

# The Citation Blueprint: Final Report

Sanjiv Shrirang Joshi, Raghu Ram Sattanapalle

Northeastern University, Boston

{joshi.sanj,sattanapalle.r}@northeastern.edu

January 26, 2025



## The Citation Blueprint

Mapping Scholarly Impact: Insights from the DBLP Citation Network

Sanjiv Shrirang Joshi | Raghu Ram Sattanapalle

Hello, everyone! Welcome to our presentation, The Citation Blueprint. Today, we will explore how the DBLP Computer Science Bibliography dataset provides a rich foundation for analyzing scholarly impact and research collaborations.

Throughout this presentation, we will highlight our methods, findings, and the new metrics we developed to evaluate author influence, moving beyond traditional approaches like the h-index.

# Agenda



Introduction & Motivation



Dataset Overview & Objectives



Topic Modeling & Clustering



Scholarly Influence Analysis



Challenges and Lessons Learned



Contributions and Future Work

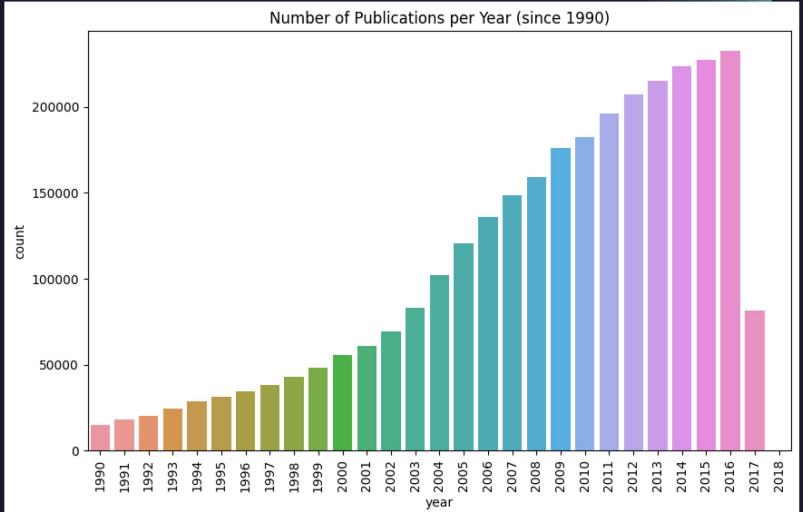
This agenda provides an overview of our presentation and highlights the key sections we'll cover.

- **Motivation and Objectives:** We'll start by discussing why this project is important and the key questions we set out to answer.
- **Dataset Overview:** Next, we'll provide details about the DBLP Citation Dataset, including its structure and significance for our analysis.
- **Topic Modeling and Clustering:** We'll dive into how we identified research topics, grouped papers into clusters, and visualized interconnections using techniques like LDA and t-SNE.
- **Scholarly Influence Analysis:** In this section, we'll explain how we measured influence using PageRank and other metrics, including a novel metric we developed.
- **Challenges and Lessons Learned:** We'll highlight some of the challenges we faced during the analysis and the insights we gained from overcoming them.
- **Contributions and Future Work:** Finally, we'll summarize the key contributions of this project and discuss potential areas for further exploration.

Let's begin by discussing the motivation behind this project and the objectives we aimed to achieve.

# Introduction & Motivation

- Citation networks are a rich source of insight into knowledge flow across disciplines.
- Current influence metrics (e.g., h-index) lack nuance in measuring impact across diverse contexts.
- Understanding research trends, influential works, and collaboration patterns is key to improving academic and industry collaboration.



- **Insight into Knowledge Flow:** Analyzing citation networks helps us see how research topics evolve, how knowledge disseminates, and how various disciplines influence one another. For instance, we can track the emergence of interdisciplinary research fields or identify how foundational studies influence applied research.
- **Gaps in Traditional Metrics:** While traditional influence metrics like the h-index are widely used, they have significant limitations. These metrics fail to capture nuanced aspects of influence, such as:
  - The role of citations from prestigious venues, which often signal higher-quality contributions.
  - The impact of interdisciplinary research, where influence spreads across multiple fields but may not result in high citation counts within any single domain.
- **Relevance to Collaboration:** Understanding research trends and identifying influential works can also enhance collaboration. This is not only important in academia, where collaboration drives innovation, but also in bridging the gap between academia and industry. By identifying influential authors and papers, institutions can form strategic partnerships to advance key areas of research.
- **Graph Explanation:** To highlight the growing importance of analyzing academic networks, we include this graph showing the exponential growth in publications over the decades. By the 2010s, the numbers soared to over 1.5 million publications, reflecting a dramatic increase in research output. This exponential growth underscores the need for advanced methods to analyze and understand this rapidly expanding academic landscape.

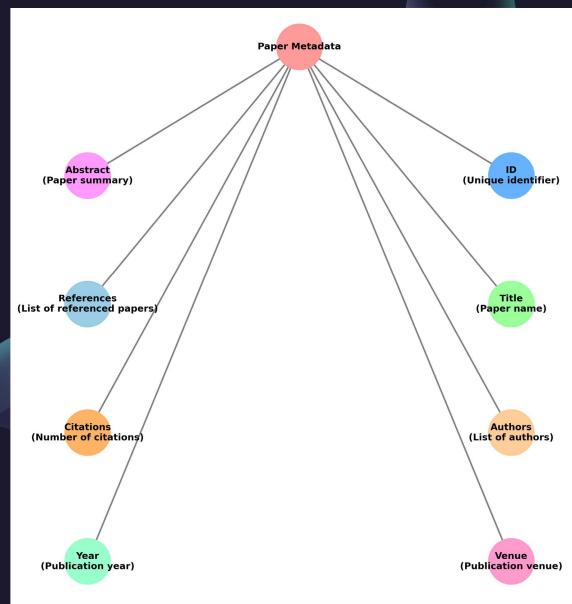
With this motivation in mind, let's look at the specific objectives we set out to achieve in this project. The next slide outlines our goals and how we planned to address the challenges presented by this rich dataset.

# Dataset Overview

## Key Dataset Information:

- **Total Papers:** 3,079,007
- **Year Range:** 1936–2017
- **Average References per Paper:** 8.17
- **Papers with No References:** 574,884
- **Papers with No Citations:** 718,250
- **Average Citations per Paper:** 35.22

This dataset provides a comprehensive view of academic contributions in computer science, enabling in-depth analysis of trends, influence, and collaborations over time.



The DBLP Citation Dataset is one of the most comprehensive sources for studying academic contributions in computer science. It spans over 80 years, from 1936 to 2017, and comprises more than 3 million papers and 25 million citations.

## Key Dataset Information:

- **Total Papers:** The dataset includes 3,079,007 papers, providing a diverse range of research contributions.
- **Year Range:** Covers the period from 1936 to 2017, offering insights into the evolution of computer science as a discipline.

**Dataset Structure:** The dataset contains detailed metadata for each paper, which makes it uniquely suited for in-depth analysis. This metadata includes:

- **Title:** The name of the research paper.
- **Authors:** A list of contributing authors.
- **Venue:** The publication venue (e.g., conference or journal).
- **Year:** The year of publication.
- **Citations:** The number of times the paper has been cited.
- **References:** The list of other papers cited by this work.
- **Abstract:** A summary of the research.

The dataset's large scale and detailed structure make it ideal for studying the evolution of academic research and identifying influential works. Now that we have a solid understanding of the dataset, let's explore the insights gained from our initial exploratory data analysis and move toward uncovering the underlying topics and trends in computer science research.

# Objectives

- **Understand Research Topics and Interconnections:**

Use techniques like LDA, clustering, and t-SNE to identify research topics and their relationships.

- **Measure Scholarly Influence:**

Implement PageRank and centrality metrics to identify key authors and papers, introducing a new metric that combines citation networks and co-authorship networks.

- **Predict Future Collaborations:**

Explore predictive techniques such as link prediction and association rule mining to uncover patterns of collaboration.

This project aims to uncover meaningful patterns in research and collaboration within the DBLP Citation Dataset. Our objectives are threefold:

- **Understand Research Topics and Interconnections:**

- Use LDA for topic modeling, clustering to group papers, and t-SNE for visualizing overlaps between topics.
- Map the research landscape and identify connections between fields.

- **Measure Scholarly Influence:**

- Implement PageRank and centrality metrics to identify influential papers and authors.
- Develop a novel metric combining citation and co-authorship networks to go beyond traditional measures like the h-index.

- **Predict Future Collaborations:**

- Use link prediction and association rule mining to identify emerging research collaborations. This is part of our high-risk analysis and we believe we made good progress on this front.

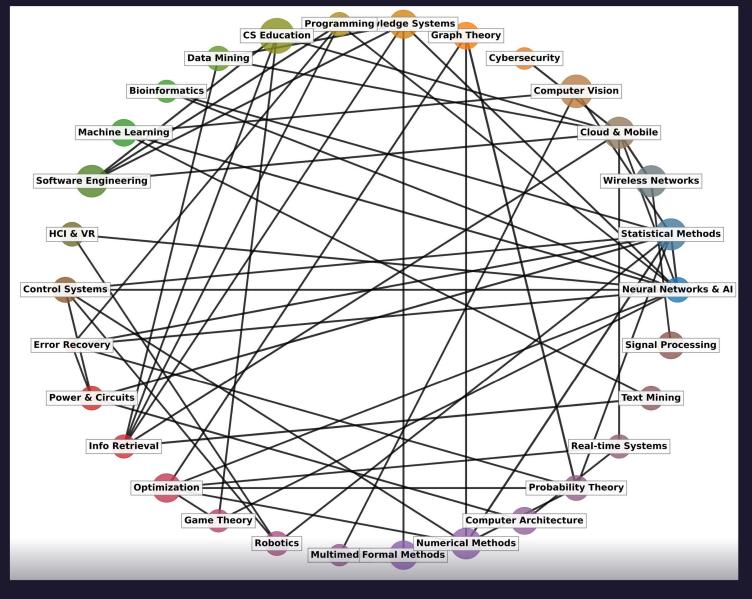
# Understanding the Building Blocks of Research

**Objective:** Identify dominant research themes as a foundation for deeper analysis.

- Applied Latent Dirichlet Allocation (LDA) on paper titles and abstracts from the DBLP dataset to extract 30 core research topics.
- Topics include "Machine Learning," "Graph Theory," "Cybersecurity," and others.

The circular graph shows:

- **Node size:** Topic prevalence in the dataset.
- **Edges:** Topic relationships or overlaps.



- This slide focuses on the first stage of our analysis: identifying key research areas within the DBLP dataset.

- **Methodology:**

- Used **LDA (Latent Dirichlet Allocation)** on the *titles and abstracts* of research papers in the DBLP dataset.
- Identified 30 dominant research topics across computer science.

- **Visual Insights (Circular Graph):**

- **Node Sizes:** Larger nodes, such as "*Software Engineering*" or "*Computer Vision*," represent more frequent research areas.
- **Edges:** Represent thematic overlaps between topics. For instance, "*Graph Theory*" is strongly connected to "*Optimization*" and "*Algorithms*" due to shared methodologies.

- **Key Observations:**

- Popular topics, like "*Neural Networks & AI*," show connections to multiple fields, highlighting their interdisciplinary nature.
- Specialized areas, such as "*Cybersecurity*," have fewer connections, reflecting their niche focus.

- **Challenges with LDA:**

- LDA effectively identifies broad themes but struggles to:
  - \* Capture detailed overlaps within subdomains.
  - \* Quantify the broader influence of each topic.
- To address these limitations, we applied **K-Means clustering** and used **t-SNE visualization** to uncover finer relationships. This sets the foundation for building a *comprehensive metric of scholarly influence*.

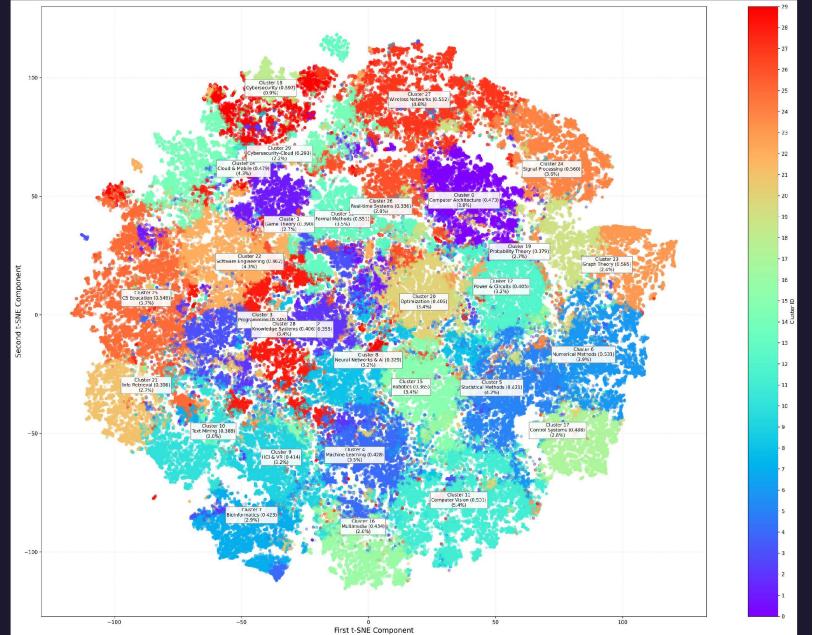
# Understanding Research Overlaps Through Clustering

**Objectives:** Explore finer relationships between research topics and clusters.

- Applied **K-Means clustering** on LDA topic distributions (30 clusters determined using silhouette and other metrics).
- Visualized clusters using **t-SNE** for a dense, interpretable 2D map.

## Insights:

- Identified overlaps between research topics (e.g., "Machine Learning" links with "Data Mining" and "Neural Networks & AI").
- Highlighted topic diversity within clusters.



After identifying topics with LDA, we applied **K-Means clustering** to group papers based on similar topic distributions. Using metrics such as the silhouette score and Calinski-Harabasz index, we determined that **30 clusters** best represented the dataset. To make these clusters interpretable, we used **t-SNE**, which mapped the high-dimensional data into a 2D space while preserving local relationships. Each cluster corresponds to a distinct research area, such as "*Machine Learning*," "*Graph Theory*," or "*Cybersecurity*." The density of points within a cluster reflects the strength of relationships among papers, while proximity between clusters indicates thematic overlaps:

- **Example 1:** "*Machine Learning*" overlaps with "*Data Mining*" and "*Neural Networks & AI*."
- **Example 2:** "*Cybersecurity*" forms a distinct cluster, reflecting its specialized focus.
- Foundational topics such as "*Optimization*" and "*Graph Theory*" connect multiple clusters, showcasing their interdisciplinary importance.

## Challenges and Progress

While clustering and t-SNE reveal valuable insights, the dense visualization highlights the complexity of academic research networks. This analysis lays the groundwork for developing a **scholarly metric** that integrates these insights to comprehensively measure influence.

Next, we'll explore **PageRank** to rank influential papers. By incorporating **weighted PageRank**, we account for both citation quality and venue prestige, providing a more nuanced measure of scholarly impact.

# Weighted PageRank: Rethinking Influence in Academic Networks

**Objective:** Rank influential papers by integrating **venue prestige** into PageRank.

- Citation Network:
  - Nodes: Papers
  - Edges: Citations
- Venue Prestige:
  - Weight citations by citing venue's prestige.
- Weighted PageRank:
  - Higher scores for citations from prestigious venues.

## Insights:

- **50% overlap in top papers.**
- **PageRank favors methodological contributions.**
- **Citations reflect interdisciplinary impact.**

Traditional metrics, such as citation counts, often overlook nuances like the quality of the publication venue. **Weighted PageRank** addresses this by assigning greater importance to citations from *prestigious venues*, providing a more nuanced measure of academic influence.

## Methodology

- Constructed a **citation network** connecting papers and their references.
- Incorporated **venue prestige** into PageRank by weighting edges based on the prestige of the citing venues.

## Insights

**Overlap:** 50% of the Top 10 papers are common across PageRank and citation-based rankings.  
**Differences:**

- **PageRank:** Highlights foundational works like "*Probabilistic Reasoning in Intelligent Systems*".
- **Citation Counts:** Emphasize interdisciplinary impact, such as "*Bowling Alone*" from social sciences.

## Takeaway

**Weighted PageRank** effectively combines *venue quality* with *citation counts*, offering a more comprehensive measure of scholarly influence.

Next, we'll examine **collaboration patterns** using *association rule mining* to further deepen this analysis.

# Uncovering Collaboration Patterns: Association Rule Mining

## Objective:

- Identify frequent collaboration patterns among prolific authors.
- Use association rules to uncover strong co-authorship relationships.

## Key Steps:

- Filtered authors with **≥50 papers** to focus on prolific collaborators.
- Applied **Apriori algorithm** to generate frequent itemsets and association rules.
- Measured relationships using **support, confidence, and lift**.

## Top Collaboration Pairs

- Irith Pomeranz & Sudhakar M. Reddy**
  - Highest support (0.000519)
  - Focus on Computer Engineering/Testing
- Makoto Takizawa & Tomoya Enokido**
  - Support: 0.000443
  - Collaboration in distributed systems and computer networks
  - Strong academic partnership with shared institutional history
- Fatos Xhafa & Leonard Barolli**
  - Support: 0.000428
  - International collaboration in Computer Networks

## Strongest Associations

- Tetsuya Oda → Leonard Barolli**
  - Confidence: 1.0
  - Lift: 1046.41
  - Indicates strong research group connection
- Nora Cuppens-Boulahia → Frédéric Cuppens**
  - Confidence: 1.0
  - Lift: 4179.86
  - Family/research partnership
- Mayank Vatsa → Richa Singh**
  - Confidence: 1.0
  - Lift: 5327.85
  - Strong research partnership

Collaboration is a fundamental aspect of academic influence. Our goal was to:

- Identify **frequent co-authorship patterns**.
- Uncover **strong associations** between authors using *association rule mining*.

## Methodology

- Focused on **prolific authors** with 50 or more publications.
- Applied the **Apriori algorithm** to extract:
  - Frequent itemsets:** Groups of authors who collaborated frequently.
  - Association rules:** Relationships between authors, quantified using *support, confidence, and lift*.

## Frequent Collaborators:

- Irith Pomeranz & Sudhakar M. Reddy:** Top collaboration in Computer Engineering.
- Makoto Takizawa & Tomoya Enokido:** Focused on Distributed Systems and Computer Networks.

## Strongest Association Rules:

- Tetsuya Oda → Leonard Barolli:** Perfect confidence (1.0) and high lift (1046.41), reflecting a strong research group connection.
- Nora Cuppens-Boulahia → Frédéric Cuppens:** Strong familial and research partnership (Lift: 4179.86).

Collaboration patterns highlight the structure of research groups, ranging from *institutional partnerships* to *cross-institutional dynamics*. These insights provide a critical layer in understanding **scholarly influence**.

While collaboration is vital, quantifying overall influence requires combining these patterns with **Weighted PageRank** and other metrics. Next, we'll discuss how these analyses integrate into our comprehensive scholarly metric.

# Scholarly Influence: Venue Prestige

- Recognized venues
- Reviewed papers
- Cited more often
- More papers, more importance
- Field of study
- Upcoming venues

venue	year	n_citations
International Journal of Computer Vision	2004	42508
conference on computer supported cooperative work	2000	34288
ACM Transactions on Intelligent Systems and Technology	2011	33016
neural information processing systems	1999	29285
Machine Learning	2001	28679
Management Information Systems Quarterly	1989	27068
Machine Learning	1995	26114
systems man and cybernetics	1985	25835

## Venue Prestige and Scholarly Influence

### Role of Venue Prestige

- Venues like ACM and IEEE ensure rigorous peer review and publish highly cited, impactful papers.
- They advance research fields through consistent, high-quality contributions.

### Citations as Indicators

- Longstanding venues (e.g., DBLP) with high citation counts demonstrate lasting impact.
- Citations reflect a venue's academic influence.

### Challenges with New Venues

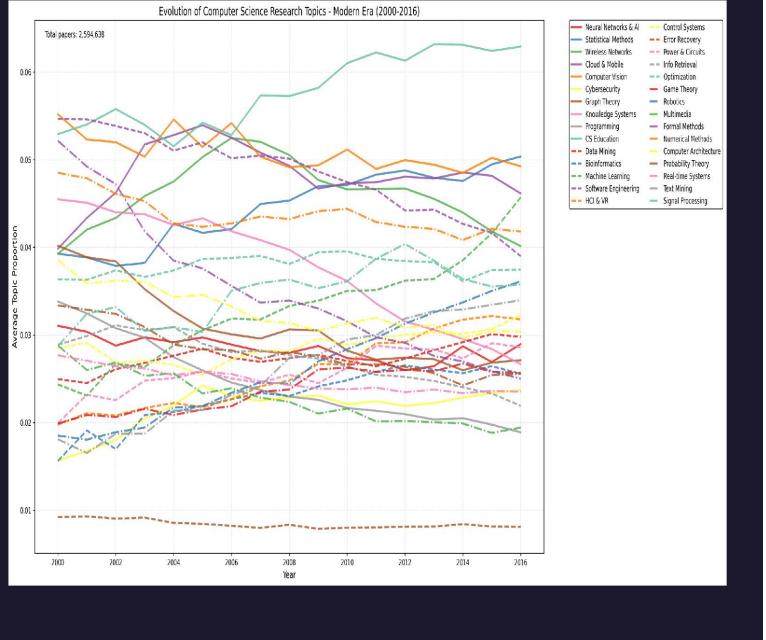
- New venues may publish strong work but lack recognition due to limited history.
- Perception challenges hinder their inclusion in traditional metrics.

### Incorporating New Venues

- Metrics should account for peer-review rigor, community engagement, and innovation.
- Inclusive systems are needed to recognize contributions across venues.

# Temporal Patterns

- More papers in growth and modern era
- Citations are more recent.
- A gauge for effectiveness of paper and author
- Interesting areas of study
- No longer prestigious
- Trends change
- Delay in citation and peaking



## Temporal Analysis of Scholarly Influence

### Evolving Citation Patterns

- Citation behavior has changed significantly since the 1980s, with papers in recent decades receiving more citations due to increased accessibility and expanding research fields.
- Dynamic fields like machine learning show rapid citation growth for impactful papers when the domain advances.

### Delayed and Long-Term Impact

- Some papers gain influence only after extended periods, as seen in historical cases like Huygens' pendulum research.
- Temporal impact highlights the need to evaluate scholarly contributions over both short and long terms.

### Shifting Research Trends

- Publication delays and shifting trends can affect citation patterns and venue prestige.
- Research domains may become obsolete, influencing the impact of papers and authors over time.

# Centrality Analysis

- Degree Centrality Metrics
- Betweenness Centrality Analysis
- Eigenvector Centrality Analysis

## Top Authors by Eigenvector Centrality

### 1. Completely Different Group:

- Led by Arthur W.Toga (0.053419)
- Much lower paper counts (22-307)
- Lower citation counts compared to other metrics

### 2. Notable Characteristics:

- Focused in medical/neuroscience field
- More specialized research community

## Top Authors by Degree Centrality

### 1. Wei Wang

- Highest degree centrality (0.002205)
- 75,765 citations
- 2,518 papers
- Indicates extensive collaboration network

### 2. Notable Pattern

- Top authors have high paper counts (>1,000)
- Citations range from 29,804 to 75,765
- Most authors have Chinese names, suggesting strong representation in the field

## Top Authors by Betweenness Centrality

### 1. Completely Different Group:

- Wei Wang (again leading)
- Highest betweenness (0.015474)
- Suggests role as key connector between research communities

### 2. Notable Characteristics:

- Focused in medical/neuroscience field
- More specialized research community

## Centrality Measures in Research Networks

### Degree Centrality

- Focuses on direct collaborations, highlighting prolific authors with extensive networks.
- Reveals the volume and breadth of an author's immediate research connections.

### Betweenness Centrality

- Identifies researchers who connect distinct research communities.
- Often correlates strongly with degree centrality, emphasizing authors bridging collaborative groups.

### Eigenvector Centrality

- Measures influence based on the prominence of an author's collaborators.
- Effectively identifies leaders within specialized research communities.

### Prioritization of Eigenvector Centrality

- Balances field-specific publication patterns, offering fair evaluation across disciplines.
- Analogous to comparing athletes from different sports, enabling equitable comparison of researchers in diverse fields.

# The new metric (APR-Index)

$$APR = 0.5 * P_{\text{coauthor\_network}}(\text{author}) + 0.5 * \sum_{\text{publications}} P_{\text{citation\_network}}(\text{publication})$$

- In good faith of the analysis.
- General weighted approach
- More in-depth meaning than a simple H-Index
- Gradient than integral
- Base implementation
- Simplify values

Author	H-Index	Our metric
David G. Lowe	50	0.549801
Chih-Jen Lin	50	0.530669
Gerard Salton	50	0.419323
Lotfi A. Zadeh	50	0.40901
David L. Donoho	50	0.3947

## The Academic Page Rank Index (APR-Index)

### Overview of the APR-Index

- Builds on gradient learning to evaluate author impact using a generalized weighted approach.
- Provides finer granularity than the discrete H-index, distinguishing authors with identical H-indices.

### Comprehensive Assessment

- Incorporates multiple features for a more holistic evaluation of academic influence.
- Aims to deliver a nuanced metric for assessing scholarly contributions.

### Current Limitations

- Requires further refinement, including extensive testing and regression analysis to validate effectiveness.
- Output values need to be made more interpretable for practical use in academic evaluations.



# Impact and Expected Outcomes

- Door to more insights
- Gauge the influence
- Deeper meaning
- Extensibility
- Ranking system
- Addition of authors and papers
- Removal of authors and papers

## Unified Analysis Approach

### Enhanced Insights

- Combines multiple metrics to provide a more accurate and granular ranking of academic impact.
- Offers unprecedented insights into author and paper influence in academia.

### Extensibility and Future Enhancements

- Designed as an extensible framework, allowing for the addition of new metrics as needed.
- Adapts to evolving academic evaluation requirements.

### Computational Trade-offs

- Updating rankings with additional data may require more computational resources.
- Modern computing power minimizes these costs, making the trade-off negligible.
- The comprehensive evaluation benefits significantly outweigh these minor costs.



# Benefits for researchers and institutions

- Similar interests
- Collaboration
- Value from the beginning
- Promote and support research areas
- See a better picture of authors and their papers

## Applications of the Metric

### Identifying Potential Collaborators

- Helps researchers find potential collaborators with shared interests and research areas.
- The weighted approach provides meaningful evaluations from a paper's initial publication.

### Institutional Use

- Institutions can identify underexplored research areas and assess authors and their work objectively.
- Facilitates informed decisions about resource allocation and support.

# Challenges and Future Work

- Just the beginning
- Address drawbacks
- Usability and reliability
- Extend APR to its fullest capability
- Adoptability
- Verifiability



## Multidisciplinary Approach to Academic Evaluation

### Future Enhancements

- Our approach is ambitious, aiming to address limitations and improve the metric's reliability through iterative enhancements.
- As the APR-Index evolves, it will more comprehensively reflect all aspects of academic impact.

### Potential for Standardization

- The goal is for the APR-Index to become a standardized tool for academic evaluation across institutions worldwide.
- Its broader adoption will ensure more consistent and meaningful assessments of academic impact.

### Validation and Practical Value

- Ultimate validation would involve predicting academic excellence, such as identifying future Fields Medal winners.
- The primary goal is to provide a practical tool that aids researchers and institutions in decision-making regarding collaboration and resource allocation.

