

# DIGITAL HUMANITIES FOR PHILOSOPHY OF MATHEMATICAL PRACTICE

# The DH4PMP LATEX-pipeline for processing arXiv sources

Henrik Kragh Sørensen

Digital Humanities for Philosophy of Mathematical Practice, Section for History and Philosophy of Science, Department of Science Education, University of Copenhagen, Denmark. henrik.kragh@ind.ku.dk.

February 8, 2022

#### Abstract

The purpose of the DH4PMP IATEX-pipeline is to provide a robust way of processing IATEX source into an abstract representation that preserves the features that interest us. This note documents the purpose, construction, and affordances of the pipeline, while Sørensen (2021) documents the present stage of running the pipeline and internal assessments of the quality of data, including descriptive statistics. Additional papers such as SÃ, rensen (2021b) report on the construction of particular datasets from the pipeline and various external validations of the data.

### 1 Introduction

The arXiv poses a rich corpus of research prepublications which can be utilized to explore philosophical questions. Some such questions might be answerable just from the posted PDF files, but the fact the vast majority of papers in the arXiv are written in LATEX (especially in the mathematics and computer science categories) offers additional possibilities (Mejia-Ramos et al., 2019). LATEX is best understood as a programming language producing typeset text and mathematics. Thus the structure of the LATEX sources provide additional structural detail to the otherwise hierarchically organized linear text. Thus to understand and exploit the structural information, we have set up a pipeline which inputs LATEX source and provides a representation of a reduced LATEX source, split into so-called *contexts*. Using this pipeline, it is possible to compare e.g. word prevalence between the *theorem context* and the *proof context*.

# 2 LATEX: structured mathematical typesetting

 $I^{A}T_{E}X$  is a further development of  $T_{E}X$  which was designed and developed by Donald Knuth and first released in 1978 (Knuth, 1986). Knuth's  $T_{E}X$  was designed to typeset complicated

1

mathematical formulas, and it holds tremendous power but is also (now) considered very difficult to master. Instead, mathematicians, computer scientists and many other academics have adopted LATEX which was written during the early 1980s by Leslie Lamport as an extension of TEX (Lamport, 1994). LATEX forms the core of a very complex and evolving software system to which new functionality is constantly added through the publication of *packages* on the opensource repository ctan.org. In addition, LATEX has been compiled and distributed for all the major operating systems, and document editors have been developed to ease the composition of text and online collaboration.<sup>1</sup>

# 2.1 Commands and environments

The power of  $\mathbb{IATEX}$  is that it is actually a *programming language* intended for producing typeset text. The central programming constructs are *commands* and (as extensions thereof) *environments*. A command (or *macro*) has a name, is prefixed by a backslash, can take arguments (either optional or mandatory), and expands to its definition through the  $\mathbb{IATEX}$  compilation process. An environment is technically a set of combined commands to *begin* and *end* the environment. These commands can be taylored to specific applications, and are often used to scope together elements that belong together, for instance in a list or a theorem declaration.

A lot more can be said about  $IaT_EX$  formatting and programming (see Lamport, 1994), but here it suffices to know that  $IaT_EX$ :

- is a widely used system for typesetting articles and books within certain academic disciplines including mathematics, computer science, and physics,
- is a programming language in which *source files* are written as ordinary text files which are *compiled* into the output (a PDF-file) by running the LATEX compiler either from a command prompt or from a dedicated editor,
- has the potential to organize the input text in a hierarchy using sectioning commands as well as in a structure using environments,
- is extremely customizable through the use of *packages* and can include input from secondary sources such as code listings or tables,
- is able to handle comments that do not produce output but can serve as notes, messages, outlines etc. between the authors.

### 2.2 Accessing arXiv data

Our raw data consists of essentially three distinct sets of files:

- 1. A JSON file containing the metadata of all articles published on the arXiv. This file is deposited from the arXiv to Kaggle and scheduled for monthly updates.<sup>2</sup>
- 2. A collection of PDF files obtained from Kaggle via a Google cloud storage.
- 3. A collection of TAR files containing the LATEX source of articles. These are obtained from the AWS<sup>3</sup> for which there is a requester fee.

<sup>&</sup>lt;sup>1</sup>The  $L^{A}T_{E}X$  software system is robust in the sense that its produced output is highly platform independent, whereas small changes in the source may actually lead to substantial differences in the output.

<sup>&</sup>lt;sup>2</sup>See https://www.kaggle.com/Cornell-University/arxiv. To download the metadata from Kaggle, you will need to create a free account.

<sup>&</sup>lt;sup>3</sup>See https://arxiv.org/help/bulk\_data\_s3.



Figure 1: The output of the a simple  ${\rm I\!AT}_{\rm E}\!X$  source, indicating various structural and hierarchical elements.

```
1 \documentclass{article}
  \usepackage[T1]{fontenc}
2
3
  \usepackage[latin1]{inputenc}
4
5 \usepackage{amsthm,amsmath}
6 \newtheorem{theorem}{Theorem}
7
  \newtheorem{lemma}{Lemma}
8
  \begin{document}
9
10
  \author{The author(s) of the paper, given by a \emph{command}}
11
  \title{The title of the paper, given by a \emph{command}}
12
  \date{The date of the compilation of the paper, given by a
13
14
     \emph{command} or defaulting to today}
15
  \maketitle
16
17
  \thispagestyle{empty}
18
19
  \begin{abstract}
20
     Here goes the abstract of my paper into an \emph{environment}.
21
22
  \end{abstract}
23
  \section*{This is a section produced by a \emph{command}}
24
25
26 Here goes some ordinary text, referred to as the \emph{outer context}.
27
28
  \subsection*{There are a number of sectioning \emph{commands}, proving
     organization and (possibly) numbering}
29
30
  \begin{theorem}
31
     This is a theorem produced by an \emph{environment}
32
33
  \end{theorem}
34
35 % Some lines can be commented out and thus produce output while still
36 % being part of the source
37
  \begin{proof}
38
     And here goes the proof, again produced by an \emph{environment}.
39
  \end{proof}
40
41
42 Again, we can have outer-context text, and equations and figures can
43 be inserted either in the outer context or within an environment.
  \begin{align*}
44
     e^{i\pi}+1
                 = 0.
45
  \end{align*}
46
47
48
  \begin{lemma}
49
     Other environments, such as lemma, corollary, definition,
     proposition etc. may either be predefined in the style class or
50
51
     added by the user.
52 \setminus end{lemma}
53
  \end{document}
54
```

Figure 2: The LATEX source responsible for producing the output in Figure 1.



Figure 3: Schematic overview of the data processing.

From Kaggle, we first obtain a JSON file with metadata on all articles in the arXiv. These are then filtered to include only articles assigned to (at least) one mathematical category math.<sup>\*</sup>.

Using <code>gsutil</code>, we can download PDF files of the articles from, say, September 2020 by using the command:<sup>4</sup>

#### gsutil cp -m -r gs://arxiv-dataset/arxiv/arxiv/2009\* .

Here, we meet the arXiv way of indicating production time: the yymm format: two digits to indicate the year (after 2000) and two digits to indicate the month (0-prefixed). If we are only interested in specific articles (say the the mathematical ones), we filter which TAR-files to download based on the JSON metadata.

At the AWS, article sources are stored in chronological collections, so given the *manifest* we identify those collections that are relevant for us, download those and extract the relevant individual article sources. Each source for an arXiv article is collected in a (gzipped) TAR-file, which will typically contain IATEX source as well as optional figures, bibliography files etc.

In our 'backend', we can prioritize downloads from AWS to focus on specific months or specific categories; ultimately it is our intention to keep a full duplicate of the arXiv on AWS.

The result of this download process is (almost) identical to the result of downloading each of these article sources individually from arXiv, but in accordance with arXiv policy, we wish to help curtail the traffic to the site. However, individual article sources can later be supplemented by download from the arXiv. The only difference between the two approaches is that articles may be *updated* which will reflect in the directly downloaded source and PDF-files but may not always be incorporated in the TAR-file.

Each arXiv article is identified by its arxivID, a number of the format yymm.ddddd where yymm is the time indication of initial upload, and the digits are the separator indicates a running index (0-prefixed to either 4 or 5 places).

We process information from the raw data sources into *pandas* dataframes through a series of import routines written in *python* and constituting the DH4PMP LATEX-pipeline.

# 2.3 Processing LATEX sources

Considered as a programming language, the typical way to process LATEX sources would be to build 1) a generic parser which could parse LATEX source into an abstract representation and 2) a

<sup>&</sup>lt;sup>4</sup>See https://www.kaggle.com/Cornell-University/arxiv/discussion/200472.



Figure 4: Processing IAT<sub>F</sub>X sources.

specific backend for producing the desired output (PDF or other formats); see Figure 4. However, because  $LAT_{EX}$  is such a complex and powerful language (it is 'Turing complete'), no complete grammar was ever defined (Beebe, 2004, p. 15). Instead, the full description of the  $LAT_{EX}$  'programming language' is the implementation of  $T_{EX}$  as well as the adjunction of particular constructs in  $LAT_{EX}$  and the full set of packages extending that even further. And since Knuth's original  $T_{EX}$  implementation is written in a presently almost obsolete programming language (*WEB*), the typical contemporary distributions are actually made by compiling these into *C* using *Web2c* and then compiling the to the desired platform. Thus, accessing the intermediate abstract representations of the typeset is difficult, and instead we are often left to implement parsers and backends for a (small) subset of the entire expressive power of  $T_{EX}$ .

A number of processors for  $T_EX$  and  $I_AT_EX$  which provide the kind of structured representations that we want have been developed, including *latexML* and *pandoc* (Miller and Ginev, 2020) which are able to convert large subsets of  $I_AT_EX$  into XML or HTML (or a number of other markup formats) which could, in turn, be parsed with relative ease. However, to have control over the entire pipeline and have it implemented in *python*, we chose to build our own pipeline on top of the module *pylatexenc* (Faist, 2019).

# 3 The DH4PMP LATEX-pipeline

The purpose of the LATEX-pipeline, that we describe in the following, is to provide a robust way of processing LATEX source into an abstract representation that preserves the features that interest us. The core of the pipeline is our adaptation of the *pylatexenc* parser (walker) and processor (Faist, 2019). However, to overcome some important shortcomings of *pylatexenc*, we have had to add an additional level of preprocessing.

The pipeline consists of the following steps:

- 1. File access (depending on application area such as arXiv or MathSciNet parsing). For arXiv parsing, we rely on our class ArxivSourceManager.
- 2. Preprocessor, working on an individual file, doing the following:
  - a) Removing comments
  - b) Extracting certain information from the source files
  - c) Resolving static conditionals
  - d) Making adjustments to the source (primarily to balance tokens) before it can be submitted to *pylatexenc*
- 3. Parser, based on *pylatexenc*, producing a nodelist and implementing persistent serialization for subsequent recall.



Figure 5: Outline of the pipeline.

4. Postprocessors (depending on application) such as context splitters or linguistic analyses.

The pipeline is implemented as inheritance of classes for steps 3–4, and attribute classes for steps 1 and 2.

# 4 Preprocessor

The preprocessor uses relatively simple tools such as regular expressions and simple *lark*-based parsers (Shinan, 2021) to prepare the invidual  $IAT_EX$ -files for processing by *pylatexenc*. In particular, the preprocessor performs the following tasks:

1. Remove comments from the  ${\rm IAT}_{\rm E}\!{\rm X}$  source



Figure 6: Inheritance structure.

- 2. Remove conditionals from the LATEX source
- 3. Extract simple macro definitions to be expanded in step 4 and process environment constructions in the LAT<sub>F</sub>X source; see Figure 7
- 4. Perform simple macro expansion
- 5. Remove math constructs from the  $LAT_EX$  source.

A central aspect of the preprocessor is circumventing the nested balancing of braces and math environments. For instance, it is legal (and frequent in the arXiv) to write

 $x\in[0,1)$ 

which contains an unbalanced bracket.<sup>5</sup> Similarly, it is legal to shorthand math environments, writing e.g.

```
\newcommand\ba{\begin{align*}},
```

which includes an unbalanced opening of a math environment, albeit within a nested (and balanced) set of braces.

Thus, to consistently process such issues, we implemented a layered number of runs of parsers (multiple parse runs) to process  $LAT_FX$  files.<sup>6</sup>

Some of the grammars, in particular parts of gobble-conditionals.lark and of gobble-math.lark are fully or partially auto-generated from other configuration files. Then, depending on the run of the preprocessor, *lark* rules can be given different priorities to match the kind of lexing that is currently going on.

# 4.1 Stages of the preprocessor

The preprocessor works through a number of *stages* (sometimes referred to as *runs* or *jobs*), which constitute the preprocessor pipeline:

<sup>&</sup>lt;sup>5</sup>Such legal strings are exactly what messes with the *pylatexenc*, giving rise to the need for a preprocessor to handle them before passing the string on to *pylatexenc*.

<sup>&</sup>lt;sup>6</sup>Obviously, we could also have opted for a different model with increasing handling of special cases. This might still be a preferable option for some specific cases, and can then be implemented as a separate run of the preprocessor.

```
(newcmd | newcmdtwo | newenv | docclass | inputfile
  elems:
1
                                                                          | ref
      | cite | includegraphics| letstmt | elem)*
   // | ref | includegraphics | cite
\mathbf{2}
3
  ?elem.999:
4
                      string
                | SEP
5
6
   //WRONG.999: "\\include" /[a-zA-Z]+/
7
8
   newcmd.100:
                    NEWCMD (MACRONAME | arg) [optarg [optarg]] (
9
      MACRONAME | exparg)
                      DEF DEFMACRONAME defopts ["="] [SEP] (
10
                  /
   //
      MACRONAME (exparg)
                | DEF DEFMACRONAME ["="] [SEP] (MACRONAME|exparg)
11
12
   //DEFOPTS.100:
                      /#\d/ | /##\d/ | /\[#\d\]/
13
  //DEFSEP.100:
                      /[\(\):\- ;,<>\/]/ / "\\@nil"
14
  //defopts.100:
                      ([DEFSEP+] DEFOPTS)* [DEFSEP+]
15
   //DEFMACRONAME.100:
                          (/ \setminus \{ \ \# \setminus \{ \setminus [] + / \} \}
16
  DEFMACRONAME.100:
                        /\\[^\{]+/
17
18
   newcmdtwo.101: NEWCMDTWO (MACRONAME|arg) [optarg [optarg [optarg]]] (
19
      MACRONAME | exparg)
20
21
  newenv.100:
                    NEWENV arg [optarg [optarg]] exparg exparg [SEP]
22
                   DOCCLASS [optarg] arg
23
  docclass.100:
24
                             INPUTFILE [SEP] (arg | CHARS+)
25 inputfile.100:
                             REF [SEP] arg | REF SEP | REF MACRONAME
26
  ref.100:
                             INCLUDEGRAPHICS [SEP] [optarg+] arg
  includegraphics.100:
27
  cite.100:
                             CITE [SEP] [optarg+] arg
28
29
30 CITE.20:
                            /\\cite[pt]?(?![a-zA-Z~])/ ["*"]
                            "\\includegraphics" ["*"]
31 INCLUDEGRAPHICS.20:
                            (/\\ref(?![a-zA-Z~])/ | /\\eqref(?![a-zA-Z~])/ |
32 REF.20:
       /\\autoref(?![a-zA-Z~])/) ["*"]
   INPUTFILE.20:
                            (//\input(?![a-zA-Z~])/ | //\include(?![a-zA-Z
33
      ~])/)
34
35
   //COMMANDTOKENS.20: REF | CITE | INCLUDEGRAPHICS | INPUTFILE
36
                    ("\\newcommand" ["x"] ["*"] | "\\providecommand" ["x"]
  NEWCMD.100:
37
      ["*"] | "\\renewcommand" ["x"] ["*"]) [SEP]
                   ("\\newcommandtwoopt") ["*"] [SEP]
   NEWCMDTWO.101:
38
                    ("\\newenvironment" ["*"] | "\\renewenvironment" ["*"])
   NEWENV.100:
39
       [SEP]
  DOCCLASS.100:
                    "\\documentclass" [SEP]
40
41
              /\\def(?![a-zA-Z~])/ [SEP]
42
  DEF.999:
43
                    LET (MACRONAME | CITE | REF) ["="] (MACRONAME | RBQ |
44
  letstmt.200:
      LBQ)
                 /\\let(?![a-zA-Z~])/ [SEP]
  LET.200:
45
46
                   (CHARS | MACRONAME | LBQ | RBQ | LBR | RBR ) +
47 ?string.10:
48 CHARS.0:
                     /[^\[\]\{\}]/
49 ?sbstring:
                      (CHARS+|MACRONAME|SEP|INPUTFILE) | arg | optarg
                 |REF|CITE
                      (CHARS+|MACRONAME|SEP|INPUTFILE|LBQ|RBQ) | arg |
50
  ?brstring:
                |REF|CITE
      newcmd
                                         9
51
52 LBQ.50: "["
                                                                                 s &
53 RBQ.50: "]"
```

gobble-comments removes comments in LATEX source, preserving the structure.

- gobble-conditionals resolves static conditionals (i.e. \iffalse and \iftrue). To process other static conditional constructs which also can contain \else and \fi (e.g. such as provided by LATEX classes and packages), it is (presently) necessary to supply a list of such conditionals to the processor. To aid in this, the processor picks out \newif constructs in the file, but these were found to account for only a small portion of the static conditionals.
- macro-extract extracts macro and environment definitions (such as arising from \newcommand or \newenvironment). For macros, it stores the literal definition for later. For for environments, it simply creates an empty environment, since we are only interested in the structure of the source, not the actually formatting of environments, including their nesting.
- macro-expansion performs literal expansion of previoulsy defined macros (picked up through macro-extract), resulting in a one-level literal expansion in each individual file.
- gobble-math removes all mathematical constructs, potentially replacing them with a customizable placeholder (not yet implemented).

The first of these can be (and is) implemented using regular expressions in *python*, but the remaining stages are all implemented using a (rather) simple LALR(1) parser constructed from EBNF-form through lark.<sup>7</sup>

Each stage of the preprocessor (after the first one) consists of a standard parser structure:

- 1. A *parser* which takes the input string and produces an AST. Here, we use the default 'contextual' parser provided in *lark*.
- 2. A transformer which makes transformations to the AST.
- 3. An *interpreter* which takes the AST and produces the output, in our case a new (valid) LATEX string with certain (irrelevant) problematic constructs processed (typically removed).

As a matter of convenience, we use the same transformer and interpreter for all stages in the pipeline.

# 5 Parser

The purpose of the parser is the more detailed processing of the  $IAT_EX$  source into a *nodelist* (an AST) that picks up on the important structural information in the sources. This is (at the moment) achieved through the customizable *pylatexenc* interface. The results of the parser are stored in a database structure for which nodelists for invidual papers can easily be retrieved.

As part of the parser, the following steps are performed:

- 1. The list of filenames is retrieved from the source TAR.
- 2. The main file is identified as one that includes a \documentclass macro, and this information was picked up in the preprocessor. If more than one file satisfies this, the paper is marked for follow-up by human intervention.
- 3. Each file in the TAR is extracted and stored based on its relative path. The contents of auxiliary files are then substituted for relevant \input commands in the main file, iterated until a stable text is reached.

<sup>&</sup>lt;sup>7</sup>It is possible that some stages can be combined, but they are split for purposes of debugging the complex process and for the possibility to assign different priorities to the same syntactic constructs at different stages.



Figure 8: Pipeline of processing LATEX source for structural and textual information.

- 4. At this point, the preprocessor is run, preparing a revised version of the expanded main file that can be processed by our customized *pylatexenc* parser. In the preprocessor, numerous structures are extracted, in particular citations to the bibliography, internal references in the document, inclusion of graphics files, and inline and displayed mathematics.
- 5. In the next step, the customized *pylatexenc* parser is run in the resulting text. As part of this phase, various environments are picked out and stored as *contexts* according to the many-to-one mapping described below.

### 6 Postprocessors

The purpose of a postprocessor is to turn the nodelist produced in the parser into formats suited for particular research questions. This could involve adding PoS (position-of-speech) tagging to identify e.g. nouns or verbs or associated clusters. The (currently) most important postprocessor splits the  $LAT_EX$  source text into *contexts* defined by a many-to-one relation between the environments used in the  $LAT_EX$  source and a set of predefined contexts (section 6.1).

# 6.1 Context splitting

The pipeline, see Figure 8 consists of a 'global' part which allows for configuration, and a 'local' part in which each article is processed individually based on the global configuration. In order to process the structure of the LATEX source, we first pick out *all* environments used in the corpus of sources by a regular expression search for the corresponding LATEX construct, i.e. \begin{. This provides us with a list of environments and their respective frequencies. This list is then processed by hand as an Excel sheet where a partial mapping is established which groups variant forms of similar environments together; for instance the construct \begin{thm} ... \end{thm} is equated with the theorem environment.

A *context* is defined through two sets of many-to-one mappings that combine environments from the LATEX sources into a small number of seperate structural contexts.

The output of this procedure is, for each  $IAT_EX$  mainfile (i.e. a file containing a \documentclass command and thus signifying a compilable  $IAT_EX$  source), the following:

Copyright © 2022 • www.dh4pmp.dk

Henrik Kragh SÄ,rensen (May 24, 2021a). The DH4PMP LaTeX-pipeline for processing arXiv sources. DH4PMP Work-in-Progress & Proof-of-Concept Series 12. Version 1. 18 pp. URL: https://www.erda.dk/vgrid/DH4PMP/wippoc/paper12.pdf.



Figure 9: Small extract from the two many-to-one maps which define contexts from environment names. The first one is intended to capture synonyms, abbreviations, and translations. The second one maps these into contexts.

- For each context, a list of lists of paragraphs. This lists are true to the sequence in which the paragraphs occur in the source.
- A list of structural elements combined with the level of nesting for the particular element.
- The default main contexts currently implemented are *definition*, *theorem*, *proof*, and *meta*. Additional contexts are also implemented, primarily to capture one-off environments such as *abstract* and *acknowledgements*. The two sets of mappings are defined in the file environments.xlsx.
- All environments which are not mapped as in Figure 9 are assigned to the *other* context, whereas the *outer* context contains all level-0 text, i.e. text not embedded in an environment. A few groups are not matched to a context; this currently applies to *algorithm*, *figure*, *table*, *diagram*, and *equation* groups.
- For each of the elements CITE\_EXTERNAL, REF\_INTERNAL, GRAPHICS\_FILE, and MATH\_DISPLAY, the occurance in the main file is substituted with a tag of the form MATH\_DISPLAY(1), which points to a list of extracted text from the original occurance. To allow reversing this extraction or studying particular elements such as displayed mathematics, the resulting data also contain lists of these elements. If one is to run linguistic analyses, it may be preferable to remove the indexes, e.g. running

re.sub(r'(MATH\_DISPLAY)\(\d+\)', g<1>, x)

on the text x.

The main contexts can be described as follows:

- Theorem contexts, in which knowledge claims are stated. These include environments such as theorem, proposition etc.
- Proof contexts, which corresponds to the proof environment.
- Meta contexts, which collect reflexive parts of the article and include environments such as example etc.
- Outer context, which captures the outer level of the source with the above-mentioned contexts removed. This can be considered the 'paratext'.

In addition to the text split into contexts, it is also possible to adjoin metadata either from the original metadata provided by arXiv (*arXiv Dataset 2020*), derived from it, or built from other studies. A full list of the variables that can be adjoined can be obtained, but the most applicable ones are summarized in Appendix A.

# 7 Building a corpus from the DH4PMP LATEX-pipeline

To extract a corpus of articles split into contexts from the DH4PMP  $LAT_EX$ -pipeline, one must form a *corpus*. In this situation, a corpus can be defined either

- 1. by specifying a list of arxivIDs to include, or
- 2. by specifying a selection criterion from the database, using the SQL commands of SQLite3, using the variables maincat or yymm (others can be added).

For instance, to build a corpus of all papers in math.GT from 2020, one would define the corpus as

#### maincat LIKE "math.GT" AND yymm LIKE "20\_\_"

where the underscores represent any character. Similarly, to define a corpus of all texts with a main mathematical category, one would specify the criterion

#### maincat LIKE "math.\_\_"

The implementation of the arXiv corpus is the class arXivCorpus derived from the abstract class Corpus. Its function export\_dataframe returns a dataframe of the corpus along with a dictionary of information about the corpus in a format that can be processed directly or included in the LATEX include file generated by the dh4pmp-paper class. The function export\_dataframe takes a number of parameters, determining which information is included in the resulting dataframe:

- contexts: A list of which contexts to include (or "all" to include them all).
- data: A list of which elements in the extracted data dictionary to include in the mappings (or "all"). These elements include REF\_INTERNAL, MATH\_DISPLAY, CITE\_EXTERNAL, and GRAPHICS\_FILE.
- metadata: A list of which metadata to include (or "all"). Presently, only yymm, maincat, mathcats can be included directly; other metadata can be added by a left-join with the JSON-metadata on arxivid.
- params: A list of which other parameters to compute (or "all"). Presently, includes lenghts, which counts the number of instances of each context, timestamp which includes the processing timestamp, and subversion which includes the version of the processing pipeline.
- minimum\_version: A string to indicate the minimum processing pipeline version to allow in the dataframe. Currently not implemented.

Each context column in the dataframe is (typically) returned in a JSON-encoded form, thus should be converted before it becomes a list (of lists):

```
df['proof'] = df['proof'].apply(json.loads)
```

The info dictionary contains information such as

- metadata:size indicating the size of the corpus as defined.
- database:size indicating the number of papers downloaded and processed which satisfy the corpus definition.
- corpus:size providing the resultant size of the corpus.
- metadata:definition containing the corpus definition (SQL).

# References

arXiv Dataset (Nov. 22, 2020). arXiv dataset and metadata of 1.7M+ scholarly papers across STEM. Cornell University. URL: https://www.kaggle.com/Cornell-University/arxiv/ metadata (visited on 01/30/2021).

Beebe, Nelson H. F. (2004). "25 Years of T<sub>E</sub>X and METAFONT. Looking Back and Looking Forward". TUG 2003 Keynote Address. *TUGboat*, vol. 25, no. 1, pp. 7–30.

Faist, Philippe (2019). pylatexenc. Version 2.8. URL: https://pylatexenc.readthedocs.io/.

Henrik Kragh SÃ, rensen (May 24, 2021a). The DH4PMP LaTeX-pipeline for processing arXiv sources. DH4PMP Work-in-Progress & Proof-of-Concept Series 12. Version 1. 18 pp. urL: https://www.erda.dk/vgrid/DH4PMP/wippoc/paper12.pdf.

Knuth, Donald E. (1986). The T<sub>F</sub>Xbook. Reading et al.: Addison-Wesley. ISBN: 978-0-201-13448-3.

- Lamport, Leslie (1994). LaTeX: A Document Preparation System. User's Guide and Reference Manual. 2nd ed. Reading (Mass.): Addison-Wesley Publishing Company.
- Mejia-Ramos, Juan Pablo et al. (2019). "Using Corpus Linguistics to Investigate Mathematical Explanation". In: Methodological Advances in Experimental Philosophy. Ed. by Eugen Fischer and Mark Curtis. London et al.: Bloomsbury Academic. Chap. 8, pp. 239–264.
- Miller, Bruce and Deyan Ginev (Nov. 17, 2020). LaTeXML. A LaTeX to XML/HTML/MathML Converter. URL: https://dlmf.nist.gov/LaTeXML/ (visited on 05/20/2021).
- SA, rensen, Henrik Kragh (May 23, 2021b). Validating the DH4PMP pipeline and corpus for studying explanations in the arXiv. DH4PMP Work-in-Progress & Proof-of-Concept Series 16.

Shinan, Erez (2021). Lark. URL: https://lark-parser.readthedocs.io/.

Sørensen, Henrik Kragh (Nov. 14, 2021). Status of the DH4PMP LaTeX-pipeline. DH4PMP Work-in-Progress & Proof-of-Concept Series 13.

www.dh4pmp.dk

Copyright © 2022

# A Metadata

abstract	:	Abstract
	-	

Catalog(s): basic

```
authors : Authors
```

Catalog(s): basic

authors parsed : Authors processed into list

Catalog(s): basic

categories : Categories separated by spaces

Catalog(s): basic

categories parsed : List of categories

Catalog(s): basic\_auto

Depends on: categories

#### **comments** : Comments on publication etc.

Catalog(s): basic

### ${\bf doi} \ : {\rm DOI} \ {\rm of} \ {\rm eventual} \ {\rm publication}$

Catalog(s): basic

 $e_Figs_pp$  : Number of figure (etc.) environments per page

Symbol:	$e(\mathrm{Figs})_{pp}$
Catalog(s):	env_auto
Depends on:	e_figure $\bullet$ e_tikzpicture $\bullet$ e_wrapfigure $\bullet$ numpages

**e\_figure** : Number of figure environments in source

Symbol:	e(figure $)$
Catalog(s):	env

#### e tikzpicture : Number of tikzpicture environments in source

Symbol:	e(tikzpicture)
Catalog(s):	env

e wrapfigure : Number of wrapfigure environments in source

Symbol:	e(wrapfigure)
Catalog(s):	env

 $\mathbf{f} \quad \mathbf{gif} \quad : \text{ Number of gif files in tar archive}$ 

Catalog(s): src

f images : Number of image files in tar archive

Catalog(s): src\_auto Depends on:  $f_{gif} \bullet f_{jpg} \bullet f_{pdf}$ 

f\_jpg : Number of jpg files in tar archive Catalog(s): src

f\_pdf : Number of pdf files in tar archive Catalog(s): src

 $\mathbf{f}$ \_tex : Number of tex files in tar archive

Catalog(s): src

 ${\bf f} \quad {\bf total} \quad : \ {\rm Total} \ {\rm number} \ {\rm of} \ {\rm files} \ {\rm in} \ {\rm tar} \ {\rm archive}$ 

Catalog(s): src

id : arXiv identification

Catalog(s):	basic $\bullet$ basic_auto $\bullet$ env $\bullet$ env_auto $\bullet$ pdf $\bullet$ src $\bullet$ src_auto
Details:	This is used as the key for combining the different metadata catalogs. Often
	renamed arxivid.

maincat : Main (first) category listed

Catalog(s): basic\_auto Depends on: categories\_parsed

mainmathcat : Main (first) math.\* category listed

Catalog(s):	basic_auto
Depends on:	categories_parsed
Details:	Notice that this variable <i>need not</i> refer to the <i>first</i> category assigned to the article, but only to the first math.* category listed.

missing files : List of files not found while parsing the source tar

Catalog(s): src

num\_missing\_files : Number of missing files in source tar

Catalog(s): src\_auto Depends on: missing\_files

**numauthors** : Number of authors

Symbol: #authors

Catalog(s): basic\_auto

Depends on: authors\_parsed

**numcats** : Number of categories

Catalog(s): basic\_auto Depends on: categories parsed

#### nummathcats : Number of math.\* categories

Catalog(s): basic\_auto Depends on: categories\_parsed

#### numpages : Number of pages in pdf file

Symbol: #pages Catalog(s): pdf

#### pdf filename :

Catalog(s): pdf

### **pdfnumwords** : Number of words in pdftext

Catalog(s): Depends on: pdftext

#### **pdftext** : Extracted text from pdf file

Catalog(s): pdf

#### $\mathbf{t\_meta}~:~ \mathrm{Extracted~meta-context~text~from~src~files}$

Catalog(s): env

t\_outer : Extracted paratext (outer-context text) from src files Catalog(s): env

### $\mathbf{t\_proof}~:$ Extracted proof-context text from src files

Catalog(s): env

#### ${\bf t}$ ${\bf theorem}$ : Extracted theorem-context text from src files

Catalog(s): env

#### texstructure : Structure of the tex document from src files

Catalog(s): env

### title : Title

Catalog(s): basic

### update\_date :

Catalog(s): basic

version : Maximal version number

Catalog(s): basic\_auto

basic	<code>abstract</code> $\bullet$ <code>authors</code> $\bullet$ <code>authors_parsed</code> $\bullet$ <code>categories</code> $\bullet$ <code>comments</code> $\bullet$ <code>doi</code> $\bullet$
	id • title • update_date • versions
basic_auto	$\texttt{categories\_parsed} \bullet \texttt{id} \bullet \texttt{maincat} \bullet \texttt{mainmathcat} \bullet \texttt{numauthors} \bullet \texttt{numcats}$
	• nummathcats • version
env	e_figure $\bullet$ e_tikzpicture $\bullet$ e_wrapfigure $\bullet$ id $\bullet$ t_meta $\bullet$ t_outer $\bullet$
	$t_proof \bullet t_theorem \bullet texstructure$
env_auto	e_Figs_pp • id
pdf	id • numpages • pdf_filename • pdftext
src	$f_gif \bullet f_jpg \bullet f_pdf \bullet f_tex \bullet f_total \bullet id \bullet missing_files$
<pre>src_auto</pre>	f_images • id • num_missing_files

Table 1: List of catalogs in the arXiv data. Catalogs with suffix auto are derived from the fundamental catalogs.

Depends on: versions

#### versions :

Catalog(s): basic

 $\mathbf{yymm}~$  : Year and month extracted from arXiv ID

Catalog(s):

Depends on: id