# High-energy physics theory

Is the theory based on only a few articles?

*By Manglai Manglai (20230012 )*

*Nov. 2, 2024, 5:04 a.m  Cy time*

## Introduction

This article explores a dataset from arXiv (arXiv.org), an open-access archive for nearly 2.4 million scholarly articles across various disciplines, including physics, mathematics, computer science, and quantitative biology. This study focuses on the **High Energy Physics Theory (HEP-TH)** section, covering papers from **January 1993 to April 2003**. The dataset includes **27,770 papers** and **352,807 citations**, representing a near-complete record of HEP-TH's early contributions.

We analyze this dataset as a **citation network**, where nodes represent papers and directed edges signify citations from one paper to another. By examining this network, we uncover patterns of influence, community structures, and citation behaviors that provide insights into the collaborative nature and research trends within the HEP-TH domain.

## Chapter 1 Loading the Data

In this section, we load and prepare the dataset using NetworkX, a Python library designed for the analysis of complex networks.

### Loading Citation Data:

The dataset consists of a list of citations where each entry connects a citing paper to a cited paper. Using NetworkX, we represent this data as a directed graph where each node is a paper, and each directed edge signifies a citation from one paper to another.

Number of papers (nodes): 27770

Number of citations (edges): 352807

## Adding Metadata:

We add additional data to each node, such as the submission date, to allow for additional analysis.

```python
import networkx as nx

import pandas as pd


# Function to load the graph from the dataset file

def load_graph(file_path):

    G = nx.DiGraph()  # Directed graph

    with open(file_path, 'r') as f:

        for line in f:

            try:

                node_from, node_to = map(int, line.strip().split())

                G.add_edge(node_from, node_to)

            except ValueError:

                # Handle any improperly formatted lines

                continue

    return G


# Path to the citation data file

file_path = 'C:/Users/mangl/Desktop/assigement/Cit-HepTh.txt'
```

# Chapter 2: Analyzing the Data

In this chapter, we analyze network metrics and community structures to reveal patterns within the HEP-TH dataset.

## 2.1 Degree Centrality Analysis

We begin by analyzing **degree centrality**, which helps us identify the most cited papers (in-degree) and the most citing papers (out-degree).

**In-degree**: Indicates the number of times a paper is cited by others.

| Top 10 most cited papers: | |
| --- | --- |
| Paper 9711200: Cited | 2414 papers |
| Paper 9802150: Cited | 1775 papers |
| Paper 9802109: Cited | 1641 papers |
| Paper 9407087: Cited | 1299 papers |
| Paper 9610043: Cited | 1199 papers |
| Paper 9510017: Cited | 1155 papers |
| Paper 9908142: Cited | 1144 papers |
| Paper 9503124: Cited | 1114 papers |
| Paper 9906064: Cited | 1032 papers |
| Paper 9408099: Cited | 1006 papers |

**Out-degree**: Indicates the number of citations a paper makes to others.

| Top 10 most cited papers: | |
| --- | --- |
| Paper 9905111: Cited | 562 papers |
| Paper 9710046: Cited | 359 papers |
| Paper 110055: Cited | 302 papers |
| Paper 210157: Cited | 289 papers |
| Paper 101126: Cited | 274 papers |
| Paper 7170: Cited | 263 papers |
| Paper 204089: Cited | 246 papers |
| Paper 201253: Cited | 226 papers |
| Paper 9809039: Cited | 216 papers |
| Paper 9802067: Cited | 214 papers |

```
 in_degrees = citation_graph.in_degree()  # Papers cited by others

out_degrees = citation_graph.out_degree()  # Papers citing others


# Get the top 10 papers that are cited the most

top_cited_papers = sorted(in_degrees, key=lambda x: x[1], reverse=True)[:10]

print("Top 10 most cited papers:")

for paper, cites in top_cited_papers:

    print(f"Paper {paper}: Cited by {cites} papers")
```

In-degree centrality (most cited papers): Papers with high in-degree centrality are influential.

| Total_in_degree | |
|---|---|
| 9711200 | 0.086931 |
| 9802150 | 0.06392 |
| 9802109 | 0.059095 |
| 9407087 | 0.046779 |
| 9610043 | 0.043178 |
| 9510017 | 0.041593 |
| 9908142 | 0.041197 |
| 9503124 | 0.040117 |
| 9906064 | 0.037164 |
| 9408099 | 0.036227 |

Out-degree centrality (most citing papers)

| Total_out_degree | |
| --- | --- |
| 9905111 | 0.020238 |
| 9710046 | 0.012928 |
| 110055 | 0.010875 |
| 210157 | 0.010407 |
| 101126 | 0.009867 |
| 7170 | 0.009471 |
| 204089 | 0.008859 |
| 201253 | 0.008139 |
| 9809039 | 0.007778 |
| 9802067 | 0.007706 |
| | |

Total degree centrality (both in and out citing papers)

| Total_out_degree | |
| --- | --- |
| 9711200 | 0.088876085 |
| 9802150 | 0.064712449 |
| 9802109 | 0.05952681 |
| 9905111 | 0.049299579 |
| 9407087 | 0.047102885 |
| 9908142 | 0.043897872 |
| 9610043 | 0.04386186 |
| 9510017 | 0.041953257 |
| 9503124 | 0.040476791 |
| 9906064 | 0.037379812 |

```
# In-degree centrality (most cited papers)

in_degree_centrality = nx.in_degree_centrality(citation_graph)

top_in_degree = sorted(in_degree_centrality.items(), key=lambda x: x[1], reverse=True)[:10]


# Out-degree centrality (most citing papers)

out_degree_centrality = nx.out_degree_centrality(citation_graph)

top_out_degree = sorted(out_degree_centrality.items(), key=lambda x: x[1], reverse=True)[:10]




degree_centrality=nx.degree_centrality(citation_graph)
```

## 2.2 PageRank Analysis

We apply the **PageRank algorithm** to identify influential papers based on both the number and quality of their citations.

```
pagerank = nx.pagerank(citation_graph)

top_pagerank = sorted(pagerank.items(), key=lambda x: x[1], reverse=True)[:10]
```

| Page Rank | |
|---|---|
| 9407087 | 0.006239 |
| 9503124 | 0.004633 |
| 9510017 | 0.004385 |
| 9402044 | 0.003935 |
| 9711200 | 0.00341 |
| 9410167 | 0.003407 |

| | |
|---|---|
| 9408099 | 0.00319 |
| 9207016 | 0.003114 |
| 9402002 | 0.002962 |
| 9610043 | 0.002753 |

Papers with high PageRank scores, such as Paper 9407087 and Paper 9503124, serve as influential hubs within the network, connecting different research areas and helping disseminate foundational ideas.

## 2.3 Community Detection

We use the Louvain method, a popular algorithm for community detection in complex networks, to identify groups of nodes that are more densely connected to each other than to the rest of the network. These groups are referred to as communities. In the context of citation networks, these communities can represent subfields or research clusters within a broader research area. By detecting communities within the network, this analysis helps identify clusters of related research topics or research groups that are more interconnected.

Community 12: 2413 papers

Community 1: 2065 papers

Community 7: 1993 papers

Community 0: 1946 papers

Community 14: 1887 papers

Community 3: 1795 papers

Community 9: 1613 papers

Community 26: 1454 papers

Community 2: 1407 papers

Community 25: 1390 papers

```
import community.community_louvain as community_louvain


# Apply Louvain community detection

partition = community_louvain.best_partition(citation_graph.to_undirected())


# Count the number of nodes in each community

community_sizes = {}

for com in set(partition.values()):

    community_sizes[com] = list(partition.values()).count(com)


# Print top communities by size

top_communities = sorted(community_sizes.items(), key=lambda x: x[1], reverse=True)[:5]

print("Top 5 communities (most papers):")

for community, size in top_communities:

    print(f"Community {community}: {size} papers")
```

## 2.4 Papers with No or Low Citations

In any citation network, certain papers are less likely to be cited, either due to limited relevance, niche research topics, or lack of visibility within the community. In our dataset:

**4,590 papers** have never been cited by any other paper in the network. These papers might represent niche topics, initial ideas that didn't gain traction, or were simply overlooked.

**14,712 papers** have been cited less than 5 times, while **19,591 papers** have been cited less than 10 times.

These statistics highlight the skewed nature of academic citations, where a few papers receive the majority of citations, while many papers receive little to no attention

## 2.5 Influence Concentration and the Top 20% of Papers

In order to assess the concentration of influence within the HEP-TH network, we analyzed the distribution of citations among the most-cited papers. Using **in-degree centrality**—which measures the number of times a paper is cited by others—we identified the **top 20% of papers** and calculated the proportion of total citations they received.

The results reveal an even stronger concentration of influence than previously observed:

- The **top 20% of papers** account for **75.66% of all citations** in the network.

```
# Calculate in-degree centrality for each paper

in_degree_centrality = nx.in_degree_centrality(citation_graph)


# Sort papers by in-degree centrality and select the top 20%

sorted_in_degree = sorted(in_degree_centrality.items(), key=lambda x: x[1],
reverse=True)

top_20_percent_cutoff = int(0.2 * len(sorted_in_degree))

top_20_percent_nodes = {node for node, _ in
sorted_in_degree[:top_20_percent_cutoff]}


# Calculate the percentage of citations received by the top 20% of papers

total_citations = sum(dict(citation_graph.in_degree()).values())

top_20_citations =
sum(dict(citation_graph.in_degree(nbunch=top_20_percent_nodes)).values())

top_20_percentage = (top_20_citations / total_citations) * 100
```
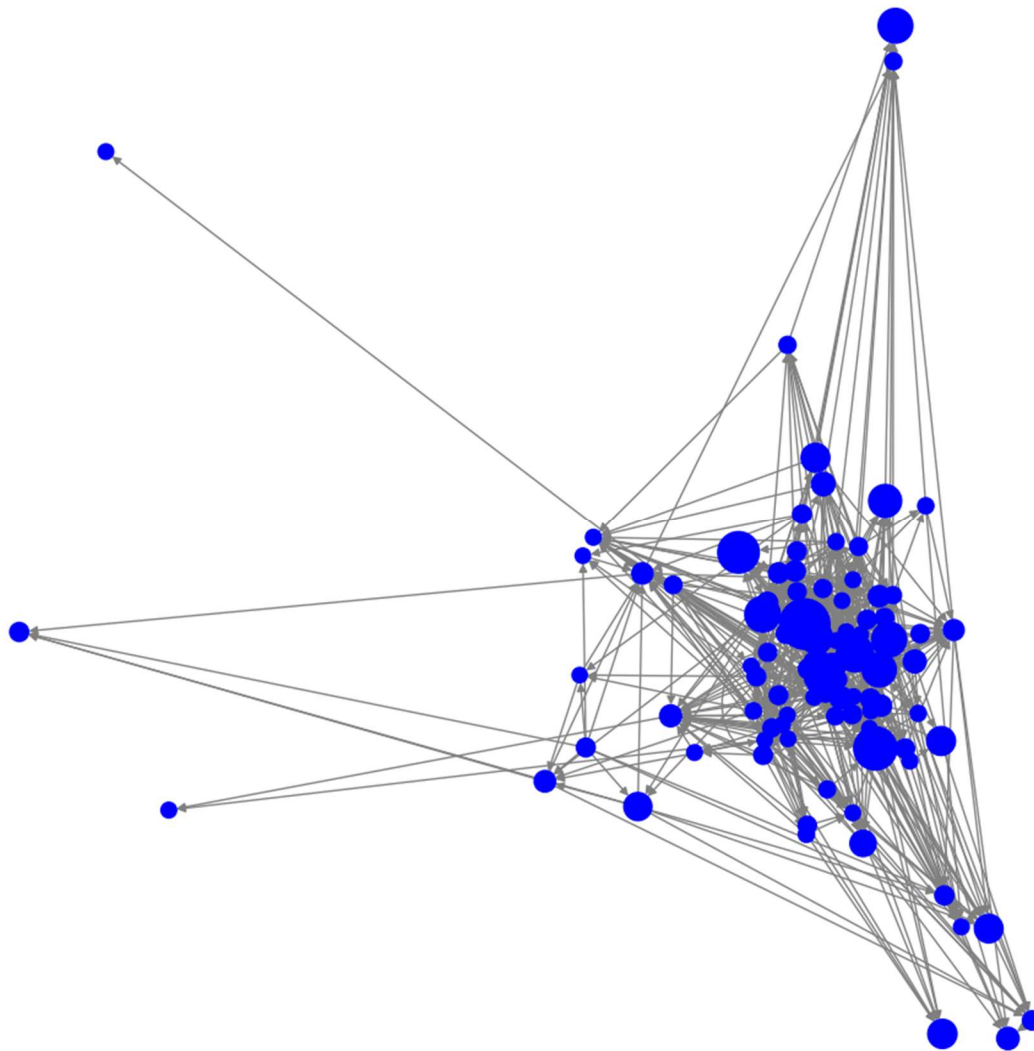
# Chapter 3: Visualizing the Results

To make our findings visually engaging, we create various plots and graphs.

## 3.1 Degree Centrality Visualization

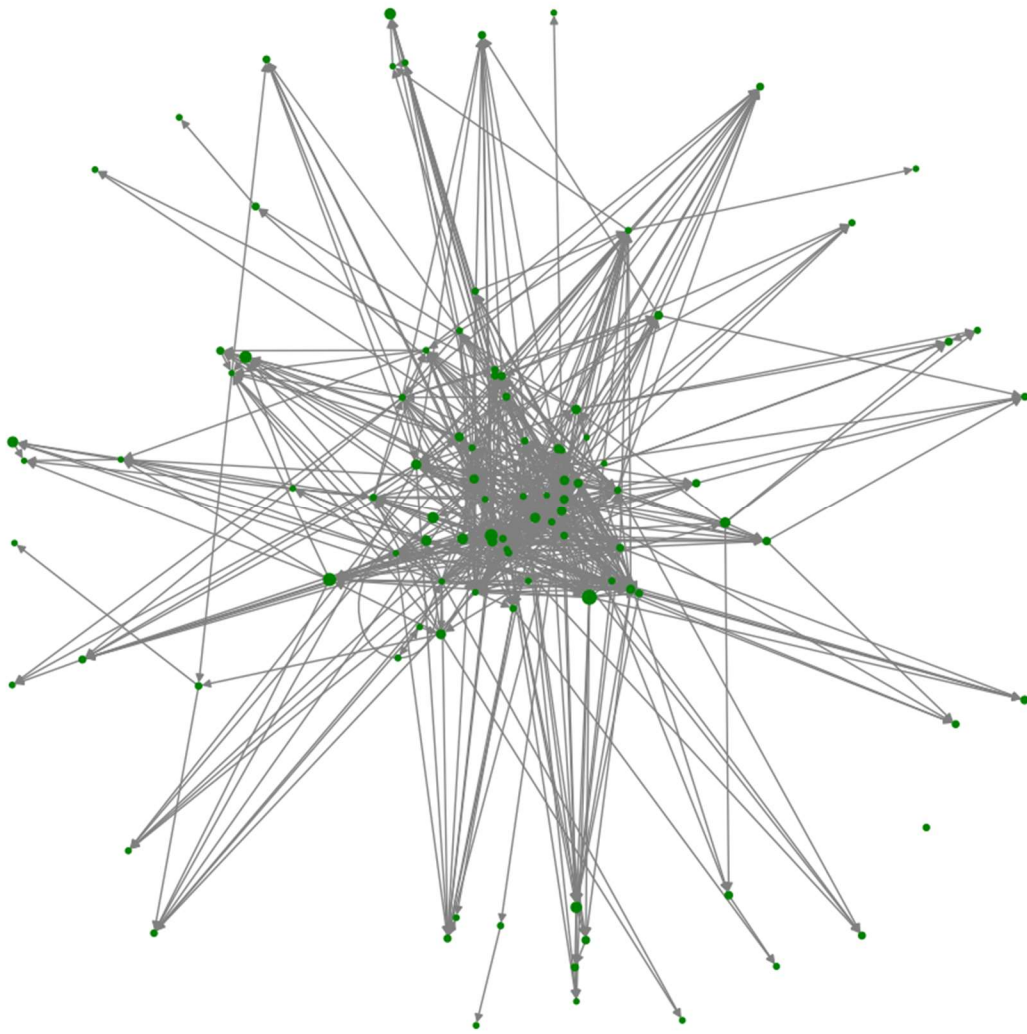We visualize the top 100 most-cited papers, with node size representing the degree
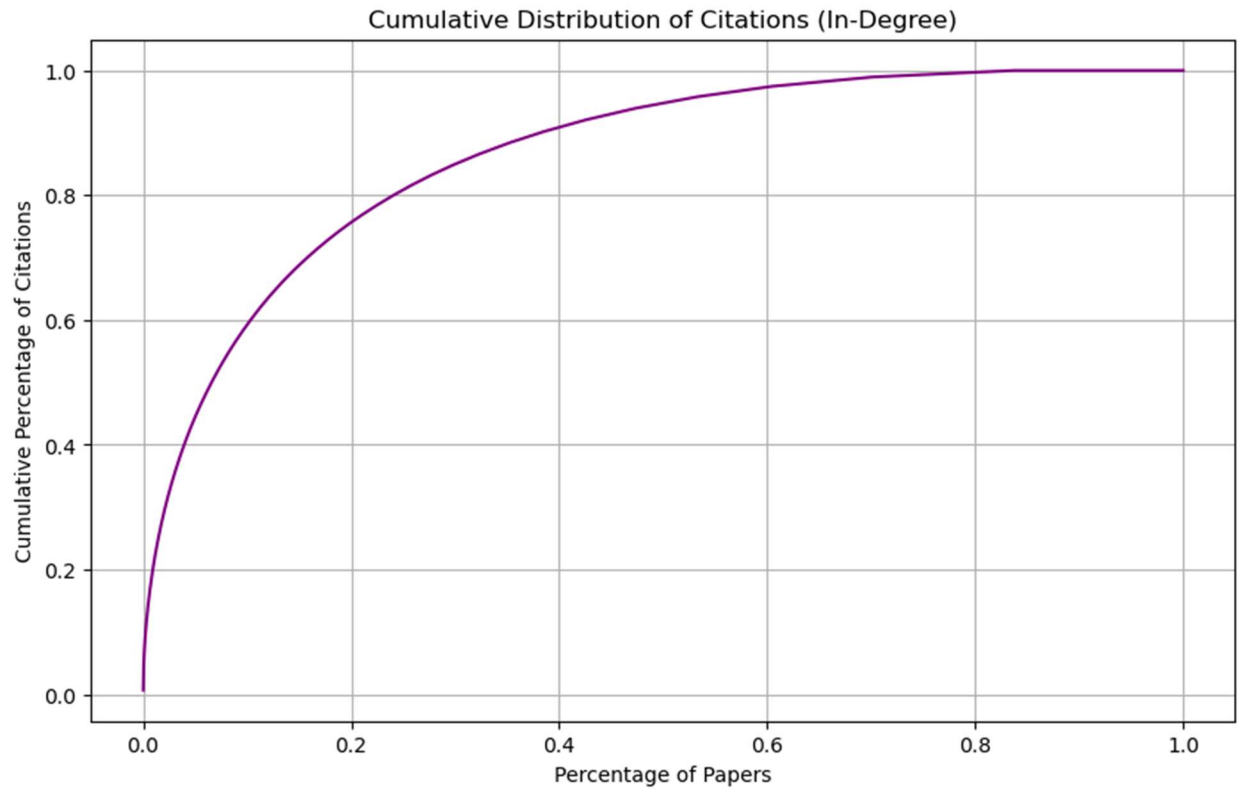centrality.

Top 100 Most Cited Papers

## 3.2 PageRank Visualization

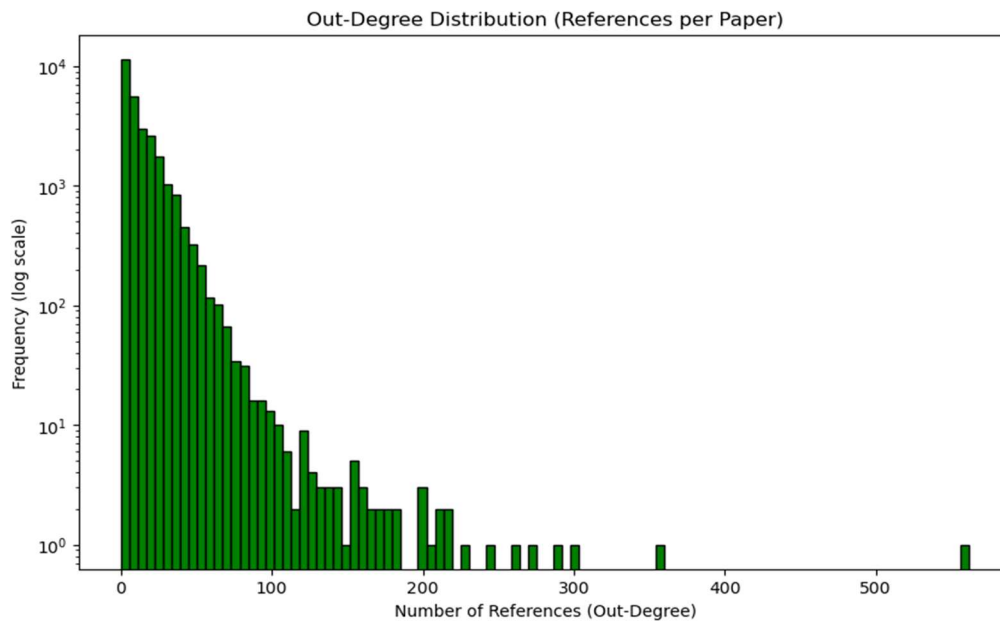This visualization shows the top papers based on PageRank scores, highlighting the most influential ones.
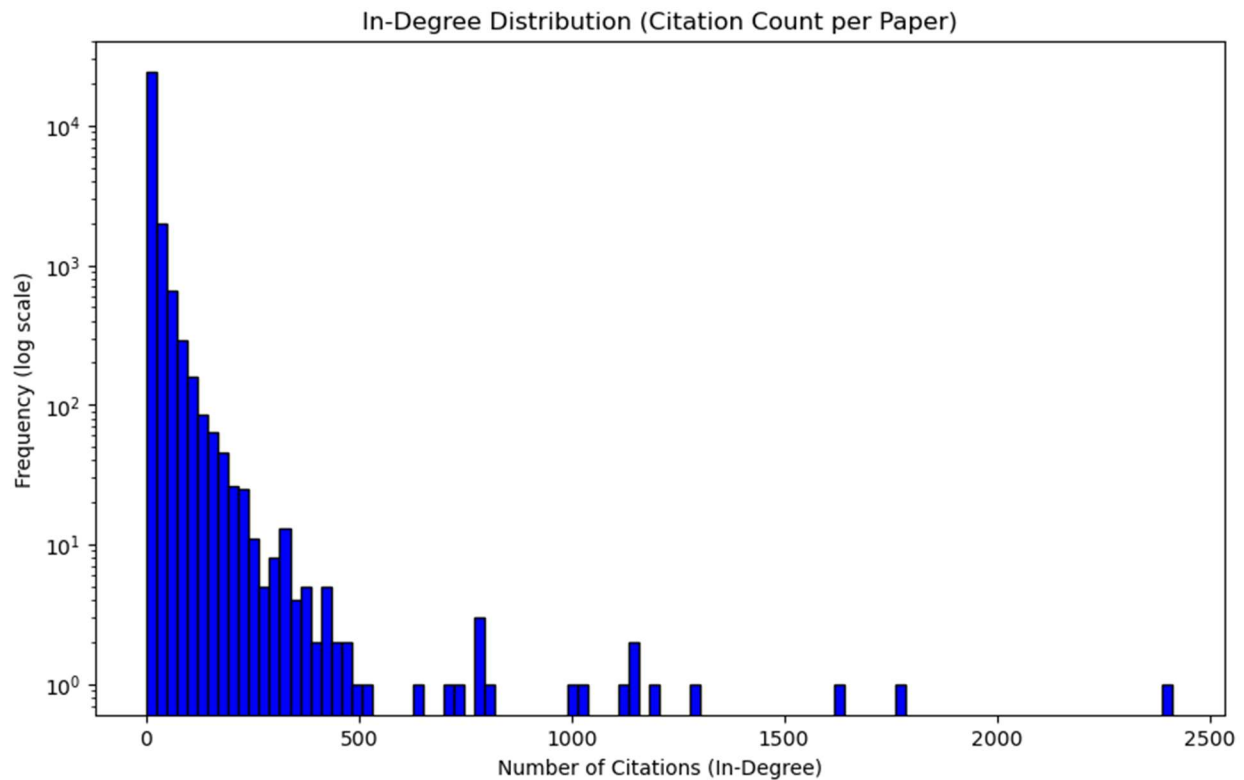
Top 100 Papers by PageRank

Cumulative Distribution of Citations (In-Degree)

## 3.3 degree

Get out-degrees (number of references each paper makes)



Out-Degree Distribution (References per Paper)

# Get in-degrees (number of citations for each paper)



In-Degree Distribution (Citation Count per Paper)

# Chapter 4: Results and Insights

This chapter synthesizes the findings from our analysis of the HEP-TH citation network. By examining key measures such as degree centrality, PageRank, community structure, and citation patterns, we reveal underlying patterns of influence, community clustering, and the overall structure of the network.

## 4.1 Influence Concentration and Core Papers

Our analysis confirms a strong concentration of influence within a small subset of highly cited papers:

- **Top 10 Most-Cited Papers**: The top 10 papers, led by Paper 9711200 (cited by 2,414 papers), demonstrate the skewed nature of citations in this network. This subset of influential papers forms the backbone of the research community, serving as foundational works that many others reference.

- **Degree Centrality Patterns**: High in-degree centrality scores indicate that these papers are not only frequently cited but also highly respected within the field. The out-degree centrality analysis further shows that while some papers contribute extensively to the foundational knowledge by citing numerous other works, others are primarily reference points and do not engage in extensive citation.

This disparity aligns with the **rich-get-richer effect** commonly observed in citation networks, where a few influential papers accumulate the majority of citations, forming a **hub-and-spoke structure** in which influence radiates outward from a small core.

## 4.2 PageRank Analysis of Influence

Using the PageRank algorithm, we identified papers that are both frequently cited and serve as connectors across different parts of the network:

- **Top Papers by PageRank**: Papers with the highest PageRank scores, such as 9407087 and 9503124, not only receive numerous citations but also link important parts of the network. These papers play a pivotal role in disseminating knowledge across the broader HEP-TH community.

- **Secondary Hubs**: The PageRank results also highlight a set of secondary hubs, which, while not as influential as the top papers, help bridge communities. These secondary hubs support the spread of ideas and facilitate the interconnectedness of research within high-energy physics theory.

## 4.3 Community Detection and Subfield Clustering

Applying the **Louvain method** for community detection revealed several distinct research communities within the HEP-TH network:

- **Major Communities**: The top five communities, each with over 1,500 papers, represent prominent subfields or research clusters within high-energy physics. For example, **Community 12** (2,413 papers) and **Community 1** (2,065 papers) show high levels of internal citation, suggesting they represent specialized areas of research with strong internal connections.

- **Community Structure and Interconnections**: While each community is highly interconnected, there are cross-community citations involving highly cited papers that serve as **bridges** between different subfields. These bridging papers indicate the presence of foundational studies with applications or influence across multiple sub-disciplines.

The Louvain method helps us visualize the clustering of research topics, where the most connected papers within communities likely drive major research themes in HEP-TH.

## 4.4 Citation Behavior: Independent and Low-Citation Papers

The citation behavior within the network reveals several notable patterns:

- **Papers with No incoming Citations**: A total of **4,590 papers** did not be cited by any other work in the network. Alternatively, some might be comprehensive surveys summarizing existing knowledge.

- **Low-Citation Papers**: **19,591 papers** cited by fewer than 10 other works. This subset likely includes specialized studies, exploratory papers, or works with a narrow focus that rely on a limited set of prior research.

These findings suggest that while a few papers serve as core connectors within the broader research community, a significant number operate independently or with minimal citations, emphasizing the hierarchical nature of the network.

## 4.5 Influence Concentration in the Top 20%

The in-depth analysis of in-degree centrality reveals a highly skewed distribution of citations within the HEP-TH network. Our findings show that the **top 20% of papers receive 75.66% of all citations**, a concentration that underscores the hierarchical nature of the citation network.

This distribution implies that:

- **Foundational Works Dominate**: A small subset of highly influential papers acts as the intellectual core of the field, receiving the vast majority of citations.

- **Hierarchical Structure of Influence**: The network's structure resembles a hub-and-spoke model, where these central hubs (the top 20% of papers) draw the majority of citations, while the remaining 80% of papers receive relatively few citations.

This pattern exemplifies the **rich-get-richer phenomenon** seen in citation networks, where highly cited papers continue to attract more citations over time, solidifying their status as cornerstone works within the HEP-TH research community. The high degree of concentration suggests that these influential papers play a critical role in shaping the direction and foundation of research in high-energy physics.

# Conclusion

The HEP-TH citation network reveals a hierarchical structure where influence is concentrated within a small number of highly cited papers. These core papers act as hubs that connect different parts of the network, driving research direction and facilitating knowledge dissemination. Community detection further highlights clustering within subfields, showing how research themes form and evolve over time. This analysis provides a comprehensive view of the citation patterns, community structure, and influence dynamics within the HEP-TH field, shedding light on the interconnected nature of academic research in high-energy physics.