

PseudoSeer: a Search Engine for Pseudocode

Levent Toksoz¹ **Mukund Srinath²** **Gang Tan¹** **C. Lee Giles^{1,2}**

The Pennsylvania State University,
¹ Computer Science and Engineering,
 University Park, Pennsylvania, 16802
 lkt5297@psu.edu mus824@psu.edu

The Pennsylvania State University,
²Information Sciences and Technology,
 University Park, Pennsylvania, 16802
 gtan@psu.edu clg20@psu.edu

Abstract

A novel pseudocode search engine is designed to facilitate efficient retrieval and search of academic papers containing pseudocode. By leveraging Elasticsearch, the system enables users to search across various facets of a paper, such as the title, abstract, author information, and LaTeX code snippets, while supporting advanced features like combined facet searches and exact-match queries for more targeted results. A description of the data acquisition process is provided, with arXiv as the primary data source, along with methods for data extraction and text-based indexing, highlighting how different data elements are stored and optimized for search. A weighted BM25-based ranking algorithm is used by the search engine, and factors considered when prioritizing search results for both single and combined facet searches are described. We explain how each facet is weighted in a combined search. Several search engine results pages are displayed. Finally, there is a brief overview of future work and potential evaluation methodology for assessing the effectiveness and performance of the search engine is described.

1 Introduction

The ever-growing volume of academic literature necessitates efficient search tools that cater to the specific needs of researchers. Within this domain, code plays a crucial role in various disciplines, serving as a foundation for replicating experiments, validating methodologies, and promoting understanding. Traditional academic search engines, however, often lack the functionality to effectively search for code within research papers. In this paper, we present PseudoSeer, a pseudocode search engine capable of searching over pseudocode in academic documents.

PseudoSeer bridges the gap between traditional text-based search and the need to locate code snippets embedded within academic publications. Pseudocode uses a language that closely resembles programming languages but focuses on readability and general concepts over syntax specifics, allowing researchers to represent algorithms and data structures in an easy-to-understand manner. A dedicated search engine capable of indexing and retrieving pseudocode offers a significant advantage by enabling targeted searches based on the inherent logic and functionality described within a paper.

This paper details the design and implementation of a pseudocode search engine built using the robust capabilities of Elasticsearch. The system empowers users to search across various facets of a research paper, including the title, abstract, author information, and LaTeX code sections, as well as perform combined searches across multiple facets. This functionality caters to diverse search scenarios. For instance, a researcher might utilize the engine to locate papers implementing a specific algorithm or explore the use of a particular data structure across various publications.

The development of this system involved several key considerations. Dataset is carefully chosen, with Arxiv serving as the primary source for academic papers containing code, ensuring both extensive coverage and diversity to handle broad range of search queries. A critical component we discuss in the paper is the indexing strategy, where decisions regarding how different data elements (text, code snippets, author information) are stored and optimized for search significantly impact retrieval

effectiveness. Another core aspect is the ranking algorithm, which determines the order in which search results are presented to the user. Factors such as performing single or combined facet searches, along with the relevance of each facet, all play a role in prioritizing the results.

Developing a pseudocode search engine presents both opportunities and challenges. On one hand, the system offers a powerful tool for code-centric literature exploration. However, complexities arise in accurately identifying and parsing code sections within academic papers. Additionally, tradeoffs exist between comprehensiveness and search performance. Indexing a vast amount of code data can be computationally expensive, while focusing solely on keywords within code snippets might miss the underlying logic and functionality. This paper addresses these challenges by outlining the chosen approaches and discussing the inherent difficulties involved. In summary, our contributions are:

- **Development of PseudoSeer:** A specialized search engine for pseudocode retrieval
- **Facet-Based Search Functionality:** Support for combined facet searches across title, abstract, references, author, and pseudocode for refined queries.
- **Exact-Match Search:** Exact-match query functionality using quotation marks.
- **Weighted Ranking Mechanism:** Implementation of a BM25-based ranking system with facet-specific weighting to enhance relevance in search results.

2 Related Work

Several pseudocode datasets have been developed to support research in various code related translation and generation tasks. Notably, the SPOC dataset Kulal et al. (2019) contains around 20,000 programs with human-authored pseudocodes and test cases, designed for translating pseudocode to C++ code. Another significant dataset is by Oda et al. (2015), which includes approximately 16,000 manually crafted pseudocodes. The dataset is used for statistical machine translation. Additionally, Zavershynskyi et al. (2018) provide a dataset of about 2,000 manually written pseudocodes for program synthesis tasks.

While these datasets are valuable for specific tasks, they represent a limited variety of pseudocodes and are hand-crafted by programmers. Moreover, their relatively small size and lack of diversity make them less suitable for use in search engines, which require more extensive and varied datasets to effectively handle a broad range of pseudocode search queries. To that end, Toksoz et al. (2024) provide a much larger dataset containing 320,000 pseudocode samples extracted from scholarly papers on arXiv, which will serve as the foundation for the dataset we use to build our pseudocode search engine.

Several search engines have tackled similar challenges in retrieving structured content, particularly in domains that involve specialized formats similar to pseudocode. For instance, Kohlhase and Şucan (2006) and Kohlhase et al. (2012) developed a math formula search engine called MathWebSearch, which structurally indexes mathematical expressions to enhance retrieval efficiency. Further, Normann and Kohlhase (2007) introduced extended normalization techniques aimed at mathematical formulas to detect not only structurally identical queries but also logically equivalent ones, regardless of the nomenclature chosen.

In addition to search engines, various efforts have focused on extracting structured information from documents, employing methods designed to retrieve elements such as mathematical equations, tables, figures, and metadata. These works demonstrate the potential of machine learning models to accurately extract structured content from PDFs, yet there remains a gap in addressing the unique challenges of generating pseudocode datasets from scholarly sources like arXiv. Recent advancements in this area include methods developed by Kardas et al. (2020), who built a machine learning pipeline for extracting results from research paper, and Nassar et al. (2022), who introduced TableFormer, a model specifically designed to capture table data. Similarly, Davila et al. (2021) present techniques for

chart extraction, aiming to improve the accessibility of visual data in academic texts. Other notable contributions, such as and Hu et al. (2005), focus on math equation detection, while Mali et al. (2020), Hou et al. (2019) and Hou et al. (2021) address the extraction of tasks, datasets, and evaluation metrics as distinct entities. Blecher et al. (2023) and tools like GRO (2008–2023) further demonstrate advancements in general metadata extraction and scientific document processing.

3 Data Collection

The pseudocode search engine utilizes data provided by Toksoz et al. (2024)¹, which is extracted from the LaTeX files of scholarly papers on arXiv, starting with year 1991 and ending in June of 2023. These files are stored across both Amazon S3 buckets and Google Cloud. To gather the LaTeX files, their extraction process scanned over 10 TB of arXiv data, recursively extracting ZIP files and matching the extracted content to the corresponding papers using arXiv identifiers. They also stored all the metadata information such as the arXiv identifier, any equations referenced by the pseudocode, and the year it was stored in arXiv. The pseudocode extraction process operates as follows: first, it identifies the `\begin{algorithm}` and `\end{algorithm}` tags within the LaTeX files of scholarly papers on arXiv. Next, it extracts the sections enclosed by these tags, recognizing them as pseudocode. In total nearly 320,000 pseudocode examples over 2.2 million scholarly papers are extracted.

The pseudocode dataset obtained from Toksoz et al. (2024) is further processed to include additional metadata information such as arXiv weblinks of the papers. Following the Reference Detection Algorithm described by Toksoz et al., references are extracted by matching labels with their associated ‘ref’ tags using regular expressions, and capturing a window of text around the referenced section. The extracted reference snippets are then further processed by shortening and cleaning the LaTeX syntax as much as possible to provide a more accessible overview and improved search experience. It should be noted that not all pseudocode references could be extracted, particularly those lacking associated reference tags or having tags that are too complex for regular expressions to handle. Additionally, papers and pseudocode that could not be parsed and indexed without errors were omitted, such as those with incomplete LaTeX code.

4 Search Interface

The pseudocode search engine² prioritizes a user-friendly interface that resembles common search engines to ensure ease of use. The landing page presents a search bar where users can enter their queries. Figures 2 and 1 illustrate these features. Unlike traditional search engines, however, this interface offers additional functionalities to refine the search based on the specific elements within a research paper. The page provides a set of radio buttons that allow users to specify which part of the paper they want to search, with options including title, abstract, author, or LaTeX code sections. Users can also select multiple buttons simultaneously to perform a combined search. By default, if no buttons are selected, a combined search across all fields is conducted. The interface also supports exact query matching by allowing users to wrap quotation marks around the desired query. This functionality caters to researchers with diverse search goals. For instance, a user aiming to locate papers implementing a specific algorithm might choose to search within the code sections, while someone interested in an author’s code contributions could utilize the author search option.

The search results page displays a list of retrieved papers, ordered by a custom ranking function discussed in the following section. Each entry on the results page provides relevant information about the retrieved paper, including the title, arXiv website, authors, and a snippet of text containing the matched keywords. The system also highlights the specific sections of the paper where the search

¹The dataset can be accessed via the *arxiv-pseudocode* repository on GitHub.

²The search engine can be accessed via the *pseudoseer* link.

terms were found (title, abstract, code snippet). This visual cue can be beneficial for users by allowing them to quickly assess if the retrieved paper aligns with their search intent. Furthermore, the interface offers options for pagination, enabling users to explore a larger set of results if the initial search retrieves a significant number of papers.

Results 1 - 5 of 1817

Approximate k-NN Graph Construction: a Generic Online Approach

Wan-Lei Zhao, Hui Wang, Chong-Wah Ngo

Website:
1804.03032

Abstract:

Nearest neighbor search and k-nearest neighbor graph construction are two fundamental issues arise from many disciplines such as multimedia information retrieval, data-mining and machine learning. They become more and more imminent given the big data emerge in various fields in recent years. In this paper, a simple but effective solution both for approximate k-nearest neighbor search and approximate k-nearest neighbor graph construction is presented. These two issues are addressed jointly in our solution. On the one hand, the approximate k-nearest neighbor graph construction is treated as a search task. Each sample along with its k-nearest neighbors are joined into the k-nearest neighbor graph by performing the nearest neighbor search sequentially on the graph under construction. On the other hand, the built k-nearest neighbor graph is used to support k-nearest neighbor search. Since the graph is built online, the dynamic update on the graph, which is not possible from most of the existing solutions, is supported. This solution is feasible for various distance measures. Its effectiveness both as k-nearest neighbor construction and k-nearest neighbor search approaches is verified across different types of data in different scales, various dimensions and under different metrics.

1. Pseudocode

LaTeX: ▼

```
\begin{algorithm}
\KwData{\textit{q}: query, $G$: \textit{k}-NN Graph, $Q$: \textit{k}-NN list of \textit{q}}
\KwResult{\textit{Q}: \textit{k}-NN list of \textit{q}}
$Q \leftarrow \emptyset$; $Flag[1 \dots n] \leftarrow 0$;
$R[1 \dots p] \leftarrow \textit{p}$ random samples;
\For {each $r \in R$} {
    $InsertQ(r, m(q, r), Q)$;
}
\For {$r \leftarrow top(Q)$} {
    \If {$Flag[r] == 0$} {
        \For {each $n_i \in G[r]$} {
            $InsertQ(n_i, m(q, n_i), Q)$;
        }
        \For {each $n_j \in \overline{G}$} {
            $InsertQ(n_j, m(q, n_j), Q)$;
        }
    }
    $Flag[r] = 1$;
}
\caption{Enhanced Hill-Climbing search (EHC)}
\label{alg:hcsearch}
\end{algorithm}
```

Figure 1: Results page obtained by using the search query 'nearest neighbor' in the abstract field

PseudoSeer{}

The screenshot shows the PseudoSeer landing page. At the top is a large title "PseudoSeer{}". Below it is a search bar with the placeholder "Enter Query Here". Underneath the search bar are five search filters: "LaTeX", "References", "Title", "Abstract", and "Authors", each preceded by a small square checkbox. Below these filters is a "Search" button. A text overlay states "Over 2.2 Million papers have been explored." At the bottom, it says "Powered by: elasticsearch" with the Elasticsearch logo.

Enter Query Here

LaTeX References Title Abstract Authors

Search

Over 2.2 Million papers have been explored.

Powered by: elasticsearch

Figure 2: Landing page with options to search in LaTeX, references, title, abstract, and authors, individually or combined.

5 Tokenization and Indexing

The effectiveness of the pseudocode search engine hinges on a robust indexing strategy that efficiently stores and retrieves information from various elements within a research paper. This section delves into the process of indexing different data types, with a particular focus on the complexities involved in handling both LaTeX pseudocode and its cleaned references.

Standard text processing techniques are employed for indexing the title, abstract, and author information of research papers. This involves tokenization, where the text is broken down into individual words or meaningful units. We tokenized each of the above mentioned fields using grammar based tokenization that works based on the Unicode text segmentation algorithm.

5.1 Indexing Code

Indexing code sections, especially those written in LaTeX, presents a unique challenge. LaTeX, a document typesetting system, utilizes a combination of plain text and markup commands to define the structure and formatting of the document. A straightforward approach involves converting the entire LaTeX code to plain text before performing tokenization. This simplifies the indexing process but can lead to loss of code-specific constructs that are essential for meaningful search results. For instance, when commands such as "\for" or "\begin{algorithm}" are removed, structure and logical flow in the pseudocode can be compromised.

In our search engine, we retain the LaTeX code itself for indexing, treating it as regular text while keeping important commands intact to preserve the structural cues within the pseudocode. This allows us to capture specific elements like "\for" loops and "\If" conditions directly, without breaking down the LaTeX into plain text and risking a loss of context and structure, especially since there are a considerable number of LaTeX pseudocode examples that are like Figure 3.

```

\For each element in the list
  \If element satisfies condition
    perform action
  \EndIf
\EndFor

```

Figure 3: Example pseudocode with for and if statements.

We also index the references surrounding each pseudocode snippet to enable keyword-based searches within the pseudocode’s descriptive context, while avoiding complex LaTeX syntax. This dual-layered indexing is particularly useful for users who may not know or prefer not to use the exact code syntax, as it allows them to locate pseudocode based on broader thematic keywords and descriptions found in the surrounding text, rather than being restricted to the code-construct-heavy language within the pseudocode itself. While this strategy requires more storage and adds complexity to the retrieval and ranking process, it enhances the search engine’s flexibility and precision for both structural and contextual searches.

6 Ranking

Another crucial factor that affect the effectiveness of the pseudocode search engine is an accurate and efficient ranking mechanism that prioritizes the most relevant results for a given query. This section details the ranking process employed by the search engine. For single field searches, the search engine utilizes the BM25 algorithm, which is based on TF-IDF (Term Frequency-Inverse Document Frequency) metric. The BM25 algorithm considers the frequency of query terms in a document and adjusts the score based on the rarity of those terms across all indexed documents, while applying a saturation effect to prevent overly frequent terms from dominating the score. The score is also normalized by the document length to ensure that more relevant results are ranked higher while accounting for natural variations in term frequency and document length.

In multi-field searches, the search engine assigns predetermined weights to each field. These weights are then incorporated into the BM25 algorithm to calculate the weighted scores. The weights are chosen based on the potential relevance of the fields in real-world use cases. For instance, results obtained by searching in the LaTeX pseudocode and authors fields are likely to be less relevant to the user than those from the abstract and title. To that end, the weights for the LaTeX pseudocode and authors fields are set to one, while the weights for other fields are set to two.

7 Evaluation

Currently, the search engine evaluation is conducted by manually inspecting the results obtained from various query types. Figure 4 illustrates a single-field search and the corresponding results displayed by the search engine. Figure 5 demonstrates a combined-field search, showing how multiple fields can be queried simultaneously with the results displayed accordingly. Similarly, Figure 6 presents an exact-match search query and its respective results as produced by the search engine. Enhanced evaluation will be a focus of future.

8 Future Work

We briefly discuss potential improvements for the search engine. Currently, the data only consists of pseudocode from papers containing `\begin{algorithm}` and `\end{algorithm}` tags. A machine learning-based solution could be implemented to detect additional patterns and extract pseudocode,

thereby expanding the dataset. Additionally, the search functionality could be enhanced by using a BERT-based model to better match user queries with the desired pseudocode. Moreover, incorporating different ranking algorithms that can dynamically adapt to both single-field and multi-field search queries could improve the precision of results based on the query context. Additional improvements could include integrating LaTeX pseudocode rendering into the interface, along with creating a more comprehensive evaluation framework to measure the effectiveness of ranking algorithms and user experience enhancements.

9 Conclusion

A novel pseudocode search engine is introduced and built with Elasticsearch. It enables efficient searching across various facets of research papers containing pseudocode. These facets include the title, abstract, author information, and LaTeX code sections, providing a comprehensive tool for locating relevant academic content. Users can choose to search within a single field or combine multiple facets for more refined results. Additionally, the ranking algorithm is designed to prioritize relevant search results based on whether the user is performing a single-field or multi-field search, while also considering the relevance of each field in multi-field searches.

PseudoSeer{}

tree

LaTeX References Title Abstract Authors

Search

Over 2.2 Million papers have been explored.

Powered by:  elasticsearch

(a) Single Field Search

Unbiased Math Word Problems Benchmark for Mitigating Solving Bias

Zhicheng Yang, Jinghui Qin, Jiaqi Chen, and Xiaodan Liang

Website:

[2205.08108](https://github.com/yangzhch6/UnbiasedMWP)

Abstract:

In this paper, we revisit the solving bias when evaluating models on current Math Word Problem (MWP) benchmarks. However, current solvers exist solving bias which consists of data bias and learning bias due to biased dataset and improper training strategy. Our experiments verify MWP solvers are easy to be biased by the biased training datasets which do not cover diverse questions for each problem narrative of all MWPs, thus a solver can only learn shallow heuristics rather than deep semantics for understanding problems. Besides, an MVP can be naturally solved by multiple equivalent equations while current datasets take only one of the equivalent equations as ground truth, forcing the model to match the labeled ground truth and ignoring other equivalent equations. Here, we first introduce a novel MWP dataset named UnbiasedMWP which is constructed by varying the grounded expressions in our collected data and annotating them with corresponding multiple new questions manually. Then, to further mitigate learning bias, we propose a Dynamic Target Selection (DTS) Strategy to dynamically select more suitable target expressions according to the longest prefix match between the current model output and candidate equivalent equations which are obtained by applying commutative law during training. The results show that our UnbiasedMWP has significantly fewer biases than its original data and other datasets, posing a promising benchmark for fairly evaluating the solvers' reasoning skills rather than matching nearest neighbors. And the solvers trained with our DTS achieve higher accuracies on multiple MWP benchmarks. The source code is available at <https://github.com/yangzhch6/UnbiasedMWP>.

1. Pseudocode

LaTeX: ▾

```
\begin{algorithm}[t]
\small{
    \caption{Equivalent Expression Tree Generation}
    \label{alg:tree_variation}
    \textbf{Function}: \textit{Variation}($\mathbf{tree}$, $\mathbf{root}$)
    \KwIn{Expression $\mathbf{tree}$: $ \mathbf{tree} $; Root node of the input $\mathbf{tree}$: $ \mathbf{root} $;}
    \KwOut{Equivalent expression list: $ \mathbf{equList} $}
    \BlankLine
    \If{$\mathbf{rm} \ $ $ \mathbf{root} $ is null}
}
```

(b) Results page for Single Field

Figure 4: Single-field search using a search query 'tree' and the search field LaTeX

PseudoSeer{}

LLM

LaTeX References Title Abstract Authors

Search

Over 2.2 Million papers have been explored.

Powered by:  elasticsearch

(a) Combined search

Learning to Generate Better Than Your LLM

Jonathan D. Chang, Kiante Brantley, Rajkumar Ramamurthy, Dipendra Misra, Wen Sun

Website:
[2306.11816](https://arxiv.org/abs/2306.11816)

Abstract:
Reinforcement learning (RL) has emerged as a powerful paradigm for fine-tuning Large Language Models (LLMs) for text generation. In particular, recent LLMs such as ChatGPT and GPT-4 can engage in fluent conversations with users after finetuning with RL. Capitalizing on key properties of text generation, we seek to investigate RL algorithms beyond general purpose algorithms like Proximal Policy Optimization (PPO). In particular, we extend RL algorithms to allow them to interact with a dynamic black-box guide LLM and propose RL with guided feedback (RLGF), a suite of RL algorithms for LLM fine-tuning. We provide two ways for the guide LLM to interact with the LLM to be optimized for maximizing rewards. The guide LLM can generate text which serves as additional starting states for the RL optimization procedure. The guide LLM can also be used to complete the partial sentences generated by the LLM that is being optimized, treating the guide LLM as an expert to imitate and surpass eventually. We experiment on the IMDB positive sentiment, CommonGen, and TL;DR summarization tasks. We show that our RL algorithms achieve higher performance than supervised learning (SL) and the RL baseline PPO, demonstrating the benefit of interaction with the guide LLM. On both CommonGen and TL;DR, we not only outperform our SL baselines but also improve upon PPO across a variety of metrics beyond the one we optimized for. Our code can be found at <https://github.com/Cornell-RL/tril>.

1. Pseudocode

LaTeX: ▾

```
\begin{algorithm}[H]
\caption{\texttt{cpi{} } }
\label{alg:finetuning_cpi}
\begin{algorithmic}[1]
\State \texttt{Input: } $\pi_{\theta}$, guide $\pi_{ref}$, iterations $T$, mini-batch size $B$ 
\For{$t$ \in [T]}
    \State Rollin with $\beta\pi_{ref} + (1-\beta)\pi^t_{\theta}$ starting from $x_t$ 
    \State Rollout with $\pi^t_{\theta}$ to collect trajectories 
    \State Update $V^{\pi^t_{\theta}}_{\phi}$ with trajectories and compute advantage $A^{\pi^t_{\theta}}$ 
    \State Update $\pi_{\theta}$ using $ppo$ loss with $A^{\pi^t_{\theta}}$ 
\EndFor
\end{algorithmic}

```

(b) Results page for combined search

Figure 5: Combined search using the search query 'LLM' across the title and abstract fields

PseudoSeer{}

"bubble sort"

LaTeX References Title Abstract Authors

Search

Over 2.2 Million papers have been explored.

Powered by:  elasticseach

(a) Exact Query Matching

Resilient Computing with Reinforcement Learning on a Dynamical System: Case Study in Sorting

Aleksandra Faust, James B. Aimone, Conrad D. James and Lydia Tapia

Website:

[1809.09261](https://arxiv.org/abs/1809.09261)

Abstract:

Robots and autonomous agents often complete goal-based tasks with limited resources, relying on imperfect models and sensor measurements. In particular, reinforcement learning (RL) and feedback control can be used to help a robot achieve a goal. Taking advantage of this body of work, this paper formulates general computation as a feedback-control problem, which allows the agent to autonomously overcome some limitations of standard procedural language programming: resilience to errors and early program termination. Our formulation considers computation to be trajectory generation in the program's variable space. The computing then becomes a sequential decision making problem, solved with reinforcement learning (RL), and analyzed with Lyapunov stability theory to assess the agent's resilience and progression to the goal. We do this through a case study on a quintessential computer science problem, array sorting. Evaluations show that our RL sorting agent makes steady progress to an asymptotically stable goal, is resilient to faulty components, and performs less array manipulations than traditional Quicksort and Bubble sort.

1. Pseudocode

LaTeX: ▾

```
\begin{algorithm}[h]
    \caption{\small{RL Sort}}
    \label{alg:rlsort}
    \begin{algorithmic}[1]
        \INPUT $\mathbf{x}_{in}$, an array to sort
        \INPUT feature vector $\mathbf{F}$
        \INPUT learned feature vector parametrization $\boldsymbol{\theta} = [\theta_1, \dots, \theta_n]^\top$ 
        \OUTPUT $\mathbf{x}_s$, sorted array
        \STATE $\mathbf{x} \leftarrow \mathbf{x}_{in}$
        \WHILE{$\boldsymbol{\theta}^\top \mathbf{F}(\mathbf{x}) < 0$} \label{ln:while}
            \STATE $(i,j) = \operatorname{argmax}_{(k,l)} \mathbf{x}_{in}[1, \dots, \mathbf{x}_{in}]^2 \cdot (\boldsymbol{\theta}^\top \mathbf{F}(\mathbf{x}))_{k,l}^2$ 
            \STATE $\mathbf{x} \leftarrow \mathbf{x} \setminus \mathbf{x}_{ij} \cup \mathbf{x}_{ij} \leftarrow \mathbf{x}_{in}[i, \dots, j]^\top \leftarrow \mathbf{x}_{in}[j, \dots, i]^\top$ \COMMENT{place $i^{th}$ element at $j^{th}$ position}
        \ENDWHILE
    \end{algorithmic}

```

(b) Results page for Exact Match

Figure 6: Exact query matching with the phrase 'bubble sort' in quotation marks, searched within the abstract field

10 Acknowledgements

The arXiv is gratefully acknowledged for providing access to documents with pseudocode and their latex versions.

References

- 2008–2023. Grobid. <https://github.com/kermitt2/grobid>.
- Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. 2023. Nougat: Neural optical understanding for academic documents.
- Kenny Davila, Srirangaraj Setlur, David Doermann, Bhargava Urala Kota, and Venu Govindaraju. 2021. Chart mining: A survey of methods for automated chart analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):3799–3819.
- Yufang Hou, Charles Jochim, Martin Gleize, Francesca Bonin, and Debasis Ganguly. 2019. Identification of tasks, datasets, evaluation metrics, and numeric scores for scientific leaderboards construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Florence, Italy, 27 July – 2 August 2019.
- Yufang Hou, Charles Jochim, Martin Gleize, Francesca Bonin, and Debasis Ganguly. 2021. Tdmsci: A specialized corpus for scientific literature entity tagging of tasks datasets and metrics. In *Proceedings of the 16th conference of the European Chapter of the Association for Computational Linguistics*, Online, 19–23 April 2021.
- Yunhua Hu, Hang Li, Yunbo Cao, Dmitriy Meyerzon, and Qinghua Zheng. 2005. Automatic extraction of titles from general documents using machine learning. In *Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries*, JCDL ’05, page 145–154, New York, NY, USA. Association for Computing Machinery.
- Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor, and Robert Stojnic. 2020. Axcell: Automatic extraction of results from machine learning papers. *CoRR*, abs/2004.14356.
- Michael Kohlhase, Bogdan A. Matican, and Corneliu C. Prodescu. 2012. Mathwebsearch 0.5 – scaling an open formula search engine. In *Intelligent Computer Mathematics*, LNAI, pages 342–357. Springer.
- Michael Kohlhase and Ioan Sucan. 2006. A search engine for mathematical formulae. In *Proceedings of Artificial Intelligence and Symbolic Computation, AISC’2006*, LNAI, pages 241–253. Springer.
- Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy Liang. 2019. Spoc: Search-based pseudocode to code. *CoRR*, abs/1906.04908.
- Parag Mali, Puneeth Kukkadapu, Mahshad Mahdavi, and Richard Zanibbi. 2020. Scanssd: Scanning single shot detector for mathematical formulas in PDF document images. *CoRR*, abs/2003.08005.
- Ahmed Nassar, Nikolaos Livathinos, Maksym Lysak, and Peter Staar. 2022. Tableformer: Table structure understanding with transformers.
- Ioana Normann and Michael Kohlhase. 2007. Extended formula normalization for ϵ -retrieval and sharing of mathematical knowledge. In *MKM/Calculemus - Towards Mechanized Mathematical Assistants*, LNAI, pages 266–279. Springer.
- Yusuke Oda, Hiroyuki Fudaba, Graham Neubig, Hideaki Hata, Sakriani Sakti, Tomoki Toda, and Satoshi Nakamura. 2015. Learning to generate pseudo-code from source code using statistical machine translation (t). In *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 574–584.
- Levent Toksoz, Gang Tan, and C. Lee Giles. 2024. Automatic pseudocode extraction at scale. In *2024 IEEE International Conference on Information Reuse and Integration for Data Science (IRI)*, pages 264–269.
- Maksym Zavershynskyi, Alex Skidanov, and Illia Polosukhin. 2018. Naps: Natural program synthesis dataset.