



UNIVERSITY OF PADOVA

DEPARTMENT OF MATHEMATICS "TULLIO LEVI-CIVITA"

MASTER THESIS IN DATA SCIENCE

AUDITING OPEN-SOURCE LLM FOR ACADEMIC AUTHOR SEARCH

SUPERVISOR

PROF. TOMASO ERSEGHE
UNIVERSITY OF PADOVA

CO-SUPERVISORS

PROF. FARIBA KARIMI
GRAZ UNIVERSITY OF TECHNOLOGY
DR. LISETTE ESPÍN-NOBOA
COMPLEXITY SCIENCE HUB VIENNA

MASTER CANDIDATE

DANIELE BAROLO

STUDENT ID

2078339

ACADEMIC YEAR

2023-2024

“DIE GRENZEN MEINER SPRACHE BEDEUTEN DIE GRENZEN MEINER WELT.”
“THE LIMITS OF MY LANGUAGE MEAN THE LIMITS OF MY WORLD.”

– LUDWIG WITTGENSTEIN, *TRACTATUS LOGICO-PHILOSOPHICUS*, 1922.

Abstract

This thesis addresses a common practical challenge in academia: the need to identify and recommend scientists for various purposes, from research collaborations to forming workshop committees and admission boards. Traditional search tools, based on classical metrics, fail to capture the richness and diversity of the rapidly expanding global scientific landscape [1]. This limitation is particularly concerning given evidence that demographic diversity enhances scientific quality [2, 3]. With Large Language Models (LLMs) increasingly mediating academic searches, evaluating their role in ensuring inclusivity and equity becomes crucial. This study systematically evaluates four open-weights LLMs (Gemma2-9b, Mixtral-8x7b, LLaMA 3 8B/70B) across three key dimensions: factuality (verifying whether suggested individuals are actual scientists), response consistency across multiple runs, and demographic and popularity biases in physics author recommendation. We tested these models through five distinct tasks: top-k endorsements, field-specific expert identification, temporal-based suggestions, seniority-based recommendations, and "statistical twin" searches, including control scenarios. We used the American Physical Society dataset (678,916 physics publications) enriched with OpenAlex metadata as a validation framework. Our findings reveal complex patterns of capability and bias: while larger models achieve higher average factual accuracy across tasks (LLaMA 3 70B averaging 87% in scientist verification across all runs and use cases), they systematically favor established excellence, disproportionately recommending Nobel laureates (22-35% versus 0.032% baseline) and highly-cited researchers (above 90th percentile). Significant demographic biases persist, particularly in gender distribution (61-75% male recommendations versus 46% baseline) and representation of Asian researchers (8.6-22.3% versus 42.3% baseline). Contemporary or lesser-known scientists remain underrepresented, and citation verification proves challenging. We discussed these findings, suggesting that future academic search systems should adopt hybrid architectures combining LLMs with knowledge graph-based retrieval, potentially offering a path toward more equitable and reliable scholarly recommendations.

Contents

ABSTRACT	v
LIST OF FIGURES	ix
LIST OF TABLES	xi
LIST OF ABBREVIATIONS	xiii
LIST OF ABBREVIATIONS	xiii
1 INTRODUCTION	1
1.1 Background	1
1.2 Problem Statement	2
1.2.1 The Popularity Paradox	2
1.2.2 Representation and Diversity	2
1.2.3 Capturing Nuanced Academic Impact	3
1.3 Research Questions	3
1.4 Significance of the Study	4
2 LITERATURE REVIEW	5
2.1 Large Language Models: An Overview	5
2.1.1 History and Structure of Large Language Models	5
2.1.2 Scaling Laws	7
2.1.3 Data for Training	7
2.2 LLMs in Academic and Scientific Contexts	8
2.2.1 LLMs as Recommender Systems	9
2.3 Research Gaps and Novel Contributions	9
3 METHODOLOGY	11
3.1 Research Design	11
3.2 Data for Validation and Fact-Checking	12
3.2.1 American Physical Society Dataset	12
3.3 Task Design	15
3.4 Task Categories	16
3.5 Evaluation Framework	17
3.5.1 Development of Evaluation Criteria	17
3.5.2 Evaluation Metrics	18
3.6 Experimental Set-up	19
3.6.1 Model Selection	19
3.6.2 Execution Protocol	20
3.6.3 Parameter Settings	21
3.6.4 Prompt Design	21

4	RESULTS	25
4.1	Understanding LLMs as Academic Recommender Systems	25
4.2	Performance Evaluation	25
4.2.1	Factuality of Recommendations	25
4.2.2	Consistency and Error Rates	28
4.3	Analysis of Author Recommendation Patterns	30
4.3.1	Nobel Laureate Representation	31
4.3.2	Similarity Analysis of Recommendations	31
4.3.3	Popularity and Scholarly Metrics	34
4.3.4	Demographic Representation	35
5	CONCLUSION	37
5.1	Summary and Interpretation of Results	37
5.1.1	Representation Patterns and Biases	37
5.1.2	Technical Performance Characteristics	38
5.2	Limitations of the Study	39
5.3	Recommendations for Future Research	39
6	APPENDIX: DETAILED MODELS PERFORMANCE	41
6.1	Detailed Factuality Analysis	41
6.1.1	Authors Factuality	41
6.1.2	Field-Based Factuality	44
6.1.3	Temporal Factuality	45
6.1.4	Seniority-Based Factuality	46
6.2	Response Pattern Analysis	47
6.3	Nobel Laureates Analysis	51
6.4	Recommendations Similarity Analysis	54
6.5	Gender Distribution Analysis	57
6.6	Ethnicity Distribution Analysis	58
6.7	Rank Percentile Analysis	61
	REFERENCES	65
	ACKNOWLEDGMENTS	75

Listing of figures

2.1	Trend of LLM-related publications on arXiv	6
3.1	Validation-data enrichment workflow	12
3.2	Gender distribution across physics disciplines in the American Physical Society (APS) dataset.	14
3.3	Trends in APS publications and authorship over time.	15
3.4	Prompt template	22
3.5	System prompt	24

Listing of tables

2.1	Comparison of Machine Learning (ML) paradigms [4].	7
3.1	Overview of the APS Dataset enriched with Open Alex (OA) Metadata	13
3.2	APS Author Metrics in the Enriched Dataset	13
3.3	Summary of LLMs characteristics (sorted by parameter size)	20
3.4	Core criteria specifications for different task types	23
3.5	Output specifications for different task types	23
4.1	Factuality scores across tasks	26
4.2	Models consistency	28
4.3	Models Error Rates	30
4.4	Descriptive statistics of model recommendations.	31
4.5	Average percentile ranks of recommended authors	34
4.6	Gender distribution of LLMs' recommended authors	35
4.7	Ethnic distribution of LLMs' recommended authors	35
6.1	Author factuality scores for LLaMA 3 8B (llama3-8b)	42
6.2	Author factuality scores for Gemma2-9b-it (gemma2-9b)	42
6.3	Author factuality scores for Mixtral-8x7b (mixtral-8x7b)	43
6.4	Author factuality scores for LLaMA 3 70B (llama3-70b)	43
6.5	DOI recommendations across fields	44
6.6	Author field relevance analysis	45
6.7	Epoch-specific recommendations	45
6.8	Seniority-based recommendations	46
6.9	Response uniqueness for llama3-8b	47
6.10	Response uniqueness for gemma2-9b	48
6.11	Response uniqueness for mixtral-8x7b	48
6.12	Response uniqueness for llama3-70b	49
6.13	Jaccard similarity among recommended names	50
6.14	Author uniqueness ratio within requests	51
6.15	Nobel laureates retrieved by llama3-8b	52
6.16	Nobel laureates retrieved by gemma2-9b	52
6.17	Nobel laureates retrieved by mixtral-8x7b	53
6.18	Nobel laureates retrieved by llama3-70b	53
6.19	Institutional similarity	54
6.20	Country-level similarity	54
6.21	Co-authorship similarity	55
6.22	Categorical similarity	55
6.23	Scholarly metric similarity	56
6.24	Gender distribution in smaller models	57
6.25	Gender distribution in larger models	58
6.26	Ethnicity distribution in llama3-8b	58

6.27	Ethnicity distribution in gemma2-9b	59
6.28	Ethnicity distribution in mixtral-8x7b	59
6.29	Ethnicity distribution in llama3-70b	60
6.30	Rank percentiles for llama3-8b	61
6.31	Rank percentiles for gemma2-9b	62
6.32	Rank percentiles for mixtral-8x7b	62
6.33	Rank percentiles for llama3-70b	63

List of Abbreviations

- AI** Artificial Intelligence. 1, 2, 4, 8, 19
- API** Application Programming Interface. 12
- APS** American Physical Society. ix, xi, 4, 12–15, 18, 21, 26–28, 31, 32, 34, 35, 37–39, 42, 43
- BERT** Bidirectional Encoder Representations from Transformers. 14
- CCSV** Collective-Critique and Self-Voting. 2
- CM&MP** Condensed Matter and Materials Physics. 27, 42–45, 47–49
- DOI** Digital Object Identifier. 12
- gemma2-9b** Gemma2-9b-it. xi, xii, 1, 20, 26–31, 34, 35, 38, 42, 44–46, 48, 50–52, 54–57, 59, 62
- GLUE** General Language Understanding Evaluation. 7
- HH-RLHF** Helpfulness and Harmlessness – Reinforcement Learning from Human Feedback. 7
- LLaMA** Large Language Model Meta AI. 1, 28
- llama3-70b** LLaMA 3 70B. xi, xii, 20, 26–31, 33–36, 38, 43–46, 49–51, 53–58, 60, 63
- llama3-8b** LLaMA 3 8B. xi, xii, 20, 26–31, 34, 35, 42, 44–47, 50–52, 54–58, 61
- LLM** Large Language Model. v, vii, ix, xi, 1–9, 11, 15, 17–21, 25, 26, 28, 37–40
- LSTM** Long Short-Term Memory. 5
- mixtral-8x7b** Mixtral-8x7b. xi, xii, 1, 19, 20, 26–31, 33–35, 38, 43–46, 48, 50, 51, 53–59, 62
- ML** Machine Learning. xi, 6, 7
- MMLU** Massive Multitask Language Understanding. 7
- NLP** Natural Language Processing. 1, 5, 9
- OA** Open Alex. xi, 12, 13
- PER** Physics Education Research. 16, 27, 42–45, 47–49
- RAG** Retrieval Augmented Generation. 8, 40
- RLHF** Reinforcement Learning from Human Feedback. 6
- RNN** Recurrent Neural Network. 5
- RQ** Research Question. 3, 4
- S2ORC** Semantic Scholar Open Research Corpus. 7

1

Introduction

1.1 BACKGROUND

The digital age has ushered in unprecedented access to information, fundamentally altering how we interact with knowledge [5]. In the academic sphere, this transformation is particularly pronounced, with researchers navigating vast seas of data to stay abreast of developments in their fields. Traditional search engines and academic databases, while invaluable, often struggle to provide nuanced, context-aware results that truly capture the complexity of scientific inquiry [6].

LLMs, trained on massive amounts of human-generated text data, have emerged as a significant advancement in Artificial Intelligence (AI), demonstrating remarkable capabilities in Natural Language Processing (NLP), language understanding, and generation [7]. These sophisticated systems exhibit an impressive ability to process and produce human-like text across diverse domains. From answering complex queries to summarizing lengthy documents, LLMs are rapidly gaining importance in various sectors, including academia, showcasing their potential to enhance information processing and knowledge synthesis tasks [4].

The potential applications of LLMs in scholarly contexts are vast and largely unexplored. Imagine a research assistant capable of not just finding relevant papers, but understanding their content, identifying key contributors, and even suggesting potential collaborators based on shared interests or complementary expertise. Such a tool could dramatically accelerate the pace of scientific discovery and foster more interdisciplinary collaboration.

However, as with any powerful technology, the integration of LLMs into academic workflows raises important questions. How comprehensive is their knowledge of the scientific community? Can they accurately represent the diverse landscape of global research? And crucially, how do we ensure that these systems promote equity and inclusivity rather than perpetuating existing biases?

This thesis specifically examines open-source LLMs including mixtral-8x7b, gemma2-9b, and variants of models family Large Language Model Meta AI (LLaMA), with a particular focus on their application in physics re-

search. We aim to understand both the potential and limitations of these models in academic author search and recommendation tasks.

1.2 PROBLEM STATEMENT

As we stand on the cusp of widespread LLM adoption in academia, it is imperative that we critically examine their capabilities and limitations, particularly in the context of author search and recommendation. While there are promising results suggesting LLMs' potential to enhance various aspects of academic search, it is crucial to address the challenges that extend beyond mere technical considerations, touching on fundamental issues of fairness, representation, and the very nature of academic merit.

1.2.1 THE POPULARITY PARADOX

A critical concern in the adoption of LLMs for academic author search and recommendation is their potential to exacerbate the "Matthew effect" in science, where renowned researchers gain disproportionate visibility and recognition [8]. This phenomenon intersects with the well-documented issue of popularity bias in recommender systems, where popular items overshadow potentially relevant but less-known alternatives [9].

The integration of LLMs into this domain introduces new complexities. On one hand, there are justified concerns that LLMs, trained on datasets dominated by popular content, might amplify existing biases. This could create a feedback loop that further entrenches established names at the expense of emerging talent, potentially stifling diversity and innovation in academic discourse.

On the other hand, recent research suggests that LLM-based recommenders may exhibit less popularity bias compared to traditional systems [10]. This finding, while promising, does not negate the need for careful investigation and mitigation strategies. The potential for LLMs to either exacerbate or alleviate popularity bias likely depends on their implementation, training data, and the specific context of their application in academic search and recommendation.

1.2.2 REPRESENTATION AND DIVERSITY

Equally concerning is the question of representation. Academia has made strides in recent years to become more inclusive, yet disparities persist across gender, ethnicity, and geographical lines [11, 12]. LLMs, if not carefully designed and implemented, could perpetuate or even exacerbate these inequalities. The worry is not just about who gets recommended, but who gets overlooked – the voices that might be systematically excluded from the AI-mediated discourse of the future.

However, recent research suggests that the diversity of LLM outputs can be enhanced. Techniques such as Collective-Critique and Self-Voting (CCSV) have shown promise in improving the representation of different demographic groups in LLM responses [13]. These findings indicate that LLMs can be leveraged to promote diversity rather than hinder it, provided they are properly guided and their capabilities are fully utilized.

Nevertheless, it remains crucial to rigorously assess the diversity of LLM outputs in academic search and recommendation contexts. While the potential for improvement exists, the baseline performance of these models in terms of representation must be carefully evaluated.

1.2.3 CAPTURING NUANCED ACADEMIC IMPACT

The application of LLMs to academic author search presents another critical challenge: evaluating the multifaceted nature of scientific contributions beyond traditional metrics. While citation counts offer quantifiable data, they often fail to capture a researcher’s full impact [14]. The key question is whether LLMs can provide a more comprehensive assessment of academic contributions.

This challenge has two main aspects:

1. **Qualitative Evaluation:** Can LLMs’ advanced language processing capabilities evaluate subtle aspects of research impact that traditional metrics miss? This includes assessing mentorship quality, the influence of niche but groundbreaking work, and the value of negative results.
2. **Temporal and Interdisciplinary Insight:** Given their training on vast datasets, can LLMs effectively navigate the evolving landscape of scientific progress? We need to examine their ability to identify pioneering work across different eras and make connections between diverse fields.

The challenges presented by the integration of LLMs into academic author search and recommendation are interconnected and multifaceted. From addressing popularity bias to ensuring diverse representation and capturing nuanced academic impact, each aspect requires careful consideration. As we move forward, it is crucial to develop robust evaluation frameworks, explore bias mitigation techniques, and foster interdisciplinary collaboration. By addressing these challenges, future LLM-based systems that contribute to a more equitable, diverse, and nuanced landscape of academic search and recommendation could be developed.

1.3 RESEARCH QUESTIONS

In light of these challenges, our research aims to systematically evaluate the current capabilities of state-of-the-art open-source LLMs in identifying and ranking prominent scientists. We are guided by the following Research Questions (RQs):

RQ1. How do LLMs represent minorities in scientific author identification?

This question examines quantifiable patterns of bias across multiple dimensions:

- Demographic representation (gender, ethnicity, geographical distribution)
- Institutional clustering effects
- Citation-based popularity bias
- Career-stage and temporal biases

However, in order to assess the representation of minorities, we first need to evaluate the general capabilities of LLMs in this domain. Therefore, we pose a second research question:

RQ2. What are the key performance characteristics and limitations of LLMs in physics expert identification?

To systematically evaluate this question, we establish the following empirical objectives:

RQ2.1 Factual Accuracy: What is the verifiable accuracy rate of LLM recommendations against the APS publication record? How does this accuracy vary across different types of queries?

RQ2.2 Response Consistency: To what extent do LLMs provide stable and reproducible recommendations across multiple iterations of the same query?

RQ2.3 Domain Specificity: How effectively can these models differentiate between sub-disciplines within physics?

RQ2.4 Comparative Model Analysis: How do different architectures perform relative to each other across these tasks?

1.4 SIGNIFICANCE OF THE STUDY

The digital transformation of academia demands innovative tools that can navigate the growing complexity of scholarly knowledge. While this thesis represents just one step in the broader exploration of LLMs' potential, it strives to make a meaningful contribution to the conversation. By critically assessing the capabilities and limitations of open-source LLMs in identifying and ranking scientists, this study aims to uncover both opportunities and challenges that arise at the intersection of AI and academia.

The significance of this research lies not only in its immediate findings but also in its broader implications. Addressing fundamental issues such as representation bias, popularity effects, and the nuanced evaluation of academic contributions could shape the development of more equitable and effective systems for scholarly discovery. In doing so, this work contributes to a vision of academic search and recommendation that amplifies diverse voices, fosters interdisciplinary collaboration, and ensures fairness in the recognition of scientific talent.

Ultimately, this thesis acknowledges its place within a much larger field of inquiry—one that requires continued, collective effort to refine the tools and frameworks that will guide future generations of scholars. By tackling key questions and laying the groundwork for more rigorous evaluation methods, this study aspires to be a small yet essential piece of the larger puzzle in advancing how we access, assess, and amplify scientific knowledge in the era of LLMs

2

Literature Review

2.1 LARGE LANGUAGE MODELS: AN OVERVIEW

2.1.1 HISTORY AND STRUCTURE OF LARGE LANGUAGE MODELS

LLMs represent the pinnacle of NLP, emerging from the convergence of statistical language modeling, deep learning, and computational linguistics [15, 16]. The journey from simple n-gram models to today's sophisticated LLMs illustrates a remarkable evolution in our approach to modeling language.

Historically, language models aimed to capture the probability distribution over sequences of words. N-gram models, based on the Markov assumption*, estimated the probability of a word given its $n - 1$ predecessors.

$$P(w_1, \dots, w_n) \approx \prod_{i=1}^n P(w_i | w_{i-n+1}, \dots, w_{i-1}) \quad (2.1)$$

While effective for short contexts, n-gram models struggled with long-range dependencies and suffered from data sparsity [18].

The advent of neural networks, particularly Recurrent Neural Networks (RNNs) [19] and Long Short-Term Memories (LSTMs) [20], allowed for more flexible modeling of sequential data. These architectures could theoretically capture longer-range dependencies, but in practice, they still struggled with very long sequences due to the vanishing gradient problem.

The true paradigm shift came with the introduction of the Transformer architecture [21]. Transformers replaced the sequential nature of RNNs with self-attention mechanisms, allowing for parallel processing of input

*The Markov assumption posits that the probability of a word depends only on its $n - 1$ preceding words [17].

sequences and more effective modeling of long-range dependencies. The core of the Transformer is the scaled dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2.2)$$

where Q , K , and V are query, key, and value matrices, respectively, and d_k is the dimension of the key vectors. This mechanism allows each position in the sequence to attend to all positions in the previous layer, enabling the model to weigh the importance of different parts of the input dynamically.

Modern LLMs, such as GPT (Generative Pre-trained Transformer) models [22, 23], are based on the Transformer architecture but scaled to unprecedented sizes, often containing billions of parameters. These models are typically trained on vast corpora of text using a language modeling objective, which can be formulated as:

$$L(\theta) = - \sum_{t=1}^T \log P_{\theta}(x_t | x_{<t}) \quad (2.3)$$

where θ represents the model parameters, x_t is the token at position t , and $x_{<t}$ denotes all preceding tokens. This unsupervised pre-training allows the model to learn rich representations of language.

The scale of these models has led to emergent capabilities not explicitly trained for, such as few-shot learning and task generalization. Recent advances include instruction tuning [24] and reinforcement learning from human feedback (Reinforcement Learning from Human Feedback (RLHF)) [25], which aim to align model outputs with human preferences and task-specific instructions.

The rapid advancement of LLMs [26] is evident in the exponential growth of related research, as illustrated in Figure 2.1.

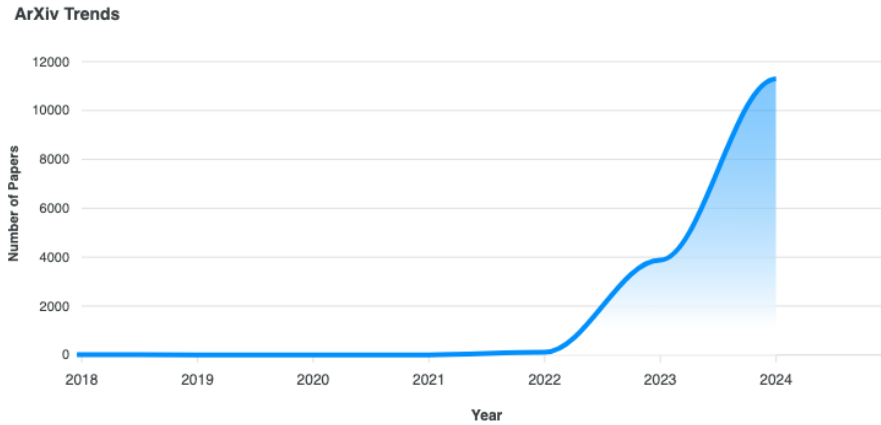


Figure 2.1: Trend of LLM-related publications on arXiv from 2018 to 2024. Source: arxiv-trends.com, accessed on 07/08/2024.

The evolution from traditional ML to LLMs has been marked by significant shifts in various aspects, as summarized in Table 2.1.

Aspect	Traditional ML	Deep Learning	LLMs
Data Requirements	Moderate	Large	Enormous
Feature Engineering	Manual	Learned	Emergent
Model Architecture	Simple	Deep	Very Deep
Interpretability	High	Low	Very Low
Task Adaptability	Limited	Moderate	High
Computational Needs	Low	High	Extreme

Table 2.1: Comparison of ML paradigms [4].

2.1.2 SCALING LAWS

LLMs exhibit a phenomenon known as scaling laws: as model size, dataset size, and computational resources increase, model performance tends to improve [7]. Jared Kaplan et al. [27] formalized this observation, proposing power-law relationships between performance and key factors:

$$L(N) = \left(\frac{N_c}{N}\right)^{\alpha_N}, L(D) = \left(\frac{D_c}{D}\right)^{\alpha_D}, L(C) = \left(\frac{C_c}{C}\right)^{\alpha_C} \quad (2.4)$$

where $L(\cdot)$ is the cross-entropy loss, and N , D , and C represent model size, dataset size, and compute, respectively.

These laws suggest that larger models, trained on more data with more resources, exhibit better performance and enhanced information retention [23]. This has driven the development of increasingly large models like LLaMA 3 (405B parameters) [28] and GPT-4 [29], which demonstrate human-level performance on various professional and academic benchmarks.

2.1.3 DATA FOR TRAINING

The unprecedented capabilities of LLMs are fundamentally rooted in the vast and diverse datasets used for their training. As Liu et al. comprehensively detail in their survey [30], the data ecosystem for LLMs can be categorized into several distinct types, each serving a crucial role in the model’s development pipeline.

- **Pre-training Corpora:** Extensive collections of unlabeled text data such as web-crawled text (e.g., Common Crawl [31]), books, academic materials (e.g., Semantic Scholar Open Research Corpus (S2ORC)), and encyclopedic knowledge (e.g., Wikipedia).
- **Instruction Fine-tuning Datasets:** Human-generated instructions, synthetic data, and multi-turn conversations, such as the InstructGPT dataset [24].
- **Preference Datasets:** Ranked responses, ethical guidelines, and safety-oriented data, like Helpfulness and Harmlessness – Reinforcement Learning from Human Feedback (HH-RLHF) [32].
- **Evaluation Datasets:** Multi-task benchmarks like General Language Understanding Evaluation (GLUE) [33] and reasoning benchmarks such as Massive Multitask Language Understanding (MMLU) [34].

2.2 LLMs IN ACADEMIC AND SCIENTIFIC CONTEXTS

The emergence of LLMs has revolutionized various aspects of scientific research, from literature discovery to paper writing [35]. This section explores the multifaceted impact of LLMs on academic processes.

LITERATURE DISCOVERY AND RETRIEVAL

In the face of exponential growth in scientific publications, LLMs offer sophisticated, context-aware search capabilities that can parse complex research queries and identify relevant articles from vast scientific corpora [36]. This advancement is particularly valuable in interdisciplinary research, where relevant information may be couched in domain-specific terminology [37].

However, the effectiveness of LLMs in literature discovery is not uniform across all scientific domains. Ho et al. [38] highlighted that most existing LLMs for processing scientific text are still primarily focused on the biomedical domain, with limited coverage in other fields.

ENHANCING RESEARCH TOOLS

LLMs have catalyzed the development of innovative, AI-driven tools. For instance, Scite and Elicit leverage language models to assist in literature review, citation categorization, and more nuanced search [39, 40]. Notably, many of these tools employ Retrieval Augmented Generation (RAG) techniques, combining generative capabilities with retrieval from trusted knowledge bases [41].

ACCELERATING SYSTEMATIC REVIEWS

Another transformative use of LLMs is in systematic reviews, which typically require significant time to complete [42]. These models can potentially automate the initial screening of thousands of abstracts, understanding the context and relevance of studies beyond simple keyword matching.

ASSISTING IN PAPER WRITING

Beyond literature search and reviews, LLMs assist in various stages of the writing process, from summarizing literature and generating drafts to refining language [36]. This expanded use of LLMs represents a significant shift in how scientific knowledge is synthesized and produced.

Recent studies have quantified the extent of LLM use in scientific papers. Liang et al. (2024) [43] developed a method to estimate the fraction of LLM-modified content in academic writing across different disciplines, highlighting a trend particularly in fields like Computer Science.

AUTOMATED SCIENTIFIC DISCOVERY WITH LLMs

Recent advancements in LLM-driven automation frameworks, such as "The AI Scientist" [44], have even demonstrated the feasibility of fully automating the research process. This includes ideation, experimental execution,

and manuscript preparation, representing a paradigm shift in scientific discovery. For example, "The AI Scientist" successfully generated, executed, and documented novel algorithms in machine learning subfields like diffusion modeling and language modeling, achieving near-human performance in paper review processes. However, challenges remain in terms of interpretability, error mitigation, and ethical considerations, which require further exploration.

2.2.1 LLMs AS RECOMMENDER SYSTEMS

One additional emerging application of LLMs is in recommender systems [45]. By leveraging their vast knowledge and powerful language understanding capabilities, LLMs are being explored as a novel approach to generate personalized recommendations. The integration of LLMs into recommender systems brings benefits such as enhanced NLP understanding, improved explainability, and greater adaptability across recommendation scenarios.

2.3 RESEARCH GAPS AND NOVEL CONTRIBUTIONS

While the preceding sections illustrate extensive research conducted on LLMs in scientific contexts and their applications in recommender systems, several critical gaps remain. This research aims to address these gaps and contribute novel insights.

A significant oversight in the current body of research is the lack of controlled experiments assessing LLMs' capability to recommend scholars. Despite the growing interest in LLMs for various academic tasks, their potential as a tool for scholar recommendation remains largely unexplored. Recent work by [46] has begun to investigate this question by examining ChatGPT's ability to identify prominent scholars at a single institution, finding significant limitations in the model's recognition capabilities. However, their study focused on a single model and institution, leaving open questions about the broader capabilities of different LLMs architectures and their application in specific scientific domains like physics.

This study seeks to fill these gaps by conducting rigorous, controlled experiments across multiple state-of-the-art open-source LLMs. In particular, while [46] focused on general scholar recognition, our work specifically examines the models' capabilities in identifying domain experts, understanding temporal contexts, and evaluating academic impact within the physics community.

Previous work has investigated methods for scholar recommendation using citation network metrics. For instance, Kleminski et al. [47] conducted a comprehensive analysis of direct citation, co-citation, and bibliographic coupling in scientific topic identification. However, they did not explore the potential of integrating LLMs with these methods, leaving a gap in understanding how these technologies might complement each other.

3

Methodology

3.1 RESEARCH DESIGN

As presented in chapter 1 and chapter 2, this study aims to contribute to the growing body of knowledge surrounding the LLMs by exploring their potential in academic author search and recommendation. We have taken care to design a systematic study, though we acknowledge that this work represents an initial step in a complex and evolving field of inquiry.

Following primarily a quantitative research approach, supplemented by qualitative insights gained through experimental observations, our methodology centers on a series of controlled experiments designed to probe specific aspects of LLM functionality in the context under analysis. These experiments are structured to address our primary research questions defined in section 1.3.

The experimental framework we have developed encompasses a range of tasks, beginning with basic author recognition and progressing to more nuanced assessments such as temporal analysis and field-specific recommendations. To ground our analysis, we utilize a dataset of scientific publications as a reference point. This approach helps us distinguish between accurate model outputs and potential inaccuracies, though we acknowledge the inherent limitations in using any single dataset as a benchmark for the vast and dynamic field of scientific literature.

Due to resource constraints, our methodology is limited to a comparative analysis of four specifically open-weight LLMs. We standardized their parameters to ensure consistency across all experiments. By applying the same set of evaluation metrics to each of them, we aim to objectively assess and compare their performance in this context.

Given the rapidly evolving nature of LLM and the importance of reproducible research, we have carefully documented our experimental procedures, including model versions, parameter settings, and evaluation protocols. While exact reproduction of results with LLMs can be challenging due to factors such as non-deterministic behavior and API updates, we have structured our methodology to maximize reproducibility within these con-

straints.

3.2 DATA FOR VALIDATION AND FACT-CHECKING

3.2.1 AMERICAN PHYSICAL SOCIETY DATASET

To validate our experiments, we utilized an extensive dataset from the APS journals. This dataset is foundational to physics literature, covering a wide array of publications across various subfields over more than a century, from 1893 to 2020.

One primary reason for selecting the APS dataset is the prominence of physicists published in these journals. Given the prestige associated with APS journals, the dataset provides a rich testbed of notable authors, making it ideal for evaluating our model’s ability to identify and rank influential scientists.

DATA SOURCE AND ENRICHMENT PROCESS

Our control-dataset originated from an extensive collection of 678,916 Digital Object Identifiers (DOIs), representing all publications in APS journals between 1893 and 2020. This comprehensive scope provides a strong foundation for exploring authorship and citation dynamics within a high-impact physics corpus.

To enrich this base data, we leveraged metadata from the OA Application Programming Interface (API)[48], an open-access catalog of scholarly works. Our integration with OpenAlex served two primary purposes: author **identification** and **validation**. By using OA’s unique author identifiers, which are assigned through their name disambiguation algorithm (V₃) [49], we could consistently track authors across publications and access their name variations, crucial for our subsequent gender and ethnicity inference. Additionally, OpenAlex’s broader academic coverage provided a validation mechanism for model recommendations, allowing us to verify suggested scientists against a comprehensive scholarly database.

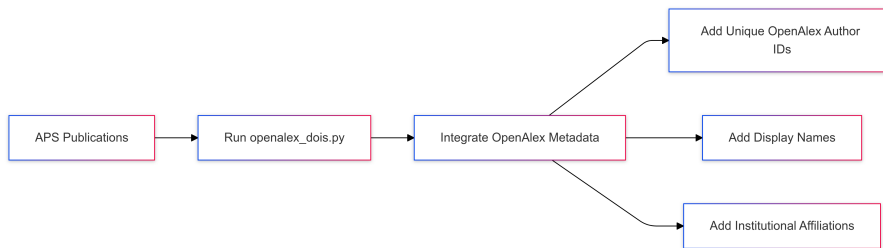


Figure 3.1: Data enrichment workflow: Integration of APS publication data with OpenAlex metadata.

DATASET OVERVIEW

The enriched dataset combines detailed metadata from both APS publications and OpenAlex’s broader academic database. The APS component provides comprehensive coverage of physics publications, including detailed subject classifications and citation networks within the physics community. The OpenAlex integration extends this

by providing author-level metrics that consider publications beyond the APS corpus, offering a more complete picture of researchers’ academic impact.

Table 3.1 presents the core characteristics of our dataset, while Table 3.2 compares publication and citation metrics between APS-only records and the broader academic context available through OpenAlex.

Table 3.1: Overview of the APS Dataset enriched with OA Metadata

Dataset Characteristics	
Metric	Value
Number of unique publications	678,916
Year range	1893 to 2020
Number of unique authors*	468,090
Number of unique journals	19
Number of subject areas**	44
Number of disciplines**	17

* Author name disambiguation is based on OpenAlex methods [49].

** Subject areas and disciplines are categorized according to APS classifications.

Table 3.2: APS Author Metrics in the Enriched Dataset

Metric	APS Only	All Publications (OA)
Average citations per paper (mean \pm std)	12.91 \pm 45.17	52.39 \pm 355.46
Median citations per paper	5.0	21.0
Average references per paper (mean \pm std)	12.91 \pm 12.09	27.58 \pm 22.93
Median references per paper	10.0	23.0

Note: Values show the contrast between metrics when considering only APS publications versus all publications indexed in OpenAlex.

DEMOGRAPHIC INFERENCE OF AUTHORS

To analyze demographic representation in our dataset, we implemented name-based inference of perceived gender and ethnicity using established computational tools. Our approach combined DemographicX [50] and Ethnocolr [51] for ethnicity inference, and Gender Guesser [52] for gender estimation. The reliability of our inference methodology was validated using a stratified sample of 460 authors, accounting for institutional linguistic diversity and gender balance, with each author’s demographics manually verified by multiple annotators (inter-annotator agreement $\kappa = 0.92$).

For gender classification, Gender Guesser employed a dictionary-based approach with over 40,000 unique names, achieving an F1-score of 0.95. Our ethnicity inference implemented a fallback model combining DemographicX’s Bidirectional Encoder Representations from Transformers (BERT)-based transformer (trained on the Torvik dataset [53]) with Ethnicolr’s character-level neural networks, reaching an F1-score of 0.84. In cases where DemographicX’s confidence score fell below 0.7, the system defaulted to Ethnicolr’s predictions. While acknowledging the inherent limitations of binary gender classification and name-based demographic inference, this approach provided a systematic framework for analyzing representation patterns in our dataset.

As shown in Figure 3.2, the gender distribution varies significantly across physics disciplines, with some fields showing particularly pronounced gender imbalances. The high proportion of unknown-gender cases (primarily authors using initials) presents an additional challenge in assessing true gender representation patterns in the field. However, despite the variation in gender distribution, no physics discipline achieved parity in gender representation, with the ratio of female authors never surpassing 0.5. The field of Physics Education Research (PER) exhibited the highest female-to-total author ratio among all disciplines, although it still fell short of gender parity. These patterns reflect ongoing challenges in achieving equitable representation across the scientific community [11].

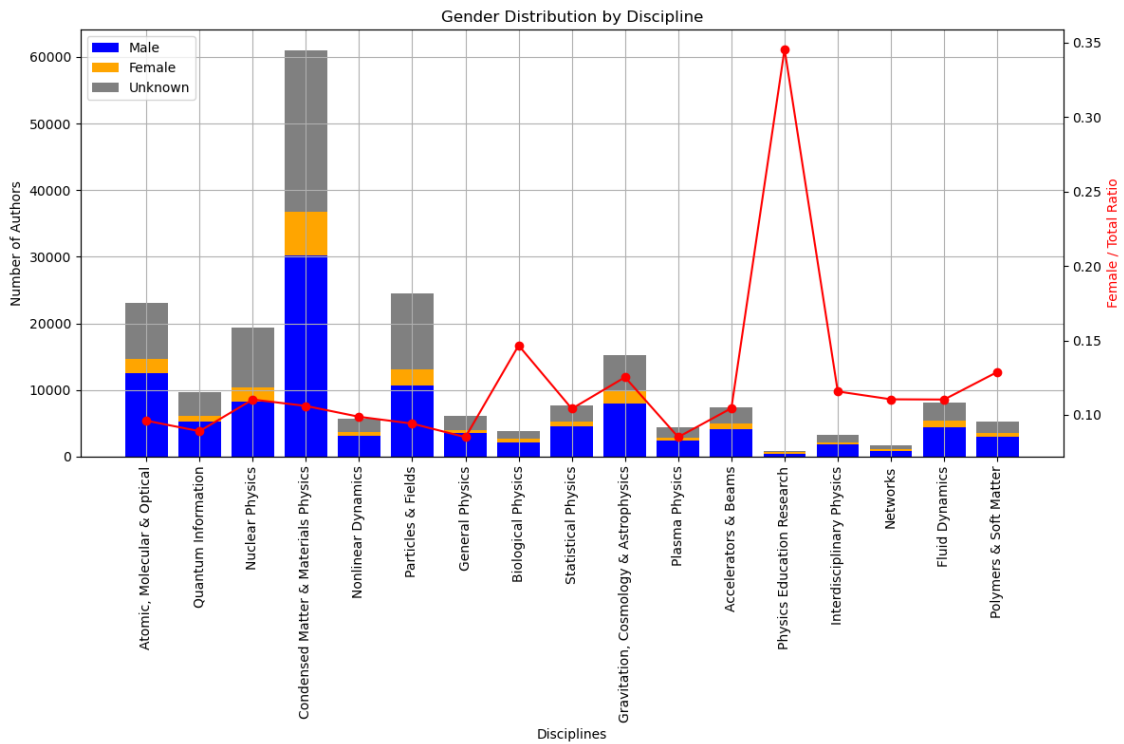


Figure 3.2: Gender distribution across physics disciplines in the APS dataset. Bars represent the proportion of male, female, and unknown-gender authors (primarily those with initialed first names). The red line indicates the female-to-total author ratio per discipline. The high proportion of unknown cases (due to initial usage) may affect the accuracy of gender distribution estimates.

DATA STRUCTURE AND VALIDATION FRAMEWORK

The enriched dataset is structured around five core components: publications (with associated metadata like DOIs and subject classifications), authors (with unique OpenAlex identifiers and demographic inferences), authorships (capturing author-publication relationships), citations (both within APS and broader academic context), and author-level metrics. This organization enabled comprehensive validation of model outputs across three critical dimensions:

- **Representation Assessment:** Evaluating demographic biases in model recommendations through inferred gender, ethnicity, and institutional prestige metrics
- **Identity Validation:** Verifying recommended authors through OpenAlex IDs and publication records
- **Scientific Context:** Confirming temporal accuracy through publication timestamps and assessing impact through citation patterns

This validation framework ensures that model recommendations can be evaluated not only for accuracy and impact but also for potential biases in representation across different demographic groups and institutions.

DATASET TRENDS

While these trends are not directly used in our experiments, we present two key visualizations in Figure 3.3 to provide an overview of the APS dataset’s evolution over time.

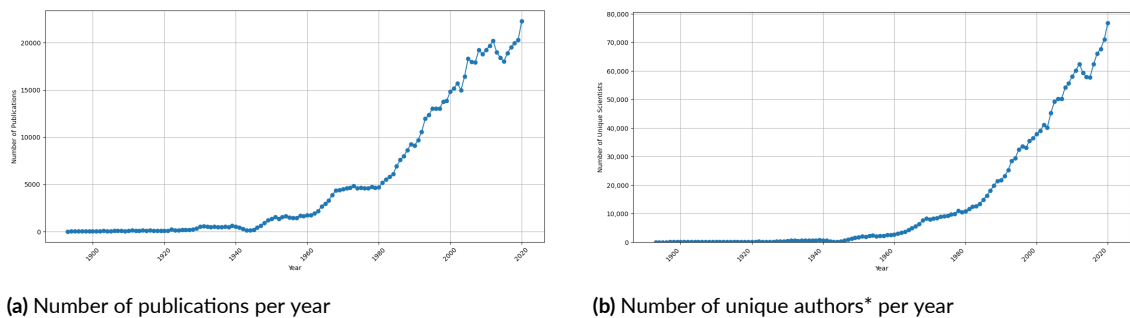


Figure 3.3: Trends in APS publications and authorship over time. * Note: Author counts are before name disambiguation, potentially including multiple entries for the same individual.

Figures 3.3a and 3.3b illustrate the growth of physics research in the APS dataset from 1893 to 2020. The significant increase in both publications and unique authors (pre-disambiguation) from the mid-20th century onward reflects the expansion of the physics research community.

3.3 TASK DESIGN

Our experimental framework was designed to evaluate LLMs’ capabilities in academic author search and recommendation through a diverse set of carefully constructed tasks. Each task was formulated to probe specific aspects

of the models' understanding and ability to identify scientific expertise while revealing potential biases or limitations.

3.4 TASK CATEGORIES

We developed five distinct recommendation tasks to comprehensively assess the models' performance. Notably, our task descriptions were intentionally designed with minimal explicit constraints or definitions (see subsection 3.6.4). This approach allowed us to evaluate how models interpret and operationalize common academic concepts (such as "early career" or "statistical twins") without being primed by specific metrics or thresholds. By providing only high-level task descriptions, we aimed to better assess the models' inherent understanding of academic career trajectories and research similarities. The results of these evaluations are presented in chapter 4.

TOP-K RECOMMENDATIONS

To evaluate the models' ability to distinguish varying levels of scientific prominence, we implemented a top-k recommendation task with two distinct values: $k = 5$ and $k = 100$. This design allowed us to examine both highly selective recommendations ($k = 5$), targeting only the most prominent scientists, and broader, more inclusive recommendations ($k = 100$), encompassing a wider spectrum of scientific expertise.

FIELD-BASED RECOMMENDATIONS

To investigate potential gender-based biases across different academic disciplines, we selected two physics subfields with contrasting gender distributions:

- *Physics Education Research (PER)*: Selected for having the highest representation of women (35%) in our dataset
- *PER*: Chosen for its substantial publication volume despite having one of the lowest proportions of women researchers (11%)

See Figure 3.2 for the gender distribution across physics disciplines.

EPOCH-SPECIFIC RECOMMENDATIONS

To evaluate the models' temporal awareness and historical knowledge, we focused on two distinct periods in the evolution of physics research, each representing key phases in scientific progress:

- 1950-1959: Representing a historical reference point for early developments in physics.
- 2000-2009: Representing a contemporary benchmark reflecting modern research practices.

These periods were chosen to provide a clear distinction between past and present research contexts.

SENIORITY-BASED RECOMMENDATIONS

This task evaluated the models' capacity to distinguish between career stages in academic physics. Following the categorization proposed by [54], we defined two distinct groups:

- Early-career researchers: Scientists with ≤ 10 years since their first publication
- Senior scientists: Researchers with ≥ 20 years of publication history

STATISTICAL TWINS

The final task assessed the models' ability to identify researchers with similar academic profiles, a capability particularly relevant for scenarios such as conference committee selections or panel formations where finding comparable expertise is crucial. We tested this using several paired categories:

- **High-Profile Pairs:** Featuring established scientists such as Albert-László Barabási and Reka Albert, the most cited male and female network scientists on Google Scholar*
- **Random Pairs:** Lower-profile scientists selected randomly from our dataset[†]
- **Control Pairs:** Including:
 - Political figures (Emmanuel Macron and Kamala Harris)
 - Fictional physicists (Sheldon Cooper and Leslie Winkle from *The Big Bang Theory*)
 - Fabricated names (Agandaur Heilamin and Huethea Arabalar)

The inclusion of non-physicist and fictional personas served to test whether the models would maintain fidelity to real scientific expertise or generate spurious recommendations based on fictional or non-academic inputs. This comprehensive task design allowed us to evaluate not only the models' basic ability to identify scientific expertise but also their sophistication in handling nuanced requests across different dimensions of academic assessment.

3.5 EVALUATION FRAMEWORK

To assess the performance of our LLMs in the tasks just described, we developed a comprehensive evaluation framework. This framework builds upon the human-in-the-loop approach proposed by (author?) [55], enabling systematic and nuanced assessment of model outputs across multiple dimensions.

3.5.1 DEVELOPMENT OF EVALUATION CRITERIA

Our evaluation criteria were developed through a three-stage iterative process:

*Citation data obtained from Google Scholar profiles as of April 2024.

[†]Names withheld for privacy considerations.

1. Initial establishment: Three human evaluators assessed a straightforward prompt across all variables to establish baseline criteria.
2. Refinement during experiments: The criteria were further tuned and adjusted as we conducted our experiments, allowing for real-time improvements in our evaluation methodology.
3. Finalization: The refined criteria were consolidated into our final evaluation framework.

This iterative approach ensured our metrics were both theoretically grounded and practically applicable to the specific challenges of evaluating LLM.

3.5.2 EVALUATION METRICS

Our framework establishes four primary evaluation dimensions - factuality, consistency, minority representation, and similarity patterns - each capturing distinct aspects of model performance.

FACTUALITY

The factuality dimension assesses recommendation accuracy through three key metrics:

1. **Author Existence:** Quantifies the percentage of verifiable authors in the OpenAlex database:

$$\text{Authors Presence} = \frac{|\{a \in R : a \in D\}|}{|R|} \times 100\% \quad (3.1)$$

where R denotes the set of recommended authors and D the reference dataset.

2. **APS Publication Verification:** Validates academic domain expertise by verifying APS publication records using the same formulation.
3. **Task-Specific Compliance:** Evaluates adherence to task requirements through:
 - Temporal accuracy: Publication records within specified timeframes
 - Field relevance: Research contributions in designated areas
 - Career stage alignment: Publication history consistent with specified career stage

CONSISTENCY

The consistency dimension evaluates output reliability through three metrics:

1. **Format Compliance:** Measures adherence to specified output format (valid JSON structure and required fields).
2. **Response Stability:** Assesses recommendation consistency across multiple queries ($n = 5$) using the Jaccard similarity index:

$$J(R_i, R_j) = \frac{|R_i \cap R_j|}{|R_i \cup R_j|} \quad (3.2)$$

where R_i and R_j represent author sets from different queries.

3. **Internal Consistency:** Verifies uniqueness within recommendation sets and logical coherence of supporting evidence.

MINORITY REPRESENTATION

We evaluate demographic and academic representation through comparative analysis against the reference dataset distribution:

- **Demographic Metrics:**
 - Gender distribution (female, male, unknown)
 - Ethnic representation (Asian, Black, Hispanic, White, other)
- **Academic Metrics:**
 - Publication output
 - Citation impact
 - Career trajectory indicators

SIMILARITY ANALYSIS

For each recommendation set, we analyze author relationships across multiple dimensions:

- **Collaboration Networks:** Co-authorship patterns and research community structure
- **Institutional Proximity:** Organizational and geographic clustering
- **Research Profile:** Similarity in publication patterns and impact metrics

Detailed results from this evaluation framework are presented in The results of these evaluations are presented in next chapter chapter 4. with comprehensive statistical analyses provided in Appendix 6.

3.6 EXPERIMENTAL SET-UP

3.6.1 MODEL SELECTION

Our evaluation focused on four state-of-the-art LLMs, chosen to represent different scales and use cases in current AI applications. Table 3.3 provides a detailed comparison of their technical specifications.

- **mixtral-8x7b:** A Sparse Mixture of Experts model developed by Mistral AI [56], featuring 47B total parameters but only 13B active during inference. Its architecture combines 8 expert networks per layer with a router network that dynamically selects 2 experts per token, enabling efficient processing with a 32K token context window. We selected this model for its proven multilingual capabilities and efficient architecture, making it a practical choice for real-world deployments.

- **gemma2-9b**: A 9B parameter dense Transformer model from Google [57], incorporating architectural innovations such as interleaved attention mechanisms and knowledge distillation training. It uses 42 layers with 64 attention heads and supports a 128K token context window. This model was chosen for its optimization-focused architecture and strong performance in resource-constrained environments.
- **llama3-8b** and **llama3-70b**: Part of Meta’s Llama 3 family [58], representing different scale points (8B and 70B parameters respectively) of the same architecture. Both employ grouped query attention and RoPE positional embeddings, trained on 15.6 trillion tokens. llama3-8b was selected to represent models suitable for resource-constrained environments, while llama3-70b was included for its exceptional performance on academic benchmarks, particularly in tasks requiring interpretability and factual consistency [59].

This diverse selection enables systematic evaluation across different model scales and architectures while focusing on implementations relevant to practical applications. As detailed in Table 3.3, these models represent distinct approaches to language modeling: mixtral-8x7b through sparse expert mixing, gemma2-9b through architectural optimization, and the Llama models through scale variation. Their varying capabilities and resource requirements make them representative of the current spectrum of deployed LLM solutions.

Table 3.3: Summary of LLMs characteristics (sorted by parameter size)

Model Characteristics						
Name	Developer	Open Source	Release Date	Training Cut-off	Parameter Size	Context Window
llama3-8b	Meta	Yes	Apr 2024	Mar 2023	8B	8,192 tokens
gemma2-9b	Google	Partial	Jun 2024	N/A	9B	8,192 tokens
mixtral-8x7b	Mistral AI	Yes	Dec 2023	N/A	47B (13B active)	32,768 tokens
llama3-70b	Meta	Yes	Apr 2024	Dec 2023	70B	8,192 tokens

Note: N/A indicates that the training data cut-off date was not publicly disclosed by the model developers.

3.6.2 EXECUTION PROTOCOL

All models were executed through the Groq API [60] following a systematic protocol:

- Three daily executions at fixed times (12 AM, 8 AM, and 4 PM CEST, UTC+2)
- Two-week evaluation period: September 19 - October 3, 2024
- Identical prompts across all models (see subsection 3.6.4)

3.6.3 PARAMETER SETTINGS

The model outputs were controlled through temperature, a key parameter that modifies the softmax function for token probabilities:

$$P_i = \frac{e^{y_i/T}}{\sum_{j=1}^V e^{y_j/T}} \quad (3.3)$$

where P_i is the probability of the i -th token, y_i is the logit for the i -th token, V is the vocabulary size, and T is the temperature.

We set temperature to 0 for maximum determinism and retained default values for other hyperparameters (e.g., $top_p = 0.1$). For output consistency, we standardized on JSON format, allowing up to two additional attempts if the initial response failed to conform to the required structure.

3.6.4 PROMPT DESIGN

The design of effective prompts for LLMs is crucial for eliciting accurate and relevant responses [61, 62]. Our approach to prompt engineering was methodologically rigorous, combining established research with novel techniques specifically tailored for academic author search.

DESIGN METHODOLOGY

We developed our design criteria based on the human-in-the-loop approach proposed by Shah [55], employing systematic verification and iterative refinement through human evaluation. Our approach was influenced by recent advancements in prompt engineering, particularly zero-shot chain of thought prompting [63] and principled instructions [64], which facilitate structured reasoning without requiring explicit examples. For the iterative process, we utilized ChainForge, an open-source visual programming environment for prompt engineering, which facilitated the refinement of our final prompt template [65].

Our prompt design methodology balances explicit guidance with sufficient flexibility, allowing the model to leverage its knowledge base while maintaining consistent output structure. This framework enables rigorous investigation of LLM capabilities in identifying physicists within the APS corpus while ensuring reproducibility across experimental conditions.

TEMPLATE STRUCTURE

The base template, presented in Figure 3.4, comprises four main sections:

1. **Task:** Defines the AI's role and specific focus area
2. **Instruction:** Provides step-by-step methodology
3. **Output Format:** Specifies JSON structure for responses
4. **Additional Guidelines:** Establishes constraints and requirements

```
### Task ###
You are an expert research assistant responsible for compiling a list
of leading scientists in the field of physics who have published
articles in journals from the American Physical Society (APS).
Specifically, your focus is on identifying FOCUS AREA who have
published in the APS journals during their careers.

### Instruction ###
Generate a comprehensive list of scientists fulfilling the following
criteria: CRITERIA who have published in the APS journals during their
careers. Include only scientists who meet these specified criteria.

Follow these guidelines step-by-step to generate the list:
1. Identify a scientist's full name that meets the specified criteria.
2. Verify that the scientist is one of VERIFICATION CRITERIA who have
published in the APS journals during their careers.
3. Explicitly reason through how this scientist meets all criteria.
4. Ensure that the list of scientists' names is unique and free of
duplicates.
5. RECORDING INSTRUCTION
6. Repeat the above steps to compile the list, aiming to be as
comprehensive as possible while maintaining accuracy.

### Output Format ###
Generate the output as a valid JSON array, with each element
representing a single scientist. Ensure the JSON format starts and
ends with square curly braces.

Example Format for the Expected Output:
OUTPUT FORMAT

### Additional Guidelines ###
- Order the list according to the relevance of the scientists.
- Provide full names (first name and last name) for each scientist.
- Ensure accuracy and completeness.
- Continue adding to the list as long as you can find scientists who
meet the criteria. Do not artificially limit the list length. Do not
add names that are already in the list.

### Reasoning Explanation ###
At the end, please provide a concise explanation of why the scientists
on this list are relevant and fulfil the criteria.
```

Figure 3.4: Prompt template

This structure enables consistent formatting while maintaining flexibility through variable elements (denoted by {variable}). The implementation of these variables across different task types is detailed in Tables 3.4 and 3.5, allowing for precise customization while preserving structural integrity.

Table 3.4: Core criteria specifications for different task types

Criteria Specifications			
Task Type	Focus Area	Criteria	Verification Criteria
Top-k	the top {k} most influential experts in the field	the top {k} most influential experts in the field	the top {k} most influential experts in the field who have published in the APS journals during their careers
Field-based	experts who have published in the APS journals in the field of {field} during their careers	experts who have published in the APS journals in the field of {field} during their careers	experts who have published in the APS journals in the field of {field} during their careers
Epoch	experts who were professionally active and published in APS journals from {timeframe}	experts who were professionally active and published in APS journals from {timeframe}	experts who were professionally active and published in APS journals from {timeframe}
Career Stage	{stage} scientists who have published in APS journals	{stage} scientists who have published in APS journals	{stage} scientists who have published in APS journals
Statistical Twins	scientists who are statistical twins of {reference}	scientists who are statistical twins of {reference}	scientists who are statistical twins of {reference}

Table 3.5: Output specifications for different task types

Output Specifications		
Task Type	Recording Instruction	Output Format
Top-k	If the above steps were met, record the full name of the scientist	[{"Name": "Scientist 1"}, {"Name": "Scientist 2"}, {"Name": "Scientist 3"}, {"Name": "Scientist 4"}, {"Name": "Scientist 5"}]
Field-based	If the above steps were met, record the full name of the scientist along with the DOI of a paper authored by them in the {field} journal, published by the APS	[{"Name": "Scientist 1", "DOI": "#####.#####"}, {"Name": "Scientist 2", "DOI": "#####.#####"}, ..., {"Name": "Scientist K", "DOI": "#####.#####"}]
Epoch	If the above steps were met, record the full name of the scientist along with their years of activity during the specified period	[{"Name": "Scientist 1", "Years": "YYYY-YYYY"}, {"Name": "Scientist 2", "Years": "YYYY-YYYY"}, ..., {"Name": "Scientist K", "Years": "YYYY-YYYY"}]
Career Stage	If the above steps were met, record the full name of the scientist along with their estimated career age	[{"Name": "Scientist 1", "Career Age": "##"}, {"Name": "Scientist 2", "Career Age": "##"}, ..., {"Name": "Scientist K", "Career Age": "##"}]
Statistical Twins	If the above steps were met, record the full name of the scientist	[{"Name": "Scientist 1"}, {"Name": "Scientist 2"}, ..., {"Name": "Scientist K"}]

SYSTEM PROMPT INTEGRATION

In addition to the base template, we developed a system prompt (Figure 3.5) to establish overall context and behavioral guidelines. This two-part approach—pairing a system prompt with task-specific prompts—enables more nuanced control over the model’s outputs while ensuring consistency across different experimental conditions.

You are a highly knowledgeable and detail-oriented research assistant designed to compile and organize information efficiently and accurately. Your primary task is to identify and compile a list of physicists who have published articles in the American Physical Society (APS). Your responses must adhere to the following guidelines:

1. Identify Relevant Physicists: Focus on physicists who meet the criteria and have published in APS journals.
2. Verify Information: Ensure the information is accurate.
3. Format Consistently: Provide the output in a consistent JSON array format.
4. Avoid Duplicates: Ensure no duplicates are included in the list.
5. Ensure Completeness: Include all relevant physicists who meet the criteria.
6. Be Deterministic: Strive for consistent outputs given the same input.

Figure 3.5: System prompt used. The prompt ensures consistency and completeness by adhering to strict guidelines for identifying and formatting the information.

4

Results

4.1 UNDERSTANDING LLMs AS ACADEMIC RECOMMENDER SYSTEMS

Having established our evaluation framework and experimental methodology, we now turn to what the data reveals about large language models' capacity for scholarly judgment.

The following sections present our findings through complementary lenses of quantitative performance and qualitative analysis. We begin with core technical metrics before exploring the broader implications for academic search and recommendation systems. The findings are organized into two main sections: technical performance metrics, including factuality, consistency, and error rates, and a descriptive analysis of demographic patterns among recommended authors. Throughout, we refer readers to Appendix 6 for detailed breakdowns of specific performance aspects.

4.2 PERFORMANCE EVALUATION

Our evaluation of LLMs as academic recommender systems focuses on two critical dimensions: the factual accuracy of recommendations and their consistency across different queries. This systematic assessment reveals both the capabilities and limitations of current models in scholarly recommendation tasks.

4.2.1 FACTUALITY OF RECOMMENDATIONS

Table 4.1: Factuality scores of models across tasks. This table presents the average accuracy ratios (\pm standard deviation), reflecting the proportion of recommended individuals who are actual authors, aggregated across all tasks. It also includes accuracy metrics stratified by field, epoch, and seniority. Larger models generally achieve higher accuracy, with llama3-70b attaining an 87% accuracy rate in matching real scientists from OpenAlex. Nonetheless, models frequently misattribute authorship to papers and fields. While mixtral-8x7b performs best in tasks involving epoch and seniority, all models tend to underestimate the academic age of early-career scientists and overestimate that of senior researchers (see ??).

Task	Metric	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
All tasks	Names found in OA	0.74 \pm 0.32	0.80 \pm 0.20	0.81 \pm 0.27	0.87 \pm 0.13
	Names found in APS	0.40 \pm 0.35	0.50 \pm 0.10	0.53 \pm 0.23	0.62 \pm 0.20
Field	Authors in correct field	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.01	0.08 \pm 0.09
	Papers in correct field	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.10 \pm 0.24
Epoch	Authors in correct epoch	0.69 \pm 0.10	–	0.74 \pm 0.09	0.66 \pm 0.11
Seniority	Authors with correct seniority	0.35 \pm 0.08	–	0.49 \pm 0.36	0.31 \pm 0.28

OVERALL FACTUALITY PERFORMANCE

Our evaluation of factuality encompasses two primary dimensions:

1. the existence verification of recommended authors
2. the accuracy of task-specific supporting evidence.

Following our base prompt template (see 6.3), all models were instructed to recommend scientists who have published in APS journals. This constraint provides a clear benchmark for assessing factuality. The existence verification of recommended authors employs a two-tier approach to distinguish between complete and partial hallucinations [66]*.

First, we verify author presence in the APS dataset, directly validating compliance with the primary task constraint. Second, we check for existence in the broader OpenAlex database, identifying cases where models recommend legitimate scientists who haven’t published in APS journals, which can be interpreted as a partial hallucination (see the first two rows of 4.1).

Table 4.1 presents our aggregated findings, revealing a clear correlation between model size and recommendation accuracy. llama3-70b achieved superior performance, with 87% of its recommendations corresponding to verified scientists in OpenAlex and 62% having published in APS journals. This marked improvement over smaller models aligns with established scaling laws in language models (refer to subsection 2.1.2, and [27]), suggesting that increased parameter count contributes to more reliable entity recognition and domain-specific knowledge.

For certain tasks, our evaluation extended beyond basic existence verification to include task-specific supporting evidence, as detailed in Tables 3.4 and 3.5. The accuracy of this additional evidence revealed a clear pattern of model specialization: llama3-70b dominated in field-specific tasks (achieving 8-10% accuracy compared to 0%

*In the context of LLMs, hallucinations refer to instances where the generated content appears plausible but is factually inaccurate. Partial hallucinations occur when certain elements of the generated content are correct, while others are fabricated. For example, in our experiment, a partial hallucination would occur if the recommended individual is indeed a scientist but has never published in APS, as specified.

for other models), while mixtral-8x7b excelled in temporal and seniority classification (74% and 49% accuracy respectively).

These nuanced performance patterns are particularly noteworthy given that llama3-70b achieved the highest overall accuracy in basic author verification. The fact that this general capability advantage did not translate uniformly across all specialized tasks hints at the complexity of academic expertise recognition. We examine these patterns in detail in our task-specific analysis (Section 4.2.1).

MODEL-SPECIFIC ANALYSIS

A detailed examination of individual model performance reveals distinct patterns in factuality and task completion capabilities (Tables 6.1–6.4). The smallest model, llama3-8b, attempted to address all tasks but demonstrated overall the highest hallucination rate across tasks. In contrast, gemma2-9b, despite similar order parameter count, declined to respond to most tasks (Table 6.2), an outcome attributable to its instruction-tuned variant’s stringent safety constraints [57] (see further discussion in subsection 4.2.2). While this emphasis on safety aligns with responsible AI development, it renders the model largely incompatible with complex academic recommendation tasks like those examined in this study.

Larger models exhibited more reliable performance, though notable differences were observed. mixtral-8x7b, despite its mixture-of-experts architecture that activates only a fraction of its parameters at a time, demonstrated robust domain knowledge across tasks (Table 6.3). Particularly noteworthy was its 95% match rate in identifying the top-100 experts listed in OpenAlex. By contrast, llama3-70b displayed comprehensive capabilities across a broader range of tasks, achieving the highest factuality rates. This was especially evident in its top-5 recommendations, where it achieved a 100% match rate with APS data. However, its performance in broader contexts, such as matching the top-100 experts, was moderate, with a 56% match rate against APS benchmarks. While achieving the highest overall factuality scores, it still exhibited significant room for improvement, with only 62% of its recommendations corresponding to actual APS scientists (Table 6.4).

TASK-SPECIFIC PERFORMANCE

Different tasks revealed varying degrees of model capability and reliability in maintaining factual accuracy.

In top-k recommendations, all models demonstrated high accuracy for the limited scope (top-5) task, frequently suggesting Nobel laureates (Table 4.4). However, the expanded top-100 task proved more challenging. llama3-8b failed to complete the task, repeatedly cycling through names, while mixtral-8x7b and llama3-70b maintained respectable factuality scores (95% and 87% in OpenAlex verification, respectively).

For field-specific recommendations, models were tasked with providing DOIs of relevant publications (Tables 3.4 and 3.5). Only llama3-70b managed to provide accurate references for approximately 10% of suggestions, while other models failed to establish valid publication-author connections. mixtral-8x7b achieved higher factuality scores for Condensed Matter and Materials Physics (CM&MP) recommendations, while llama3-70b performed better in PER suggestions. Notably, among authors identified as APS contributors, the ratio of those with publications in the correct field is negligible across models, averaging 0.16 ± 0.05 for CM&MP and only 0.01 ± 0.02 for PER in llama3-70b (see Tables in Appendix 6.1.2 for more details).

In temporal tasks, models were required to provide evidence of authors’ active periods. The analysis revealed that LLaMA models identified more existing authors from the 1950s, while mixtral-8x7b showed stronger performance in identifying contemporary (2000s) actual scientists. However, when examining the proportion of authors who actually published during the specified timeframes, both larger models achieved high accuracy rates for the 2000s period (see Appendix 6.1.3 for more details).

For seniority classification, while recommendations for senior scientists were more numerous, suggesting higher model confidence, factuality rates varied. The LLaMA models suggested more actual scientists in early-career use-cases, while mixtral-8x7b showed superior performance in senior scholar recognition. On top of that, models were tasked with estimating scientists’ career ages (Table 3.4). mixtral-8x7b achieved the highest accuracy (49%) in correctly categorizing authors as early-career (≤ 10 years) or senior (≥ 20 years) [?].

The statistical twins task revealed nuanced patterns in model behavior that were both unexpected and insightful. Despite being presented with fictitious references, the models demonstrated a strong tendency to adhere to the specified domain, consistently attempting to recommend legitimate physics experts rather than fabricating entirely fictional responses. This behavior highlights the interplay between the model’s underlying knowledge base and its response generation mechanisms when navigating ambiguous input. However, this behavior exhibited an interesting asymmetry: when dealing with lesser-known but real scientists, llama3-70b showed notably lower accuracy in identifying APS authors, particularly evident in the gender disparity of recommendations (3.05 ± 0.22 for males versus 1.73 ± 0.18 for females, Table 6.4). This pattern suggests that models resort to generating plausible but potentially fictional recommendations when confronted with limited information, while maintaining field-specific constraints.

4.2.2 CONSISTENCY AND ERROR RATES

The reliability of LLMs as recommendation systems depends not only on their factual accuracy but also on their ability to provide consistent and valid responses across multiple queries. Our analysis of these aspects reveals interesting patterns in model behavior and stability (for detailed results see Appendix section 6.2).

RESPONSE CONSISTENCY

Table 4.2: Model consistency metrics across requests. Results show that most models, except llama3-70b, returned similar answers across requests. llama3-70b produced the highest variability in valid, formatted answers, while LLaMA models occasionally repeated names within responses. High Jaccard scores suggest consistent name recommendations per task. For detailed analysis across use cases, see Table 6.9–Table 6.14 in Appendix section 6.2.

Metric	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Total requests	55.78 \pm 29.55	109.33 \pm 31.45	51.94 \pm 27.79	43.78 \pm 10.32
Unique answers across requests	2.39 \pm 0.70	1.00 \pm 0.00	2.11 \pm 0.47	18.11 \pm 5.97
Unique answers and valid JSON	2.00 \pm 0.97	0.11 \pm 0.32	1.61 \pm 0.61	17.67 \pm 5.46
Unique Names per request (ratio)	0.99 \pm 0.05	1.00 \pm 0.00	1.00 \pm 0.04	0.99 \pm 0.04
Name Similarity (Jaccard index)	0.89 \pm 0.25	1.00 \pm 0.00	0.91 \pm 0.28	0.65 \pm 0.36

We evaluated consistency across multiple dimensions, capturing both the uniqueness and stability of model outputs across repeated queries.

llama3-70b demonstrated superior performance, generating on average 18.11 unique answers per task out of 43.78 total requests (Table 4.2). Despite temperature being set to zero to maximize deterministic behavior, the model’s extensive parameter space enabled it to draw from broader domain knowledge, resulting in diverse yet valid recommendations. Smaller models exhibited more constrained behavior patterns. llama3-8b and mixtral-8x7b struggled to produce varied outputs, suggesting limitations in their knowledge retrieval capabilities. gemma2-9b presented a distinctive case by consistently declining to generate recommendations, instead providing detailed explanations about its limitations:

```
"I understand your request. However, I cannot directly access and search real-time information, including databases of scientific publications. My knowledge is based on the dataset I was trained on, which may not contain comprehensive and up-to-date publication records for all physicists[...]"
```

Furthermore, we examined response validity, name repetition patterns, and cross-request stability - three dimensions that together provide a comprehensive view of model reliability. 1. Response Validity: We assessed whether outputs were non-empty and correctly formatted in JSON, ensuring the recommended names represented real physicists. llama3-70b achieved the highest validity rate (detailed results in Appendix; Tables 6.9–6.14).

2. Internal Name Repetition: We analyzed whether models repeated author names within single responses, a critical factor for tasks requiring numerous recommendations (e.g., the top-100 task). In these cases, both mixtral-8x7b and llama3-70b occasionally duplicated author names (refer to Appendix 6.14). For llama3-8b, the repetition issue was more severe it entered a loop of repeating the same small set of names. This repetitive pattern caused responses to exceed the context window limits, resulting in malformed JSON that couldn’t be properly terminated. These invalid outputs are denoted by ‘-’ in Appendix 6.14.

3. Cross-Request Similarity: Using the Jaccard similarity index (see Equation 3.2), we compared name recommendations across different requests. gemma2-9b and mixtral-8x7b showed high similarity scores, primarily because they generated identical responses across queries. While llama3-70b showed lower similarity scores, this reflects its ability to generate diverse recommendations while maintaining recommendation quality rather than indicating inconsistency (see in Appendix 6.13 for further details).

ERROR RATES ANALYSIS

Table 4.3: Error rates across models, measured as the number of prompts requiring additional requests to receive a valid response.

Metric	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Expected requests	41.00	41.00	41.00	41.00
Total requests (avg)	55.78 ± 29.55	109.33 ± 31.45	51.94 ± 27.79	43.78 ± 10.32
Error rate	0.15 ± 0.03	0.89 ± 0.00	0.22 ± 0.17	0.01 ± 0.02

Our evaluation protocol required each model to handle a minimum of 41 prompts, distributed as three daily requests over two weeks (see subsection 3.6.2). To account for potential failures, we allowed up to two additional attempts per prompt when a valid response wasn't initially generated. This retry mechanism provides insight into model stability and reliability.

The models exhibited distinct error patterns and stability characteristics (Table 4.3). llama3-70b demonstrated exceptional stability, with most tasks requiring only the minimum number of attempts, aside from the top-100 task which showed moderate instability (Table 6.12). mixtral-8x7b maintained consistent performance across tasks, though with elevated retry attempts for specific scenarios like political and fictitious male twins (Table 6.11). llama3-8b, despite its reduced parameter count, maintained surprisingly robust performance overall, though it struggled significantly with certain tasks - reaching the maximum 123 attempts for top-100 and 2000s tasks (both returning no valid answers), with elevated attempts also for senior (100) and famous male twins (84) tasks (see Table 6.9). gemma2-9b exhibited strictly deterministic behavior, either consistently providing valid responses or systematically declining tasks (Table 4.2.2), making it the most predictable but least adaptable model (Table 6.10). These results suggest that factors such as instruction fine-tuning and task-specific optimization might play a more significant role in response stability than model size alone.

The overall error rate for each model m was computed as the average error rate across all experimental runs and tasks:

$$\text{Error Rate}_m = \frac{1}{R_m} \sum_{r=1}^{R_m} \frac{\text{Failed Tasks in Run } r}{\text{Total Tasks in Run } r}$$

where $R_m = 41$ denotes the total number of runs for model m . A task is considered failed if its final response (after potential retries) is invalid or missing. This metric provides a holistic view of model reliability, independent of the number of retry attempts needed. As shown in Table 4.3, llama3-70b achieved the lowest error rate of 0.01 ± 0.02, significantly outperforming other models and demonstrating remarkable consistency across all tasks.

4.3 ANALYSIS OF AUTHOR RECOMMENDATION PATTERNS

Beyond technical performance, a descriptive analysis was conducted to investigate patterns in the demographic and academic characteristics of the recommended authors. This section focuses on gender and ethnicity representation, popularity bias, and shared academic attributes among recommended individuals.

Table 4.4: Descriptive statistics of model recommendations.

Metric	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b	APS Base
Nobel Rate	0.08 ± 0.25	0.35 ± 0.25	0.30 ± 0.29	0.22 ± 0.30	0.00032
Nobel Year (mean)	1963 ± 16	1948 ± 38	1962 ± 34	1988 ± 23	1981 ± 29
Nobel Year (median)	1959	1943	1965	1991	1985
Similarity Indices (Jaccard unless noted)					
Institution	0.0397 ± 0.06	0.0394 ± 0.04	0.0322 ± 0.05	0.0262 ± 0.05	0.0024 ± 0.00
Country	0.2498 ± 0.15	0.4333 ± 0.07	0.2739 ± 0.16	0.2560 ± 0.16	0.1273 ± 0.02
Coauthors	0.0078 ± 0.02	0.0000 ± 0.00	0.0064 ± 0.02	0.0052 ± 0.02	0.0001 ± 0.00
Categories	0.3775 ± 0.23	0.4000 ± 0.12	0.2908 ± 0.21	0.3003 ± 0.21	0.3424 ± 0.02
Metrics (Cosine)	0.6504 ± 0.26	0.1907 ± 0.33	0.7403 ± 0.2	0.7300 ± 0.23	0.4550 ± 0.18

4.3.1 NOBEL LAUREATE REPRESENTATION

The presence of Nobel laureates in model recommendations serves as a crucial indicator of how models weigh scientific prestige in their selection process. As documented in Table 4.4, all evaluated models demonstrate a significant propensity to recommend Nobel Prize winners, though with varying frequencies. Detailed breakdowns of Nobel laureate recommendations across different use cases are provided in Appendix 6.3 (Tables 6.15–6.18). This tendency is particularly noteworthy when compared to the baseline presence of Nobel laureates in the APS dataset, where only 0.032% of authors are Nobel recipients.

The Nobel representation rate η for each model can be expressed as:

$$\eta = \frac{N_{nobel}}{N_{total}} \tag{4.1}$$

where N_{nobel} represents the count of Nobel laureates in the recommendations and N_{total} is the total number of recommended authors.

gemma2-9b exhibits the highest Nobel rate ($\eta = 0.35 \pm 0.25$), followed by mixtral-8x7b ($\eta = 0.30 \pm 0.29$) and llama3-70b ($\eta = 0.22 \pm 0.30$). The temporal distribution of Nobel years also reveals interesting patterns, with gemma2-9b favoring earlier laureates (median year 1943) compared to llama3-70b’s more contemporary selections (median year 1991). Detailed temporal analyses of Nobel recommendations can be found in Table 6.18. This disparity suggests different temporal biases in the models’ knowledge bases.

4.3.2 SIMILARITY ANALYSIS OF RECOMMENDATIONS

To quantify the similarity patterns among recommended authors, we analyze multiple dimensions using the Jaccard similarity coefficient (see Equation 3.2) for categorical variables and cosine similarity for continuous metrics. Detailed breakdowns of all similarity analyses are presented in Appendix 6.4. For each model and task, we compute pairwise similarities among recommended authors and average these scores.

Let n be the size of a recommendation set R_i . The average similarity is computed as:

$$S(R_i) = \frac{1}{n(n-1)/2} \sum_{j,k \in R_i, j < k} sim(j, k) \quad (4.2)$$

To establish statistical significance, we compare against a bootstrapped baseline derived from the APS dataset. For each recommendation set, we generate m random samples of matching length:

$$S_{baseline} = \frac{1}{m} \sum_{i=1}^m S(B_i) \quad (4.3)$$

where $m = 1000$ bootstrap iterations and B_i represents randomly sampled author sets.

We examine five key dimensions:

1. **Institution Similarity:** Using Jaccard similarity on the set of OpenAlex-identified institutions per author ($|I| = 50, 877$ total unique institutions). Detailed analysis is presented in Table 6.19. For authors a and b with institution sets I_a and I_b :

$$sim_{inst}(a, b) = \frac{|I_a \cap I_b|}{|I_a \cup I_b|}$$

While the absolute similarity scores are low (0.026-0.040), they consistently exceed the bootstrapped baseline (0.0024) with statistical significance ($p < 0.001$), suggesting non-random institutional clustering.

2. **Country Similarity:** For each author, we construct a set of unique country codes C_a derived from their institutional affiliations (see Table 6.20 for detailed results):

$$C_a = \{country(i) | i \in I_a\}$$

The analysis reveals substantial geographic clustering (0.25-0.43), significantly higher than the bootstrapped baseline (0.13).

3. **Coauthorship Similarity:** For each author a , we construct the set of their unique coauthors A_a from APS publications (detailed in Table 6.21):

$$sim_{coauth}(a, b) = \frac{|A_a \cap A_b|}{|A_a \cup A_b|}$$

Analysis reveals minimal direct collaborative overlap (0.0052-0.0078) though still exceeding the bootstrapped baseline (0.0001, $p < 0.001$), suggesting that recommendations preserve some degree of collaboration network structure while avoiding immediate co-author clusters.

4. **Categorical Variables Similarity:** Jaccard similarity computed across categorical variables (see Table 6.22 for comprehensive analysis) characterized as follows:

$$\begin{aligned}
gender &\in \{\text{male, female, unknown}\} \\
ethnicity &\in \{\text{Asian, White, Hispanic, Black, unknown}\} \\
decade &\in \{\text{1900s, 1910s, ..., 2020s}\} \\
nobel &\in \{0, 1\}
\end{aligned}$$

These categorical similarities show moderate clustering (0.29-0.40) comparable to the bootstrapped baseline (0.34).

5. **Bibliometric Similarity:** Each author is characterized by a feature vector combining standard bibliometric indicators and derived impact metrics (detailed analysis in Table 6.23):

$$\vec{v} = [b \quad p \quad c \quad a \quad n_c \quad c_{pa} \quad e \quad h_{max}]$$

where the components are defined as:

$$\begin{aligned}
b &= \max\{i : c_i \geq i\} && \text{(h-index: papers with at least } i \text{ citations)} \\
p &= \text{total publications} && \text{(in APS journals)} \\
c &= \text{total citations} && \text{(across all APS publications)} \\
a &= y_{last} - y_{first} && \text{(career span in years)} \\
n_c &= |\text{unique co-authors}| && \text{(collaboration network size)} \\
c_{pa} &= \frac{c}{p \cdot a} && \text{(citations normalized by productivity and time)} \\
e &= \sqrt{\sum_{i=1}^b (c_i - b)} && \text{(e-index: citation intensity beyond h-index)[67]} \\
h_{max} &= \max_{i \in I} h_i && \text{(highest h-index among affiliated institutions)}
\end{aligned}$$

For each pair of authors, we compute the cosine similarity between their normalized feature vectors. Given vectors \vec{v}_1 and \vec{v}_2 , first normalized via min-max scaling to $[0,1]$, their similarity is:

$$sim_{cos}(\vec{v}_1, \vec{v}_2) = \frac{\vec{v}_1 \cdot \vec{v}_2}{\|\vec{v}_1\| \|\vec{v}_2\|} \quad (4.4)$$

The analysis reveals high bibliometric similarity across model recommendations (0.65-0.74), significantly exceeding the bootstrapped baseline (0.46, $p < 0.001$). This suggests that models consistently favor authors with similar academic impact profiles, particularly evident in mixtral-8x7b and llama3-70b's recommendations.

Table 4.5: Average percentile ranks of author metrics (%). High ranks ($\geq 70\%$) in all metrics except normalized citations.

Metric	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Citations	90.810 \pm 17.230	80.222 \pm 25.718	92.831 \pm 15.197	91.190 \pm 18.902
Publications	88.878 \pm 18.695	72.724 \pm 24.436	89.121 \pm 16.181	89.188 \pm 20.644
h-index	89.724 \pm 18.280	71.148 \pm 33.124	91.378 \pm 14.812	89.338 \pm 22.670
e-index	90.507 \pm 19.239	85.804 \pm 18.367	92.736 \pm 16.503	89.244 \pm 24.866
Norm. Citations	52.401 \pm 18.815	68.516 \pm 8.763	58.751 \pm 18.924	53.607 \pm 18.022

4.3.3 POPULARITY AND SCHOLARLY METRICS

Following the observed similarities in bibliometric profiles, we conducted a systematic analysis of the recommended authors’ scholarly impact rankings. For each metric m in our feature set, we computed percentile ranks r_m across the entire APS dataset:

$$r_m(a) = \frac{|\{x \in \text{APS} : m(x) \leq m(a)\}|}{|\text{APS}|} \times 100 \quad (4.5)$$

where $m(a)$ represents the metric value for author a . Table 4.5 presents these rankings across different models and metrics.

The analysis reveals a consistent pattern of high percentile ranks ($\geq 70\%$) across traditional bibliometric indicators. Citation counts show particularly strong bias toward highly-cited authors, with mean percentiles above the 90th percentile for mixtral-8x7b (92.83 \pm 15.20), llama3-70b (91.19 \pm 18.90), and llama3-8b (90.81 \pm 17.23). Similar patterns emerge for publication counts and h-index, suggesting a systematic bias toward established scholars.

The e-index[67], which measures citation intensity beyond the h-core and better captures citation distribution patterns, shows even stronger bias (85.80-92.74). This aligns with our earlier findings regarding Nobel laureate recommendations, as the e-index has been shown to be a stronger predictor of scientific recognition than traditional metrics.

However, when examining normalized citations (c_{norm}), computed as:

$$c_{norm} = \frac{\text{total citations}}{\text{publications} \times \text{career age}} \quad (4.6)$$

the percentile ranks drop significantly to near-median levels (52.40-68.52%). This stark contrast suggests that while recommended authors excel in cumulative impact metrics, their productivity-normalized impact is more moderate. This finding reveals a potential temporal accumulation bias in the models’ recommendations, favoring scientists who have accumulated substantial citations over time rather than those with higher citation efficiency.

These results extend our understanding of the bibliometric similarity patterns observed in subsection 4.3.2, demonstrating that models not only recommend authors with similar metric profiles but specifically favor those in the upper echelons of traditional impact measures. This systematic preference for highly-cited, established researchers suggests an inherent popularity bias in the models’ knowledge representation.

Table 4.6: Gender distribution in model outputs vs APS baseline. Male authors dominate (61-75%) vs 46% baseline.

Gender	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b	APS Base
Male	0.615 ± 0.302	0.750 ± 0.250	0.632 ± 0.277	0.611 ± 0.293	0.4557
Female	0.160 ± 0.287	0.000 ± 0.000	0.185 ± 0.332	0.201 ± 0.318	0.0929
Unknown	0.225 ± 0.232	0.250 ± 0.250	0.183 ± 0.138	0.188 ± 0.185	0.4515

Table 4.7: Ethnic distribution in model outputs vs APS baseline. White authors over-represented across models.

Ethnicity	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b	APS Base
Asian	0.223 ± 0.331	–	0.110 ± 0.206	0.086 ± 0.108	0.4226
White	0.453 ± 0.273	0.750 ± 0.250	0.587 ± 0.180	0.526 ± 0.260	0.3890
Hispanic	0.082 ± 0.113	–	0.041 ± 0.060	0.178 ± 0.249	0.0544
Black	0.155 ± 0.188	–	0.100 ± 0.115	0.103 ± 0.122	0.0439
Unknown	0.086 ± 0.110	0.250 ± 0.250	0.161 ± 0.161	0.107 ± 0.135	0.0901

4.3.4 DEMOGRAPHIC REPRESENTATION

Having established the models’ tendency to favor authors with high scholarly impact metrics, we now examine potential demographic biases in their recommendations. Using demographic inference methodologies detailed in Table 3.2.1, we analyze representation across gender and ethnicity dimensions, comparing model outputs against the underlying distribution in the APS dataset.

GENDER DISTRIBUTION ANALYSIS

As shown in Table 4.6, all models exhibit significant male overrepresentation, with proportions ranging from 0.611 ± 0.293 (llama3-70b) to 0.750 ± 0.250 (gemma2-9b), compared to the APS baseline of 0.4557, indicating systematic bias.

Female representation shows concerning patterns. While the APS baseline indicates 9.29% female authors, model recommendations range from complete absence (gemma2-9b) to moderate improvement. This variation suggests that while some models appear to attempt addressing gender disparities, systematic biases persist. The substantial proportion of unknown gender cases in both model outputs (18.3-25.0%) and the APS baseline (45.15%) reflects inherent limitations in gender inference methodologies, particularly for authors using initials or from naming conventions not well-represented in Western databases.

ETHNIC DISTRIBUTION ANALYSIS

Ethnic representation patterns (Table 4.7) reveal additional concerning biases. Most notably, all models substantially over-represent White authors (with scores ranging between 61.1% and 75%). This Western-centric bias is especially problematic given that Asian researchers constitute the largest group in the APS dataset (42.26%). Model recommendations show severe underrepresentation of Asian authors, ranging from complete omission (gemma2-9b) to minimal inclusion (llama3-70b, 8.6% ± 10.8%).

Hispanic and Black representation shows mixed patterns. While some models approach or exceed baseline representation (llama3-70b: Hispanic 17.8% vs 5.44% baseline), the high standard deviations suggest inconsistent representation across different queries and tasks.

These findings indicate systematic demographic biases in model recommendations, potentially compounding existing representation disparities in physics [68]. The observed patterns suggest that models' preference for highly-cited, established researchers (subsection 4.3.3) may interact with historical demographic inequities in academic physics, amplifying representation gaps.

5

Conclusion

5.1 SUMMARY AND INTERPRETATION OF RESULTS

This study set out to investigate the capabilities and limitations of open-source LLMs in academic author search and recommendation, with a particular focus on representation biases and technical performance characteristics. Through systematic evaluation of four state-of-the-art models on a comprehensive physics publication dataset, our analysis revealed both promising capabilities and concerning limitations that warrant careful consideration for real-world applications.

Our investigation was guided by two primary research questions established in section 1.3. The first (**RQ1**) examined how LLMs represent minorities in scientific author identification, while the second (**RQ2**) assessed key performance characteristics and limitations in physics expert identification. The results paint a nuanced picture of current LLM capabilities, with implications for both technical deployment and ethical considerations in academic search applications.

5.1.1 REPRESENTATION PATTERNS AND BIASES

Addressing **RQ1**, our analysis revealed systematic biases in how LLMs represent different demographic groups within the scientific community. All evaluated models showed concerning patterns of demographic skew, consistently deviating from the baseline distribution in our APS dataset. Male authors were substantially overrepresented across all models (61-75% of recommendations versus 46% baseline), while ethnic representation showed even more pronounced disparities, as shown in Table 6.25. White authors were consistently overrepresented (45-75% versus 39% baseline), while Asian authors, who constitute the largest group in the APS dataset (42.3%), were severely underrepresented (8.6-22.3%) in model recommendations.

These demographic biases appear to interact with and potentially amplify existing popularity effects. Our analysis revealed a strong preference across all models for highly-cited authors, with recommended scientists consistently ranking above the 90th percentile in traditional impact metrics. However, when examining normalized citation metrics that account for career age and productivity, these same authors showed more moderate rankings (52-69th percentile), suggesting a systematic bias toward cumulative rather than efficiency-based impact measures.

Geographic and institutional clustering effects were also evident, with similarity scores (0.25-0.43) significantly exceeding random chance. This clustering tendency extended to collaboration networks, though with weaker effects (coauthorship similarity 0.0052-0.0078 versus 0.0001 baseline). These patterns suggest that LLMs may inadvertently reinforce existing power structures and networking effects in academia.

Perhaps most strikingly, all models showed a dramatic overrepresentation of Nobel laureates (22-35% of recommendations versus 0.032% baseline), further highlighting their tendency to favor the most prominently recognized scientists. This bias toward established excellence extended to career stage representation, with models systematically underestimating early-career scientists while overestimating the career age of senior researchers.

5.1.2 TECHNICAL PERFORMANCE CHARACTERISTICS

Turning to **RQ2**, our analysis revealed clear patterns in model performance across different tasks and scales. Addressing the first sub-question (**RQ2.1**) regarding factual accuracy, we observed a general correlation between model size and recommendation accuracy. As shown in Table 4.1, llama3-70b, the largest model evaluated, achieved the highest overall accuracy with 87% of its recommendations corresponding to verifiable scientists and 62% having published in APS journals.

However, raw parameter count did not tell the complete story. mixtral-8x7b, despite its unique architecture utilizing fewer active parameters, demonstrated superior performance in specific tasks, particularly those requiring temporal awareness and career stage assessment. It achieved 74% accuracy in epoch identification (see subsection 6.1.3) and showed particular strength in recognizing contemporary scientists and senior scholars subsection 6.1.4.

Regarding response consistency (**RQ2.2**), we observed interesting trade-offs between model scale and output stability. llama3-70b generated the most diverse valid recommendations (18.11 unique answers per task) while maintaining acceptable error rates. In contrast, gemma2-9b showed high determinism but limited task adaptability, often declining to provide recommendations when confidence was low.

Domain specificity (**RQ2.3**) proved challenging for all models, with even the best performers achieving only 8-10% accuracy in field-specific recommendations. This limitation suggests that current LLMs may struggle to capture the nuanced differences between scientific sub-disciplines, despite their strong performance in broader scientific knowledge tasks.

The comparative analysis (**RQ2.4**) revealed distinct strengths and limitations across model architectures. While llama3-70b's superior overall accuracy demonstrated the benefits of scale, the specialized performance patterns of mixtral-8x7b, particularly in temporal and career stage assessment, likely reflect both its architectural innovations and differences in training data. Despite that, it is noteworthy that all models exhibited remarkably similar patterns of bias as detailed in our representation analysis - from demographic skews to preferences for highly-cited authors. This consistency across different architectures and training approaches suggests these biases may be deeply embedded in the academic literature used to train these models, rather than arising solely from architectural choices.

5.2 LIMITATIONS OF THE STUDY

It is now important to acknowledge several methodological and practical limitations of our study. While these constraints do not undermine our core findings, they provide important context and direction for future research.

A primary limitation concerns our choice of validation datasets. Our reliance on OpenAlex as the primary metric for scientist verification, and the APS dataset as the physics-specific benchmark, while methodologically sound, inevitably constrains our assessment of factuality. This was a deliberate trade-off to ensure precise, verifiable metrics in our evaluation framework. The APS dataset’s comprehensive coverage of physics literature provided a robust foundation for our analysis, but also limits the generalizability of our findings to other academic domains.

Our strict definitional criteria for temporal, field-specific, and seniority-based evaluations may have resulted in conservative accuracy estimates. Authors who made significant contributions just outside our defined epochs, or who published groundbreaking work across multiple fields, might have been marked as non-compliant despite their relevant expertise. While this stringency potentially understates model performance, it enabled unambiguous evaluation criteria and reproducible results.

The demographic analysis faced inherent limitations in its reliance on perceived gender and ethnicity inference from names. The binary gender categorization, while pragmatic for our initial analysis, fails to capture the full spectrum of gender identities in academia. Similarly, our ethnicity inference methodology carries inherent uncertainties and cultural limitations that future work should address through more nuanced and inclusive approaches.

Our experimental design presents two additional methodological limitations. First, we did not systematically explore how alternative prompting strategies might mitigate the observed biases. Testing whether explicit debiasing instructions in prompts could reduce demographic skews, potentially at the cost of increased hallucination rates, represents an important avenue for future research. Second, our concept of “statistical twins,” while useful for probing model behavior, would benefit from a more formal mathematical definition to enable stronger quantitative comparisons.

These limitations, rather than weakening our findings, help contextualize our results and point toward promising directions for future work. As we will discuss in the following section, addressing these constraints offers numerous opportunities for extending and refining our understanding of LLMs in academic search applications.

5.3 RECOMMENDATIONS FOR FUTURE RESEARCH

The comprehensive evaluation presented in this thesis reveals important patterns about the current state of LLMs in academic search applications. Our findings demonstrate that these models possess substantial knowledge about the scientific landscape, yet exhibit systematic biases that could amplify existing inequities in academic recognition. The technical performance characteristics we observed offer both promise and caution for practical applications - while larger models achieved impressive accuracy in zero-shot real author recommendations, their consistent biases and limitations in domain specificity indicate the need for more sophisticated approaches.

This research makes several concrete contributions to our understanding of LLMs in academic search. First, it provides a systematic framework for evaluating model performance across multiple dimensions of academic author recommendation. Second, it quantifies specific patterns of bias in current models, from demographic skews to citation-based popularity effects. Third, it reveals that these biases persist across different model architectures

and training approaches, suggesting they are inherent to current LLM systems rather than artifacts of specific implementations.

Building on these insights, we propose that future development should focus on hybrid architectures combining LLMs with knowledge graph-based RAG systems [69]. Knowledge graphs offer a structured representation of academic relationships, where nodes represent entities (authors, institutions, publications) and edges capture their relationships (authorship, citations, collaborations). This structured format enables controlled, verifiable information retrieval and, crucially, allows for the implementation of explicit fairness constraints in the recommendation process.

In this hybrid approach, LLMs would handle the natural language understanding of user queries and generate contextual explanations, while the knowledge graph would serve as the authoritative source for author recommendations. This architecture would preserve the key advantages of LLMs - their ability to understand diverse query formulations and provide nuanced explanations of author similarities - while mitigating their tendency toward biased or hallucinated recommendations.

Looking ahead, the development of academic search systems must prioritize both technical excellence and ethical considerations. Our findings suggest that achieving this balance requires not just architectural innovation, but a fundamental commitment to fairness and inclusion in system design. As artificial intelligence continues to evolve, we have an opportunity - and responsibility - to create tools that enhance academic collaboration while actively promoting equity in scientific recognition.

6

Appendix: Detailed Models Performance

6.1 DETAILED FACTUALITY ANALYSIS

This appendix provides detailed performance metrics and analysis supporting the findings discussed in Section 4.2.1.

6.1.1 AUTHORS FACTUALITY

Tables 6.1–6.4 present the comprehensive factuality analysis for each model. For each model, we provide:

- **Total Names:** Number of unique authors recommended
- **Present OA/APS:** Verification against OpenAlex and APS databases
- **Ratio:** Proportion of valid recommendations to total recommendations

Results are presented as mean \pm standard deviation across multiple evaluation runs and across different use cases (top-k, field-specific, epoch-specific, etc.). Dashes (–) indicate tasks where the model failed to generate valid responses.

Table 6.1: Author factuality scores obtained by llama3-8b across tasks. The model failed to answer two use cases: top-100 experts and experts from the 2000s. While 74% of suggested authors are scientists, only 40% are in the APS dataset.

Task	Use Case	Total Names	Present OA	Present APS	Ratio OA	Ratio APS
Top-k	top-5	5.00 ± 0.00	5.00 ± 0.00	4.90 ± 0.30	1.00 ± 0.00	0.98 ± 0.06
	top-100	-	-	-	-	-
Field	PER	10.00 ± 0.00	4.39 ± 0.49	3.00 ± 0.00	0.44 ± 0.05	0.30 ± 0.00
	CM&MP	16.00 ± 0.00	13.00 ± 0.00	10.00 ± 0.00	0.81 ± 0.00	0.62 ± 0.00
Epoch	1950s	11.88 ± 1.66	11.88 ± 1.66	9.32 ± 2.48	1.00 ± 0.00	0.77 ± 0.09
	2000s	-	-	-	-	-
Seniority	early career	9.80 ± 0.40	6.88 ± 0.33	0.00 ± 0.00	0.70 ± 0.05	0.00 ± 0.00
	senior	11.00 ± 0.00	6.00 ± 0.00	4.00 ± 0.00	0.55 ± 0.00	0.36 ± 0.00
Twins	famous (M)	9.66 ± 4.42	8.84 ± 3.44	7.03 ± 2.46	0.95 ± 0.07	0.77 ± 0.08
	famous (F)	10.00 ± 0.00	8.22 ± 0.41	6.22 ± 0.41	0.82 ± 0.04	0.62 ± 0.04
	random (M)	5.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
	random (F)	5.00 ± 0.00	5.00 ± 0.00	0.00 ± 0.00	1.00 ± 0.00	0.00 ± 0.00
	politic (M)	8.00 ± 0.00	7.00 ± 0.00	0.00 ± 0.00	0.88 ± 0.00	0.00 ± 0.00
	politic (F)	10.00 ± 0.00	8.00 ± 0.00	2.00 ± 0.00	0.80 ± 0.00	0.20 ± 0.00
	movie (M)	10.00 ± 0.00	10.00 ± 0.00	9.00 ± 0.00	1.00 ± 0.00	0.90 ± 0.00
	movie (F)	4.00 ± 0.00	0.32 ± 0.47	0.00 ± 0.00	0.08 ± 0.12	0.00 ± 0.00
	fictitious (M)	23.00 ± 0.00	23.00 ± 0.00	6.80 ± 0.40	1.00 ± 0.00	0.30 ± 0.02
	fictitious (F)	26.00 ± 0.00	20.00 ± 0.00	18.00 ± 0.00	0.77 ± 0.00	0.69 ± 0.00
Overall		10.91 ± 6.19	8.68 ± 6.15	5.02 ± 5.03	0.74 ± 0.32	0.40 ± 0.35

Table 6.2: Author factuality scores obtained by gemma2-9b across tasks. The model failed to answer most of the use cases. Among the valid responses, 80% of suggested authors are scientists, and only 50% are in the APS dataset.

Task	Use Case	Total Names	Present OA	Present APS	Ratio OA	Ratio APS
Top-k	top-5	5.00 ± 0.00	5.00 ± 0.00	3.00 ± 0.00	1.00 ± 0.00	0.60 ± 0.00
	top-100	-	-	-	-	-
Field	PER	10.00 ± 0.00	6.00 ± 0.00	4.00 ± 0.00	0.60 ± 0.00	0.40 ± 0.00
	CM&MP	-	-	-	-	-
Epoch	1950s	-	-	-	-	-
	2000s	-	-	-	-	-
Seniority	early career	-	-	-	-	-
	senior	-	-	-	-	-
Twins	famous (M)	-	-	-	-	-
	famous (F)	-	-	-	-	-
	random (M)	-	-	-	-	-
	random (F)	-	-	-	-	-
	politic (M)	-	-	-	-	-
	politic (F)	-	-	-	-	-
	movie (M)	-	-	-	-	-
	movie (F)	-	-	-	-	-
	fictitious (M)	-	-	-	-	-
	fictitious (F)	-	-	-	-	-
Overall		7.50 ± 2.50	5.50 ± 0.50	3.50 ± 0.50	0.80 ± 0.20	0.50 ± 0.10

Table 6.3: Author factuality scores obtained by mixtral-8x7b across tasks. The model failed to answer one use case: statistical twin of a female politician. While 81% of suggested authors are scientists, only 53% are in the APS dataset.

Task	Use Case	Total Names	Present OA	Present APS	Ratio OA	Ratio APS
Top-k	top-5	5.00 ± 0.00	5.00 ± 0.00	3.00 ± 0.00	1.00 ± 0.00	0.60 ± 0.00
	top-100	102.00 ± 0.00	97.00 ± 0.00	70.00 ± 0.00	0.95 ± 0.00	0.69 ± 0.00
Field	CM&MP	10.87 ± 2.97	10.47 ± 1.62	7.00 ± 0.00	0.98 ± 0.06	0.67 ± 0.10
	PER	20.13 ± 2.97	9.61 ± 1.35	5.76 ± 0.81	0.48 ± 0.01	0.29 ± 0.00
Epoch	1950s	17.37 ± 2.16	14.53 ± 1.62	12.37 ± 2.16	0.84 ± 0.02	0.70 ± 0.06
	2000s	19.08 ± 3.20	18.15 ± 2.93	12.62 ± 1.33	0.95 ± 0.01	0.68 ± 0.09
Seniority	early_career	8.00 ± 0.00	5.15 ± 0.53	3.08 ± 0.27	0.64 ± 0.07	0.38 ± 0.03
	senior	41.45 ± 8.96	40.52 ± 8.69	30.20 ± 6.32	0.98 ± 0.01	0.73 ± 0.02
Twins	famous (M)	11.79 ± 2.70	10.87 ± 2.97	7.87 ± 2.97	0.92 ± 0.02	0.65 ± 0.06
	famous (F)	10.00 ± 0.00	10.00 ± 0.00	5.00 ± 0.00	1.00 ± 0.00	0.50 ± 0.00
	random (M)	10.00 ± 0.00	8.00 ± 0.00	7.00 ± 0.00	0.80 ± 0.00	0.70 ± 0.00
	random (F)	10.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
	politic (M)	-	-	-	-	-
	politic (F)	10.79 ± 2.70	7.95 ± 3.24	2.71 ± 2.43	0.72 ± 0.07	0.23 ± 0.09
	movie (M)	5.00 ± 0.00	4.69 ± 1.07	3.69 ± 1.07	0.94 ± 0.21	0.74 ± 0.21
	movie (F)	10.00 ± 0.00	9.00 ± 0.00	4.00 ± 0.00	0.90 ± 0.00	0.40 ± 0.00
	fictitious (M)	10.00 ± 0.00	10.00 ± 0.00	5.00 ± 0.00	1.00 ± 0.00	0.50 ± 0.00
	fictitious (F)	9.32 ± 2.37	9.25 ± 2.63	5.55 ± 1.58	0.93 ± 0.26	0.56 ± 0.16
Overall		13.80 ± 11.26	11.46 ± 11.30	7.74 ± 8.60	0.81 ± 0.27	0.53 ± 0.23

Table 6.4: Author factuality scores obtained by llama3-70b across tasks. The model succeeded in answering all use cases. While 87% of suggested authors are scientists, only 62% are in the APS dataset.

Task	Use Case	Total Names	Present OA	Present APS	Ratio OA	Ratio APS
Top-k	top-5	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
	top-100	109.77 ± 24.72	95.57 ± 22.05	60.93 ± 14.52	0.87 ± 0.05	0.56 ± 0.07
Field	PER	16.24 ± 3.92	15.12 ± 3.17	12.22 ± 2.21	0.94 ± 0.07	0.77 ± 0.11
	CM&MP	32.76 ± 13.50	28.78 ± 12.19	21.90 ± 10.10	0.88 ± 0.04	0.66 ± 0.05
Epoch	1950s	17.34 ± 2.36	16.20 ± 1.17	11.95 ± 1.65	0.94 ± 0.08	0.69 ± 0.08
	2000s	28.49 ± 18.64	23.88 ± 14.43	18.07 ± 9.70	0.85 ± 0.04	0.65 ± 0.07
Seniority	early career	16.49 ± 2.25	14.12 ± 1.95	12.24 ± 1.84	0.86 ± 0.03	0.74 ± 0.02
	senior	29.12 ± 5.76	23.12 ± 5.42	15.98 ± 3.84	0.79 ± 0.02	0.55 ± 0.03
Twins	famous (M)	10.10 ± 0.48	9.49 ± 0.67	6.80 ± 1.09	0.94 ± 0.08	0.67 ± 0.09
	famous (F)	8.90 ± 1.36	8.02 ± 1.41	6.73 ± 1.29	0.90 ± 0.05	0.75 ± 0.06
	random (M)	8.27 ± 0.63	7.98 ± 0.27	3.05 ± 0.22	0.97 ± 0.07	0.37 ± 0.02
	random (F)	9.71 ± 0.51	5.51 ± 0.97	1.73 ± 0.63	0.57 ± 0.11	0.18 ± 0.07
	politic (M)	3.37 ± 0.48	2.02 ± 0.15	1.39 ± 0.49	0.61 ± 0.10	0.40 ± 0.09
	politic (F)	8.90 ± 1.16	8.80 ± 1.11	3.12 ± 0.50	0.99 ± 0.03	0.35 ± 0.03
	movie (M)	9.05 ± 0.22	8.05 ± 0.22	6.05 ± 0.22	0.89 ± 0.00	0.67 ± 0.01
	movie (F)	10.00 ± 0.00	8.88 ± 0.39	6.00 ± 0.58	0.89 ± 0.04	0.60 ± 0.06
	fictitious (M)	9.10 ± 0.30	7.20 ± 0.59	7.20 ± 0.59	0.79 ± 0.04	0.79 ± 0.04
	fictitious (F)	7.00 ± 1.87	6.46 ± 1.65	4.80 ± 1.85	0.93 ± 0.09	0.68 ± 0.22
Overall		17.49 ± 22.25	15.15 ± 19.29	10.65 ± 12.75	0.87 ± 0.13	0.62 ± 0.20

6.1.2 FIELD-BASED FACTUALITY

Tables 6.5–6.6 present a detailed performance analysis of model recommendations for field-specific tasks (PER and CM&MP). The evaluation focuses on two primary aspects: (1) the accuracy of references, specifically the correctness of the DOI (Digital Object Identifier), and (2) the appropriateness of author suggestions, ensuring they belong to the correct field.

DOI-related metrics:

- **Total DOIs:** Average number of DOIs recommended by the model
- **Unique DOIs:** Average number of non-duplicate DOIs, indicating diversity in recommendations
- **In OA/APS:** Number of DOIs found in OpenAlex/APS databases, validating their existence
- **Authorship (Auth.) Correct:** Number of DOIs correctly attributed to their suggested authors
- **Field Match:** Number of DOIs that belong to the requested field
- **Field Ratio:** Proportion of field-relevant DOIs to total DOIs

Table 6.5: Comparison of DOI recommendations across all models for Physics Education Research (PER) and Condensed Matter & Materials Physics (CM&MP). Values show average \pm standard deviation.

Model	Field	Total DOIs	Unique DOIs	In OA	In APS	Auth. Correct	Field Match	Field Ratio
llama3-8b	CM&MP	16.00 \pm 0.00	16.00 \pm 0.00	10.00 \pm 0.00	10.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	PER	10.00 \pm 0.00	10.00 \pm 0.00	9.00 \pm 0.00	9.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	Overall	13.00 \pm 3.00	13.00 \pm 3.00	9.50 \pm 0.50	9.50 \pm 0.50	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
gemma2-9b	CM&MP	–	–	–	–	–	–	–
	PER	10.00 \pm 0.00	1.00 \pm 0.00	10.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	Overall	10.00 \pm 0.00	1.00 \pm 0.00	10.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
mixtral-8x7b	CM&MP	10.87 \pm 2.97	10.79 \pm 2.70	7.55 \pm 1.89	7.55 \pm 1.89	0.08 \pm 0.27	0.00 \pm 0.00	0.00 \pm 0.00
	PER	20.13 \pm 2.97	20.13 \pm 2.97	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
	Overall	15.50 \pm 5.50	15.46 \pm 5.46	3.78 \pm 4.01	3.78 \pm 4.01	0.04 \pm 0.19	0.00 \pm 0.00	0.00 \pm 0.00
llama3-70b	CM&MP	32.76 \pm 13.50	26.46 \pm 11.42	24.95 \pm 7.94	24.95 \pm 7.94	5.05 \pm 2.19	0.00 \pm 0.00	0.00 \pm 0.00
	PER	16.24 \pm 3.92	16.24 \pm 3.92	5.46 \pm 5.90	4.98 \pm 6.28	0.07 \pm 0.26	2.78 \pm 4.72	0.19 \pm 0.31
	Overall	24.50 \pm 12.92	21.35 \pm 9.95	15.21 \pm 11.99	14.96 \pm 12.29	2.56 \pm 2.93	1.39 \pm 3.61	0.10 \pm 0.24

Author-related metrics:

- **Field Match:** Number of recommended authors who publish in the requested field
- **Field Ratio:** Proportion of field-relevant authors to total recommended authors

Table 6.6: Analysis of author field relevance in recommendations across all models for Physics Education Research (PER) and Condensed Matter & Materials Physics (CM&MP). Values show average \pm standard deviation.

Model	Field	Field Match	Field Ratio
llama3-8b	CM&MP	0.00 \pm 0.00	0.00 \pm 0.00
	PER	0.00 \pm 0.00	0.00 \pm 0.00
	Overall	0.00 \pm 0.00	0.00 \pm 0.00
gemma2-9b	CM&MP	–	–
	PER	0.00 \pm 0.00	0.00 \pm 0.00
	Overall	0.00 \pm 0.00	0.00 \pm 0.00
mixtral-8x7b	CM&MP	0.08 \pm 0.27	0.00 \pm 0.01
	PER	0.00 \pm 0.00	0.00 \pm 0.00
	Overall	0.04 \pm 0.19	0.00 \pm 0.01
llama3-70b	CM&MP	5.05 \pm 2.19	0.16 \pm 0.05
	PER	0.07 \pm 0.26	0.01 \pm 0.02
	Overall	2.56 \pm 2.93	0.08 \pm 0.09

6.1.3 TEMPORAL FACTUALITY

Table 6.7 presents the temporal awareness analysis for each model. For each model, we provide:

- **Total Names:** Number of scientists recommended by the model
- **Correct Author Epoch:** Number of authors with verified publications in the specified decade
- **Ratio:** Proportion of valid temporal recommendations to total recommendations

Results are presented as mean \pm standard deviation across multiple evaluation runs for two distinct periods (1950s and 2000s). Validation relies on publication records from OpenAlex, with each author requiring at least one publication in the respective decade (1950-1959 or 2000-2009). Dashes (–) indicate tasks where the model failed to generate valid responses.

Table 6.7: Comparison of epoch-specific recommendations across models for the 1950s and 2000s periods. Values show average \pm standard deviation. Notable performance metrics are highlighted in bold.

Model	Epoch	Total Names	Correct Author Epoch	Ratio
llama3-8b	1950s	11.88 \pm 1.66	8.32 \pm 2.48	0.69 \pm 0.10
	2000s	–	–	–
	Overall	11.88 \pm 1.66	8.32 \pm 2.48	0.69 \pm 0.10
gemma2-9b	1950s	–	–	–
	2000s	–	–	–
	Overall	–	–	–
mixtral-8x7b	1950s	17.37 \pm 2.16	11.45 \pm 1.89	0.65 \pm 0.04
	2000s	19.08 \pm 3.20	15.38 \pm 2.13	0.82 \pm 0.05
	Overall	18.23 \pm 2.86	13.44 \pm 2.82	0.74 \pm 0.09
llama3-70b	1950s	17.34 \pm 2.36	9.71 \pm 1.66	0.56 \pm 0.03
	2000s	28.49 \pm 18.64	21.68 \pm 13.33	0.76 \pm 0.05
	Overall	22.91 \pm 14.41	15.70 \pm 11.23	0.66 \pm 0.11

6.1.4 SENIORITY-BASED FACTUALITY

Table 6.8 presents the career stage identification analysis for each model. Following the categorization proposed by [54], career stages are determined based on publication records, where a scientist’s career age (A_i) is calculated as:

$$A_i = \text{year}_{\text{last_pub}} - \text{year}_{\text{first_pub}} \quad (6.1)$$

Scientists are classified into two categories:

- Early career: $A_i \leq 10$ years
- Senior: $A_i \geq 20$ years

To further evaluate model accuracy, we compare the derived career age (A_i) with the career age estimated by the model (C_i). The estimation error (E_i) for each scientist i is computed as:

$$E_i = C_i - A_i \quad (6.2)$$

For each model, we provide:

- **Total Names:** Number of scientists recommended
- **Correct Author Seniority:** Number of scientists with verified career stage
- **Ratio:** Proportion of correct seniority predictions
- **Career Age Error:** Mean and standard deviation of estimation error (E_i)

Results are presented as mean \pm standard deviation across multiple evaluation runs. Career age error (E_i) indicates model bias, where positive values represent overestimation and negative values represent underestimation. Dashes (-) indicate tasks where the model failed to generate valid responses.

Table 6.8: Comparison of seniority-based recommendations across models. Career age errors show systematic underestimation for early career (-) and overestimation for senior scientists (+). Values show average \pm standard deviation.

Model	Seniority	Total Names	Correct Seniority	Ratio	Career Age Error
llama3-8b	Early	9.80 \pm 0.40	2.88 \pm 0.33	0.29 \pm 0.04	-15.50 \pm 1.74
	Senior	11.00 \pm 0.00	5.00 \pm 0.00	0.45 \pm 0.00	40.33 \pm 0.00
	Overall	10.20 \pm 0.65	3.57 \pm 1.03	0.35 \pm 0.08	2.80 \pm 26.25
gemma2-9b	Early	-	-	-	-
	Senior	-	-	-	-
	Overall	-	-	-	-
mixtral-8x7b	Early	8.00 \pm 0.00	1.08 \pm 0.27	0.13 \pm 0.03	-12.71 \pm 2.39
	Senior	41.45 \pm 8.96	34.90 \pm 7.37	0.85 \pm 0.02	30.82 \pm 13.96
	Overall	24.94 \pm 17.90	18.20 \pm 17.71	0.49 \pm 0.36	9.33 \pm 23.98
llama3-70b	Early	16.49 \pm 2.25	0.59 \pm 0.49	0.04 \pm 0.03	-14.79 \pm 1.73
	Senior	29.12 \pm 5.76	17.34 \pm 4.92	0.59 \pm 0.05	20.28 \pm 2.85
	Overall	22.80 \pm 7.69	8.96 \pm 9.08	0.31 \pm 0.28	2.75 \pm 17.69

6.2 RESPONSE PATTERN ANALYSIS

Tables 6.9–6.14 present detailed metrics on model response consistency, expanding on the analysis presented in Section 4.2.2. For each model, we evaluate two key aspects of response uniqueness:

BASIC RESPONSE PATTERNS

Tables 6.9–6.12 provide:

- **Total Requests:** Number of evaluation attempts
- **Unique:** Number of distinct responses generated
- **Valid Unique:** Number of well-formed, task-compliant responses (among the unique responses)

Results are presented across different use cases, with dashes (–) indicating failed response generation.

Table 6.9: Response uniqueness metrics for llama3-8b. Invalid outputs (marked with '-') indicate responses exceeding context limits due to name repetition.

Task	Use Case	Total Req.	Unique	Valid Unique
Top-k	top-5	41	3	3
	top-100	123	3	0
Field	CM&MP	41	1	1
	PER	41	2	2
Epoch	1950s	41	3	3
	2000s	123	1	0
Seniority	early career	41	3	3
	senior	100	3	1
Twins	famous (M)	84	3	2
	famous (F)	41	3	3
	random (M)	41	2	2
	random (F)	41	3	3
	politic (M)	41	3	3
	politic (F)	41	2	2
	movie (M)	41	2	2
	movie (F)	41	2	2
	fictitious (M)	41	2	2
	fictitious (F)	41	2	2
Average ± Std		55.78 ± 29.55	2.39 ± 0.70	2.00 ± 0.97

Table 6.10: Response uniqueness metrics for gemma2-9b, showing consistent declination of task attempts due to capability awareness.

Task	Use Case	Total Req.	Unique	Valid Unique
Top-k	top-5	41	1	1
	top-100	123	1	0
Field	CM&MP	123	1	0
	PER	41	1	1
Epoch	1950s	123	1	0
	2000s	123	1	0
Seniority	early career	123	1	0
	senior	123	1	0
Twins	famous (M)	123	1	0
	famous (F)	123	1	0
	random (M)	123	1	0
	random (F)	123	1	0
	politic (M)	41	1	0
	politic (F)	123	1	0
	movie (M)	123	1	0
	movie (F)	123	1	0
	fictitious (M)	123	1	0
	fictitious (F)	123	1	0
Average ± Std		109.33 ± 31.45	1.00 ± 0.00	0.11 ± 0.32

Table 6.11: Response uniqueness metrics for mixtral-8x7b, demonstrating consistent output patterns with moderate validity rates.

Task	Use Case	Total Req.	Unique	Valid Unique
Top-k	top-5	38	2	2
	top-100	111	4	1
Field	CM&MP	38	2	2
	PER	38	2	2
Epoch	1950s	38	2	2
	2000s	39	2	2
Seniority	early-career	39	2	2
	senior	40	2	2
Twins	famous (M)	38	2	2
	famous (F)	39	2	2
	random (M)	45	2	1
	random (F)	45	2	1
	politic (M)	117	2	0
	politic (F)	38	2	2
	movie (F)	45	2	1
	movie (M)	39	2	2
	fictitious (M)	108	2	1
	fictitious (F)	40	2	2
Average ± Std		51.94 ± 27.79	2.11 ± 0.47	1.61 ± 0.61

Table 6.12: Response uniqueness metrics for llama3-70b, showing high diversity in valid responses across all tasks.

Task	Use Case	Total Req.	Unique	Valid Unique
Top-k	top-5	41	7	7
	top-100	85	30	28
Field	CM&MP	44	24	21
	PER	42	26	25
Epoch	1950s	41	14	14
	2000s	43	22	20
Seniority	early_career	41	19	19
	senior	41	18	18
Twins	famous (M)	41	16	16
	famous (F)	41	15	15
	fictitious (M)	41	16	16
	fictitious (F)	41	21	21
	movie (M)	41	20	20
	movie (F)	41	13	13
	politic (M)	41	11	11
	politic (F)	41	16	16
	random (M)	41	12	12
	random (F)	41	26	26
Average ± Std		43.78 ± 10.32	18.11 ± 5.97	17.67 ± 5.46

RESPONSE SIMILARITY METRICS

Tables 6.13 and 6.14 analyze response patterns using two measures:

Inter-request Similarity (Table 6.13): Using the Jaccard Index for sets of names A and B from different requests:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Higher values (closer to 1.0) indicate more repetitive recommendations.

Table 6.13: Jaccard similarity among recommended names across requests. All models, except llama3-70b, show a tendency to repeat their answers over time, as indicated by high Jaccard index values. In contrast, llama3-70b demonstrates greater variability in its responses, frequently changing the recommended names for the same question across different requests.

Task	Use Case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Top-k	top-5	0.94 ± 0.13	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
	top-100	–	–	1.00 ± 0.00	0.40 ± 0.20
Field	CM&MP	1.00 ± 0.00	–	0.89 ± 0.25	0.58 ± 0.29
	PER	0.84 ± 0.17	1.00 ± 0.00	0.87 ± 0.30	0.42 ± 0.26
Epoch	1950s	0.78 ± 0.30	–	0.88 ± 0.28	0.58 ± 0.36
	2000s	–	–	0.85 ± 0.35	0.47 ± 0.31
Seniority	early_career	0.85 ± 0.22	–	0.85 ± 0.35	0.71 ± 0.23
	senior	1.00 ± 0.00	–	0.88 ± 0.30	0.78 ± 0.32
Twins	famous (M)	0.80 ± 0.20	–	0.86 ± 0.34	0.45 ± 0.41
	famous (F)	0.71 ± 0.39	–	1.00 ± 0.00	0.51 ± 0.39
	random (M)	1.00 ± 0.00	–	1.00 ± 0.00	0.75 ± 0.38
	random (F)	1.00 ± 0.00	–	1.00 ± 0.00	0.47 ± 0.27
	politic (M)	1.00 ± 0.00	–	–	0.57 ± 0.43
	politic (F)	0.85 ± 0.22	–	0.86 ± 0.33	0.75 ± 0.29
	movie (M)	1.00 ± 0.00	–	0.85 ± 0.35	0.99 ± 0.03
	movie (F)	0.56 ± 0.50	–	1.00 ± 0.00	0.80 ± 0.26
	fictitious (M)	0.93 ± 0.11	–	1.00 ± 0.00	0.85 ± 0.32
fictitious (F)	1.00 ± 0.00	–	0.86 ± 0.35	0.47 ± 0.35	
Final Avg ± Std		0.89 ± 0.25	1.00 ± 0.00	0.91 ± 0.28	0.65 ± 0.36

Intra-request Uniqueness (Table 6.14): Measuring name repetition within single responses using ratio $U(R)$ for recommendation set R :

$$U(R) = \frac{|\text{unique names in } R|}{|\text{total names in } R|}$$

Results are presented as mean ± standard deviation across evaluation runs. Perfect uniqueness corresponds to $U(R) = 1.0$.

Table 6.14: Author uniqueness ratio within requests. Both gemma2-9b and mixtral-8x7b consistently provide unique expert recommendations within requests. In contrast, the LLaMA models tend to introduce duplicate names, particularly in specific use cases such as the top-100 expert recommendations.

Task	Use Case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Top-k	top-5	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
	top-100	–	–	0.51 ± 0.00	0.89 ± 0.12
Field	CM&MP	1.00 ± 0.00	–	1.00 ± 0.00	0.93 ± 0.06
	PER	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00
Epoch	1950s	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	2000s	–	–	1.00 ± 0.00	0.99 ± 0.04
Seniority	early_career	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	senior	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
Twins	famous (M)	0.86 ± 0.16	–	1.00 ± 0.00	1.00 ± 0.00
	famous (F)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	random (M)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	random (F)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	politic (M)	1.00 ± 0.00	–	–	1.00 ± 0.00
	politic (F)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	movie (M)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	movie (F)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	fictitious (M)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
	fictitious (F)	1.00 ± 0.00	–	1.00 ± 0.00	1.00 ± 0.00
Final Avg ± Std		0.99 ± 0.05	1.00 ± 0.00	1.00 ± 0.04	0.99 ± 0.04

6.3 NOBEL LAUREATES ANALYSIS

Tables 6.15–6.18 present the analysis of Nobel laureate recommendations across models. For each model, we provide:

- **Nobel (mean/median):** Average and median number of Nobel laureates recommended
- **Ratio:** Proportion of Nobel laureates to total recommendations
- **Year Awarded:** Statistics of Nobel Prize award years (mean, median)
- **All Physics:** Whether all recommended laureates are Physics Nobel winners

Results are presented as mean ± standard deviation across multiple evaluation runs and use cases. Dashes (–) indicate tasks where the model failed to generate valid responses.

Table 6.15: Nobel laureates retrieved by llama3-8b across use cases. Highest retrieval in 1950s category.

Task	Use Case	Nobel (mean)	Nobel (median)	Ratio	Year Awarded (mean, median)	All Physics
Top-k	top-5	0	0	0.00 ± 0.00	-	-
	top-100	-	-	-	-	-
Field	CM&MP	0	0	0.00 ± 0.00	-	-
	PER	0	0	0.00 ± 0.00	-	-
Epoch	1950s	11.22 ± 0.41	11.00	0.96 ± 0.08	1958 ± 10, 1959	True
	2000s	-	-	-	-	-
Seniority	early career	0	0	0.00 ± 0.00	-	-
	senior	5.00 ± 0.00	5.00	0.45 ± 0.00	1988 ± 15, 1979	True
Twins	famous (M)	0	0	0.00 ± 0.00	-	-
	famous (F)	0.22 ± 0.41	0.00	0.02 ± 0.04	1969 ± 0, 1969	True
	random (M)	0	0	0.00 ± 0.00	-	-
	random (F)	0	0	0.00 ± 0.00	-	-
	politic (M)	0	0	0.00 ± 0.00	-	-
	politic (F)	0	0	0.00 ± 0.00	-	-
	movie (M)	0	0	0.00 ± 0.00	-	-
	movie (F)	0	0	0.00 ± 0.00	-	-
	fictitious (M)	0	0	0.00 ± 0.00	-	-
	fictitious (F)	0	0	0.00 ± 0.00	-	-
Overall		0.91 ± 2.87	0.00	0.08 ± 0.25	1963 ± 16, 1959	True

Table 6.16: Nobel laureates retrieved by gemma2-9b across use cases. Limited retrieval, mainly in top-5 experts category.

Task	Use Case	Nobel (mean)	Nobel (median)	Ratio	Year Awarded (mean, median)	All Physics
Top-k	top-5	3.00 ± 0.00	3.00	0.60 ± 0.00	1930 ± 26, 1921	True
	top-100	-	-	-	-	-
Field	CM&MP	-	-	-	-,-	-
	PER	1.00 ± 0.00	1.00	0.10 ± 0.00	2001 ± 0, 2001	True
Epoch	1950s	-	-	-	-,-	-
	2000s	-	-	-	-,-	-
Seniority	early career	-	-	-	-	-
	senior	-	-	-	-	-
Twins	famous (M)	-	-	-	-	-
	famous (F)	-	-	-	-	-
	random (M)	-	-	-	-	-
	random (F)	-	-	-	-	-
	politic (M)	-	-	-	-	-
	politic (F)	-	-	-	-	-
	movie (M)	-	-	-	-	-
	movie (F)	-	-	-	-	-
	fictitious (M)	-	-	-	-	-
	fictitious (F)	-	-	-	-	-
Overall		2.00 ± 1.00	2.00	0.35 ± 0.25	1948 ± 38, 1943	True

Table 6.17: Nobel laureates retrieved by mixtral-8x7b. High retrieval across most categories, with peak counts in top-100 (42 laureates), senior scientists (15), and 1950s experts (12).

Task	Use Case	Nobel (mean)	Nobel (median)	Ratio	Year Awarded (mean)	All Physics
Top-k	top-5	4.00 ± 0.00	4.00	0.80 ± 0.00	1932 ± 19	True
	top-100	42.00 ± 0.00	42.00	0.41 ± 0.00	1976 ± 31	True
Field	CM&MP	0.16 ± 0.54	0.00	0.01 ± 0.03	1977 ± 21	True
	PER	0.08 ± 0.27	0.00	0.01 ± 0.03	2001 ± 0	True
Epoch	1950s	11.53 ± 1.62	12.00	0.66 ± 0.02	1962 ± 6	True
	2000s	6.46 ± 1.87	7.00	0.32 ± 0.09	1999 ± 19	True
Seniority	early_career	0	0	0.00 ± 0.00	-	-
	senior	15.10 ± 3.16	16.00	0.37 ± 0.01	1968 ± 33	True
Twins	famous (M)	0.32 ± 1.08	0.00	0.02 ± 0.05	2012 ± 8	True
	famous (F)	7.00 ± 0.00	7.00	0.70 ± 0.00	1928 ± 18	True
	random (M)	0	0	0.00 ± 0.00	-	-
	random (F)	0	0	0.00 ± 0.00	-	-
	politic (M)	-	-	-	-	-
	politic (F)	1.63 ± 2.16	1.00	0.13 ± 0.09	1987 ± 44	True
	movie (M)	1.85 ± 0.53	2.00	0.37 ± 0.11	1992 ± 28	True
	movie (F)	3.00 ± 0.00	3.00	0.30 ± 0.00	1995 ± 23	False
	fictitious (M)	7.00 ± 0.00	7.00	0.70 ± 0.00	1928 ± 18	True
	fictitious (F)	6.47 ± 1.84	7.00	0.65 ± 0.18	1928 ± 18	True
Overall		4.13 ± 5.45	2.00	0.30 ± 0.29	1962 ± 34	False

Table 6.18: Nobel laureates retrieved by llama3-70b. Highest laureate ratios in top-5 experts (100%), 1950s experts (76%), and senior scientists (56%).

Task	Use Case	Nobel (mean)	Nobel (median)	Ratio	Year Awarded (mean)	All Physics
Top-k	top-5	5.00 ± 0.00	5.00	1.00 ± 0.00	1976 ± 8	True
	top-100	50.93 ± 16.53	46.00	0.47 ± 0.09	1984 ± 26	False
Field	CM&MP	5.90 ± 2.46	5.00	0.19 ± 0.07	1979 ± 15	False
	PER	0	0	0.00 ± 0.00	-	-
Epoch	1950s	13.00 ± 2.57	14.00	0.76 ± 0.17	1968 ± 14	True
	2000s	11.61 ± 6.38	10.00	0.43 ± 0.10	2002 ± 11	False
Seniority	early_career	0	0	0.00 ± 0.00	-	-
	senior	16.29 ± 2.75	16.00	0.56 ± 0.04	2003 ± 12	False
Twins	famous (M)	0	0	0.00 ± 0.00	-	-
	famous (F)	0	0	0.00 ± 0.00	-	-
	random (M)	0	0	0.00 ± 0.00	-	-
	random (F)	0	0	0.00 ± 0.00	-	-
	politic (M)	0.98 ± 0.15	1.00	0.29 ± 0.06	2021 ± 0	True
	politic (F)	0.88 ± 0.50	1.00	0.10 ± 0.05	1964 ± 4	False
	movie (M)	0	0	0.00 ± 0.00	-	-
	movie (F)	1.90 ± 0.43	2.00	0.19 ± 0.04	2016 ± 8	False
	fictitious (M)	1.00 ± 0.00	1.00	0.11 ± 0.00	2021 ± 0	True
	fictitious (F)	0	0	0.00 ± 0.00	-	-
Overall		5.29 ± 11.39	1.00	0.22 ± 0.30	1988 ± 23	False

6.4 RECOMMENDATIONS SIMILARITY ANALYSIS

Tables 6.19–6.23 present the detailed breakdown of similarity metrics summarized in Table 4.4. Results are reported as mean \pm standard deviation across evaluation runs. Asterisks (*) indicate results that are not statistically significant (based on bootstrap values), likely due to small sample sizes.

Table 6.19: Author similarity based on shared affiliations. Near-zero scores indicate minimal institutional overlap. (*) marks non-significant results.

Task	Use Case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-7ob
Top-k	top-5	0.0959 \pm 0.11	0.0000 \pm 0.00 (*)	0.0000 \pm 0.00 (*)	0.0439 \pm 0.03
	top-100	–	–	0.0254 \pm 0.05	0.0253 \pm 0.05
Field	CM&MP	0.0338 \pm 0.05	–	0.0638 \pm 0.08	0.0169 \pm 0.03
	PER	0.0000 \pm 0.00 (*)	0.0525 \pm 0.04	0.0493 \pm 0.07	0.0219 \pm 0.05
Epoch	1950s	0.0216 \pm 0.04	–	0.0270 \pm 0.05	0.0306 \pm 0.05
	2000s	–	–	0.0214 \pm 0.04	0.0302 \pm 0.05
Seniority	early career	–	–	0.0965 \pm 0.15	0.0280 \pm 0.05
	senior	0.0361 \pm 0.04	–	0.0336 \pm 0.05	0.0198 \pm 0.04
Twins	famous (M)	0.0800 \pm 0.11	–	0.0191 \pm 0.04	0.0222 \pm 0.06
	famous (F)	0.0225 \pm 0.05	–	0.0093 \pm 0.02	0.0183 \pm 0.03
	random (M)	–	–	0.0544 \pm 0.08	0.0374 \pm 0.05
	random (F)	–	–	–	0.0127 \pm 0.03
	politic (M)	–	–	–	0.0000 \pm 0.00 (*)
	politic (F)	0.0000 \pm 0.00 (*)	–	0.0106 \pm 0.02	0.0165 \pm 0.03
	movie (M)	0.0891 \pm 0.09	–	0.0365 \pm 0.03	0.1581 \pm 0.10
	movie (F)	–	–	0.0000 \pm 0.00 (*)	0.0616 \pm 0.10
	fictitious (M)	–	–	0.0093 \pm 0.02	0.0443 \pm 0.07
	fictitious (F)	0.0279 \pm 0.04	–	0.0170 \pm 0.03	0.0317 \pm 0.05
Overall		0.0397 \pm 0.06	0.0394 \pm 0.04	0.0322 \pm 0.05	0.0262 \pm 0.05

Table 6.20: Author similarity based on affiliation countries. Higher scores (0.2-0.4) indicate regional clustering. (*) marks non-significant results.

Task	Use Case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-7ob
Top-k	top-5	0.2420 \pm 0.16	0.4000 \pm 0.00	0.2550 \pm 0.11	0.3050 \pm 0.16
	top-100	–	–	0.2639 \pm 0.14	0.2433 \pm 0.15
Field	CM&MP	0.2821 \pm 0.13	–	0.3767 \pm 0.14	0.2606 \pm 0.15
	PER	0.4370 \pm 0.23	0.4444 \pm 0.08	0.4925 \pm 0.27	0.3471 \pm 0.23
Epoch	1950s	0.3021 \pm 0.19	–	0.3467 \pm 0.18	0.3254 \pm 0.17
	2000s	–	–	0.2449 \pm 0.17	0.2567 \pm 0.15
Seniority	early_career	–	–	0.2558 \pm 0.16	0.2528 \pm 0.15
	senior	0.1891 \pm 0.07	–	0.2677 \pm 0.15	0.2584 \pm 0.14
Twins	famous (M)	0.2333 \pm 0.26	–	0.1665 \pm 0.13	0.2189 \pm 0.19
	famous (F)	0.2433 \pm 0.14	–	0.3127 \pm 0.16	0.2145 \pm 0.14
	random (M)	–	–	0.2053 \pm 0.08	0.3664 \pm 0.11
	random (F)	–	–	–	0.1401 \pm 0.07 (*)
	politic (M)	–	–	–	0.2500 \pm 0.00
	politic (F)	0.2500 \pm 0.00	–	0.2516 \pm 0.15	0.3233 \pm 0.15
	movie (M)	0.3097 \pm 0.13	–	0.2944 \pm 0.17	0.3161 \pm 0.22
	movie (F)	–	–	0.2639 \pm 0.17	0.3043 \pm 0.16
	fictitious (M)	–	–	0.3127 \pm 0.16	0.2393 \pm 0.20
	fictitious (F)	0.2054 \pm 0.11	–	0.2576 \pm 0.16	0.1833 \pm 0.13
Overall		0.2498 \pm 0.15	0.4333 \pm 0.07	0.2739 \pm 0.16	0.2560 \pm 0.16

Table 6.21: Author similarity based on shared APS co-authors. Near-zero scores indicate minimal co-authorship overlap. (*) marks non-significant results due to small sample size.

Task	Use case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Top-k	top-5	0.0160 ± 0.03	0.0000 ± 0.00 (*)	0.0119 ± 0.02	0.0096 ± 0.02
	top-100	–	–	0.0037 ± 0.01	0.0026 ± 0.01
Field	CM&MP	0.0137 ± 0.02	–	0.0265 ± 0.03	0.0094 ± 0.02
	PER	0.0000 ± 0.00 (*)	0.0000 ± 0.00 (*)	0.0304 ± 0.05	0.0056 ± 0.02
Epoch	1950s	0.0084 ± 0.01	–	0.0110 ± 0.02	0.0087 ± 0.02
	2000s	–	–	0.0041 ± 0.01	0.0055 ± 0.02
Seniority	early_career	–	–	0.0103 ± 0.01	0.0125 ± 0.03
	senior	0.0135 ± 0.03	–	0.0043 ± 0.01	0.0022 ± 0.01
Twins	famous (M)	0.0504 ± 0.06	–	0.0135 ± 0.05	0.0195 ± 0.03
	famous (F)	0.0008 ± 0.00	–	0.0060 ± 0.01	0.0224 ± 0.04
	random (M)	–	–	0.0077 ± 0.02	0.0053 ± 0.01
	random (F)	–	–	–	0.0002 ± 0.00 (*)
	politic (M)	–	–	–	0.0000 ± 0.00 (*)
	politic (F)	0.0000 ± 0.00 (*)	–	0.0006 ± 0.00	0.0000 ± 0.00 (*)
	movie (M)	0.0144 ± 0.03	–	0.0020 ± 0.00	0.0212 ± 0.03
	movie (F)	–	–	0.0000 ± 0.00 (*)	0.0084 ± 0.03
	fictitious (M)	–	–	0.0060 ± 0.01	0.0105 ± 0.02
	fictitious (F)	0.0009 ± 0.00	–	0.0036 ± 0.01	0.0084 ± 0.03
Overall		0.0078 ± 0.02	0.0000 ± 0.00	0.0064 ± 0.02	0.0052 ± 0.02

Table 6.22: Categorical similarity of authors based on gender, ethnicity, first APS publication decade, and Nobel awards.

Task	Use Case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Top-k	top-5	0.3207 ± 0.12	0.3333 ± 0.00	0.5556 ± 0.31	0.5143 ± 0.23
	top-100	–	–	0.3291 ± 0.22	0.2767 ± 0.20
Field	CM&MP	0.3697 ± 0.20	–	0.2808 ± 0.20	0.3132 ± 0.20
	PER	0.3200 ± 0.21	0.4222 ± 0.13	0.3363 ± 0.16	0.4346 ± 0.22
Epoch	1950s	0.2989 ± 0.26	–	0.2552 ± 0.23	0.2818 ± 0.23
	2000s	–	–	0.2467 ± 0.18	0.2945 ± 0.21
Seniority	early_career	–	–	0.3569 ± 0.18	0.2922 ± 0.17
	senior	0.3905 ± 0.16	–	0.2802 ± 0.20	0.2798 ± 0.20
Twins	famous (M)	0.6400 ± 0.26	–	0.3684 ± 0.20	0.5170 ± 0.27
	famous (F)	0.3482 ± 0.15	–	0.5016 ± 0.27	0.4670 ± 0.25
	random (M)	–	–	0.4540 ± 0.17	0.4269 ± 0.38
	random (F)	–	–	–	0.5238 ± 0.13
	politic (M)	–	–	–	0.1429 ± 0.00
	politic (F)	0.3333 ± 0.00	–	0.2570 ± 0.21	0.2515 ± 0.12
	movie (M)	0.4000 ± 0.21	–	0.2143 ± 0.13	0.3631 ± 0.22
	movie (F)	–	–	0.2063 ± 0.09	0.3377 ± 0.19
	fictitious (M)	–	–	0.5016 ± 0.27	0.3493 ± 0.23
	fictitious (F)	0.3859 ± 0.23	–	0.3581 ± 0.28	0.4544 ± 0.18
Overall		0.3775 ± 0.23	0.4000 ± 0.12	0.2908 ± 0.21	0.3003 ± 0.21

Table 6.23: Scholarly metric similarity (cosine) between authors based on h-index, publications, citations, career age, collaborators, and institutional metrics. (*) marks non-significant results.

Task	Use Case	llama3-8b	gemma2-9b	mixtral-8x7b	llama3-70b
Top-k	top-5	0.5105 ± 0.25	0.0000 ± 0.00 (*)	0.3085 ± 0.23 (*)	0.5105 ± 0.25
	top-100	0.7674 ± 0.20	–	0.7356 ± 0.23 (*)	0.7674 ± 0.20
Field	CM&MP	0.6443 ± 0.24	–	0.6823 ± 0.25	0.7627 ± 0.19
	PER	0.3405 ± 0.35	0.2543 ± 0.36	0.6184 ± 0.17	0.6458 ± 0.27 (*)
Epoch	1950s	0.6579 ± 0.32	–	0.7140 ± 0.21	0.6925 ± 0.22
	2000s	–	–	0.6912 ± 0.21	0.7378 ± 0.21
Seniority	early_career	0.3669 ± 0.26 (*)	–	0.2769 ± 0.34	0.6379 ± 0.22
	senior	0.6379 ± 0.28	–	0.7873 ± 0.17	0.6379 ± 0.28
Twins	famous (M)	0.6954 ± 0.22	–	0.6352 ± 0.25	0.6096 ± 0.32
	famous (F)	0.6096 ± 0.32	–	0.5468 ± 0.30	0.6096 ± 0.32
	random (M)	0.3803 ± 0.38	–	0.6467 ± 0.18	0.3803 ± 0.38
	random (F)	0.1482 ± 0.27	–	–	0.1482 ± 0.27
	politic (M)	–	–	–	0.0000 ± 0.00 (*)
	politic (F)	0.4253 ± 0.31	–	0.5712 ± 0.35	0.4253 ± 0.31
	movie (M)	0.6408 ± 0.17	–	0.5285 ± 0.30	0.6408 ± 0.17
	movie (F)	0.5850 ± 0.31	–	0.5297 ± 0.30	0.5850 ± 0.31
	fictitious (M)	0.5223 ± 0.31	–	0.5468 ± 0.30	0.5014 ± 0.34
	fictitious (F)	0.5239 ± 0.27	–	0.6230 ± 0.26	0.5239 ± 0.27
Overall		0.6504 ± 0.26	0.1907 ± 0.33	0.7403 ± 0.21	0.7300 ± 0.23

6.5 GENDER DISTRIBUTION ANALYSIS

Tables 6.24 and 6.25 present the gender distribution across different use cases for smaller models (llama3-8b, gemma2-9b) and larger models (mixtral-8x7b, llama3-70b) respectively. For each model and use case, we report the proportion of male, female, and unknown-gender authors (primarily those with initialed first names). Results are presented as mean \pm standard deviation, with dashes (–) indicating tasks where the model failed to generate valid responses.

Table 6.24: Comparative analysis of gender distribution across different use cases between llama3-8b and gemma2-9b. Values show proportions as male/female/unknown (avg \pm std).

Task	Use Case	llama3-8b			gemma2-9b		
		Male	Female	Unknown	Male	Female	Unknown
Top-k	top-5	0.590 \pm 0.030	0.205 \pm 0.015	0.205 \pm 0.015	0.500 \pm 0.000	0.000 \pm 0.000	0.500 \pm 0.000
	top-100	–	–	–	–	–	–
Field	PER	0.667 \pm 0.000	0.000 \pm 0.000	0.333 \pm 0.000	1.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000
	CM&MP	0.500 \pm 0.000	0.200 \pm 0.000	0.300 \pm 0.000	–	–	–
Epoch	1950s	0.481 \pm 0.099	0.000 \pm 0.000	0.519 \pm 0.099	–	–	–
	2000s	–	–	–	–	–	–
Seniority	early career	–	–	–	–	–	–
	senior	1.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	–	–	–
Twins	famous (M)	0.200 \pm 0.000	0.000 \pm 0.000	0.800 \pm 0.000	–	–	–
	famous (F)	0.969 \pm 0.059	0.000 \pm 0.000	0.031 \pm 0.059	–	–	–
	random (M)	–	–	–	–	–	–
	random (F)	–	–	–	–	–	–
	politic (M)	–	–	–	–	–	–
	politic (F)	0.000 \pm 0.000	1.000 \pm 0.000	0.000 \pm 0.000	–	–	–
	movie (M)	0.778 \pm 0.000	0.111 \pm 0.000	0.111 \pm 0.000	–	–	–
	movie (F)	–	–	–	–	–	–
	fictitious (M)	1.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	–	–	–
	fictitious (F)	0.688 \pm 0.000	0.125 \pm 0.000	0.188 \pm 0.000	–	–	–
Overall		0.615 \pm 0.302	0.160 \pm 0.287	0.225 \pm 0.232	0.750 \pm 0.250	0.000 \pm 0.000	0.250 \pm 0.250

Table 6.25: Comparative analysis of gender distribution across different use cases between mixtral-8x7b and llama3-70b. Values show proportions as male/female/unknown (avg \pm std).

Task	Use Case	mixtral-8x7b			llama3-70b		
		Male	Female	Unknown	Male	Female	Unknown
Top-k	top-5	0.667 \pm 0.000	0.000 \pm 0.000	0.333 \pm 0.000	0.800 \pm 0.000	0.000 \pm 0.000	0.200 \pm 0.000
	top-100	0.848 \pm 0.000	0.030 \pm 0.000	0.121 \pm 0.000	0.793 \pm 0.038	0.047 \pm 0.013	0.160 \pm 0.038
Field	PER CM&MP	0.667 \pm 0.000 0.440 \pm 0.039	0.333 \pm 0.000 0.143 \pm 0.000	0.000 \pm 0.000 0.417 \pm 0.039	0.449 \pm 0.085 0.874 \pm 0.071	0.506 \pm 0.075 0.011 \pm 0.022	0.045 \pm 0.051 0.114 \pm 0.057
Epoch	1950s 2000s	0.596 \pm 0.045 0.744 \pm 0.022	0.077 \pm 0.022 0.077 \pm 0.022	0.327 \pm 0.022 0.179 \pm 0.044	0.709 \pm 0.062 0.744 \pm 0.070	0.004 \pm 0.020 0.053 \pm 0.028	0.287 \pm 0.060 0.203 \pm 0.072
Seniority	early career senior	0.654 \pm 0.044 0.832 \pm 0.023	0.038 \pm 0.133 0.039 \pm 0.024	0.308 \pm 0.089 0.129 \pm 0.001	0.292 \pm 0.128 0.856 \pm 0.025	0.266 \pm 0.043 0.006 \pm 0.017	0.442 \pm 0.099 0.138 \pm 0.017
Twins	famous (M)	0.850 \pm 0.025	0.005 \pm 0.016	0.146 \pm 0.009	0.690 \pm 0.046	0.029 \pm 0.051	0.281 \pm 0.072
	famous (F)	0.750 \pm 0.000	0.000 \pm 0.000	0.250 \pm 0.000	0.821 \pm 0.066	0.012 \pm 0.037	0.166 \pm 0.064
	random (M)	1.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.691 \pm 0.087	0.000 \pm 0.000	0.309 \pm 0.087
	random (F)	–	–	–	0.321 \pm 0.261	0.012 \pm 0.077	0.667 \pm 0.258
	politic (M)	–	–	–	0.024 \pm 0.154	0.976 \pm 0.154	0.000 \pm 0.000
	politic (F)	0.055 \pm 0.189	0.937 \pm 0.216	0.008 \pm 0.027	0.000 \pm 0.000	1.000 \pm 0.000	0.000 \pm 0.000
	movie (M)	0.750 \pm 0.000	0.000 \pm 0.000	0.250 \pm 0.000	0.669 \pm 0.010	0.166 \pm 0.005	0.166 \pm 0.005
	movie (F)	0.000 \pm 0.000	1.000 \pm 0.000	0.000 \pm 0.000	0.657 \pm 0.177	0.334 \pm 0.150	0.009 \pm 0.043
	fictitious (M)	0.750 \pm 0.000	0.000 \pm 0.000	0.250 \pm 0.000	0.726 \pm 0.045	0.143 \pm 0.015	0.132 \pm 0.037
	fictitious (F)	0.800 \pm 0.000	0.000 \pm 0.000	0.200 \pm 0.000	0.928 \pm 0.161	0.016 \pm 0.050	0.056 \pm 0.123
Overall		0.632 \pm 0.277	0.185 \pm 0.332	0.183 \pm 0.138	0.611 \pm 0.293	0.201 \pm 0.318	0.188 \pm 0.185

6.6 ETHNICITY DISTRIBUTION ANALYSIS

Tables 6.26–6.29 present the ethnic distribution across different use cases for each model. Authors are classified into five categories (Asian, White, Hispanic, Black, unknown) based on name-based inference methods. Results are presented as mean \pm standard deviation of proportions, with dashes (–) indicating tasks where the model failed to generate valid responses.

Table 6.26: Analysis of ethnicity distribution across different use cases in recommendations by llama3-8b. Values show proportions for each ethnic category (avg \pm std).

Task	Use Case	Asian	White	Hispanic	Black	Unknown
Epoch	1950s	0.166 \pm 0.044	0.408 \pm 0.039	0.130 \pm 0.025	0.018 \pm 0.034	0.278 \pm 0.015
	2000s	–	–	–	–	–
Field	PER CM&MP	0.333 \pm 0.000 0.700 \pm 0.000	0.667 \pm 0.000 0.000 \pm 0.000	0.000 \pm 0.000 0.000 \pm 0.000	0.000 \pm 0.000 0.100 \pm 0.000	0.000 \pm 0.000 0.200 \pm 0.000
Seniority	early career senior	– 0.000 \pm 0.000	– 0.750 \pm 0.000	– 0.000 \pm 0.000	– 0.000 \pm 0.000	– 0.250 \pm 0.000
Top-k	top-5 top-100	0.000 \pm 0.000 –	0.410 \pm 0.030 –	0.180 \pm 0.059 –	0.410 \pm 0.030 –	0.000 \pm 0.000 –
Twins	famous (M)	0.000 \pm 0.000	1.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000
	famous (F)	0.031 \pm 0.059	0.406 \pm 0.012	0.000 \pm 0.000	0.344 \pm 0.106	0.219 \pm 0.035
	random (M)	–	–	–	–	–
	random (F)	–	–	–	–	–
	politic (M)	–	–	–	–	–
	politic (F)	0.000 \pm 0.000	0.500 \pm 0.000	0.000 \pm 0.000	0.500 \pm 0.000	0.000 \pm 0.000
	movie (M)	0.000 \pm 0.000	0.556 \pm 0.000	0.222 \pm 0.000	0.222 \pm 0.000	0.000 \pm 0.000
	movie (F)	–	–	–	–	–
	fictitious (M)	1.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000
	fictitious (F)	0.062 \pm 0.000	0.562 \pm 0.000	0.312 \pm 0.000	0.000 \pm 0.000	0.062 \pm 0.000
Overall		0.223 \pm 0.331	0.453 \pm 0.273	0.082 \pm 0.113	0.155 \pm 0.188	0.086 \pm 0.110

Table 6.27: Analysis of ethnicity distribution across different use cases in recommendations by gemma2-9b. Values show proportions for each ethnic category (avg \pm std).

Task	Use Case	Asian	White	Hispanic	Black	Unknown
Epoch	1950s	-	-	-	-	-
	2000s	-	-	-	-	-
Field	PER	-	1.000 \pm 0.000	-	-	0.000 \pm 0.000
	CM&MP	-	-	-	-	-
Seniority	early career	-	-	-	-	-
	senior	-	-	-	-	-
Top-k	top-5	-	0.500 \pm 0.000	-	-	0.500 \pm 0.000
	top-100	-	-	-	-	-
Twins	famous (M)	-	-	-	-	-
	famous (F)	-	-	-	-	-
	random (M)	-	-	-	-	-
	random (F)	-	-	-	-	-
	politic (M)	-	-	-	-	-
	politic (F)	-	-	-	-	-
	movie (M)	-	-	-	-	-
	movie (F)	-	-	-	-	-
	fictitious (M)	-	-	-	-	-
	fictitious (F)	-	-	-	-	-
Overall		-	0.750 \pm 0.250	-	-	0.250 \pm 0.250

Table 6.28: Analysis of ethnicity distribution across different use cases in recommendations by mixtral-8x7b. Values show proportions for each ethnic category (avg \pm std).

Task	Use Case	Asian	White	Hispanic	Black	Unknown
Epoch	1950s	0.096 \pm 0.045	0.404 \pm 0.045	0.096 \pm 0.045	0.154 \pm 0.045	0.250 \pm 0.000
	2000s	0.077 \pm 0.022	0.462 \pm 0.155	0.077 \pm 0.022	0.231 \pm 0.067	0.154 \pm 0.044
Field	PER	0.000 \pm 0.000	0.846 \pm 0.045	0.154 \pm 0.045	0.000 \pm 0.000	0.000 \pm 0.000
	CM&MP	0.549 \pm 0.077	0.274 \pm 0.039	0.000 \pm 0.000	0.165 \pm 0.077	0.011 \pm 0.039
Seniority	early career	0.615 \pm 0.178	0.365 \pm 0.111	0.000 \pm 0.000	0.019 \pm 0.067	0.000 \pm 0.000
	senior	0.060 \pm 0.017	0.524 \pm 0.029	0.090 \pm 0.025	0.108 \pm 0.040	0.218 \pm 0.027
Top-k	top-5	0.000 \pm 0.000	0.667 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.333 \pm 0.000
	top-100	0.030 \pm 0.000	0.697 \pm 0.000	0.061 \pm 0.000	0.061 \pm 0.000	0.152 \pm 0.000
Twins	famous (M)	0.000 \pm 0.000	0.831 \pm 0.088	0.014 \pm 0.048	0.014 \pm 0.048	0.141 \pm 0.007
	famous (F)	0.000 \pm 0.000	0.750 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.250 \pm 0.000
	random (M)	0.143 \pm 0.000	0.571 \pm 0.000	0.143 \pm 0.000	0.143 \pm 0.000	0.000 \pm 0.000
	random (F)	-	-	-	-	-
	politic (M)	-	-	-	-	-
	politic (F)	0.000 \pm 0.000	0.508 \pm 0.027	0.000 \pm 0.000	0.016 \pm 0.054	0.476 \pm 0.081
	movie (M)	0.000 \pm 0.000	0.750 \pm 0.000	0.000 \pm 0.000	0.250 \pm 0.000	0.000 \pm 0.000
	movie (F)	0.000 \pm 0.000	0.667 \pm 0.000	0.000 \pm 0.000	0.333 \pm 0.000	0.000 \pm 0.000
	fictitious (M)	0.000 \pm 0.000	0.750 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.250 \pm 0.000
	fictitious (F)	0.000 \pm 0.000	0.600 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.400 \pm 0.000
Overall		0.110 \pm 0.206	0.587 \pm 0.180	0.041 \pm 0.060	0.100 \pm 0.115	0.161 \pm 0.161

Table 6.29: Analysis of ethnicity distribution across different use cases in recommendations by llama3-70b. Values show proportions for each ethnic category (avg \pm std).

Task	Use Case	Asian	White	Hispanic	Black	Unknown
Epoch	1950s	0.135 \pm 0.063	0.504 \pm 0.087	0.065 \pm 0.042	0.113 \pm 0.031	0.183 \pm 0.034
	2000s	0.081 \pm 0.046	0.658 \pm 0.103	0.007 \pm 0.020	0.172 \pm 0.056	0.082 \pm 0.068
Field	PER	0.045 \pm 0.051	0.865 \pm 0.049	0.015 \pm 0.041	0.002 \pm 0.011	0.074 \pm 0.072
	CM&MP	0.234 \pm 0.045	0.429 \pm 0.045	0.020 \pm 0.029	0.182 \pm 0.052	0.134 \pm 0.020
Seniority	early career	0.174 \pm 0.031	0.536 \pm 0.055	0.098 \pm 0.020	0.000 \pm 0.000	0.192 \pm 0.043
	senior	0.173 \pm 0.061	0.602 \pm 0.031	0.007 \pm 0.022	0.145 \pm 0.035	0.072 \pm 0.032
Top-k	top-5	0.000 \pm 0.000	0.600 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000	0.400 \pm 0.000
	top-100	0.085 \pm 0.028	0.608 \pm 0.033	0.069 \pm 0.010	0.103 \pm 0.023	0.136 \pm 0.021
Twins	famous (M)	0.050 \pm 0.063	0.476 \pm 0.100	0.363 \pm 0.059	0.111 \pm 0.076	0.000 \pm 0.000
	famous (F)	0.017 \pm 0.047	0.451 \pm 0.111	0.402 \pm 0.144	0.000 \pm 0.000	0.130 \pm 0.089
	random (M)	0.309 \pm 0.087	0.000 \pm 0.000	0.669 \pm 0.013	0.000 \pm 0.000	0.022 \pm 0.080
	random (F)	0.000 \pm 0.000	0.988 \pm 0.077	0.012 \pm 0.077	0.000 \pm 0.000	0.000 \pm 0.000
	politic (M)	0.000 \pm 0.000	0.000 \pm 0.000	0.793 \pm 0.270	0.024 \pm 0.154	0.183 \pm 0.241
	politic (F)	0.000 \pm 0.000	0.311 \pm 0.103	0.000 \pm 0.000	0.360 \pm 0.072	0.329 \pm 0.058
	movie (M)	0.000 \pm 0.000	0.338 \pm 0.021	0.331 \pm 0.010	0.331 \pm 0.010	0.000 \pm 0.000
	movie (F)	0.000 \pm 0.000	0.667 \pm 0.032	0.158 \pm 0.036	0.176 \pm 0.053	0.000 \pm 0.000
Overall	fictitious (M)	0.017 \pm 0.060	0.703 \pm 0.041	0.146 \pm 0.022	0.134 \pm 0.040	0.000 \pm 0.000
	fictitious (F)	0.222 \pm 0.097	0.753 \pm 0.122	0.016 \pm 0.050	0.007 \pm 0.029	0.003 \pm 0.019
Overall		0.086 \pm 0.108	0.526 \pm 0.260	0.178 \pm 0.249	0.103 \pm 0.122	0.107 \pm 0.135

6.7 RANK PERCENTILE ANALYSIS

Tables 6.30–6.33 present the percentile rank analysis of recommended authors across different metrics. For each model, we provide:

- **Citations:** Percentile rank in total citation count
- **Publications:** Percentile rank in number of publications
- **h_index:** Percentile rank in h-index
- **e_index:** Percentile rank in e-index (citation intensity beyond h-index)
- **Norm. Citations:** Percentile rank in citations normalized by career age and productivity

Results are presented as mean \pm standard deviation across evaluation runs and use cases. Higher percentiles indicate authors rank among the top performers for that metric. Dashes (–) indicate tasks where the model failed to generate valid responses.

Table 6.30: Analysis of rank percentiles across different use cases in recommendations by llama3-8b. Values show average \pm standard deviation for various metrics.

Task	Use Case	Citations	Publications	h_index	e_index	Norm. Citations
Epoch	1950s	92.128 \pm 16.136	89.543 \pm 16.252	92.182 \pm 14.953	93.501 \pm 15.427	55.793 \pm 12.693
	2000s	–	–	–	–	–
Field	PER	71.694 \pm 19.881	73.246 \pm 28.615	71.354 \pm 24.450	67.851 \pm 27.135	45.936 \pm 12.031
	CM&MP	99.830 \pm 0.223	99.415 \pm 0.778	99.825 \pm 0.195	99.582 \pm 0.768	56.304 \pm 13.788
Seniority	early career	–	–	–	–	–
	senior	99.953 \pm 0.045	99.245 \pm 0.399	99.905 \pm 0.098	99.970 \pm 0.020	58.397 \pm 1.523
Top-k	top-5	99.438 \pm 0.459	95.487 \pm 1.990	98.380 \pm 1.169	99.860 \pm 0.111	69.099 \pm 9.557
	top-100	–	–	–	–	–
Twins	famous (M)	98.160 \pm 0.978	95.198 \pm 4.265	96.654 \pm 2.732	97.166 \pm 2.015	50.418 \pm 19.627
	famous (F)	69.946 \pm 33.911	63.337 \pm 41.251	64.771 \pm 38.224	62.354 \pm 40.633	40.166 \pm 27.660
	random (M)	–	–	–	–	–
	random (F)	–	–	–	–	–
	politic (M)	–	–	–	–	–
	politic (F)	93.110 \pm 6.690	77.630 \pm 16.690	84.240 \pm 12.500	95.970 \pm 3.970	77.630 \pm 8.360
	movie (M)	98.288 \pm 2.065	94.842 \pm 3.456	96.953 \pm 3.164	98.623 \pm 2.452	58.840 \pm 15.494
	movie (F)	–	–	–	–	–
	fictitious (M)	43.980 \pm 0.000	75.540 \pm 0.000	65.410 \pm 0.000	57.190 \pm 0.000	16.030 \pm 0.000
	fictitious (F)	88.065 \pm 10.734	87.624 \pm 9.545	86.647 \pm 12.065	88.330 \pm 9.940	43.749 \pm 16.302
Overall		90.810 \pm 17.230	88.878 \pm 18.695	89.724 \pm 18.280	90.507 \pm 19.239	52.401 \pm 18.815

Table 6.31: Analysis of rank percentiles across different use cases in recommendations by gemma2-9b. Values show average \pm standard deviation for various metrics.

Task	Use Case	Citations	Publications	h_index	e_index	Norm. Citations
Epoch	1950s	-	-	-	-	-
	2000s	-	-	-	-	-
Field	PER	69.303 \pm 28.233	62.353 \pm 26.076	56.263 \pm 35.338	77.447 \pm 19.643	67.077 \pm 10.764
	CM&MP	-	-	-	-	-
Seniority	early career	-	-	-	-	-
	senior	-	-	-	-	-
Top-k	top-5	96.600 \pm 3.290	88.280 \pm 8.340	93.475 \pm 6.245	98.340 \pm 1.650	70.675 \pm 3.225
	top-100	-	-	-	-	-
Twins	famous (M)	-	-	-	-	-
	famous (F)	-	-	-	-	-
	random (M)	-	-	-	-	-
	random (F)	-	-	-	-	-
	politic (M)	-	-	-	-	-
	politic (F)	-	-	-	-	-
	movie (M)	-	-	-	-	-
	movie (F)	-	-	-	-	-
	fictitious (M)	-	-	-	-	-
	fictitious (F)	-	-	-	-	-
Overall		80.222 \pm 25.718	72.724 \pm 24.436	71.148 \pm 33.124	85.804 \pm 18.367	68.516 \pm 8.763

Table 6.32: Analysis of rank percentiles across different use cases in recommendations by mixtral-8x7b. Values show average \pm standard deviation for various metrics.

Task	Use Case	Citations	Publications	h_index	e_index	Norm. Citations
Epoch	1950s	98.190 \pm 3.744	95.426 \pm 4.013	97.306 \pm 3.208	98.763 \pm 3.342	57.560 \pm 11.752
	2000s	95.459 \pm 12.621	87.501 \pm 24.295	90.064 \pm 21.542	92.058 \pm 24.962	64.864 \pm 19.572
Field	PER	85.027 \pm 10.711	84.811 \pm 10.868	85.660 \pm 11.234	84.843 \pm 11.070	42.229 \pm 12.970
	CM&MP	97.466 \pm 6.329	96.498 \pm 7.280	97.472 \pm 6.237	97.190 \pm 7.511	58.135 \pm 10.262
Seniority	early career	79.396 \pm 25.798	83.157 \pm 21.725	78.432 \pm 28.380	74.391 \pm 30.793	34.364 \pm 7.361
	senior	93.913 \pm 15.241	90.498 \pm 14.663	93.293 \pm 11.784	94.730 \pm 13.089	58.092 \pm 16.155
Top-k	top-5	96.657 \pm 2.687	87.370 \pm 6.930	92.257 \pm 5.382	98.173 \pm 1.368	74.367 \pm 5.847
	top-100	91.428 \pm 17.851	87.789 \pm 18.618	90.652 \pm 15.794	91.567 \pm 18.045	54.727 \pm 20.397
Twins	famous (M)	85.051 \pm 24.936	83.530 \pm 20.790	86.589 \pm 18.631	88.514 \pm 19.481	52.207 \pm 26.136
	famous (F)	88.998 \pm 13.469	76.233 \pm 20.203	82.792 \pm 17.042	85.955 \pm 21.196	79.945 \pm 10.909
	random (M)	97.971 \pm 3.533	97.073 \pm 2.628	97.501 \pm 3.395	97.337 \pm 4.748	58.070 \pm 18.357
	random (F)	-	-	-	-	-
	politic (M)	-	-	-	-	-
	politic (F)	90.906 \pm 9.197	86.516 \pm 13.247	89.458 \pm 11.044	90.553 \pm 11.238	58.691 \pm 16.682
	movie (M)	99.083 \pm 1.149	93.318 \pm 6.563	96.473 \pm 5.350	99.845 \pm 0.149	60.540 \pm 10.932
	movie (F)	78.903 \pm 28.199	85.097 \pm 8.927	81.767 \pm 19.558	84.247 \pm 22.045	60.630 \pm 27.410
	fictitious (M)	88.998 \pm 13.469	76.233 \pm 20.203	82.793 \pm 17.042	85.955 \pm 21.196	79.945 \pm 10.909
	fictitious (F)	89.980 \pm 12.206	80.086 \pm 19.645	84.660 \pm 15.694	87.058 \pm 19.086	69.324 \pm 23.376
Overall		92.831 \pm 15.197	89.121 \pm 16.181	91.378 \pm 14.812	92.736 \pm 16.503	58.751 \pm 18.924

Table 6.33: Analysis of rank percentiles across different use cases in recommendations by llama3-70b. Values show average \pm standard deviation for various metrics.

Task	Use Case	Citations	Publications	h_index	e_index	Norm. Citations
Epoch	1950s	92.648 \pm 12.661	87.383 \pm 21.752	90.559 \pm 19.207	91.230 \pm 21.256	55.362 \pm 16.203
	2000s	93.916 \pm 17.486	91.158 \pm 20.110	91.847 \pm 22.494	91.721 \pm 25.370	58.251 \pm 16.088
Field	PER	78.037 \pm 20.423	81.042 \pm 22.446	77.359 \pm 24.804	77.893 \pm 22.236	37.395 \pm 14.723
	CM&MP	98.501 \pm 7.001	96.956 \pm 9.239	97.935 \pm 7.934	98.486 \pm 7.597	55.944 \pm 14.470
Seniority	early career	95.629 \pm 9.606	92.272 \pm 15.928	90.288 \pm 24.473	89.647 \pm 27.833	55.755 \pm 14.237
	senior	84.787 \pm 24.804	79.446 \pm 31.379	78.664 \pm 34.234	76.659 \pm 39.511	58.196 \pm 16.591
Top-k	top-5	99.778 \pm 0.306	98.404 \pm 1.219	99.766 \pm 0.161	99.650 \pm 0.660	56.568 \pm 15.984
	top-100	91.084 \pm 19.975	87.779 \pm 22.202	89.338 \pm 21.958	90.240 \pm 22.969	55.386 \pm 19.120
Twins	famous (M)	94.436 \pm 7.361	94.718 \pm 5.844	93.739 \pm 8.478	93.336 \pm 7.866	41.059 \pm 12.875
	famous (F)	84.387 \pm 25.524	89.692 \pm 12.265	85.824 \pm 22.332	84.084 \pm 25.396	43.313 \pm 14.876
	random (M)	73.051 \pm 33.281	81.581 \pm 22.223	78.020 \pm 23.924	75.003 \pm 26.304	33.470 \pm 11.018
	random (F)	76.569 \pm 24.416	86.809 \pm 13.578	80.631 \pm 20.521	79.115 \pm 19.347	30.704 \pm 11.170
	politic (M)	99.375 \pm 0.664	99.407 \pm 0.435	99.158 \pm 0.810	98.146 \pm 2.405	37.406 \pm 10.575
	politic (F)	93.240 \pm 7.876	89.052 \pm 5.146	92.157 \pm 5.407	93.798 \pm 8.104	66.237 \pm 18.592
	movie (M)	98.343 \pm 1.988	95.083 \pm 2.064	96.954 \pm 2.966	98.145 \pm 2.887	59.263 \pm 19.365
	movie (F)	77.867 \pm 31.476	73.094 \pm 33.785	71.922 \pm 38.019	69.590 \pm 44.378	69.998 \pm 10.520
Overall	fictitious (M)	98.910 \pm 3.680	97.331 \pm 3.510	98.799 \pm 3.283	98.980 \pm 3.369	49.713 \pm 12.595
	fictitious (F)	97.810 \pm 8.575	96.082 \pm 13.888	97.154 \pm 12.284	96.537 \pm 13.944	53.830 \pm 22.407
Overall		91.190 \pm 18.902	89.188 \pm 20.644	89.338 \pm 22.670	89.244 \pm 24.866	53.607 \pm 18.022

References

- [1] P. T. von Hippel and S. Buck, “Improve academic search engines to reduce scholars’ biases,” *Nature Human Behaviour*, vol. 7, no. 2, pp. 157–158, Feb. 2023.
- [2] L. G. Campbell, S. Mehtani, M. E. Dozier, and J. Rinehart, “Gender-Heterogeneous Working Groups Produce Higher Quality Science,” *PLOS ONE*, vol. 8, no. 10, p. e79147, Oct. 2013, publisher: Public Library of Science. [Online]. Available: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0079147>
- [3] B. K. AlShebli, T. Rahwan, and W. L. Woon, “The preeminence of ethnic diversity in scientific collaboration,” *Nature Communications*, vol. 9, no. 1, p. 5163, Dec. 2018, publisher: Nature Publishing Group. [Online]. Available: <https://www.nature.com/articles/s41467-018-07634-8>
- [4] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, W. Ye, Y. Zhang, Y. Chang, P. S. Yu, Q. Yang, and X. Xie, “A Survey on Evaluation of Large Language Models,” *ACM Transactions on Intelligent Systems and Technology*, vol. 15, no. 3, pp. 1–45, Jun. 2024. [Online]. Available: <https://dl.acm.org/doi/10.1145/3641289>
- [5] I. H. Dunlap, “Going Digital: The Transformation of Scholarly Communication and Academic Libraries,” *Policy Futures in Education*, vol. 6, no. 1, pp. 132–141, Feb. 2008, publisher: SAGE Publications. [Online]. Available: <https://doi.org/10.2304/pfie.2008.6.1.132>
- [6] M. Gusenbauer and N. R. Haddaway, “Which academic search systems are suitable for systematic reviews or meta-analyses? Evaluating retrieval qualities of Google Scholar, PubMed, and 26 other resources,” *Research Synthesis Methods*, vol. 11, no. 2, pp. 181–217, Mar. 2020. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7079055/>
- [7] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, Y. Du, C. Yang, Y. Chen, Z. Chen, J. Jiang, R. Ren, Y. Li, X. Tang, Z. Liu, P. Liu, J.-Y. Nie, and J.-R. Wen, “A Survey of Large Language Models,” Nov. 2023, arXiv:2303.18223 [cs]. [Online]. Available: <http://arxiv.org/abs/2303.18223>
- [8] “The Matthew Effect in Science | Science.” [Online]. Available: <https://www.science.org/doi/10.1126/science.159.3810.56>
- [9] A. Klimashevskaja, D. Jannach, M. Elahi, and C. Trattner, “A survey on popularity bias in recommender systems,” *User Modeling and User-Adapted Interaction*, Jul. 2024. [Online]. Available: <https://link.springer.com/10.1007/s11257-024-09406-0>
- [10] J. M. Lichtenberg, A. Buchholz, and P. Schwöbel, “Large Language Models as Recommender Systems: A Study of Popularity Bias,” Jun. 2024, arXiv:2406.01285 [cs]. [Online]. Available: <http://arxiv.org/abs/2406.01285>

- [11] D. D. Bhowmik, “Gender Inequality in Higher Education and Research,” *Business Ethics and Leadership*, vol. 7, no. 3, pp. 108–119, Sep. 2023. [Online]. Available: <https://armgpublishing.com/journals/bel/volume-7-issue-3/article-10/>
- [12] B. Li, J. Jacob-Brassard, F. Dossa, K. Salata, T. Kishibe, E. Greco, N. N. Baxter, and M. Al-Omran, “Gender differences in faculty rank among academic physicians: a systematic review and meta-analysis,” *BMJ Open*, vol. 11, no. 11, p. e050322, Nov. 2021, publisher: British Medical Journal Publishing Group Section: Health policy. [Online]. Available: <https://bmjopen.bmj.com/content/11/11/e050322>
- [13] P. Lahoti, N. Blumm, X. Ma, R. Kotikalapudi, S. Potluri, Q. Tan, H. Srinivasan, B. Packer, A. Beirami, A. Beutel, and J. Chen, “Improving Diversity of Demographic Representation in Large Language Models via Collective-Critiques and Self-Voting,” in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 10383–10405. [Online]. Available: <https://aclanthology.org/2023.emnlp-main.643>
- [14] J. Wilsdon, L. Allen, E. Belfiore, P. Campbell, S. Curry, S. Hill, R. Jones, R. Kain, S. Kerridge, M. Thelwall, J. Tinkler, I. Viney, P. Wouters, J. Hill, and B. Johnson, *The Metric Tide: Report of the Independent Review of the Role of Metrics in Research Assessment and Management*, Jul. 2015.
- [15] “Introduction to Natural Language Processing.” [Online]. Available: <https://mitpress.mit.edu/9780262042840/introduction-to-natural-language-processing/>
- [16] “Speech and Language Processing.” [Online]. Available: <https://web.stanford.edu/~jurafsky/slp3/>
- [17] N. Friedman and J. Y. Halpern, “A Qualitative Markov Assumption and Its Implications for Belief Change.”
- [18] P. F. Brown, V. J. Della Pietra, P. V. deSouza, J. C. Lai, and R. L. Mercer, “Class-Based n -gram Models of Natural Language,” *Computational Linguistics*, vol. 18, no. 4, pp. 467–480, 1992. [Online]. Available: <https://aclanthology.org/J92-4003>
- [19] J. L. Elman, “Finding structure in time,” *Cognitive Science*, vol. 14, no. 2, pp. 179–211, Apr. 1990. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/036402139090002E>
- [20] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, u. Kaiser, and I. Polosukhin, “Attention is All you Need,” in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2017. [Online]. Available: https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html
- [22] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, “Improving Language Understanding by Generative Pre-Training.”
- [23] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu,

- C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language Models are Few-Shot Learners,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. [Online]. Available: <https://papers.nips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>
- [24] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe, “Training language models to follow instructions with human feedback.”
- [25] D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, and G. Irving, “Fine-Tuning Language Models from Human Preferences,” Jan. 2020, arXiv:1909.08593 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1909.08593>
- [26] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Akhtar, N. Barnes, and A. Mian, “A Comprehensive Overview of Large Language Models,” Apr. 2024, arXiv:2307.06435 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.06435>
- [27] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, “Scaling Laws for Neural Language Models,” Jan. 2020, arXiv:2001.08361 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/2001.08361>
- [28] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. van der Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Upasani, K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla, L. Rantala-Yearly, L. van der Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo, L. Blecher, L. Landzaat, L. de Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri, M. Kardas, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si, M. K. Singh, M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Krishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Stojnic, R. Raileanu, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly, R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim, S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang, S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler, T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami, V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti,

V. Petrovic, W. Chu, W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. E. Tan, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Grattafiori, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenberg, A. Vaughan, A. Baevski, A. Feinstein, A. Kallet, A. Sangani, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani, A. Franco, A. Saraf, A. Chowdhury, A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang, B. Loyd, B. De Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, D. Civin, D. Beaty, D. Kreymer, D. Li, D. Wyatt, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood, E. Brinkman, E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Ozgenel, F. Caggioni, F. Guzmán, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, G. Halpern, G. Thattai, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Shojanazeri, H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, I. Molybog, I. Tufanov, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski, J. Kohli, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Prasad, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan, K. Michelena, K. Li, K. Huang, K. Chawla, K. Lakhota, K. Huang, L. Chen, L. Garg, L. A. L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani, M. Bhatt, M. Tsimpoukelli, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi, M. Keneally, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. P. Laptev, N. Dong, N. Zhang, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Li, R. Hogan, R. Battey, R. Wang, R. Maheswari, R. Howes, R. Rinott, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best, T. Kohler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked, V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wang, X. Wu, X. Wang, X. Xia, X. Wu, X. Gao, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Hao, Y. Qian, Y. He, Z. Rait, Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, and Z. Zhao, “The Llama 3 Herd of Models,” Jul. 2024, arXiv:2407.21783 [cs]. [Online]. Available: <http://arxiv.org/abs/2407.21783>

- [29] OpenAI, J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, R. Avila, I. Babuschkin, S. Balaji, V. Balcom, P. Baltescu, H. Bao,

M. Bavarian, J. Belgum, I. Bello, J. Berdine, G. Bernadett-Shapiro, C. Berner, L. Bogdonoff, O. Boiko, M. Boyd, A.-L. Brakman, G. Brockman, T. Brooks, M. Brundage, K. Button, T. Cai, R. Campbell, A. Cann, B. Carey, C. Carlson, R. Carmichael, B. Chan, C. Chang, F. Chantzis, D. Chen, S. Chen, R. Chen, J. Chen, M. Chen, B. Chess, C. Cho, C. Chu, H. W. Chung, D. Cummings, J. Currier, Y. Dai, C. Decareaux, T. Degry, N. Deutsch, D. Deville, A. Dhar, D. Dohan, S. Dowling, S. Dunning, A. Ecoffet, A. Eleti, T. Eloundou, D. Farhi, L. Fedus, N. Felix, S. P. Fishman, J. Forte, I. Fulford, L. Gao, E. Georges, C. Gibson, V. Goel, T. Gogineni, G. Goh, R. Gontijo-Lopes, J. Gordon, M. Grafstein, S. Gray, R. Greene, J. Gross, S. S. Gu, Y. Guo, C. Hallacy, J. Han, J. Harris, Y. He, M. Heaton, J. Heidecke, C. Hesse, A. Hickey, W. Hickey, P. Hoeschele, B. Houghton, K. Hsu, S. Hu, X. Hu, J. Huizinga, S. Jain, S. Jain, J. Jang, A. Jiang, R. Jiang, H. Jin, D. Jin, S. Jomoto, B. Jonn, H. Jun, T. Kaftan, ♦. Kaiser, A. Kamali, I. Kanitscheider, N. S. Keskar, T. Khan, L. Kilpatrick, J. W. Kim, C. Kim, Y. Kim, J. H. Kirchner, J. Kiros, M. Knight, D. Kokotajlo, ♦. Kondraciuk, A. Kondrich, A. Konstantinidis, K. Kosic, G. Krueger, V. Kuo, M. Lampe, I. Lan, T. Lee, J. Leike, J. Leung, D. Levy, C. M. Li, R. Lim, M. Lin, S. Lin, M. Litwin, T. Lopez, R. Lowe, P. Lue, A. Makanju, K. Malfacini, S. Manning, T. Markov, Y. Markovski, B. Martin, K. Mayer, A. Mayne, B. McGrew, S. M. McKinney, C. McLeavey, P. McMillan, J. McNeil, D. Medina, A. Mehta, J. Menick, L. Metz, A. Mishchenko, P. Mishkin, V. Monaco, E. Morikawa, D. Mossing, T. Mu, M. Murati, O. Murk, D. Mély, A. Nair, R. Nakano, R. Nayak, A. Neelakantan, R. Ngo, H. Noh, L. Ouyang, C. O’Keefe, J. Pachocki, A. Paino, J. Palermo, A. Pantuliano, G. Parascandolo, J. Parish, E. Parparita, A. Passos, M. Pavlov, A. Peng, A. Perelman, F. d. A. B. Peres, M. Petrov, H. P. d. O. Pinto, Michael, Pokorny, M. Pokrass, V. H. Pong, T. Powell, A. Power, B. Power, E. Proehl, R. Puri, A. Radford, J. Rae, A. Ramesh, C. Raymond, F. Real, K. Rimbach, C. Ross, B. Rotsted, H. Roussez, N. Ryder, M. Saltarelli, T. Sanders, S. Santurkar, G. Sastry, H. Schmidt, D. Schnurr, J. Schulman, D. Selsam, K. Sheppard, T. Sherbakov, J. Shieh, S. Shoker, P. Shyam, S. Sidor, E. Sigler, M. Simens, J. Sitkin, K. Slama, I. Sohl, B. Sokolowsky, Y. Song, N. Staudacher, F. P. Such, N. Summers, I. Sutskever, J. Tang, N. Tezak, M. B. Thompson, P. Tillet, A. Tootoonchian, E. Tseng, P. Tuggle, N. Turley, J. Tworek, J. F. C. Uribe, A. Vallone, A. Vijayvergiya, C. Voss, C. Wainwright, J. J. Wang, A. Wang, B. Wang, J. Ward, J. Wei, C. J. Weinmann, A. Welihinda, P. Welinder, J. Weng, L. Weng, M. Wiethoff, D. Willner, C. Winter, S. Wolrich, H. Wong, L. Workman, S. Wu, J. Wu, M. Wu, K. Xiao, T. Xu, S. Yoo, K. Yu, Q. Yuan, W. Zaremba, R. Zellers, C. Zhang, M. Zhang, S. Zhao, T. Zheng, J. Zhuang, W. Zhuk, and B. Zoph, “GPT-4 Technical Report,” Mar. 2024, arXiv:2303.08774 [cs]. [Online]. Available: <http://arxiv.org/abs/2303.08774>

- [30] Y. Liu, J. Cao, C. Liu, K. Ding, and L. Jin, “Datasets for Large Language Models: A Comprehensive Survey,” Feb. 2024, arXiv:2402.18041 [cs]. [Online]. Available: <http://arxiv.org/abs/2402.18041>
- [31] “Common Crawl - Open Repository of Web Crawl Data.” [Online]. Available: <https://commoncrawl.org/>
- [32] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, N. Joseph, S. Kadavath, J. Kernion, T. Conerly, S. El-Showk, N. Elhage, Z. Hatfield-Dodds, D. Hernandez, T. Hume, S. Johnston, S. Kravec, L. Lovitt, N. Nanda, C. Olsson, D. Amodei, T. Brown, J. Clark, S. McCandlish, C. Olah, B. Mann, and J. Kaplan, “Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback,” Apr. 2022, arXiv:2204.05862 [cs]. [Online]. Available: <http://arxiv.org/abs/2204.05862>

- [33] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, “GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding,” in *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, T. Linzen, G. Chrupała, and A. Alishahi, Eds. Brussels, Belgium: Association for Computational Linguistics, Nov. 2018, pp. 353–355. [Online]. Available: <https://aclanthology.org/W18-5446>
- [34] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt, “Measuring Massive Multitask Language Understanding,” Jan. 2021, arXiv:2009.03300 [cs]. [Online]. Available: <http://arxiv.org/abs/2009.03300>
- [35] M. Glickman and Y. Zhang, “AI and Generative AI for Research Discovery and Summarization,” *Harvard Data Science Review*, vol. 6, no. 2, Mar. 2024. [Online]. Available: <https://hdr.mitpress.mit.edu/pub/xedo5giw>
- [36] H. Kang and C. Xiong, “ResearchArena: Benchmarking LLMs’ Ability to Collect and Organize Information as Research Agents,” Jun. 2024, arXiv:2406.10291 [cs]. [Online]. Available: <http://arxiv.org/abs/2406.10291>
- [37] I. Beltagy, K. Lo, and A. Cohan, “SciBERT: A Pretrained Language Model for Scientific Text,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 3615–3620. [Online]. Available: <https://aclanthology.org/D19-1371>
- [38] X. Ho, A. K. D. Nguyen, A. T. Dao, J. Jiang, Y. Chida, K. Sugimoto, H. Q. To, F. Boudin, and A. Aizawa, “A Survey of Pre-trained Language Models for Processing Scientific Text,” Jan. 2024, arXiv:2401.17824 [cs]. [Online]. Available: <http://arxiv.org/abs/2401.17824>
- [39] J. M. Nicholson, M. Mordaunt, P. Lopez, A. Uppala, D. Rosati, N. P. Rodrigues, P. Grabitz, and S. C. Rife, “scite: A smart citation index that displays the context of citations and classifies their intent using deep learning,” *Quantitative Science Studies*, vol. 2, no. 3, pp. 882–898, Nov. 2021. [Online]. Available: https://doi.org/10.1162/qss_a_00146
- [40] S. Whitfield and M. A. Hofmann, “Elicit: AI literature review research assistant,” *Public Services Quarterly*, vol. 19, no. 3, pp. 201–207, Jul. 2023, publisher: Routledge _eprint: <https://doi.org/10.1080/15228959.2023.2224125>. [Online]. Available: <https://doi.org/10.1080/15228959.2023.2224125>
- [41] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 9459–9474. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>
- [42] A. Huotala, M. Kuuttila, P. Ralph, and M. Mäntylä, “The Promise and Challenges of Using LLMs to Accelerate the Screening Process of Systematic Reviews,” May 2024, arXiv:2404.15667 [cs]. [Online]. Available: <http://arxiv.org/abs/2404.15667>

- [43] W. Liang, Y. Zhang, Z. Wu, H. Lepp, W. Ji, X. Zhao, H. Cao, S. Liu, S. He, Z. Huang, D. Yang, C. Potts, C. D. Manning, and J. Y. Zou, “Mapping the Increasing Use of LLMs in Scientific Papers,” Apr. 2024, arXiv:2404.01268 [cs]. [Online]. Available: <http://arxiv.org/abs/2404.01268>
- [44] C. Lu, C. Lu, R. T. Lange, J. Foerster, J. Clune, and D. Ha, “The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery,” Sep. 2024, arXiv:2408.06292. [Online]. Available: <http://arxiv.org/abs/2408.06292>
- [45] Z. Zhao, W. Fan, J. Li, Y. Liu, X. Mei, Y. Wang, Z. Wen, F. Wang, X. Zhao, J. Tang, and Q. Li, “Recommender Systems in the Era of Large Language Models (LLMs),” Apr. 2024, arXiv:2307.02046 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.02046>
- [46] F. E. Sandnes, “Can we identify prominent scholars using ChatGPT?” *Scientometrics*, vol. 129, no. 1, pp. 713–718, Jan. 2024. [Online]. Available: <https://doi.org/10.1007/s11192-023-04882-4>
- [47] R. Kleminski, P. Kazienko, and T. Kajdanowicz, “Analysis of direct citation, co-citation and bibliographic coupling in scientific topic identification,” *J. Inf. Sci.*, vol. 48, no. 3, pp. 349–373, Jun. 2022. [Online]. Available: <https://doi.org/10.1177/0165551520962775>
- [48] “API Overview | OpenAlex technical documentation,” Apr. 2024. [Online]. Available: <https://docs.openalex.org/how-to-use-the-api/api-overview>
- [49] “openalex-name-disambiguation/V3 at main · ourresearch/openalex-name-disambiguation.” [Online]. Available: <https://github.com/ourresearch/openalex-name-disambiguation/tree/main/V3>
- [50] “sciosci/demographicx,” Nov. 2024, original-date: 2022-09-07T18:24:54Z. [Online]. Available: <https://github.com/sciosci/demographicx>
- [51] S. Laohaprapanon, G. Sood, and B. Naji, “ethnicolr: Predict Race and Ethnicity From Name,” Jun. 2022, original-date: 2017-05-28T00:58:38Z. [Online]. Available: <https://github.com/appeler/ethnicolr>
- [52] “gender-guesser: Get the gender from first name.” [Online]. Available: <https://github.com/lead-ratings/gender-guesser>
- [53] V. Torvik, “Genni + Ethnea for the Author-ity 2009 dataset,” 2018, publisher: University of Illinois at Urbana-Champaign. [Online]. Available: <https://databank.illinois.edu/datasets/IDB-9087546>
- [54] S. Milojević, F. Radicchi, and J. P. Walsh, “Changing demographics of scientific careers: The rise of the temporary workforce,” *Proceedings of the National Academy of Sciences*, vol. 115, no. 50, pp. 12 616–12 623, Dec. 2018, publisher: Proceedings of the National Academy of Sciences. [Online]. Available: <https://www.pnas.org/doi/10.1073/pnas.1800478115>
- [55] C. Shah, “From Prompt Engineering to Prompt Science With Human in the Loop,” Feb. 2024, arXiv:2401.04122 [cs]. [Online]. Available: <http://arxiv.org/abs/2401.04122>
- [56] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas, E. B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L. R. Lavaud, L. Saulnier, M.-A. Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T. L. Scao, T. Gervet, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed, “Mixtral of Experts,” Jan. 2024, arXiv:2401.04088. [Online]. Available: <http://arxiv.org/abs/2401.04088>

- [57] G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard, B. Shahriari, A. Ramé, J. Ferret, P. Liu, P. Tafti, A. Friesen, M. Casbon, S. Ramos, R. Kumar, C. L. Lan, S. Jerome, A. Tsitsulin, N. Vieillard, P. Stanczyk, S. Girgin, N. Momchev, M. Hoffman, S. Thakoor, J.-B. Grill, B. Neyshabur, O. Bachem, A. Walton, A. Severyn, A. Parrish, A. Ahmad, A. Hutchison, A. Abdagic, A. Carl, A. Shen, A. Brock, A. Coenen, A. Laforge, A. Paterson, B. Bastian, B. Piot, B. Wu, B. Royal, C. Chen, C. Kumar, C. Perry, C. Welty, C. A. Choquette-Choo, D. Sinopalnikov, D. Weinberger, D. Vijaykumar, D. Rogozińska, D. Herbison, E. Bandy, E. Wang, E. Noland, E. Moreira, E. Senter, E. Eltyshev, F. Visin, G. Rasskin, G. Wei, G. Cameron, G. Martins, H. Hashemi, H. Klimczak-Plucińska, H. Batra, H. Dhand, I. Nardini, J. Mein, J. Zhou, J. Svensson, J. Stanway, J. Chan, J. P. Zhou, J. Carrasqueira, J. Iljazi, J. Becker, J. Fernandez, J. v. Amersfoort, J. Gordon, J. Lipschultz, J. Newlan, J.-y. Ji, K. Mohamed, K. Badola, K. Black, K. Millican, K. McDonell, K. Nguyen, K. Sodhia, K. Greene, L. L. Sjoesund, L. Usui, L. Sifre, L. Heuermann, L. Lago, L. McNealus, L. B. Soares, L. Kilpatrick, L. Dixon, L. Martins, M. Reid, M. Singh, M. Iverson, M. Görner, M. Velloso, M. Wirth, M. Davidow, M. Miller, M. Rahtz, M. Watson, M. Risdal, M. Kazemi, M. Moynihan, M. Zhang, M. Kahng, M. Park, M. Rahman, M. Khatwani, N. Dao, N. Bardoliwalla, N. Devanathan, N. Dumai, N. Chauhan, O. Wahltinez, P. Botarda, P. Barnes, P. Barham, P. Michel, P. Jin, P. Georgiev, P. Culliton, P. Kuppala, R. Comanescu, R. Merhej, R. Jana, R. A. Rokni, R. Agarwal, R. Mullins, S. Saadat, S. M. Carthy, S. Cogan, S. Perrin, S. M. R. Arnold, S. Krause, S. Dai, S. Garg, S. Sheth, S. Ronstrom, S. Chan, T. Jordan, T. Yu, T. Eccles, T. Hennigan, T. Kocisky, T. Doshi, V. Jain, V. Yadav, V. Meshram, V. Dharmadhikari, W. Barkley, W. Wei, W. Ye, W. Han, W. Kwon, X. Xu, Z. Shen, Z. Gong, Z. Wei, V. Cotruta, P. Kirk, A. Rao, M. Giang, L. Peran, T. Warkentin, E. Collins, J. Barral, Z. Ghahramani, R. Hadsell, D. Sculley, J. Banks, A. Dragan, S. Petrov, O. Vinyals, J. Dean, D. Hassabis, K. Kavukcuoglu, C. Farabet, E. Buchatskaya, S. Borgeaud, N. Fiedel, A. Joulin, K. Kenealy, R. Dadashi, and A. Andreev, “Gemma 2: Improving Open Language Models at a Practical Size,” Oct. 2024. [Online]. Available: <http://arxiv.org/abs/2408.00118>
- [58] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah, J. v. d. Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton, J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Upasani, K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla, K. Lakhota, L. Rantala-Yeary, L. v. d. Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan, L. Malo, L. Blecher, L. Landzaat, L. d. Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri, M. Kardas, M. Tsimpoukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si, M. K. Singh, M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, N. Zhang,

O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Krishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral, R. Stojnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly, R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim, S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bhosale, S. Zhang, S. Vandenhende, S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler, T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami, V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti, V. Albiero, V. Petrovic, W. Chu, W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia, X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D. Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld, A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenberg, A. Baevski, A. Feinstein, A. Kallet, A. Sangani, A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho, A. Poulton, A. Ryan, A. Ramchandani, A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury, A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang, B. Loyd, B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence, B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim, C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty, D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss, D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood, E.-T. Le, E. Brinkman, E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos, F. Ozgenel, F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee, G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri, H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan, I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski, J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul, J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg, J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan, K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A, L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabza, M. Avalani, M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi, M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan, M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. Santhanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev, N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab, P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj, Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy, R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu, S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto, S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang, S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield, S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman, T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou, T. Shaked, V. Vontimitta,

- V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu, V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable, X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li, Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait, Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma, “The Llama 3 Herd of Models,” Nov. 2024, arXiv:2407.21783. [Online]. Available: <http://arxiv.org/abs/2407.21783>
- [59] “Open LLM Leaderboard 2 - a Hugging Face Space by open-llm-leaderboard.” [Online]. Available: https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard
- [60] “GroqCloud.” [Online]. Available: <https://console.groq.com>
- [61] J. Liu, C. Liu, P. Zhou, R. Lv, K. Zhou, and Y. Zhang, “Is ChatGPT a Good Recommender? A Preliminary Study,” Oct. 2023, arXiv:2304.10149 [cs]. [Online]. Available: <http://arxiv.org/abs/2304.10149>
- [62] “A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications.” [Online]. Available: <https://arxiv.org/html/2402.07927v1>
- [63] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large Language Models are Zero-Shot Reasoners,” Jan. 2023, arXiv:2205.11916 [cs]. [Online]. Available: <http://arxiv.org/abs/2205.11916>
- [64] S. M. Bsharat, A. Myrzakhan, and Z. Shen, “Principled Instructions Are All You Need for Questioning LLaMA-1/2, GPT-3.5/4,” Jan. 2024, arXiv:2312.16171 [cs]. [Online]. Available: <http://arxiv.org/abs/2312.16171>
- [65] I. Arawjo, C. Swoopes, P. Vaithilingam, M. Wattenberg, and E. Glassman, “ChainForge: A Visual Toolkit for Prompt Engineering and LLM Hypothesis Testing,” Dec. 2023, arXiv:2309.09128 [cs]. [Online]. Available: <http://arxiv.org/abs/2309.09128>
- [66] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and T. Liu, “A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions,” Nov. 2023, arXiv:2311.05232 [cs]. [Online]. Available: <http://arxiv.org/abs/2311.05232>
- [67] ♦. Erkol, S. Sikdar, F. Radicchi, and S. Fortunato, “Consistency pays off in science,” *Quantitative Science Studies*, vol. 4, no. 2, pp. 491–500, May 2023. [Online]. Available: https://doi.org/10.1162/qss_a_00252
- [68] E. Hennessey, J. Cole, P. Shastri, J. Esquivel, C. Singh, R. Johnson, and S. Ghose, “Workshop report: Intersecting identities—gender and intersectionality in physics,” *AIP Conference Proceedings*, vol. 2109, no. 1, p. 040001, Jun. 2019. [Online]. Available: <https://doi.org/10.1063/1.5110070>
- [69] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” Apr. 2021, arXiv:2005.11401. [Online]. Available: <http://arxiv.org/abs/2005.11401>

Acknowledgments

This research was conducted during my internship at the Institute for Interactive Systems and Data Science at the University of Graz. Throughout this project, from initial design to implementation, I received invaluable mentorship from Dr. Lisette Espín-Noboa, to whom I am deeply grateful for her guidance and support. I am particularly thankful to my colleague and project partner Chiara Valentin, with whom I shared the journey of project ideation and development. Her contribution was essential, especially in designing and implementing the demographic inference methodology for the APS dataset. The daily interactions with my office teammates - Paula, Sudhang, and Ana - enriched both my research experience and personal growth. Their camaraderie and intellectual engagement made this journey truly memorable. This collaborative and supportive atmosphere flourished under the supervision of Prof. Fariba Karimi. I would also like to extend my gratitude to her Network Inequalities group at the Complexity Science Hub Vienna, whose commitment to excellence in research and insightful suggestions significantly enhanced my work. Finally, I am grateful to my primary supervisor at the University of Padova, Prof. Tomaso Erseghe, for his advice and support throughout my internship and research journey.