



# Research Paper Classifier: Identifying Publishable Papers and Target Conferences



Team ALPHA

Nitin & Devansh

# Introduction to the Problem: Challenges in Research Paper Evaluation

## Subjective Criteria

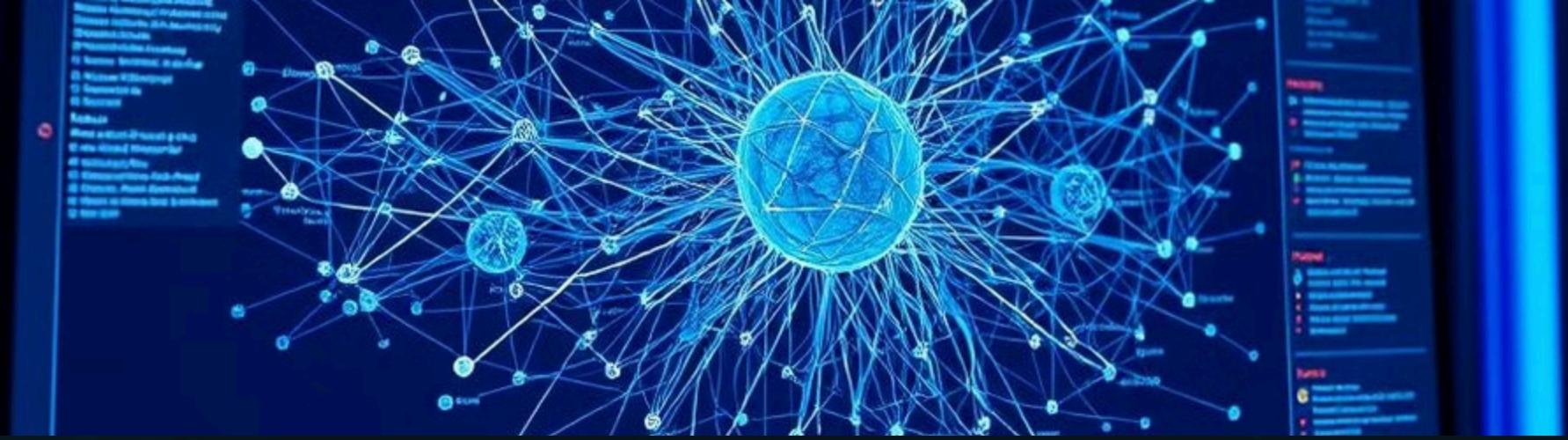
Determining publishability often relies on subjective criteria, making consistent evaluation difficult. For example, different reviewers might prioritize novelty differently, leading to variations in assessment. One reviewer might value a highly novel approach even with minor methodological flaws, while another might prioritize methodological rigor above all else, even if the novelty is less groundbreaking. This subjectivity introduces inconsistencies into the review process.

## Limited Resources

Reviewers may have limited time and expertise, leading to inconsistent and potentially biased evaluations. Reviewers are often juggling multiple responsibilities, and the time allocated to reviewing a single paper might be limited. This time constraint can lead to superficial reviews, overlooking critical details or nuances in the research. Furthermore, expertise is another crucial factor; if a reviewer isn't fully familiar with the specific field of the submitted paper, their evaluation might be less accurate and more prone to biases.

## Large Volume of Papers

The growing number of research papers necessitates a more efficient way to handle the evaluation process. The sheer volume of papers submitted to conferences and journals is increasing exponentially every year. This overwhelming number of submissions makes manual review processes slow, inefficient, and potentially unsustainable. A more automated or assisted system is needed to help filter and prioritize papers, ensuring that the best research receives the attention it deserves without overwhelming human reviewers.



# Methodology: Leveraging Large Language Models and Retrieval-Augmented Generation



## LLM-based Analysis

We utilize a large language model (LLM) to analyze the text content of research papers, extracting key insights and identifying potential weaknesses. This involves several steps: first, the paper is tokenized and embedded into a vector space. Then, the LLM processes this embedding to identify key themes, arguments, and methodological approaches. Finally, it analyzes these elements to assess the paper's strengths and weaknesses, identifying potential areas for improvement or expansion.



## RAG Implementation

Retrieval-Augmented Generation (RAG) enhances the model's accuracy by retrieving relevant information from a vast corpus of scientific literature. This corpus is pre-processed and indexed to enable efficient information retrieval. The RAG module searches the corpus for relevant parts from the reference data. These parts are then provided as context to the LLM, improving its understanding and enabling it to generate more accurate and comprehensive analyses.

# Preparing the Data for Analysis

## Dataset used

Our dataset, provided by the Khadagpur Data Science Group, consisted of 150 research papers, encompassing diverse research areas including computer science, machine learning, and natural language processing. These papers were collected from various conference proceedings, ensuring representation across different publication types. To guide initial model development and establish a reliable evaluation benchmark, 15 of these papers were manually labeled. This labeling assigned labels indicating both publishability and target conference suitability.

## Embeddings Generation

Embeddings were generated for both the research papers in our dataset and the reference data used by the RAG module. We employed FlagEmbedding from HuggingFace to generate high-quality, semantically meaningful vector representations of the text. The choice of embedding dimensions and model selection were determined through experimentation and evaluation on a held-out set of the manually labeled papers.

# Model Architecture: LLM API Integration and RAG Implementation

1

## LLM API Integration

We integrated a pre-trained LLM API of Gemini 2.0 Flash to perform natural language processing task. It is used to classify the paper as Publishable or Non-Publishable. It also generates the name of conference which the paper can be published in, and a rationale for the decision.

2

## RAG Implementation

RAG allows the model to retrieve relevant information from a knowledge base of scientific papers, enhancing its accuracy and contextual understanding. It saves vector embeddings of the reference papers, which are then compared with the embeddings of input paper to pull out the most relevant information.

3

## How the Model Works

The RAG implementation first analyzes and selects most relevant reference data, which is added to the main prompt. This prompt is then sent to LLM API for generating the necessary output. The output consists of 3 parts:

1. Yes/No : Is the paper publishable?
2. Conference(s) [If Yes] : If the paper is publishable, in what conference(s) is it publishable?
3. Rationale : The reason for these choices.

# Conference Recommendation: Matching Papers to Suitable Venues



## Relevance

The model accurately suggests conferences fitting the paper's topic. This is achieved by analyzing the keywords, abstract, and overall content of the research paper to identify relevant themes and research areas. The model then cross-references these themes against a database of conferences known for focusing on those areas.



## Coverage

The system suggests many conferences across different research areas, giving authors varied choices. This broad coverage ensures that even highly specialized papers have a chance to find a suitable venue. The model accounts for interdisciplinary research, recognizing that a paper might be relevant to multiple conference communities.



## Specificity

Recommendations consider the paper's specific contributions and novelty, suggesting appropriate venues based on theoretical or applied focus. The model assesses the paper's originality, the level of technical detail, and the overall impact of its findings. This enables it to distinguish between conferences focused on theoretical advancements and others centered on practical applications.



## Impact Factor

The model incorporates information about the impact factor and prestige of each suggested conference. This provides authors with a broader understanding of the potential reach and visibility of their work, helping them make an informed decision about where to submit.



# Conference Recommendation: Matching Papers to Suitable Venues



## Tapping Into a Wealth of Knowledge

Our model leveraged a vast and diverse dataset encompassing published and unpublished research papers from leading conferences including CVPR, EMNLP, KDD, NeurIPS, and TMLR. The diversity of this dataset was critical in enhancing the model's ability to identify key patterns indicative of high-quality, publishable research. This broad scope verifies that the model can differentiate nuanced distinctions between papers suitable for different conferences. The inclusion of unpublished papers provided valuable additional context and allowed the model to learn even more about the subtle differences in writing styles and research approaches that distinguish one conference from another.

# Future Enhancements and Potential Applications

1

## Multi-Lingual Support

Expanding the model to support multiple languages will broaden its accessibility and applicability to a global research community. This will allow researchers who are not native English speakers to benefit from the model's capabilities, increasing the overall impact of our work and fostering greater inclusivity in academic publishing.

2

## Real-Time Feedback

Providing real-time feedback to researchers during the writing process can significantly improve the quality of research papers before submission. This feature could be integrated into writing platforms or as a standalone tool, allowing researchers to receive immediate guidance on areas such as clarity, coherence, and alignment with conference guidelines. Such a system could also provide suggestions for improving the paper's structure and argumentation.

3

## Automated Review Process

Integrating the model into the peer review process can revolutionize the efficiency and transparency of academic evaluation. By automating the initial screening and providing preliminary feedback, the model could free up reviewers' time, allowing them to focus on more in-depth analysis and critical evaluation. This could lead to a faster and more efficient review cycle, accelerating the publication process and reducing the burden on reviewers.

4

## Improved Conference Matching

Enhancements could include a more sophisticated algorithm for matching papers to suitable venues. A more refined matching system would minimize mismatches and optimize the likelihood of successful publication in top-tier conferences.

5

## Integration with Citation Management

Integrating the classifier with popular citation management tools could automate the process of identifying relevant papers and conferences. Researchers could leverage this functionality to efficiently explore the research landscape and find optimal publication venues aligned with their work, enhancing the efficiency of their research workflow.

# Conclusion

Our research paper classifier represents a significant advancement in the evaluation and dissemination of research papers. By automating the process of identifying publishable papers and suggesting suitable conferences, this model promises to streamline the academic workflow and improve the efficiency of the publication process. Its potential impact extends beyond individual researchers to the broader research community, fostering greater collaboration and knowledge sharing. Future enhancements will focus on refining the model's accuracy and expanding its functionality to encompass multilingual support, real-time feedback for authors, and seamless integration with existing citation management tools. We believe that this technology can fundamentally transform how research is evaluated and disseminated, ultimately accelerating the pace of scientific discovery and knowledge advancement.

