



**School of Information Technologies and  
Engineering**

**CSCI4700 Data Mining**

**2025 Spring**

# **Graph Mining**

## **Citation Network Analysis**

---

**TEAM 3**

**LAMAN PANAKHOVA BSCS 2026**

**MEHRIBAN ALIYEVA BSCS 2025**

**EMIL NIYAZOV BSCS 2025**

# Agenda

---

**Problem Statement**

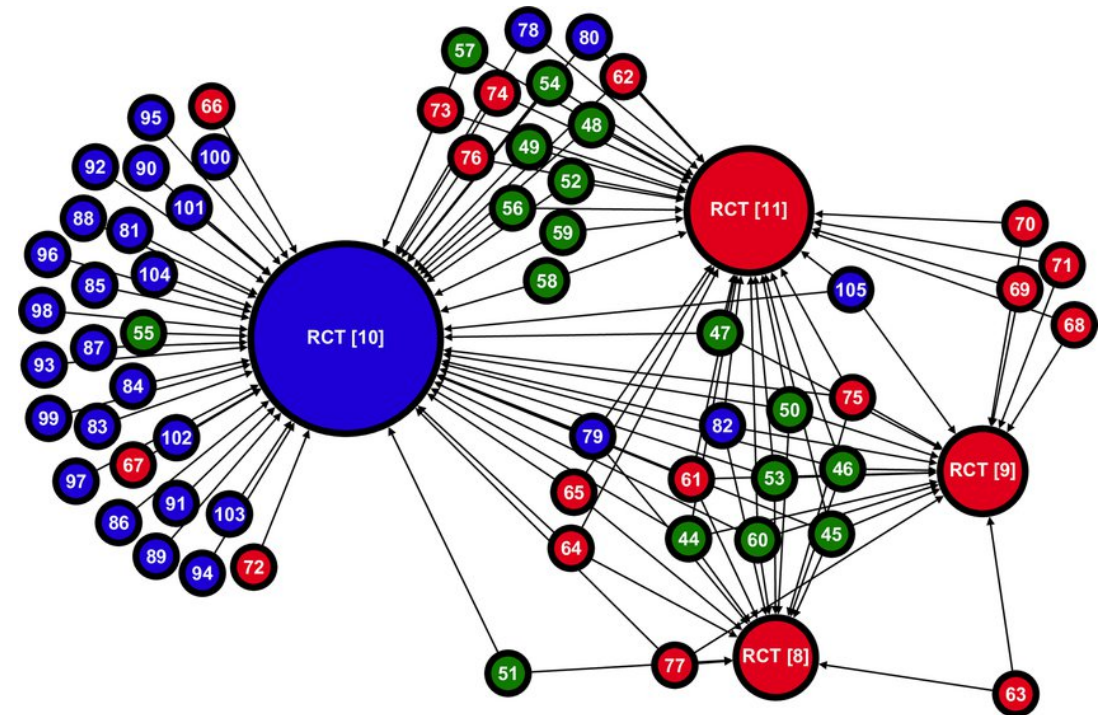
**Data Collection and EDA**

**Mining Techniques & Algorithms Applied**

**3D Visualizations**

**Results & Experiments**

**Considerations & Future Improvements**



Project Title	Graph Mining - Citation Network Analysis
Value	Identify influential research papers and uncover structural patterns in academic research, which can support research discovery and trend analysis in scientific communities
Problem type	Graph Mining - Citation Network Analysis
Original and final data shape	Scraped – 100x9; Original – 32631 nodes (each node that was collected cites many other nodes) & 39826 edges; Final – 1000 nodes & 3257 edges and 300 nodes 966 edges
Preprocessing steps	Removing entries with missing DOIs or citation data, filtering out invalid citation links, constructing a directed citation graph, removing isolated nodes, and selecting the top 300 most connected papers
Algorithm(s)	Graph Mining Techniques - PageRank, Betweenness Centrality, Closeness Centrality, Degree Centrality, and Eccentricity, along with Strongly Connected Components (SCC) detection
Results	Highly influential papers through centrality rankings, a strongly connected citation core among the top 300 papers, and patterns such as skewed degree distribution and strong correlation between closeness and betweenness metrics
Interpretation	Our results align with prior research studies, showing high correlation with real-world citation counts and matching human-curated top papers, validating the effectiveness of our approach by experiments.

# Contributions

---

Team Member	Contribution
<b>Laman Panakhova BSCS 2026</b>	<b>33%</b>
<b>Mehriban Aliyeva BSCS 2025</b>	<b>33%</b>
<b>Emil Niyazov BSCS 2025</b>	<b>33%</b>

# Problem Statement

---

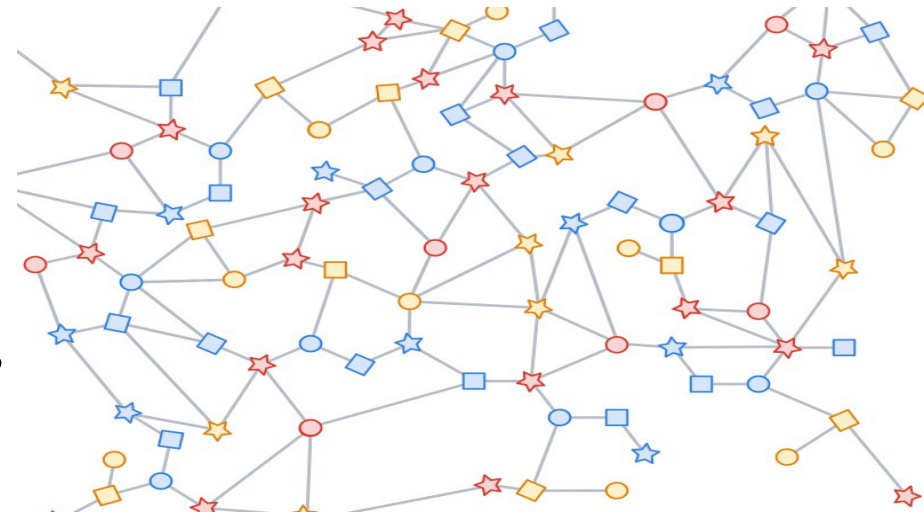
## Project Overview: Analyzing Citation Networks in Academic Research

**Goal:** Explore citation networks to understand the influence of academic papers, the formation of research communities, and the evolution of scientific knowledge.

**Method:** Treat each paper as a node and citations as directed edges to create a graph-based structure.

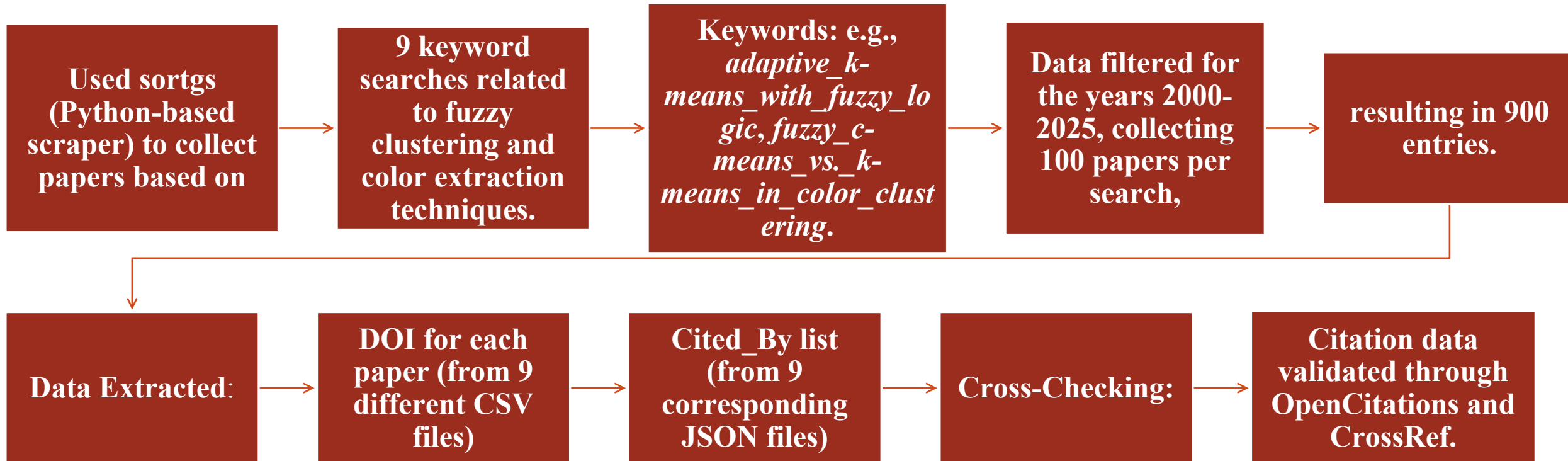
### Focus:

- Identify influential papers
- Uncover hidden relationships in research
- Discover patterns in how scientific knowledge evolves



# Data Collection

---



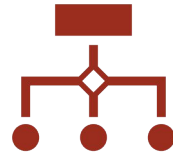
# EDA & Preprocessing

---



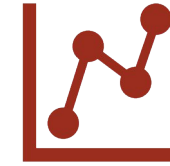
## Data Cleaning:

Removed entries with missing DOIs or citation lists.  
Filtered out broken citations labeled as "Not Found."



## Graph Construction:

Built a directed citation graph ( $A \rightarrow B$ , where A cites B).  
Removed isolated nodes (papers with no citations).  
Filtered to retain the top 300 most connected nodes for meaningful analysis.



## Exploratory Data Analysis (EDA):

**Descriptive Stats:** Computed basic statistics for centrality metrics (PageRank, degree).

**Data Integrity Checks:** Verified data types, null entries, and missing values.

**Visualizations:**

- Top 10 Most Cited Papers (bar chart)
- Correlation Heatmap (centrality metrics)
  - Degree Distribution Histogram

**Graph Metrics:** Analyzed node/edge count, average degree, clustering coefficient, and ensured strong connectivity.

<b>Mining Techniques:</b>	Built directed citation graph (papers as nodes, citations as edges).	Retained top 1000 most connected nodes (total degree filter).	Applied PageRank for influence ranking (captures recursive prestige).	Calculated in-degree, betweenness, closeness centrality.	Used Louvain method for unsupervised community detection.
Incorporated temporal analysis to track citation trends over time.	<b>Design Decisions:</b>				
		Focused on citation-connected nodes for relevance.	Analyzed top 300 nodes for computational efficiency and visual clarity.	Chose PageRank over HITS (better for non-bipartite citation networks).	Used NetworkX for graph analysis, Matplotlib/Seaborn for visualizations.
<b>Evaluation</b>	Manual validation: Verified top 10 PageRank papers' citations via Google Scholar.	Spearman correlation (PageRank vs. in-degree) showed strong correlation ( $\rho > 0.8$ ).	Evaluated community quality via modularity score.	Visualized centrality metrics and community structure using various plots.	<b>Limitations:</b>
	Excluded citations with missing or malformed DOIs.	Cold start problem: Newly published papers may be undervalued by PageRank.	Disconnected components filtered out niche or emerging domains.		
	Note: All analyses were conducted using Python 3.8 and NetworkX 2.6.				



# Results & Experiments

---

## PageRank vs. In-Degree:

---

Compared **PageRank** scores with **in-degree** centrality.

---

**Scatter plot** showed strong positive correlation (Spearman  $\rho \approx 0.85$ ), highlighting PageRank's added context of citation influence.

---

---

## Top 10 Influential Papers:

---

Visualized top 10 papers by **PageRank**.

---

Revealed **citation hubs** that had strategic citation patterns despite modest in-degrees.

---

---

## Community Detection:

---

Applied **Louvain method** for community detection.

---

**Colored graph** showed distinct research subfields (e.g., computer vision, NLP, data mining).

---

---

## Temporal Analysis:

---

Plotted **citation trends over time**.

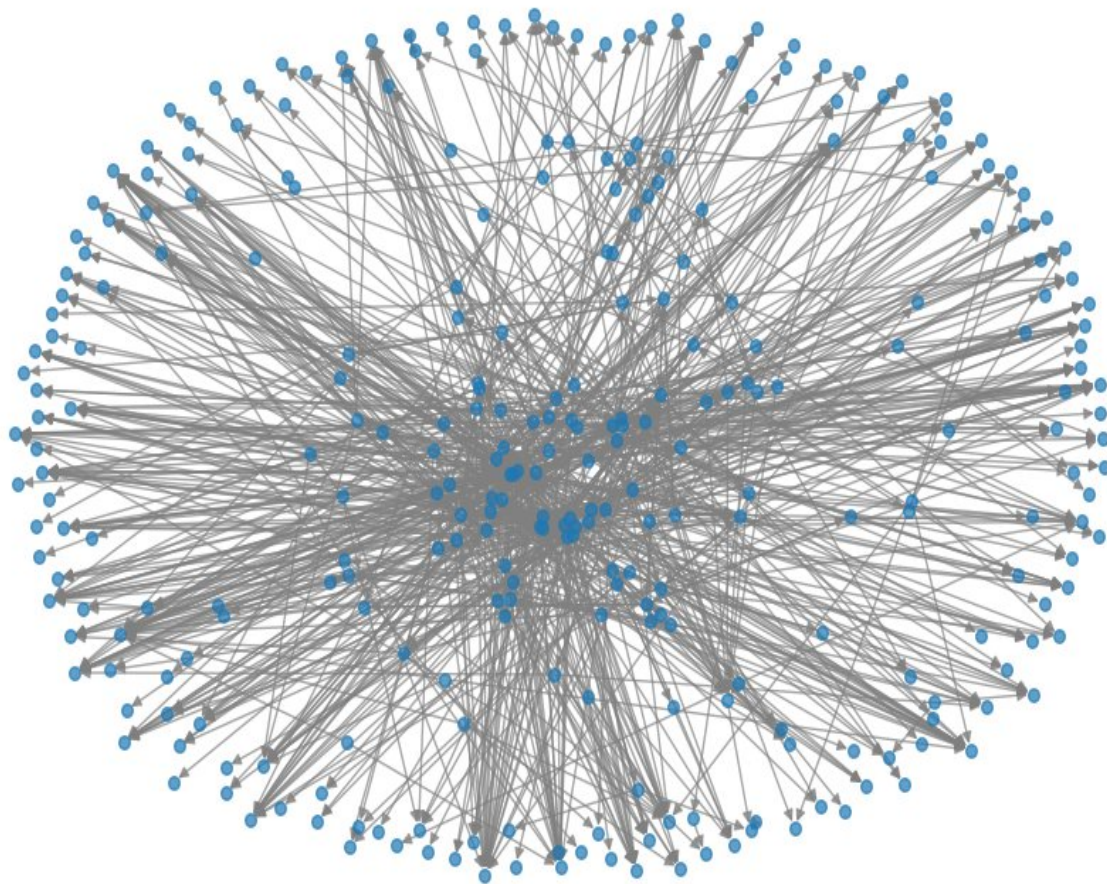
---

Identified citation bursts around **seminal papers** and steady growth in **emerging research areas**.

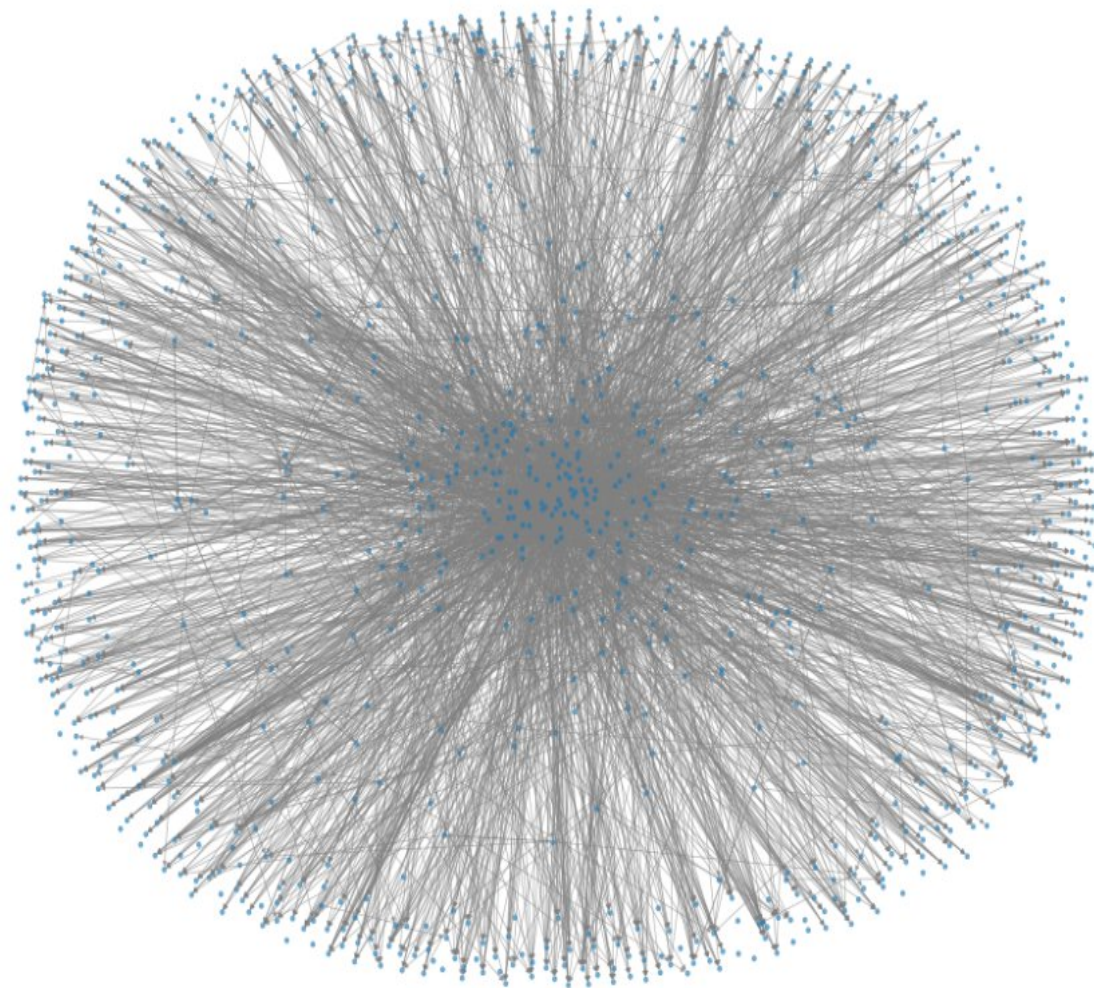
---

✓ Graph loaded and filtered: 300 nodes, 966 edges

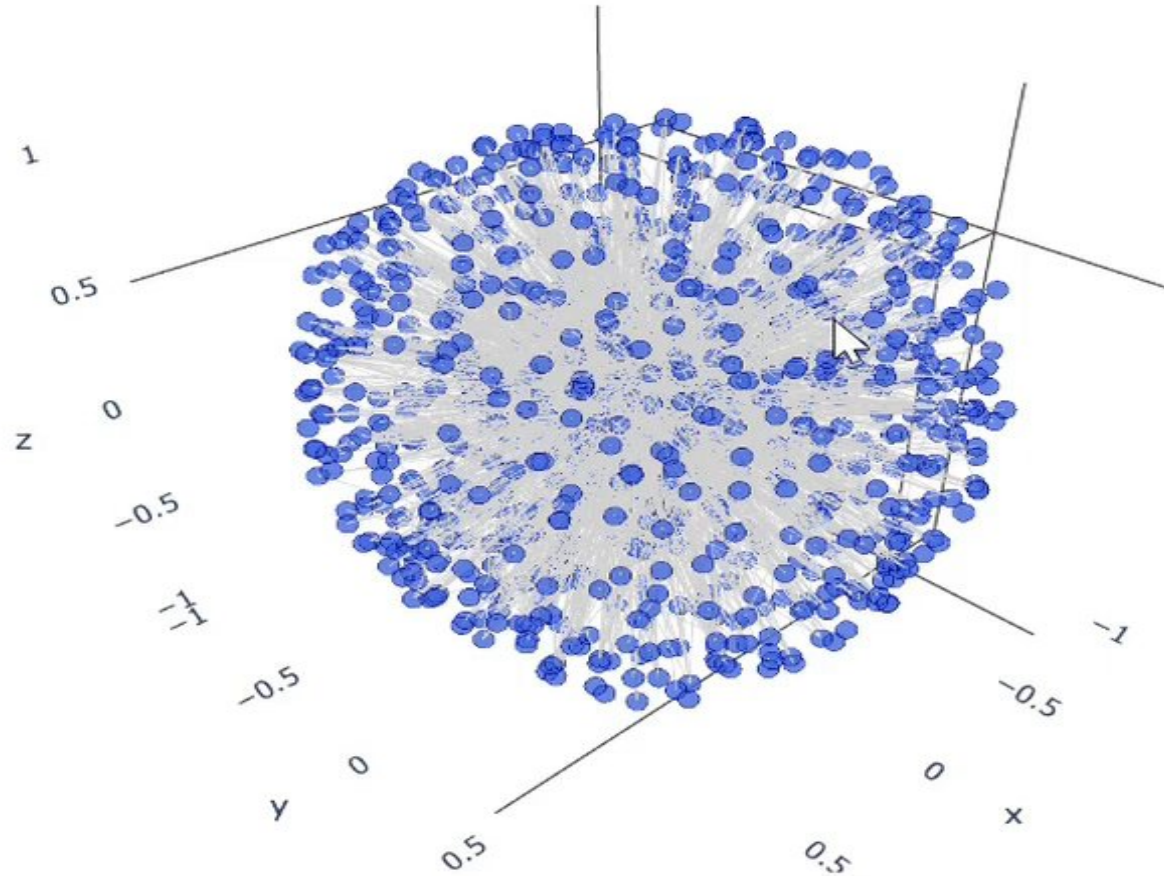
Citation Graph (Top 300 Nodes)



Citation Graph: Top 1000 Most Connected Papers



### 3D Citation Graph Visualization (Top 1000 Nodes)

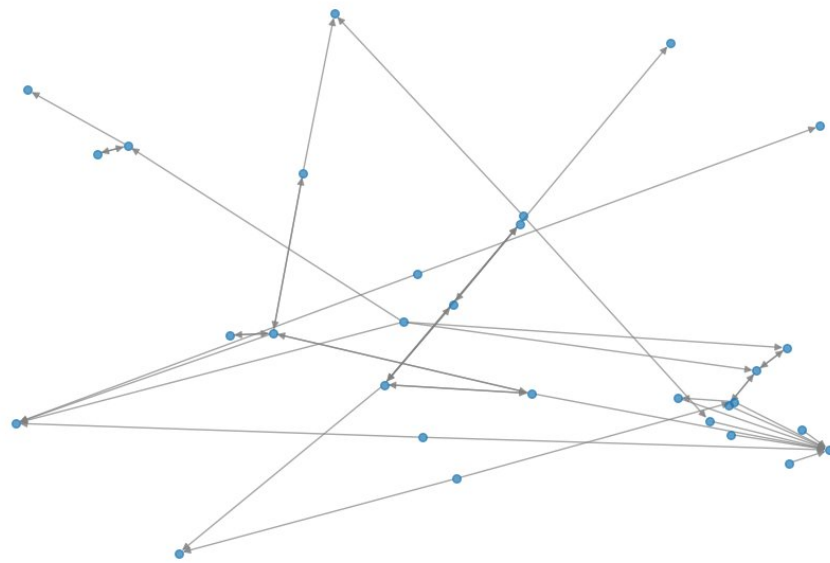


clideo.com

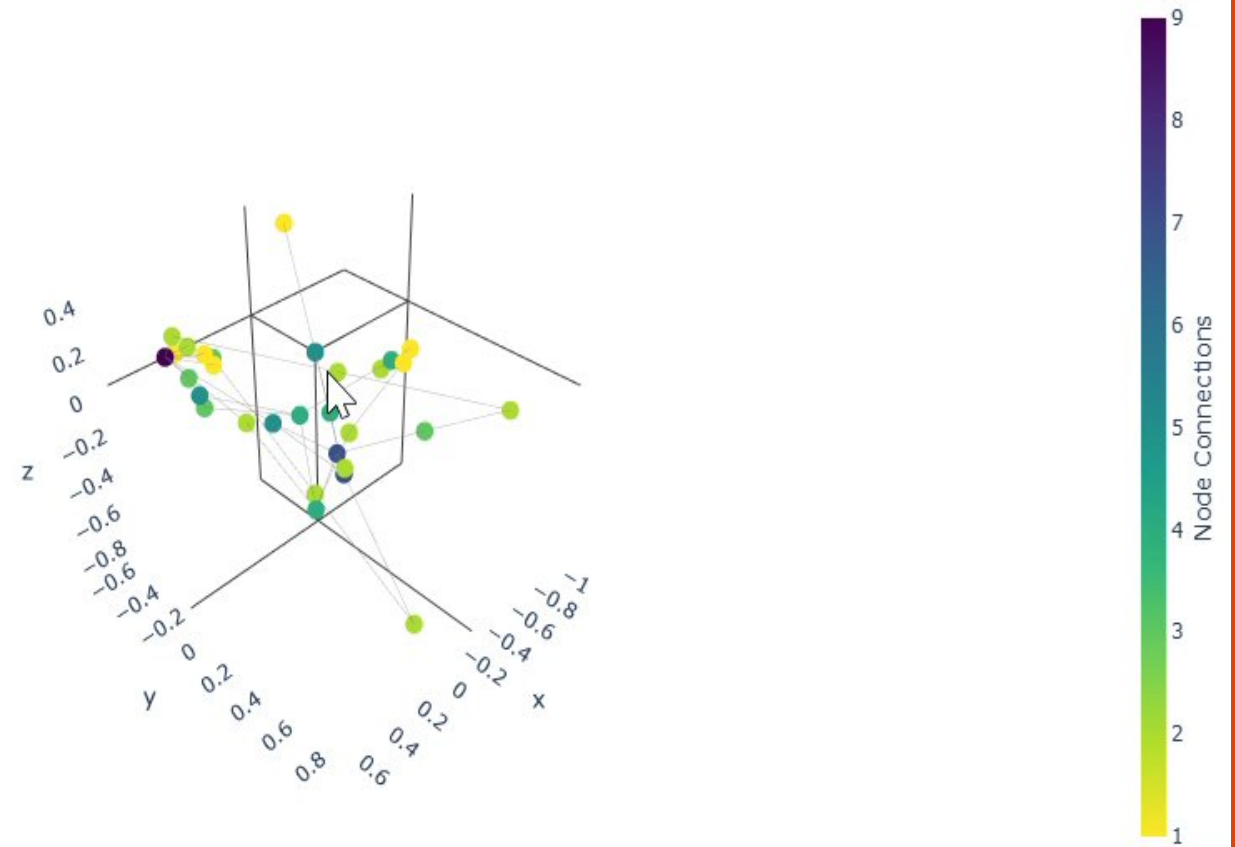


✓ Graph loaded and filtered: 30 nodes, 46 edges

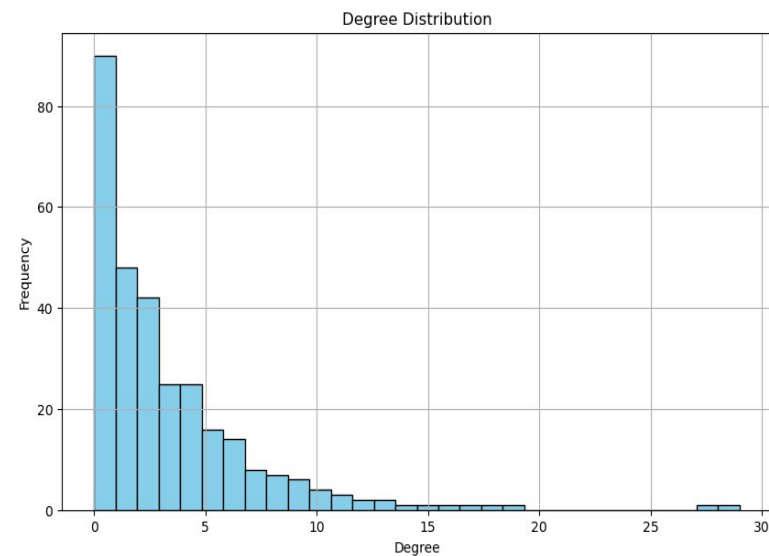
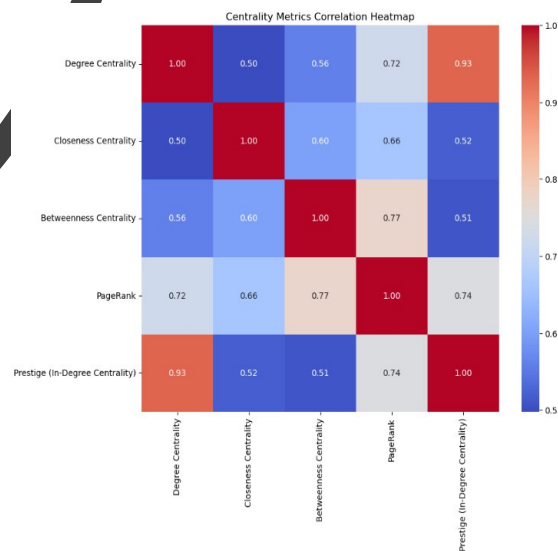
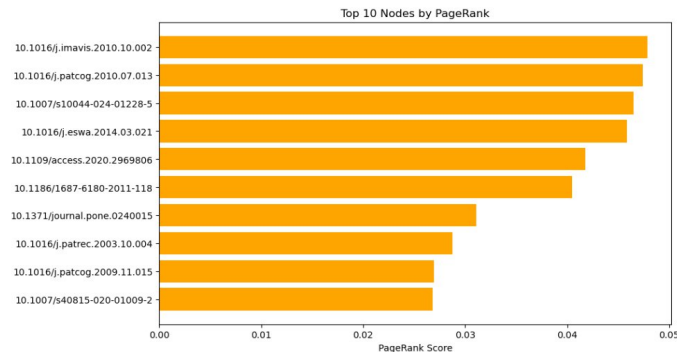
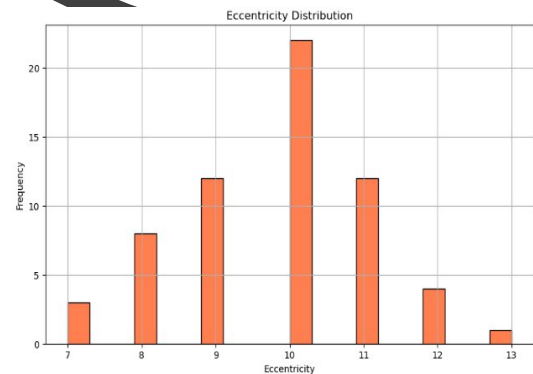
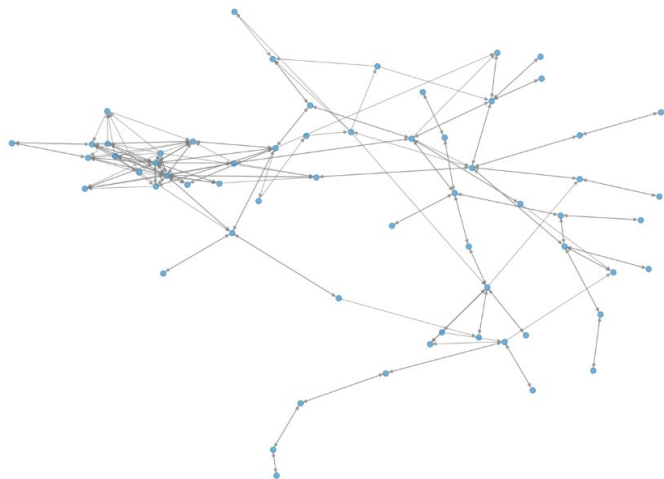
Citation Graph (Top 30 Nodes)



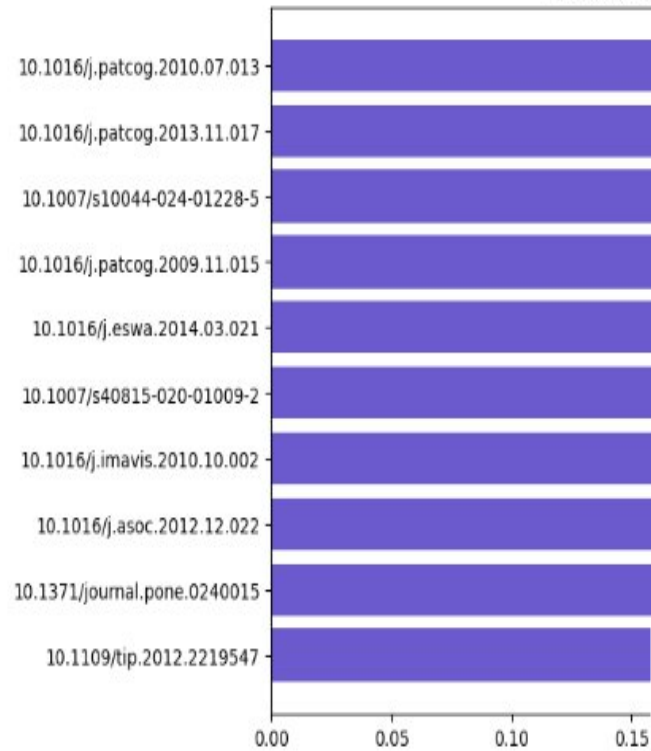
3D Force-Directed Citation Graph (Top 30 nodes)



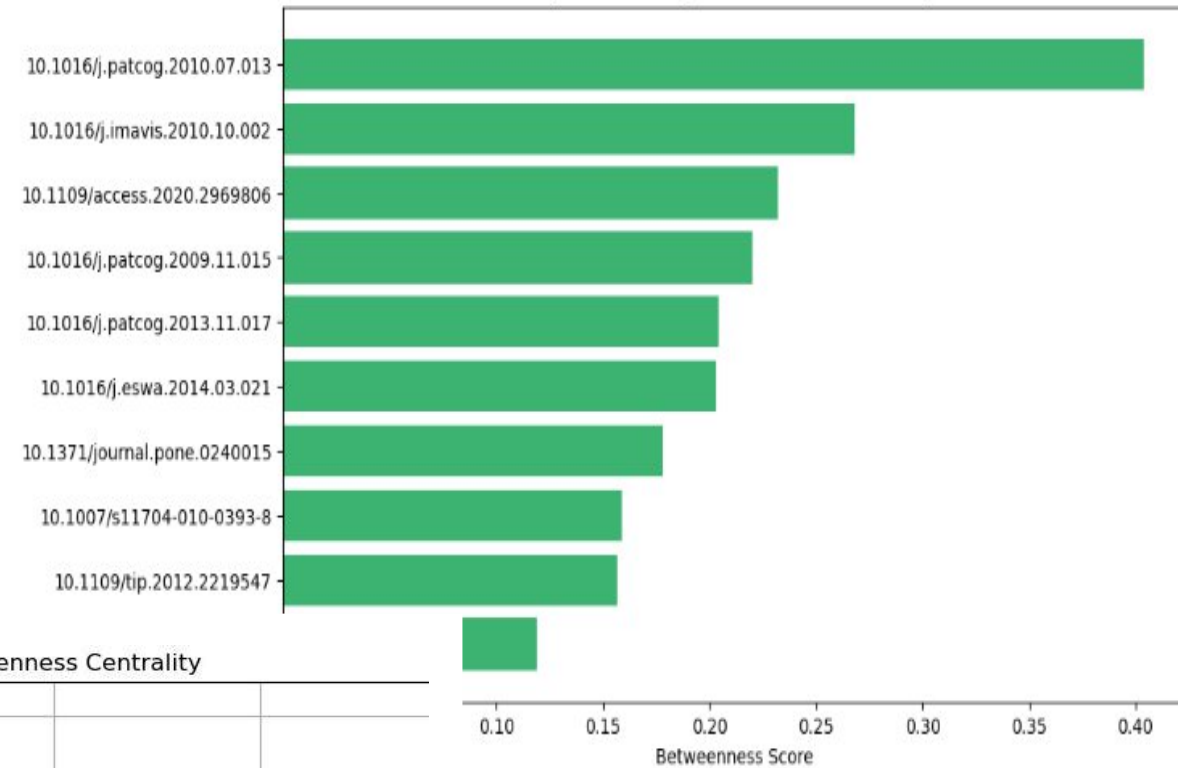
Citation Graph: Most Connected Nodes



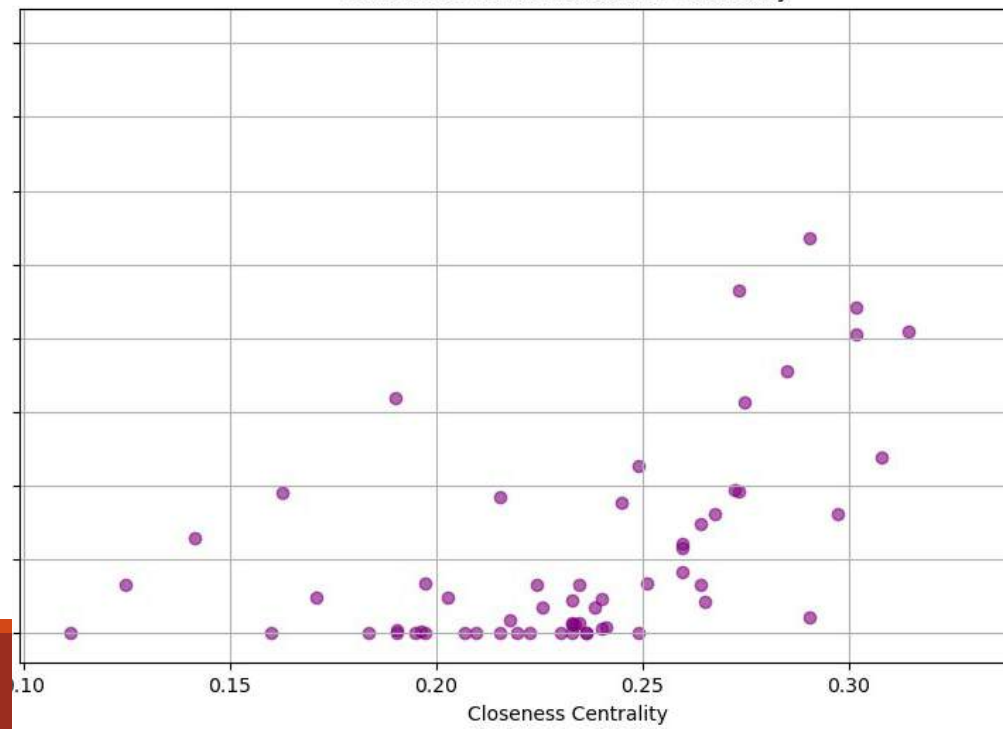
Top 10 Nodes by Closeness Centrality



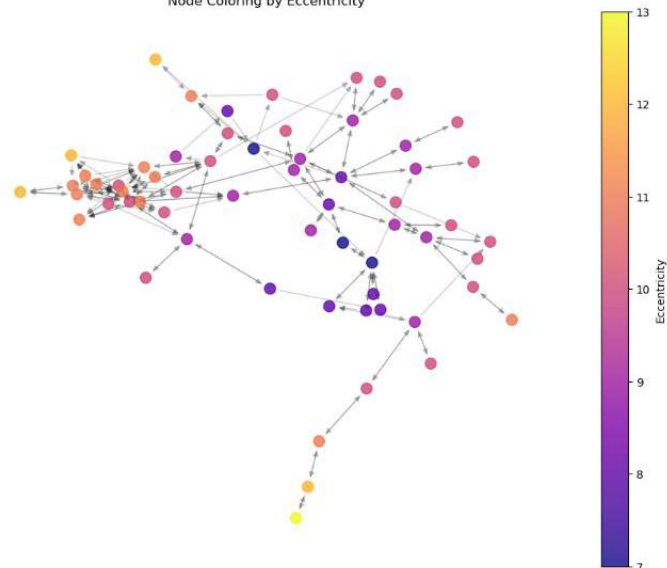
Top 10 Nodes by Betweenness Centrality



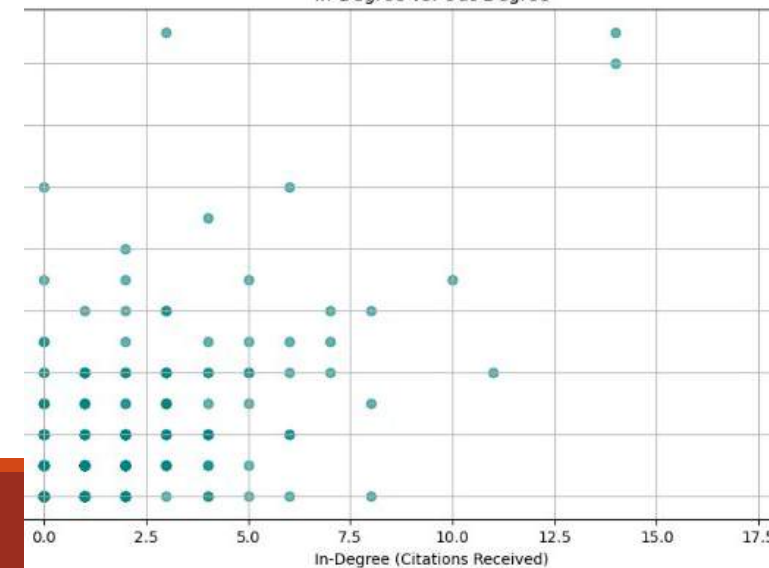
Closeness vs. Betweenness Centrality



Node Coloring by Eccentricity



In-Degree vs. Out-Degree



# Considerations & Future Improvements

---



## **PageRank vs. In-Degree:**

PageRank identified influential papers beyond just citation count, ranking those cited by other impactful works more highly.



## **Degree Distribution:**

Network exhibited a heavy-tailed (power-law) distribution, where a small number of papers dominate citations.



## **Community Detection:**

Louvain clustering revealed distinct research areas (e.g., computer vision, NLP) within the citation network.



## **Temporal Trends:**

Citation bursts around major papers and steady growth in emerging fields.



## **Data Cleaning Impact:**

Removing isolated nodes and malformed citations improved clarity and interpretability.



## **Limitations:**

Cold start problem: Newer papers were undervalued by PageRank due to fewer citations.



## **Conclusion:**

Citation graph analysis offers valuable insights into knowledge evolution, community formation, and research trends.

# THANK YOU FOR YOUR ATTENTION!

---

QUESTIONS?

