



HAL
open science

TheoremView: A Framework for Extracting Theorem-Like Environments from Raw PDFs

Shrey Mishra, Neil Sharma, Antoine Gauquier, Pierre Senellart

► To cite this version:

Shrey Mishra, Neil Sharma, Antoine Gauquier, Pierre Senellart. TheoremView: A Framework for Extracting Theorem-Like Environments from Raw PDFs. ECIR (European Conference on Information Retrieval), Apr 2025, Lucca, Italy. pp.6. hal-04894570

HAL Id: hal-04894570

<https://inria.hal.science/hal-04894570v1>

Submitted on 17 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

TheoremView: A Framework for Extracting Theorem-Like Environments from Raw PDFs

Shrey Mishra¹[0009–0004–2357–9593], Neil Sharma²[0009–0003–8770–7374],
Antoine Gauquier¹[0009–0005–9573–6364], and
Pierre Senellart^{1,3}[0000–0002–7909–5369]

¹ DI ENS, ENS, CNRS, PSL University, Inria, Paris, France
shrey.mishra@ens.psl.eu, antoine.gauquier@ens.psl.eu,
pierre@senellart.com

² Malaviya National Institute of Technology Jaipur, India
neil.sharma3000@gmail.com

³ Institut Universitaire de France (IUF)

Abstract. This paper presents TheoremView, a novel framework for extracting proofs and theorems from raw PDF scientific papers without requiring \LaTeX source files. Our approach combines three modalities (**font**, **text**, and **vision**) with sequential modeling to capture long-term dependencies and layout information. By eliminating OCR preprocessing, TheoremView reduces computational overhead for real-time applications while providing robust automated theorem extraction. Our framework is publicly available at <https://theoremkb.org/demo>, with a demonstration video at <https://theoremkb.org/video>.

Keywords: Theorem extraction · Multimodal machine learning · Document analysis · Natural language processing

1 Introduction

1.1 Motivation for Theorem Extraction

In contemporary scientific research, articles are primarily published as PDFs, and many search engines index entire papers instead of specific scientific results. This paper contributes to TheoremKB [6], a project focused on building a knowledge base of mathematical results across different fields of science. The objective is to improve the accessibility of relevant information for researchers, allowing for more effective retrieval and utilization of scientific knowledge. In particular, TheoremKB aims at streamlining the retrieval of specific proofs and theorems, allowing quick access to targeted mathematical results compared to traditional full-text search engines.

Having such a knowledge base can significantly impact the way researchers find information [6]. Some advantages of this system include:

1. **Enhanced Accessibility:** TheoremKB streamlines the retrieval of specific proofs and theorems, allowing for quick access to targeted mathematical results, in contrast to traditional search engines that index full-text papers.

2. **Facilitated Knowledge Discovery:** TheoremKB assists researchers in uncovering connections between disparate mathematical results and their applications, thereby enhancing the exploration of specific results.
3. **Identification of Theorem Interdependencies:** TheoremKB helps determine which theorems are used in the proofs of others, which is essential for assessing the impact of errors in foundational results.
4. **Support for Automated Reasoning:** TheoremKB provides a foundation for developing AI systems capable of automated theorem proving, promoting innovative approaches to mathematical problem-solving.

1.2 Prior work on Theorems and Proofs Extraction

Previous attempts to address this task include the work presented in [8], which focused on initial explorations of extraction from PDFs framed as object detection and text classification problems. This approach utilized font visuals and text modalities but operated only at the text-line level. Subsequent research, such as [7], refined the methodology by incorporating contextual information surrounding paragraphs and employing multimodal systems to unify the extraction model. The TheoremView framework offers a user interface to visualize the results extracted by various models in an end-to-end system that directly takes PDFs as input and displays the extracted results. It is designed modularly, allowing users to select which model to utilize for extraction, thereby leveraging different modalities that highlight the strengths and weaknesses of each approach. This flexibility enables users to run models on low-compute hardware, such as systems without GPU instances, for inference. The primary objective of this paper is to present an easy-to-use interface that facilitates preprocessing and inference in a modular manner.

2 Methodology

We propose a modular approach to extract raw information from PDFs. We utilize **Grobid**⁴ [3] and **pdfalto**⁵ to convert the documents into valid XML formats. The XML data generated by Grobid organizes the content into paragraphs, while that from pdfalto provides segments content into text lines along with associated font information. We then employ a merging script to correlate the font information with each paragraph extracted by Grobid. This process yields a CSV file structured by paragraphs, where each row includes the spatial location of the paragraph on the page (indicating the page number as well as vertical and horizontal coordinates), the textual content extracted from Grobid, and the font information used within those paragraphs. For a schematic diagram of the data pipeline refer to Fig. 1.

Once the information is stored in CSV format, we process the font information using an **LSTM** [1] model, where each font is encoded as a unique token

⁴ <https://github.com/kermitt2/grobid>

⁵ <https://github.com/kermitt2/pdfalto>

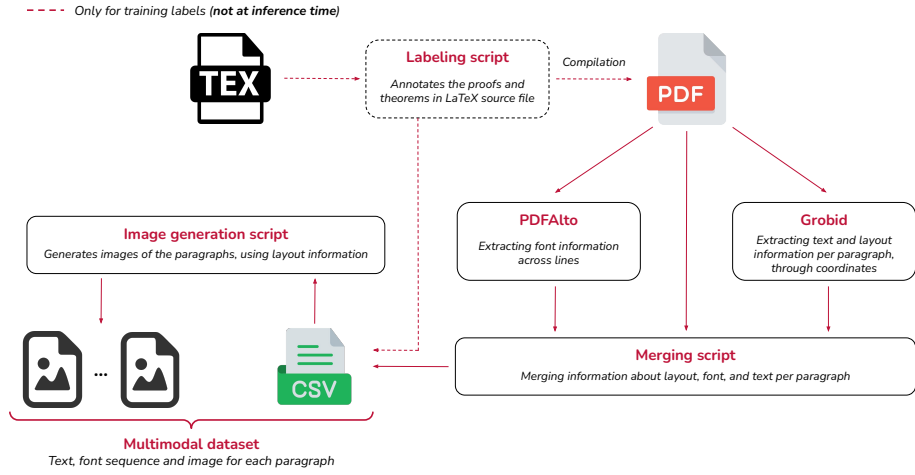


Fig. 1: Data pipeline for extracting and processing information from PDFs

to train the network. Simultaneously, we utilize the bitmap image rendering of each paragraph to train an **EfficientNetV2** model [10]. Additionally, we employ a **RoBERTa**-like language model which is pretrained from scratch [5], on a scientific corpus, to make predictions based on the text modality. Subsequently, we integrate all three trained models into a unified multimodal architecture, freezing the weights of each modality backbone and adding additional layers to capture intermodality interactions through mechanisms like Gated Multimodal Units (GMU) [9] or cross-modality attention similar to ViLBERT [4] that capture intermodality dependencies. Refer to Fig. 2 for a summarized view of the model inference pipeline and to [7, 5] for a detailed presentation of the architecture.

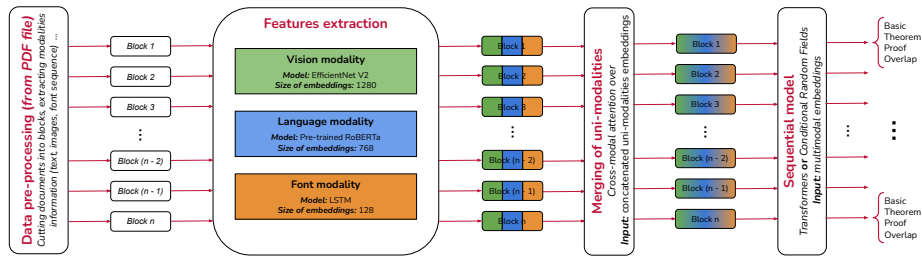


Fig. 2: Model inference pipeline (adding the sequential paragraph component)

With a set of base features extracted from either the unimodal or multimodal approaches, we generate features for all paragraphs within the PDF. This pro-

cess incorporates normalized page information, normalized coordinates of each paragraph, and paragraph embeddings derived from the raw features just before the softmax layer. To capture sequential information across multiple paragraphs, we train a Conditional Random Field (CRF) [2] or Transformer layer on top of the extracted features. The goal is to utilize relative information to contextualize each paragraph and accurately determine its label. Our model categorizes paragraphs into four major classes: (1) **Proof**, (2) **Theorem-like**, (3) **Basic** (neither proof nor theorem), and (4) an **Overlap** reject class that arises from preprocessing discrepancies.

3 Demonstration Scenario

The TheoremView demo interface, built using Streamlit, follows a modular architecture with distinct functional components. The frontend allows users to upload PDFs or select from cached samples for metadata processing using the Grobid and pdftalto tools, with results stored as CSV files. Users can run various ML models (unimodal or multimodal) through a pipeline that calculates processing time, generates bounding boxes, creates cropped images of theorems/proofs, and produces analytical graphs. The system implements efficient caching using pickle files for frequently accessed PDFs and ML model results.

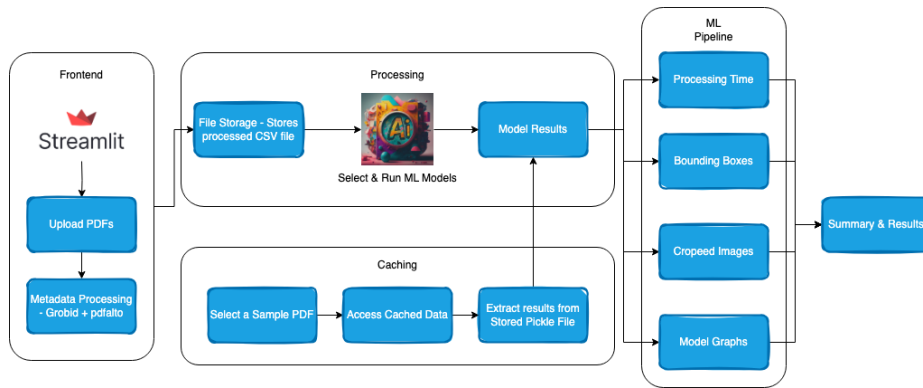
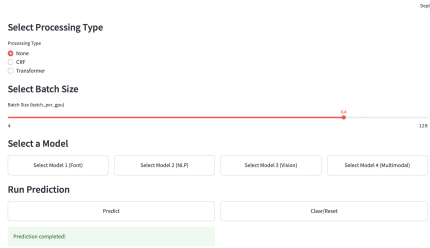


Fig. 3: System architecture of the various UI components

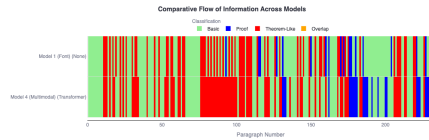
The UI of the demo is organized into several components (for an overview see Fig. 3), each serving a specific function:

1. **Upload and Process:** Users can upload PDFs or select from cached examples. The system processes PDFs using Grobid and pdftalto, converts pages to bitmap images, and merges XML outputs to generate a preprocessed `data.csv` file.

- Predict and Preview:** Users can select unimodal or multimodal base models, with optional sequential processing using CRF, Transformer, or none, offering 12 possible combinations. Results can be previewed or downloaded (see Fig. 4).
- Summary and Statistics:** Provides a breakdown of inference time for current and cached runs, enabling comparative analysis.



(a) Model selection interface



(b) Flow of information across paragraphs

Observe that taking expectations with respect to a uniform \tilde{R} on both sides in the conclusion of Lemma 4.5, we get that next-block hardness in relative entropy is equal to the sum of next-block inaccessible relative entropy and the expectation of the error term coming from the rejection sampling procedure. The following lemma upper bounds this expectation.

Lemma 4.6. *Let \tilde{G} be an online m -block generator, and let $L_i \stackrel{\text{def}}{=} 2^{|\tilde{G}_i|}$ be the size of the codomain of \tilde{G}_i , $i \in [m]$. Then for all $i \in [m]$, $r_{<i} \in \text{Supp}(\tilde{R}_{<i})$ and uniform \tilde{R}_i :*

$$\mathbb{E}_{y_i \sim \tilde{G}_i(r_{<i}, \tilde{R}_i)} \left[\log \frac{1}{\Pr[\tilde{Y}_i = y_i | Y_i = y_i, \tilde{R}_{<i} = r_{<i}]} \right] \leq \log \left(1 + \frac{L_i - 1}{T} \right).$$

Proof of Lemma 4.6. By definition of Sim^{GJ} , we have

$$\Pr[\tilde{Y}_i = y_i | Y_i = y_i, \tilde{R}_{<i} = r_{<i}] = 1 - \left(1 - \Pr[\tilde{G}_i(r_{<i}, \tilde{R}_i) = y_i] \right)^{L_i}$$

(c) PDF rendering with predictions

Classification Statistics

TOTAL PARAGRAPHS	224	100	50	0	74
		44.64%	22.32%	0.00%	33.04%

Extractions Page 1 of 19

Definition 1.1 (next-block pseudentropy, informal). Let n be a security parameter, and $X = (X_1, \dots, X_n)$ be a random variable distributed on strings of length $\text{poly}(n)$. We say that X has next-block pseudentropy at least k if there is a random variable $Z = (Z_1, \dots, Z_n)$, jointly distributed with X , such that:

Definition 1.1 (next-block pseudentropy, formal). Let n be a security parameter, and $X = (X_1, \dots, X_n)$ be a random variable distributed on strings of length $\text{poly}(n)$. We say that X has next-block pseudentropy at least k if there is a random variable $Z = (Z_1, \dots, Z_n)$, jointly distributed with X , such that:

1. For all $i = 1, \dots, m$, $(X_1, \dots, X_{i-1}, X_i)$ is computationally indistinguishable from $(X_1, \dots, X_{i-1}, Z_i)$.

Probability: 0.9254

Probability: 0.6836

(d) Extraction results

Fig. 4: User interface elements for model selection and predictions

Acknowledgments. This work was funded in part by the French government under management of Agence Nationale de la Recherche as part of the “Investissements d’avenir” program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute). This work was also made possible through HPC resources of IDRIS granted under allocation 2020-AD011012097 made by GENCI (Jean Zay supercomputer).

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Comput.* **9**(8), 1735–1780 (1997). <https://doi.org/10.1162/NECO.1997.9.8.1735>
2. Lafferty, J.D., McCallum, A., Pereira, F.C.N.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In: Brodley, C.E., Danyluk, A.P. (eds.) *Proceedings of the Eighteenth International Conference on Machine Learning (ICML 2001)*, Williams College, Williamstown, MA, USA, June 28 - July 1, 2001. pp. 282–289. Morgan Kaufmann (2001)
3. Lopez, P.: GROBID: combining automatic bibliographic data recognition and term extraction for scholarship publications. In: Agosti, M., Borbinha, J., Kapidakis, S., Papatheodorou, C., Tsakonas, G. (eds.) *Research and Advanced Technology for Digital Libraries, 13th European Conference, ECDL 2009, Corfu, Greece, September 27 - October 2, 2009. Proceedings. Lecture Notes in Computer Science*, vol. 5714, pp. 473–474. Springer (2009). https://doi.org/10.1007/978-3-642-04346-8_62
4. Lu, J., Batra, D., Parikh, D., Lee, S.: ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E.B., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. pp. 13–23 (2019), <https://proceedings.neurips.cc/paper/2019/hash/c74d97b01eae257e44aa9d5bade97baf-Abstract.html>
5. Mishra, S.: Multimodal Extraction of Proofs and Theorems from the Scientific Literature. (Extraction multimodale de preuves et de théorèmes à partir de la littérature scientifique). Ph.D. thesis, Paris Sciences et Lettres University, France (2024), <https://tel.archives-ouvertes.fr/tel-04665528>
6. Mishra, S., Brihmouche, Y., Delemazure, T., Gauquier, A., Senellart, P.: First steps in building a knowledge base of mathematical results. In: Ghosal, T., Singh, A., Waard, A., Mayr, P., Naik, A., Weller, O., Lee, Y., Shen, S., Qin, Y. (eds.) *Proceedings of the Fourth Workshop on Scholarly Document Processing (SDP 2024)*. pp. 165–174. Association for Computational Linguistics, Bangkok, Thailand (Aug 2024), <https://aclanthology.org/2024.sdp-1.16/>
7. Mishra, S., Gauquier, A., Senellart, P.: Modular multimodal machine learning for extraction of theorems and proofs in long scientific documents. In: *JCDL ’24: ACM/IEEE Joint Conference on Digital Libraries, Hong Kong, China, December 16–20, 2024*. ACM, Hong Kong, China (2024). <https://doi.org/https://doi.org/10.1145/3677389.3702540>
8. Mishra, S., Pluinage, L., Senellart, P.: Towards extraction of theorems and proofs in scholarly articles. In: Healy, P., Bilauca, M., Bonnici, A. (eds.) *DocEng ’21: ACM Symposium on Document Engineering 2021, Limerick, Ireland, August 24-27, 2021*. pp. 25:1–25:4. ACM (2021). <https://doi.org/10.1145/3469096.3475059>
9. Ovalle, J.E.A., Solari, T., Montes-y-Gómez, A., González, F.A.: Gated multimodal networks. *Neural Comput. Appl.* **32**(14), 10209–10228 (2020). <https://doi.org/10.1007/S00521-019-04559-1>
10. Tan, M., Le, Q.V.: EfficientNetV2: Smaller models and faster training. In: Meila, M., Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event. Proceedings of Machine Learning Research*, vol. 139, pp. 10096–10106. PMLR (2021)