG/JPG, TEX math, or all. Both MathML 2.0 and 3.0 are supported.

### 3.4 More features

1. Depends on LATEX3 methodologies and paradigms.
2. DOI link fetching and checking
3. Crossmark, ORCID, FundRef linking
4. Linking external objects such as Genbank accession numbers, PDB, CTGOV, OMIM, etc.
5. Source editing with track-changed source, PDF and XML at proof/final stages

6. Author proofing with Neptune
7. Tooltips in PDF
8. Technical support for TEX authors
9. NLM validator to check against NLM Article Publishing, CrossRef Deposit Schema and PMC style checker
10. Supports pdfLATEX, ConTEXt and X LATEX for PDF creation. However, for standard-compliant PDF generation, pdfLATEX is used for the time being.

### 4 Summary

Figure 3 shows the main page of TEXFolio.

For more details and screenshots, please visit `https://texfolio.org`.

◇ Rishikesan Nair T., Rajagopal C.V., Radhakrishnan C.V.
STM Document Engineering Pvt. Ltd., River Valley Campus, Mepukada, Malayinkil, Trivandrum 695571, India
rishi (at) stmdocs.in
http://stmdocs.com