

# R $\chi$ iv-Maker: An Automated Template Engine for Streamlined Scientific Publications

Bruno M. Saraiva<sup>1</sup>✉, Guillaume Jaquemet<sup>2,3,4</sup>✉, and Ricardo Henriques<sup>1,5</sup>✉

<sup>1</sup>Instituto de Tecnologia Química e Biológica António Xavier, Universidade Nova de Lisboa, Oeiras, Portugal

<sup>2</sup>Faculty of Science and Engineering, Cell Biology, Åbo Akademi University, Turku, Finland

<sup>3</sup>InFLAMES Research Flagship Center, University of Turku, Turku, Finland

<sup>4</sup>Turku Bioscience Centre, University of Turku and Åbo Akademi University, Turku, Finland

<sup>5</sup>UCL Laboratory for Molecular Cell Biology, University College London, London, United Kingdom

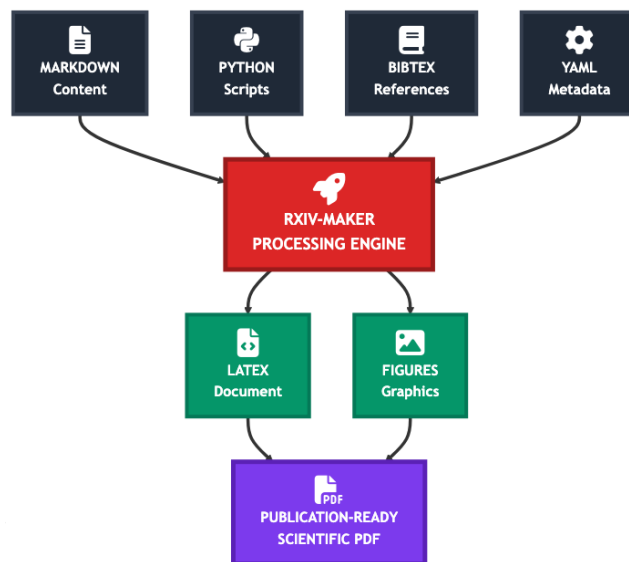
Modern scientific publishing has moved towards rapid dissemination through preprint servers, putting greater demands on researchers for preparing and quality-checking manuscripts. We introduce RXiv-Maker, a comprehensive system native to Github that simplifies scientific writing through markdown-based authoring with automated LaTeX conversion. It's specifically designed to help produce preprints for curation in arXiv, bioRxiv, and medRxiv. The system lets researchers write in the familiar and lightweight markdown syntax while generating publication-quality documents automatically. RXiv-Maker offers flexible compilation strategies, including cloud-based GitHub Actions, interactive Google Colab notebooks, and reproducible local builds via Docker containerisation, ensuring consistent environments and eliminating dependency conflicts. The framework inherently supports reproducible research by enabling programmatic figure generation using Python libraries and script-based diagramming with Mermaid.js. This self-documenting article, created entirely within the framework, shows how this markdown-centric workflow transforms scientific communication into an efficient, collaborative, and transparent process, empowering researchers to focus on content while upholding rigorous standards of quality and reproducibility.

article template | scientific publishing | preprints

Correspondence: (B. M. Saraiva) [b.saraiva@itqb.unl.pt](mailto:b.saraiva@itqb.unl.pt); (G. Jaquemet) [guillaume.jacquemet@abo.fi](mailto:guillaume.jacquemet@abo.fi); (R. Henriques) [ricardo.henriques@itqb.unl.pt](mailto:ricardo.henriques@itqb.unl.pt)

## Main

Today's scientific landscape is marked by the swift sharing of research findings, a trend largely driven by the exponential growth of preprint servers like arXiv, bioRxiv, and medRxiv (1–3). This shift, while accelerating scientific discovery and enhancing research transparency (4, 5), puts a significant new burden on researchers to produce polished manuscripts frequently and efficiently. Traditional tools and workflows for scientific writing, often reliant on proprietary word-processing software, are poorly suited to this new reality. They pose major challenges for version control, collaborative authoring, and ensuring the computational reproducibility of the final document (6, 7). To tackle these systemic issues, we've developed RXiv-Maker, a Github-native framework designed to streamline the entire scientific writing and publication pipeline. The system is built on the principle of using markdown, a simple and intuitive plain-text formatting syntax, as the primary authoring language. This approach fundamentally separates the scientific content from its



**Fig. 1. The RXiv-Maker Diagram.** The system integrates Markdown content, YAML metadata, Python scripts, and bibliography files through a processing engine. This engine leverages Docker, GitHub Actions, and LaTeX to produce a publication-ready scientific article, demonstrating a fully automated and reproducible pipeline.

final presentation, allowing researchers to focus on the substance of their work. The plain-text nature of markdown files makes them ideally suited for management with version control systems like Git. This integration provides an unparalleled level of transparency and traceability for the evolution of a manuscript. Every change, every contribution, and every revision can be precisely tracked, attributed, and, if necessary, reverted. Furthermore, it resolves the convoluted and error-prone process of merging changes from multiple collaborators, a common bottleneck in conventional workflows that often leads to duplicated effort and loss of information. At the heart of the RXiv-Maker philosophy is the pursuit of genuine reproducibility, a fundamental aspect of scientific integrity that extends beyond the experimental data to encompass the entire publication process (8–10). Our framework embodies this principle by enabling the programmatic generation of figures and tables, following best practices for computational research (11). Rather than manually inserting static image files, which obscures the connection between the data and its visualisation, our system promotes the use of scripting languages like Python, with its powerful data visualisation libraries such as Matplotlib (12) and Seaborn (13). Figures are generated directly from the source data and anal-

ysis scripts during the manuscript compilation process. This creates an unbreakable, auditable chain from raw data to the final figure. If the underlying data is updated or the analysis is refined, all affected figures are automatically regenerated, ensuring complete consistency and eliminating the possibility of outdated visuals persisting in the manuscript. This dynamic approach transforms figures from mere illustrations into reproducible scientific artefacts. The system also integrates Mermaid.js (14), a tool for generating complex diagrams and flowcharts from a simple, text-based syntax. This is particularly valuable for creating clear, version-controlled illustrations of experimental workflows, conceptual models, and algorithms, which are essential for communicating complex ideas in many scientific disciplines. RXiv-Maker, therefore, offers a holistic solution that treats the manuscript not as a static document but as the executable output of the research itself.

RXiv-Maker addresses these requirements through a markdown-centric authoring system that automatically translates familiar markdown syntax into professional LaTeX documents. Built upon the established HenriquesLab bioRxiv template (15), the system extends capabilities through automated processing pipelines, integrated figure generation, and flexible deployment strategies. The architecture, detailed in Fig. 1 and comprehensively illustrated in Fig. 2, provides automated figure generation for statistical visualisation, integrated Mermaid diagram creation, and robust build automation through containerised environments. The technical details of the figure generation system are described in [Supp. Note 3](#).

Using the RXiv-Maker framework results in a highly efficient and robust workflow for producing professional-quality scientific papers. The system's primary output is a fully typeset PDF document, as seen in the article you're currently reading, which was generated entirely using this process. The markdown source files are automatically converted into a structured LaTeX document, then compiled to produce a PDF with a clean, academic layout, proper pagination, and high-resolution figures. Bibliographic management is handled seamlessly through integration with a standard BibTeX file. The system automatically processes this file to generate correctly formatted in-text citations and a comprehensive bibliography section according to a specified citation style. This automation eliminates the tedious and error-prone task of manually formatting references.

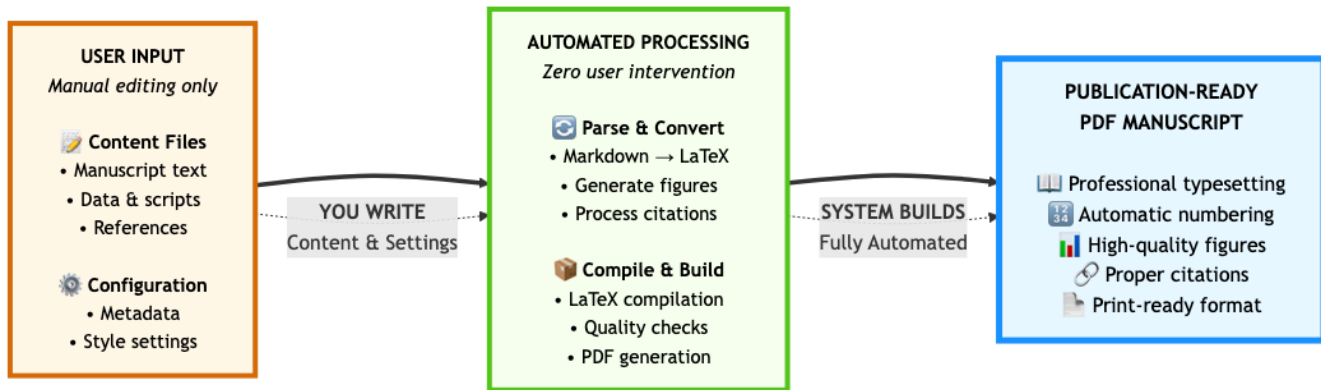
Its utility is further highlighted by the system's flexible deployment options, which cater to a broad range of user preferences and technical environments. We have successfully validated three distinct compilation pathways. Firstly, the cloud-based GitHub Actions workflow offers a fully automated, hands-off experience. Every time code is pushed to the repository, the action is triggered, building a Docker container with all necessary dependencies, compiling the manuscript, and releasing the resulting PDF as a downloadable artefact. This continuous integration pipeline ensures a current, correctly compiled version of the manuscript is always available, serving as a top-notch quality control mechanism for

collaborative projects. Secondly, for users who prefer an interactive, web-based environment, the system can be deployed in a Google Colab notebook. This method removes all local software requirements, allowing users to compile the manuscript simply by executing cells within the notebook, making the framework exceptionally accessible to those with limited command-line experience.

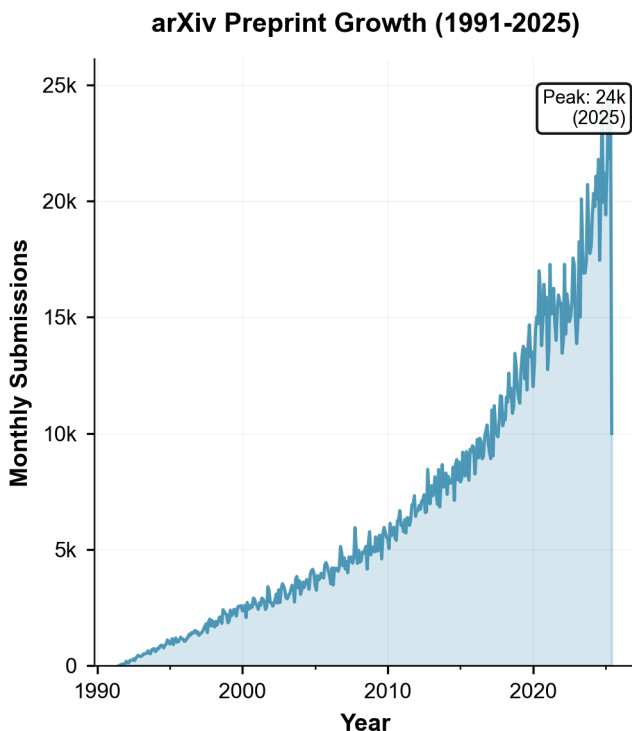
Thirdly, for local development, the Docker-based approach provides a perfectly reproducible compilation environment on any host machine. By running a simple command, a user can start the build process within a self-contained Docker container that encapsulates the exact versions of LaTeX and all other required system libraries. This completely eliminates the 'works on my machine' problem, guaranteeing that the output is identical byte-for-byte regardless of the user's local operating system or software configuration (16). The integration of programmatic figure generation was also validated, supporting interactive computational environments like Jupyter notebooks (17). Python scripts placed within the designated directory were automatically executed during compilation. These scripts loaded data, performed analyses, and generated visualisations, which were then saved as image files and seamlessly included in the final PDF. Similarly, Mermaid.js diagrams embedded within the markdown source were correctly rendered into SVG images and incorporated into the document. This programmatic integration demonstrates a closed loop of reproducibility, where the final manuscript serves as a verifiable and self-contained record of the research findings and their presentation.

Our RXiv-Maker system marks a major step forward in how scientific manuscripts are prepared. By using plain-text markdown and robust open-source tools like Docker, we've created a workflow that boosts efficiency and promotes best practices for reproducibility and collaboration. The main benefit of our framework is that it removes technical complexity from the author's hands. Scientists can focus on the research content and narrative, using simple and widely understood markdown syntax, while the system handles the complex and often frustrating aspects of typesetting, reference management, and dependency control. This approach embraces literate programming principles (18), creating documents that seamlessly blend narrative text with executable code. This is a significant departure from traditional workflows, where researchers often have to act as amateur typesetters, spending hours wrestling with the formatting intricacies of word processors or raw Latex. The integration with Git provides a robust platform for collaborative writing (19, 20), far superior to the chaotic exchange of files via email. It enables transparent attribution, conflict-free merging of contributions, and a complete, auditable history of the manuscript's development.

Within the larger context of the open science movement, RXiv-Maker acts as a practical tool for turning principles into reality. By focusing on automated figure generation and a fully containerised, reproducible build process, the path from raw data to final publication becomes transparent and verifiable. This directly tackles the 'reproducibility crisis' by



**Fig. 2. RXiv-Maker Workflow: User Input vs. Automated Processing.** The framework clearly separates user responsibilities (content creation and configuration) from automated processes (parsing, conversion, compilation, and output generation). Users only need to write content and set preferences, while the system handles all technical aspects of manuscript preparation automatically, ensuring a streamlined workflow from markdown input to publication-ready PDF output.



**Fig. 3. The growth of preprint submissions on the arXiv server from 1991 to 2025.** The data, sourced from arXiv's public statistics, is plotted using a Python script integrated into our RXiv-Maker pipeline. This demonstrates the system's capacity for reproducible, data-driven figure generation directly within the publication workflow.

making it easy to create publications that are not just reports of research, but are themselves computationally reproducible artefacts. This aligns with a growing consensus that the publication itself should be accompanied by the code and data needed to reproduce its findings (8). Our framework makes this possible by design, treating the manuscript and the code to generate it as two sides of the same coin. While other platforms and tools for scientific writing exist, including sophisticated environments like DL4MicEverywhere (21), RXiv-Maker stands out through its simplicity, flexibility, and tight integration with the GitHub ecosystem, which is already the de facto standard for collaborative software development and

is increasingly used for scientific projects.

Although RXiv-Maker has its advantages, we acknowledge certain limitations and areas for future improvement. The current system is mainly designed to produce PDF outputs using LaTeX. While this is the standard for many scientific disciplines, future versions could support other output formats, such as HTML for web-native articles, potentially leveraging universal document converters like Pandoc (22). Additionally, concerns about reporting quality in preprints (23) suggest opportunities for integrating automated quality checks. Additionally, although the system is designed to be accessible, researchers new to Git and markdown may encounter an initial learning curve. To address this, we plan to develop more comprehensive documentation and tutorial materials. Future work will also focus on deeper integration with data analysis environments like Jupyter notebooks, allowing for a more seamless transition from exploratory analysis to manuscript-ready figures. We could also explore integrating automated tools for checking style, grammar, and scientific rigour, further enhancing the system's role as a comprehensive quality control platform. Ultimately, RXiv-Maker is a contribution to a more open, efficient, and reproducible future for scientific communication, providing a powerful and accessible tool for modern researchers.

#### MANUSCRIPT PREPARATION

This manuscript was prepared using RXiv-Maker version 1.11.0.

#### DATA AVAILABILITY

Arxiv monthly submission data used in this article is available at [https://arxiv.org/stats/monthly\\_submissions](https://arxiv.org/stats/monthly_submissions). The source code and data for the figures in this article are available at <https://github.com/henriques/rxiv-maker>.

#### CODE AVAILABILITY

The RXiv-Maker computational framework is available at <https://github.com/henriques/rxiv-maker>. All source code is under an MIT License.

#### AUTHOR CONTRIBUTIONS











Both Bruno M. Saraiva, Guillaume Jacquemet and Ricardo Henriques conceived the project and designed the framework. All authors contributed to writing and reviewing the manuscript.

#### ACKNOWLEDGEMENTS

B.S. and R.H. acknowledge support from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 101001332) (to R.H.) and funding from the European Union through the Horizon Europe program (AI4LIFE project with grant agreement 101057970-AI4LIFE and RT-SuperES project with grant agreement 101099654-RTSuperES to R.H.). Funded by the European Union. However, the views and

opinions expressed are those of the authors only and do not necessarily reflect those of the European Union. Neither the European Union nor the granting authority can be held responsible for them. This work was also supported by a European Molecular Biology Organization (EMBO) installation grant (EMBO-2020-IG-4734 to R.H.), a Chan Zuckerberg Initiative Visual Proteomics Grant (vpi-0000000044 with <https://doi.org/10.37921/743590vtudfp> to R.H.) and a Chan Zuckerberg Initiative Essential Open Source Software for Science (EOSS6-0000000260). This study was supported by the Academy of Finland (no. 338537 to G.J.), the Sigrid Juselius Foundation (to G.J.), the Cancer Society of Finland (Syöpäjärjestöt, to G.J.) and the Solutions for Health strategic funding to Åbo Akademi University (to G.J.). This research was supported by InFLAMES Flagship Program of the Academy of Finland (decision no. 337531).

## EXTENDED AUTHOR INFORMATION

- **Bruno M. Saraiva:**  
 0000-0002-9151-5477;  Bruno\_MSaraiva;  bruno-saraiva
- **Guillaume Jaquemets:**  
 0000-0002-9286-920X;  guijacquemet;  guijacquemet.bsky.social
- **Ricardo Henriques:**  
 0000-0002-1234-5678;  HenriquesLab;  henriqueslab.bsky.social;  
 ricardo-henriques

## Bibliography

- Jeremy M Berg, Narinder Bhalla, Philip E Bourne, M Cacchione, D Egnor, J Flombaum, E Ghanem, M Ginsburg, B Goldstein, L Gorodetski, et al. Preprint servers. *Science*, 352 (6288):899, 2016. doi: 10.1126/science.aaf9133.
- Richard J Abdill and Ran Blekhan. The growth of biorxiv preprints and the implications for preprint discovery. *PLoS Biology*, 17(4):e3000269, 2019. doi: 10.1371/journal.pbio.3000269.
- Nicholas Fraser, Fakhri Momeni, Philipp Mayr, and Isabella Peters. The relationship between biorxiv preprints, citations and altmetrics. *Quantitative Science Studies*, 2(2):618–638, 2021. doi: 10.1162/qss\_a\_00043.
- Ronald D Vale. Accelerating scientific publication in biology. *Proceedings of the National Academy of Sciences*, 116(52):26222–26229, 2019. doi: 10.1073/pnas.1915668116.
- Jonathan P Tennant, François Waldner, Damien C Jacques, Paola Masuzzo, Lauren B Collister, and Chris HJ Hartgerink. The academic, economic and societal impacts of open access: an evidence-based review. *F1000Research*, 5:632, 2016. doi: 10.12688/f1000research.8460.3.
- Vincent Larivière, Cassidy R Sugimoto, Benoît Macaluso, Staša Milojević, and Blaise Cronin. arXiv.org and the dissemination of scholarly information: The case of computer science. *Journal of the Association for Information Science and Technology*, 65(4):844–848, 2014. doi: 10.1002/asi.23081.
- Jerson L da Silva. The challenge of being a scientist in the 21st century. *Anais da Academia Brasileira de Ciências*, 94, 2022. doi: 10.1590/0001-376520220211396.
- David L Donoho. An invitation to reproducible computational research. *The American Statistician*, 64(1):1–2, 2010. doi: 10.1198/tast.2010.09132.
- Geir Kjetil Sandve, Anton Nekrutenko, James Taylor, and Eivind Hovig. Ten simple rules for reproducible computational research. *PLoS Computational Biology*, 9(10):e1003285, 2013. doi: 10.1371/journal.pcbi.1003285.
- Greg Wilson, Dhavide A Aruliah, C Titus Brown, Neil P Chue Hong, Matt Davis, Richard T Guy, Steven HD Haddock, Kathryn D Huff, Ian M Mitchell, Mark D Plumbley, et al. Best practices for scientific computing. *PLoS Biology*, 12(1):e1001745, 2014. doi: 10.1371/journal.pbio.1001745.
- Dominik Scherer, Lars Bendix, Stephan Greiner, and Julian Lück. A survey on the state of research data reproducibility in computer science. In *2020 IEEE/ACM 42nd International Conference on Software Engineering Workshops (ICSEW)*, pages 310–317. IEEE, 2020. doi: 10.1109/ICSEW50701.2020.00052.
- John D Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007. doi: 10.1109/MCSE.2007.55.
- Michael L Waskom. seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60):3021, 2021. doi: 10.21105/joss.03021.
- Mermaid Team. Mermaid: Generation of diagrams and flowcharts from text in a similar manner as markdown, 2023. Accessed: 2024-12-01.
- Ricardo Henriques. Henriques biorxiv template, 2015. Overleaf LaTeX template. Accessed: 2025-06-16.
- Carl Boettiger. An introduction to docker for reproducible research. *ACM SIGOPS Operating Systems Review*, 49(1):71–79, 2015. doi: 10.1145/2723872.2723882.
- Thomas Kluyver, Benjamin Ragan-Kelley, Fernando Pérez, Brian E Granger, Matthias Bussonnier, Jonathan Frederic, Kyle Kelley, Jessica B Hamrick, Jason Grout, Sylvain Corlay, et al. Jupyter notebooks—a publishing format for reproducible computational workflows. In *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, pages 87–90. IOS Press, 2016. doi: 10.3233/978-1-61499-649-1-87.
- Donald E Knuth. Literate programming. *The Computer Journal*, 27(2):97–111, 1984. doi: 10.1093/comjnl/27.2.97.
- Karthik Ram. Git can facilitate greater reproducibility and increased transparency in science. *Source Code for Biology and Medicine*, 8(1):7, 2013. doi: 10.1186/1751-0473-8-7.
- Yasset Perez-Riverol, Laurent Gatto, Rui Wang, Timo Sachsenberg, Julian Uszkoreit, Felipe da Veiga Leprevost, Christian Fufezan, Tobias Tertent, Stephen J Eglen, Daniel S Katz, et al. Ten simple rules for taking advantage of git and github. *PLoS Computational Biology*, 12(7):e1004947, 2016. doi: 10.1371/journal.pcbi.1004947.
- Ivan Hidalgo-Cenalmor, Joanna W Pylvänäinen, Mariana G Ferreira, Craig T Russell, Alon Saguy, Ignacio Arganda-Carreras, Guillaume Jacquemet, and Ricardo Henriques.

- DI4miceverywhere: deep learning for microscopy made flexible, shareable and reproducible. *Nature Methods*, 21(5):804–810, 2024. doi: 10.1038/s41592-024-02295-6.
- John MacFarlane. Pandoc: A universal document converter. *Journal of Open Source Software*, 7(76):4642, 2022. doi: 10.21105/joss.04642.
- Jamie J Kirkham, Christopher M Penfold, Fiona Murphy, and John PA Ioannidis. The quality of reporting in preprints. *bioRxiv*, page 451822, 2018. doi: 10.1101/451822.
- Overleaf. Overleaf: Real-time collaborative writing and publishing tools with integrated pdf preview, 2024. Cloud-based LaTeX editor.
- Posit PBC. Quarto: An open-source scientific and technical publishing system, 2024. Multi-language scientific publishing system.
- Christopher Rackauckas, Rowan Cockett Laursen, Lindsey Heagy Thirumalai, et al. Myst-nb: Authoring and publishing computational documents with myst markdown. *Journal of Open Source Software*, 6(61):3053, 2021. doi: 10.21105/joss.03053.
- Rowan Cockett, Lindsey Heagy Thirumalai, Seogi Kang, Sia Hejazi, Bane Consummes, Bas Couwenberg, Andreas Colliander, Chris Phang, Vladimir Keleshev, Remi Rampin, et al. Executable books: Building publication-quality documents from computational material. In *Proceedings of the 19th Python in Science Conference*, pages 79–88, 2020. doi: 10.25080/Majors-342d178e-00b.
- Typst GmbH. Typst: A new markup-based typesetting system, 2024. Modern typesetting system designed for scientific documents.
- Yihui Xie. bookdown: Authoring books and technical documents with r markdown. 2016.

## Methods

**.1. The RXiv-Maker Processing Pipeline: A Deterministic Automated Workflow.** The RXiv-Maker framework implements a deterministic, multi-stage processing pipeline to convert manuscript source files into a compiled Portable Document Format (PDF) document. The entire workflow is orchestrated by a central Makefile that defines a directed acyclic graph (DAG) of dependencies, ensuring a consistent and reproducible order of operations for every build. This Makefile serves as the high-level algorithmic specification for the tool, defining distinct pathways for local and containerized compilation to enhance both accessibility for general users and control for developers.

The primary build process, typically invoked via the `make pdf` target for local compilation, proceeds through five distinct, logically ordered stages:

**Environment Setup (setup):** The initial stage prepares the build environment by creating the designated output directory (`output/`) and its necessary subdirectories (e.g., `output/Figures/`). This ensures a clean and predictable workspace for all subsequent build artifacts.

**Programmatic Content Generation (figures-conditional):** This stage manages the creation of dynamic content, primarily figures. A key methodological feature is its conditional execution logic. The system inspects the `FIGURES/` directory for source files (e.g., `.py` for Python scripts, `.mmd` for Mermaid diagrams) and checks for the existence of corresponding output files (e.g., `.pdf`, `.png`). The figure generation script (`src/py/commands/generate_figures.py`) is executed only if these output files are absent or if the user explicitly forces regeneration by setting the `FORCE_FIGURES=true` flag. This caching-like behavior represents a critical optimization, balancing the need for reproducibility with computational efficiency by avoiding redundant processing of unchanged assets.

**Core Content Conversion (generate):** This is the central conversion step, where the main Python script (`src/py/commands/generate_preprint.py`) is executed. This script parses the primary manuscript markdown file (`MANUSCRIPT/01_MAIN.md`), extracts meta-



data from the configuration file (MANUSCRIPT/00\_CONFIG.yml), and synthesizes these inputs into a master LaTeX file (MANUSCRIPT.tex) within the output directory.

**Asset Aggregation** (copy-files): Following the generation of the primary LaTeX file, this stage gathers all necessary dependencies for compilation into the output/ directory. This includes LaTeX style files (.cls, .bst, .sty) from src/tex/style/, the project's bibliography file (.bib), and all generated figures from the FIGURES/ directory. This aggregation creates a self-contained environment for the final typesetting stage.

**Final Typesetting** (pdf): The final stage executes the LaTeX compilation sequence within the output/ directory. A standard, robust sequence of pdflatex→bibtex→pdflatex→pdflatex is used. This multi-pass process ensures that all cross-references (for figures, tables, and equations) and bibliographic citations are correctly resolved, resulting in a polished, publication-ready PDF document. A timeout wrapper is applied to each command to prevent build processes from hanging indefinitely.

For users without a local LaTeX installation, the framework provides an equivalent containerized pathway via the make easy-build target, which is an alias for docker-build. This target leverages a pre-configured Docker environment to execute the same fundamental pipeline, guaranteeing bit-for-bit reproducibility of the final output across any host system. This dual-pathway design is a core methodological choice that significantly enhances the tool's accessibility and scientific utility.

The Makefile defines several key automation targets that orchestrate the build process: setup creates the output directory structure, figures-conditional manages programmatic figure generation, generate performs the core markdown-to-LaTeX conversion, copy-files aggregates all necessary assets, build ensures readiness for typesetting, pdf executes the LaTeX compilation sequence, and easy-build provides a containerized build pathway for users without local LaTeX installations.

**.2. The Markdown-to-LaTeX Conversion Engine.** The transformation of the user's markdown manuscript into a structured LaTeX document is not performed by a generic, off-the-shelf converter. Instead, RXiv-Maker employs a custom, multi-pass conversion engine written in Python (src/py/converters/md2tex.py). This engine is designed to recognize and process a specific "extended academic Markdown" syntax, which adds essential scientific publishing features to the standard markdown specification. The precise behavior of this engine is formally documented and validated by a comprehensive suite of unit tests, particularly those found in tests/unit/test\_md2tex.py.

The conversion process can be understood as a pipeline of specialized processing functions, each responsible for a specific syntactic element. This modular architecture allows for robust and maintainable code, where complex transformations are broken down into discrete, testable steps. The key stages of this conversion pipeline include:

**Code Block Protection:** The first pass identifies all fenced

code blocks (e.g., python... ) and replaces them with unique, protected placeholders. This crucial step ensures that the content within these blocks is treated as literal text and is not subjected to any subsequent markdown-to-LaTeX conversion. This preserves code syntax, YAML examples, and other verbatim content without corruption.

**Specialized Element Conversion:** The engine then applies a series of dedicated converters to the remaining text:

- *Citation Processing* (convert\_citations\_to\_latex): This function identifies two citation patterns. Single, in-text citations like @smith2023 are converted to \citesmith2023. Bracketed, multi-citation groups like [@smith2023;@jones2022] are converted to the consolidated LaTeX command \citesmith2023,jones2022.
- *Figure Processing* (convert\_figures\_to\_latex): This processor handles the complex syntax for embedding figures. It recognizes attributes for labels (\#fig:id), width (width="0.8"), and LaTeX placement (tex\_position="!ht"), translating them into the appropriate \beginfigure, \includegraphics, and \label commands.
- *Table Processing* (convert\_tables\_to\_latex): This handles GitHub-flavored markdown tables, converting them into LaTeX tabular environments. It also supports extended attributes for rotation (rotate=90) and column width, wrapping the table in \rotatebox or tabularx environments as needed.
- *Supplementary Note Processing* (process\_supplementary\_notes): A custom syntax, #snote:id Title., is used for creating structured supplementary notes. This processor converts these blocks into formatted LaTeX subsections with automatic numbering and labeling, a feature essential for organizing supplementary information.

**Standard Markdown Conversion:** After the specialized elements are handled, the engine applies standard markdown conversion rules for basic formatting, such as transforming **bold** to \textbf, *italic* to \textit, and headers (\#, \#\#) to LaTeX sectioning commands (\section, \subsection).

**Placeholder Restoration:** Finally, the protected placeholders for code blocks are restored, inserting the original verbatim content into the appropriate LaTeX environment for syntax highlighting.

This multi-pass, protection-first approach ensures that the extended academic syntax is processed correctly while preventing the accidental conversion of literal content within code blocks. The conversion engine supports several custom syntax extensions including bracketed citations, figure placement attributes, cross-references, rotated tables, language-specific code blocks, and structured supplementary notes.

**.3. Programmatic Content Generation for Reproducible Science.** A cornerstone of the RXiv-Maker methodology is its direct support for programmatic content generation, which transforms figures and diagrams from static, manually-created assets into dynamic, reproducible outputs of the scientific analysis itself. This capability is not an add-on but a core, algorithmically-defined part of the build pipeline, ensuring that the final publication is a verifiable and auditable record of the research process.

The mechanism is orchestrated by the `figures` and `figures-conditional` targets in the Makefile. The process follows a clear algorithm:

**Source Identification:** During a build, the system scans the `MANUSCRIPT/FIGURES/` directory for files with recognized executable extensions, primarily `.py` for Python scripts and `.mmd` for Mermaid diagram definitions.

**Conditional Execution:** As described previously, the `figures-conditional` logic determines whether a script needs to be executed. This check for pre-existing output files serves as an effective caching mechanism, significantly speeding up subsequent builds where figure-generating code or data has not changed.

**Interpreter Dispatch:** If execution is required, the system dispatches the source file to the appropriate interpreter:

- *Python Scripts (.py):* Files ending in `.py` are executed using the project's configured Python interpreter (`$(PYTHON_CMD)`). These scripts are expected to perform data loading, analysis, and plotting (e.g., using libraries like Matplotlib and Seaborn, as specified in `pyproject.toml`), and save the resulting visualizations as image files (e.g., `.pdf`, `.png`) back into the `FIGURES/` directory.
- *Mermaid Diagrams (.mmd):* Files ending in `.mmd` are processed using the Mermaid CLI. This converts the declarative, text-based diagram syntax into vector (`.svg`, `.pdf`) and raster (`.png`) graphics. This allows complex flowcharts and architectural diagrams to be version-controlled and automatically rendered as part of the manuscript.

**Integration:** Once the output files are generated, the markdown-to-LaTeX conversion engine recognizes the corresponding `!\href...Caption` tags in the manuscript's markdown source and embeds the newly created figures into the final document.

This entire process ensures a tight coupling between the data, the code that analyzes it, and the final visual representation in the publication. It elevates the manuscript from a static report to a dynamic, executable artifact, which is a significant step towards achieving the ideals of fully reproducible computational science.

**.4. Environment Encapsulation for Guaranteed Consistency.** Reproducibility in computational science is critically dependent on the stability of the software environment. To address this, RXiv-Maker provides a hierarchical strategy for environment management, ensuring that the tool behaves iden-

tically regardless of the host system. This is a deliberate methodological choice to eliminate the common "works on my machine" problem.

**Dependency Pinning:** At the most fundamental level, the `pyproject.toml` file explicitly lists all Python dependencies with version specifiers (e.g., `matplotlib>=3.7.0`, `ruff>=0.8.0`). This provides a baseline level of reproducibility for any installation, as it instructs the package manager to use versions known to be compatible with the tool.

**Virtual Environments:** For local development, the Makefile provides a `venv` target to create an isolated Python virtual environment. This prevents conflicts with system-wide packages and ensures that the project's dependencies are self-contained.

**Containerization with Docker:** The highest level of reproducibility is achieved through Docker. The project includes a multi-stage Dockerfile (`src/docker/Dockerfile`) that encapsulates the entire software stack—the specific operating system, LaTeX distribution, Python version, and all system and Python libraries—into a single, portable container image. The `make easy-build` command executes the entire build pipeline within this container, guaranteeing that the process is identical for every user on any platform that supports Docker. This approach provides bit-for-bit reproducibility of the final PDF output.

The sophistication of this approach is further demonstrated by the `.dockerignore` file, which is meticulously crafted to optimize the Docker build process. By excluding build artifacts (`output/`, `build/`), local caches (`__pycache__/`, `.pytest_cache/`), and version control directories (`.git/`), it minimizes the size of the build context sent to the Docker daemon, enhancing the efficiency and speed of the reproducible build process.

**.5. The Quality Assurance and Validation Framework.** In computational science, the software testing framework serves the same role as experimental controls and validation assays in empirical science. It provides the evidence that the computational instrument—the software—is behaving as specified. RXiv-Maker employs a formal, multi-level testing strategy to validate its functionality, which should be described as the validation protocol for the method.

The testing framework, managed via Makefile targets and configured in `pyproject.toml`, is structured into three distinct layers:

**Unit Tests (`tests/unit/`):** This suite forms the foundation of the validation strategy. Each test validates a single, discrete component or conversion rule in isolation. These tests function as precise positive and negative controls. For example, `test_convert_bold_text` in `test_md2tex.py` validates that the `**text**` to `\textbf{text}` conversion works correctly (a positive control). Conversely, tests for markdown inside backticks validate that this same conversion does not occur inside a code block, confirming the specificity of the rule (a negative control). This suite provides granular, verifiable evidence that the core algorithms of the conversion engine are implemented correctly.

**Integration Tests** (`tests/integration/`): This layer validates the successful orchestration of multiple components. Tests like `test_end_to_end_with_citations` verify that the entire pipeline—from parsing a markdown file containing citations and figure references, to generating a LaTeX file with the correct `\cite` and `\ref` commands—functions as a cohesive whole. These tests are analogous to validating a complete experimental protocol from sample preparation to final measurement, ensuring that the interfaces between different modules are working correctly.

**Platform Tests** (`tests/integration/test_platform_integration.py`): This layer validates the tool's functionality in its target deployment environments. These tests verify the behavior of the Docker-based workflows, ensuring that containerized builds produce the expected results across different platforms and architectures.

In addition to this validation suite, the project enforces a proactive quality control pipeline using pre-commit hooks, configured in `.pre-commit-config.yaml`. Before any code change can be committed to the version control system, a series of automated checks are executed. These include code formatting (`ruff-format`), linting (`ruff`), type checking (`mypy`), and even spell-checking of configuration files (`typos`). This automated pipeline serves as a continuous, preventative quality assurance method, ensuring a high standard of code quality, consistency, and correctness throughout the development process.

The project's tooling choices also reflect a deliberate methodology aimed at quality and performance. The configuration explicitly adopts modern, high-performance tools. For instance, the traditional Python-based tools `black`, `isort`, and `flake8` have been replaced by the single, Rust-based tool `Ruff`, reflecting a move towards faster alternatives while maintaining established tools like `mypy` for type checking. This hybrid approach demonstrates a sophisticated, evidence-based methodology for toolchain construction, prioritizing performance and developer efficiency while maintaining best-in-class functionality.

## Supplementary Information

# **R** $\chi$ iv-Maker: An Automated Template Engine for Streamlined Scientific Publications



Markdown Element	LaTeX Equivalent	Description
<b>**bold text**</b>	<code>\textbf{bold text}</code>	Bold formatting for emphasis
<b>*italic text*</b>	<code>\textit{italic text}</code>	Italic formatting for emphasis
<b># Header 1</b>	<code>\section{Header 1}</code>	Top-level section heading
<b>## Header 2</b>	<code>\subsection{Header 2}</code>	Second-level section heading
<b>### Header 3</b>	<code>\subsubsection{Header 3}</code>	Third-level section heading
<code>citation</code>	<code>\cite{citation}</code>	Single citation reference
<code>[cite1;cite2]</code>	<code>\cite{cite1,cite2}</code>	Multiple citation references
<code>fig:label</code>	<code>\ref{fig:label}</code>	Figure cross-reference
Image with attributes	<code>\begin{figure}...\end{figure}</code>	Figure with attributes (old format)
Image with caption	<code>\begin{figure}...\end{figure}</code>	Figure with separate caption (new format)
<code>- list item</code>	<code>\begin{itemize}\item...\end{itemize}</code>	Unordered list
<code>1. list item</code>	<code>\begin{enumerate}\item...\end{enumerate}</code>	Ordered list
<code>[link text](url)</code>	<code>\href{url}{link text}</code>	Hyperlink with custom text
<code>https://example.com</code>	<code>\url{https://example.com}</code>	Bare URL
<code>&lt;!-- comment --&gt;</code>	<code>% comment</code>	Comments (converted to LaTeX style)
Markdown table	<code>\begin{table}...\end{table}</code>	Table with automatic formatting
<code>&lt;newpage&gt;</code>	<code>\newpage</code>	Manual page break control
<code>&lt;clearpage&gt;</code>	<code>\clearpage</code>	Page break with float clearing

**Sup. Table S1. RXiv-Maker Markdown Syntax Overview.** Comprehensive mapping of markdown elements to their LaTeX equivalents, demonstrating the automated translation system that enables researchers to write in familiar markdown syntax whilst producing professional LaTeX output.

Deployment Method	Environment	Dependencies	Collaboration	Ease of Use	Reproducibility
<b>Docker Local</b>	Local machine	Docker only	Git-based	High	Perfect
<b>GitHub Actions</b>	Cloud CI/CD	None (cloud)	Automatic	Very High	Perfect
<b>Google Colab</b>	Web browser	None (cloud)	Shared notebooks	Very High	High
<b>Local Python</b>	Local machine	Python + LaTeX	Git-based	Medium	Good
<b>Manual LaTeX</b>	Local machine	Full LaTeX suite	Git-based	Low	Variable

**Sup. Table S2. RXiv-Maker Deployment Strategies.** Comparison of available compilation methods, highlighting the flexibility of the framework in accommodating different user preferences and technical environments whilst maintaining consistent output quality.

Format	Input Extension	Processing Method	Output Formats	Quality	Use Case
<b>Mermaid Diagrams</b>	<code>.mmd</code>	Mermaid CLI	SVG, PNG, PDF	Vector/Raster	Flowcharts, architectures
<b>Python Figures</b>	<code>.py</code>	Script execution	PNG, PDF, SVG	Publication	Data visualisation
<b>Static Images</b>	<code>.png, .jpg, .svg</code>	Direct inclusion	Same format	Original	Photographs, logos
<b>LaTeX Graphics</b>	<code>.tex, .tikz</code>	LaTeX compilation	PDF	Vector	Mathematical diagrams
<b>Data Files</b>	<code>.csv, .json, .xlsx</code>	Python processing	Via scripts	Computed	Raw data integration

**Sup. Table S3. Supported Figure Generation Methods.** Comprehensive overview of the framework's figure processing capabilities, demonstrating support for both static and dynamic content generation with emphasis on reproducible computational graphics.

Tool	Type	Markdown	Primary Use Case	Key Strengths	Open Source
<b>RXiv-Maker</b>	Pipeline	Excellent	Preprint servers	Containerized builds, automated workflows	Yes
<b>Overleaf</b> (24)	Web Editor	Limited	Academic publishing	Real-time collaboration, rich templates	Freemium
<b>Quarto</b> (25)	Publisher	Native	Multi-format publishing	Polyglot support, multiple outputs	Yes
<b>Pandoc</b> (22)	Converter	Excellent	Format conversion	Universal format support, extensible	Yes
<b>MyST-Parser</b> (26)	Extension	Native	Technical documentation	Sphinx ecosystem, rich directives	Yes
<b>Jupyter Book</b> (27)	Publisher	Native	Educational content	Interactive content, executable books	Yes
<b>Typst</b> (28)	Typesetter	Good	Modern typesetting	Fast compilation, modern syntax	Yes
<b>Bookdown</b> (29)	Publisher	R Markdown	Academic books	Cross-references, multiple formats	Yes
<b>Direct LaTeX</b>	Typesetter	Limited	Traditional publishing	Ultimate control, established workflows	Yes

**Sup. Table S4. Comprehensive Comparison of Manuscript Preparation Tools.** This comparison provides an exhaustive overview of available tools for scientific manuscript preparation, positioning each within the broader ecosystem of academic publishing workflows. RXiv-Maker is designed as a specialized solution optimizing for preprint server submissions, complementing rather than replacing established tools like Overleaf for general LaTeX collaboration or Quarto for multi-format publishing. The comparison emphasizes that different tools excel in different contexts: Overleaf dominates collaborative LaTeX editing, Quarto excels at multi-format computational publishing, and RXiv-Maker streamlines the specific workflow of preparing reproducible preprints for arXiv, bioRxiv, and medRxiv submission.

**Supp. Note 1: Architectural Philosophy and Project Organisation..** The RXiv-Maker framework embodies a carefully considered architectural philosophy that prioritises clarity, maintainability, and computational reproducibility through systematic organisation of project components. The system’s file structure reflects established software engineering principles whilst accommodating the specific requirements of scientific manuscript preparation. This organisational schema segregates content, configuration, and computational elements into distinct hierarchical domains, thereby facilitating both human comprehension and automated processing.

The primary manuscript content resides within the MANUSCRIPT directory, which houses the core intellectual contribution in easily accessible formats. This directory contains the YAML configuration file (00\_CONFIG.yml) that centralises all metadata including authorship details, institutional affiliations, and document properties, thereby enabling programmatic manipulation of manuscript attributes without requiring modifications to the narrative content. The numbered markdown files (01\_MAIN.md, 02\_SUPPLEMENTARY\_INFO.md) contain the substantive text, with the numerical prefixing ensuring logical processing order whilst maintaining intuitive organisation for collaborative authoring. The BibTeX references file (03\_REFERENCES.bib) provides standardised bibliographic management, ensuring consistent citation formatting across the entire document. Figure sources and data are organised within dedicated subdirectories (FIGURES/, TABLES/) that maintain clear separation between content types whilst enabling automated discovery during the compilation process.

The src directory encompasses the computational infrastructure that transforms markdown source into publication-ready output. This separation ensures that the technical implementation remains distinct from the scientific content, facilitating maintenance and updates to the processing pipeline without affecting the manuscript itself. The modular structure within src reflects software engineering best practices, with specialised processors for different content types that can be independently developed and tested. The output directory serves as the compilation workspace where intermediate files and final products are generated, preventing contamination of source materials with temporary compilation artefacts whilst providing transparency into the conversion process.

**Supp. Note 2: Comparative Analysis with Alternative Scientific Authoring Platforms..** Within the broader landscape of scientific authoring tools, RXiv-Maker occupies a distinctive position that reflects careful consideration of the trade-offs between functionality and simplicity. Platforms such as Overleaf (24) have revolutionised collaborative LaTeX authoring by providing sophisticated web-based environments with real-time collaboration features, comprehensive template libraries, and integrated compilation services. These systems excel in scenarios requiring complex document structures, advanced typesetting control, and seamless multi-author workflows. The platform’s strength lies in its ability to democratise LaTeX authoring by providing a familiar word-processor-like interface whilst maintaining the typographical excellence of LaTeX output.

Similarly, Quarto (25) represents a powerful framework for scientific and technical publishing that supports multiple programming languages, diverse output formats, and sophisticated computational document features. Its versatility enables researchers to create documents that seamlessly integrate narrative text with executable code, supporting formats ranging from HTML web pages to PDF documents and interactive presentations. Quarto’s strength lies in its comprehensive approach to scientific communication, enabling complex multi-format publishing workflows across various scientific domains.

Pandoc (22), as a universal document converter, provides exceptional flexibility in transforming content between numerous formats. Its strength lies in its ability to serve as a foundation for custom publishing workflows, enabling researchers to develop bespoke solutions for specific requirements. However, this flexibility comes at the cost of increased complexity in configuration and setup.

**Supp. Note 3: Programmatic Figure Generation and Computational Reproducibility..** The technical architecture underlying RXiv-Maker’s figure generation capabilities demonstrates how automated processing pipelines can maintain transparent connections between source data and final visualisations whilst ensuring computational reproducibility. The system supports two primary methodologies for figure creation: Mermaid diagram processing and Python-based data visualisation, each addressing distinct requirements within the scientific publishing workflow.

Mermaid diagram processing leverages the Mermaid CLI to convert text-based diagram specifications into publication-ready graphics. This approach enables version-controlled diagram creation where complex flowcharts, system architectures, and conceptual models can be specified using intuitive syntax and automatically rendered into multiple output formats. The system generates SVG, PNG, and PDF variants to accommodate different compilation requirements whilst maintaining vector quality where appropriate. This automation eliminates the manual effort traditionally required for diagram creation and updates, whilst ensuring that modifications to diagram specifications are immediately reflected in the final document.

Python figure generation represents a more sophisticated approach to computational reproducibility, where analytical scripts are executed during document compilation to generate figures directly from source data. This integration ensures that visualisations remain synchronised with the underlying datasets and analytical methods, eliminating the possibility of outdated or inconsistent graphics persisting in the manuscript. The system executes Python scripts within the compilation environment, automatically detecting generated image files and incorporating them into the document structure. This approach transforms figures from static illustrations into dynamic, reproducible computational artefacts that enhance the scientific rigour of the publication.

**Supp. Note 4: Markdown-to-LaTeX Conversion Architecture and Processing Pipeline..** The markdown-to-LaTeX conversion architecture demonstrates how specialised processors can handle complex document transformations whilst maintaining code modularity and testability. The system employs dedicated processors for figures, tables, citations, and other content types, each implementing specific transformation rules that preserve semantic meaning whilst ensuring typographical excellence. This modular approach enables independent development and testing of conversion components, facilitating maintenance and enhancement of the framework's capabilities.

Figure processing supports multiple syntax variants to accommodate different authoring preferences, including the new format where images are followed by attribute blocks and captions, the attributed format with inline specifications, and simple format for basic inclusions. The core conversion function implements a multi-pass approach that protects literal content during transformation, processes each figure format through dedicated functions, and restores protected content after processing. This sophisticated content protection mechanism ensures that code examples and other literal content are preserved during transformation, proving essential for technical manuscripts.

Table processing handles GitHub Flavored Markdown tables with LaTeX-specific enhancements such as rotation capabilities and sophisticated cross-referencing systems. The conversion system supports both legacy and modern caption formats, enabling authors to specify table properties including width detection for double-column layouts, rotation angles for landscape orientation, and identifier extraction for cross-referencing. The table cell formatting function implements context-aware processing that preserves markdown syntax within examples whilst properly escaping special characters and converting emphasis markers to appropriate LaTeX commands.

Reference processing demonstrates how automated systems can enhance document quality whilst reducing authoring burden. The framework automatically converts markdown-style references into appropriate LaTeX cross-references, ensuring consistent formatting and enabling LaTeX's sophisticated reference management capabilities. This automation extends to bibliographic citations, where the system integrates seamlessly with BibTeX workflows to provide professional citation formatting without requiring authors to master LaTeX citation syntax.

**Supp. Note 5: Reproducibility Features and Version Control Integration..** The RXiv-Maker framework incorporates reproducibility as a fundamental design principle rather than an afterthought, implementing features that ensure complete traceability from source data to final publication. The system's integration with Git version control provides comprehensive tracking of all components necessary for manuscript generation, including content files, configuration parameters, processing scripts, and even the framework code itself. This approach ensures that every aspect of the publication process can be reproduced, verified, and audited.

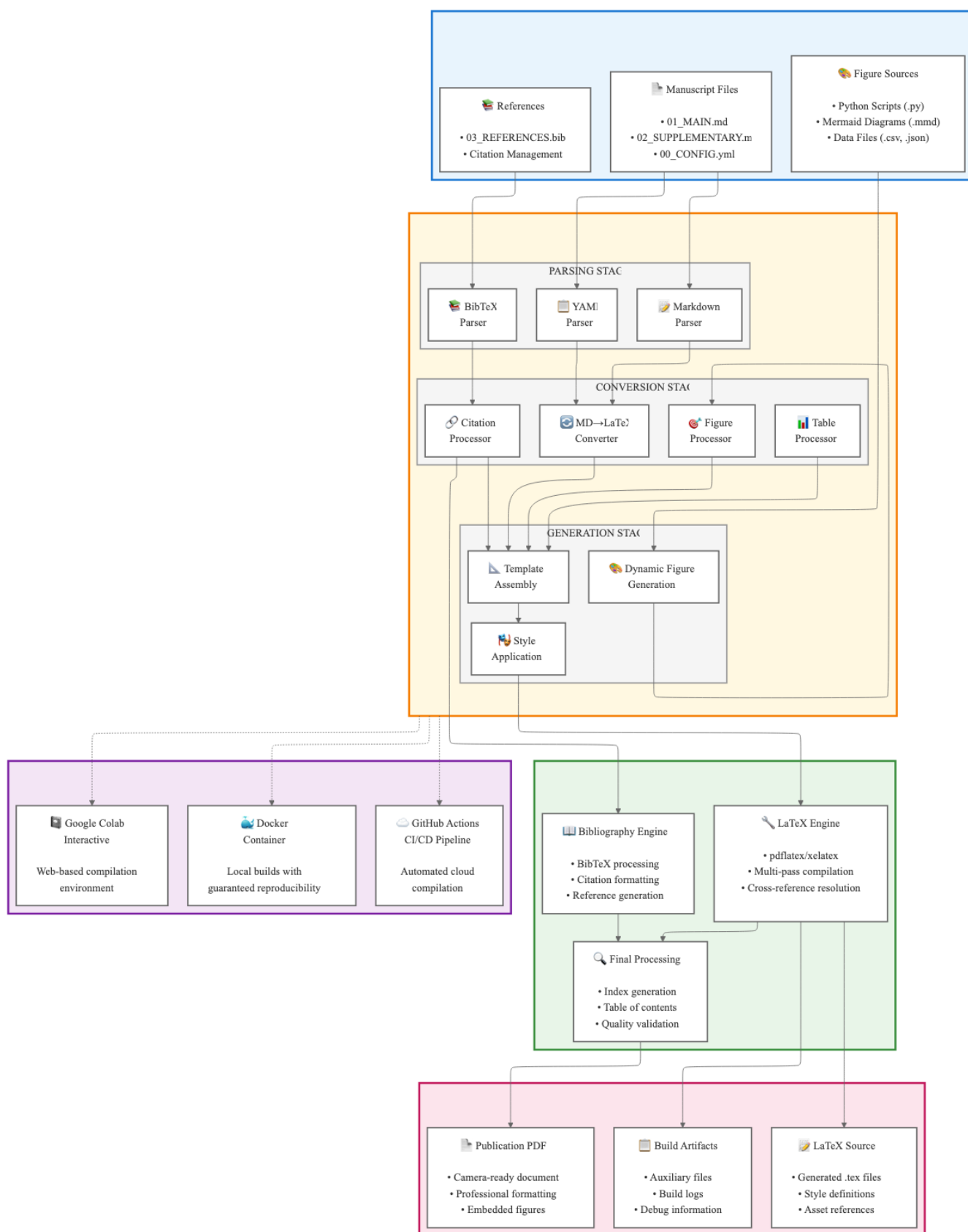
The containerised compilation environment, implemented through Docker, provides perfect isolation and reproducibility of the software environment. By encapsulating the exact versions of LaTeX, Python libraries, and system dependencies within a container image, the framework eliminates the common "works on my machine" problem that plagues many scientific computing workflows. This containerisation extends beyond mere convenience to serve as a critical component of scientific integrity, ensuring that the same input always produces identical output regardless of the host system configuration.

The framework's programmatic approach to figure generation creates an auditable chain from raw data to final visualisation. Python scripts that generate figures are version-controlled alongside the manuscript content, enabling complete reconstruction of all visual elements from source data. This approach contrasts sharply with traditional workflows where figures are created separately and inserted as static images, potentially leading to inconsistencies when data is updated or analysis methods are refined.

**Supp. Note 6: Template Customisation and Advanced Styling Options..** The RXiv-Maker framework provides extensive customisation capabilities through its LaTeX template system, enabling researchers to adapt the visual presentation to meet specific publication requirements whilst maintaining the simplicity of the markdown authoring experience. The template architecture separates content from presentation through a sophisticated class file (`rxiv_maker_style.cls`) that encapsulates all formatting decisions, typography choices, and layout specifications.

The YAML configuration system enables fine-grained control over document properties including author information formatting, institutional affiliation handling, and abstract presentation. Advanced users can modify template parameters to adjust margins, typography, colour schemes, and sectioning styles without requiring direct LaTeX modifications. The framework supports customisation of citation styles through configurable BibTeX style files, enabling compliance with specific journal requirements or institutional guidelines.

For institutions requiring consistent branding or specific formatting requirements, the framework provides extension points that enable custom style development whilst maintaining compatibility with the core processing pipeline. This extensibility ensures that RXiv-Maker can adapt to diverse institutional requirements without compromising its fundamental commitment to simplicity and ease of use.



**Sup. Fig. S1. Detailed System Architecture and Processing Layers.** Comprehensive technical diagram showing the complete RXiv-Maker architecture, including input layer organisation, processing engine components (parsers, converters, generators), compilation infrastructure, output generation, and deployment methodology integration. This figure illustrates the modular design that enables independent development and testing of system components.