

# CMS/CS/EE 144 - Homework 1 Analysis Report

## AI Use Declaration

I manually implemented the initial version of the crawler. I used AI assistance to optimize the crawler by parallelizing it across 5 workers and to help visualize the network structure. For the analysis, I used AI to identify the correct NetworkX functions for calculating diameter and clustering coefficients, as well as to clean histograms and CCDFs. Finally, AI was used to help format this L<sup>A</sup>T<sub>E</sub>X report.

### 1.1.1 Code & Execution

**Libraries:** The crawler uses `networkx` for graph management, `matplotlib` for plotting, and `concurrent.futures` for parallel execution.

#### How to run:

1. Run the crawler: `python crawler.py`  
This generates the graph, saves it to `caltech_web_graph.pkl`, and creates the network visualization.
2. Run the analysis: `python analysis.py`  
This calculates metrics and generates the histograms/CCDFs.

```
 1 from fetcher3 import fetch_links
 2 import networkx as nx
 3 import matplotlib.pyplot as plt
 4 import concurrent.futures
 5 import pickle
 6
 7 start_url = "https://www.caltech.edu"
 8 num_nodes_limit = 2000
 9 MAX_WORKERS = 5 # Polite parallelism (approx 5 concurrent requests
 10 )
11 G = nx.DiGraph()
12 G.add_node(start_url)
13
14 to_visit = [start_url]
15 visited = set()
16
17 with concurrent.futures.ThreadPoolExecutor(max_workers=MAX_WORKERS)
18     as executor:
19         while len(visited) < num_nodes_limit and to_visit:
```

```

19     # 1. Select a batch of URLs to visit (BFS)
20     batch = []
21     while len(batch) < MAX_WORKERS and to_visit:
22         url = to_visit.pop(0)
23         if url not in visited:
24             batch.append(url)
25             visited.add(url)
26
27     if not batch:
28         break
29
30     # 2. Parallel Fetch
31     future_to_url = {executor.submit(fetch_links, url): url for
32                      url in batch}
33
34     # 3. Process results as they complete
35     for future in concurrent.futures.as_completed(future_to_url):
36
37         current_url = future_to_url[future]
38         try:
39             links = future.result()
40             if links:
41                 for link in links:
42                     if "caltech.edu" in link:
43                         # Add edge (safe: main thread does this
44
45                         G.add_edge(current_url, link)
46                         # Add to queue if not seen (approximate
47                         # check, set handles exact)
48                         if link not in visited:
49                             to_visit.append(link)
50
51         except Exception as e:
52             pass # safely ignore fetch errors
53             print(f"Crawled: {len(visited)} | Nodes: {G.
54             number_of_nodes()} | Queue: {len(to_visit)}", end='\r')
55
56 unvisited_nodes = [node for node in G.nodes() if node not in
57 visited]
58 G.remove_nodes_from(unvisited_nodes)
59
60 print(f"\nFinal Graph: {G.number_of_nodes()} nodes, {G.
61             number_of_edges()} edges")
62
63 # Save the graph
64 with open('caltech_web_graph.pkl', 'wb') as f:
65     pickle.dump(G, f)
66 print("Graph saved to caltech_web_graph.pkl")
67
68 # plot the graph
69 plt.figure(figsize=(10, 10))
70 nx.draw(G,
71         pos=nx.spring_layout(G, k=0.15, iterations=20), # k
72         controls spacing
73         node_size=10,
74         width=0.1,
75         arrowsize=5,
76         with_labels=False,
77

```

```

68     node_color='blue',
69     edge_color='gray',
70     alpha=0.6)
71 plt.title("Caltech Web Crawl")
72 plt.axis('off')
73 plt.savefig('caltech_web_graph.png', dpi=500)

```

Listing 1: Crawler Implementation

```

1 import pickle
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 import os
7
8 # Graph loading moved to analyze_dataset function
9
10 def plot_histogram(data, title, xlabel, ylabel, filename):
11     plt.figure()
12     plt.hist(data, bins=range(min(data), max(data) + 2, 5), align='left', rwidth=0.8)
13     plt.yscale('log')
14     plt.title(title)
15     plt.xlabel(xlabel)
16     plt.ylabel(ylabel)
17     plt.savefig(filename)
18     plt.close()
19
20 def plot_ccdf(data, title, xlabel, ylabel, filename):
21     # Filter out zeros for log-log plot (reciprocal of infinite is
22     # 0, log(0) is undefined)
23     # or just shift them if 0 is meaningful. For degrees, 0 is
24     # possible.
25     # Usually we plot data > 0 for power law checks.
26     data = np.array(data)
27     data = data[data > 0]
28
29     if len(data) == 0:
30         print(f"Warning: No positive data for {title}")
31         return
32
33     data_sorted = np.sort(data)
34     # Calculate CCDF: P(X >= x)
35     # rank i (0-based) corresponds to N-i elements >= x
36     # y = (N - i) / N
37     n = len(data_sorted)
38     y = np.arange(n, 0, -1) / n
39
40     plt.figure()
41     plt.loglog(data_sorted, y, marker='.', linestyle='none')
42     plt.title(title)
43     plt.xlabel(xlabel)
44     plt.ylabel(ylabel)
45     plt.grid(True, which="both", ls="--", alpha=0.2)
46     plt.savefig(filename)
47     plt.close()

```

```

47 def analyze_dataset(pickle_file, label, file_prefix=''):
48     print(f"\n--- Analyzing {label} ({pickle_file}) ---")
49
50     try:
51         with open(pickle_file, 'rb') as f:
52             G = pickle.load(f)
53             print(f"Graph loaded with {G.number_of_nodes()} nodes and {G.number_of_edges()} edges.")
54     except FileNotFoundError:
55         print(f"File {pickle_file} not found. Skipping.")
56         return
57
58     # 1.3 Histograms & 1.4 CCDF
59     out_degrees = [G.out_degree(n) for n in G.nodes()]
60     in_degrees = [G.in_degree(n) for n in G.nodes()]
61
62     prefix = file_prefix
63
64     plot_histogram(out_degrees, f"Out-Degree Distribution ({label})",
65                   "Out-Degree", "Frequency", f"{prefix}hist_out_degree.png")
66     plot_histogram(in_degrees, f"In-Degree Distribution ({label})",
67                   "In-Degree", "Frequency", f"{prefix}hist_in_degree.png")
68
69     plot_ccdf(out_degrees, f"CCDF of Out-Degree ({label})", "Out-Degree",
70                "P(X >= x)", f"{prefix}ccdf_out_degree.png")
71     plot_ccdf(in_degrees, f"CCDF of In-Degree ({label})", "In-Degree",
72                "P(X >= x)", f"{prefix}ccdf_in_degree.png")
73
74     # 1.5 Clustering Coefficients - Treat as undirected
75     G_undirected = G.to_undirected()
76
77     # Handle disconnected components if necessary for diameter
78     if not nx.is_connected(G_undirected):
79         largest_cc = max(nx.connected_components(G_undirected), key=len)
80         G_comp = G_undirected.subgraph(largest_cc)
81         print(f"Graph is not connected. Using largest component ({len(largest_cc)} nodes) for diameter.")
82     else:
83         G_comp = G_undirected
84
85     avg_clustering = nx.average_clustering(G_undirected)
86     overall_clustering = nx.transitivity(G_undirected)
87
88     print(f"Average Clustering Coefficient: {avg_clustering}")
89     print(f"Overall Clustering Coefficient: {overall_clustering}")
90
91     # 1.6 Diameter
92     try:
93         diameter = nx.diameter(G_comp)
94         avg_diameter = nx.average_shortest_path_length(G_comp)
95         print(f"Maximal Diameter: {diameter}")
96         print(f"Average Diameter: {avg_diameter}")
97     except Exception as e:
98         print(f"Could not calculate diameter: {e}")
99         diameter = "N/A"
100        avg_diameter = "N/A"

```

```

97
98     # Write results
99     mode = 'a'
100    with open("analysis_results.txt", mode) as f:
101        f.write(f"\n--- {label} ---\n")
102        f.write(f"Nodes: {G.number_of_nodes()}\n")
103        f.write(f"Edges: {G.number_of_edges()}\n")
104        f.write(f"Average Clustering Coefficient: {avg_clustering:.4f}\n")
105        f.write(f"Overall Clustering Coefficient: {overall_clustering:.4f}\n")
106        f.write(f"Maximal Diameter: {diameter}\n")
107        f.write(f"Average Diameter: {avg_diameter if isinstance(avg_diameter, str) else f'{avg_diameter:.4f'}}\n")
108
109    # Clear results file
110    with open("analysis_results.txt", "w") as f:
111        f.write("Analysis Report Summary\n")
112
113    # Run analysis for both
114    # Parallel (original names maintained by empty prefix)
115    analyze_dataset('caltech_web_graph.pkl', 'Parallel', '')
116    # Sequential (prefixed names)
117    analyze_dataset('caltech_web_graph_np.pkl', 'Sequential', 'seq_')

```

Listing 2: Analysis Implementation

### 1.1.2 Selection Policy

We utilized a **Breadth-First Search (BFS)** selection policy. The crawler maintains a FIFO (First-In-First-Out) queue of URLs.

#### Strengths:

- **Trap Avoidance:** BFS is less susceptible to "infinite depth" traps (e.g., infinite calendar next-page links) compared to Depth-First Search (DFS).
- **Relevance:** It prioritizes nodes closer to the root (*caltech.edu*), which are typically more significant than deeply nested pages.
- **Parallelism Friendly:** The level-by-level exploration naturally supports our parallel batch processing approach.

#### Weaknesses:

- **Memory Overhead:** The frontier queue can grow large quickly as it stores all discovered-but-not-visited links.

### 1.1.3 Network Visualization

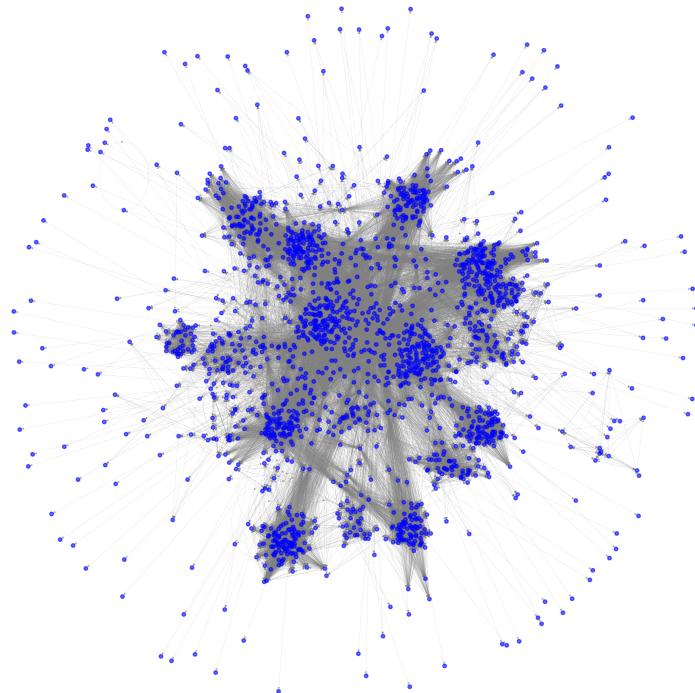


Figure 1: Visualization of the crawled Caltech web graph (Parallel).

### 1.1.4 Degree Distributions

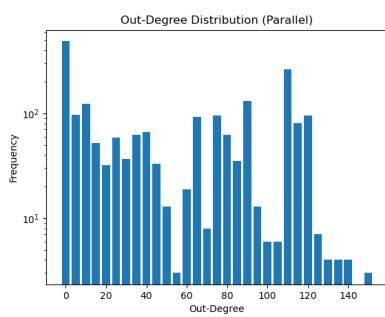


Figure 2: Out-Degree Histogram

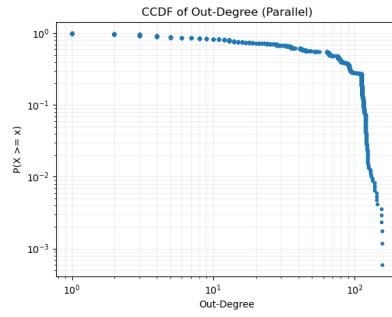


Figure 3: Out-Degree CCDF (Log-Log)

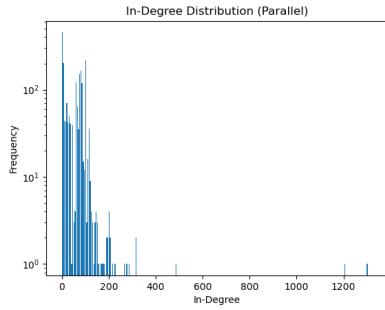


Figure 4: In-Degree Histogram

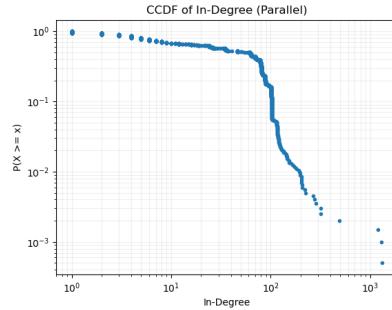


Figure 5: In-Degree CCDF (Log-Log)

### 1.1.5 Clustering Coefficients

The following values were calculated treating edges as undirected:

- **Average Clustering Coefficient:** 0.7867
- **Overall Clustering Coefficient:** 0.5905

### 1.1.6 Diameter Metrics

Calculated on the largest connected component of the undirected graph:

- **Maximal Diameter:** 4
- **Average Diameter:** 2.4183

### 1.1.7 Comparison with Collaboration Networks

**Task:** A comparison of the degree distribution, the clustering coefficients, and the diameters with those in Problem 2. Can you describe the similarities and differences regarding the "universal" properties between collaboration network and the web graph?

- Similarities:
- Differences:

### 1.1.8 Sequential Analysis (Retry)

The Cumulative Distribution Functions (CDFs) from the initial parallel crawl did not exhibit a clear power law behavior, which was unexpected for a web graph. Suspecting that the parallel execution might have influenced the sampling or traversal order non-deterministically, I implemented a sequential version of the crawler to verify the results.

```
1 from fetcher3 import fetch_links
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 import pickle
5
6 start_url = "https://www.caltech.edu"
7 num_nodes_limit = 2000
8
9 G = nx.DiGraph()
10 G.add_node(start_url)
11
12 to_visit = [start_url]
13 visited = set()
14
15 while len(visited) < num_nodes_limit and to_visit:
16     current_url = to_visit.pop(0)
17
18     if current_url in visited:
19         continue
20
21     visited.add(current_url)
22
23     try:
24         links = fetch_links(current_url)
25         if links:
26             for link in links:
27                 if "caltech.edu" in link:
28                     G.add_edge(current_url, link)
29                     if link not in visited:
30                         to_visit.append(link)
31     except Exception:
32         pass # Ignore fetch errors
33
34 print(f"Crawled: {len(visited)} | Nodes: {G.number_of_nodes()}"
| Queue: {len(to_visit)}", end='\r')
```

```

35
36 # Clean up nodes that weren't visited
37 unvisited_nodes = [node for node in G.nodes() if node not in
38     visited]
39 G.remove_nodes_from(unvisited_nodes)
40 print(f"\nFinal Graph: {G.number_of_nodes()} nodes, {G.
41     number_of_edges()} edges")
42
43 # Save the graph with a different name
44 with open('caltech_web_graph_np.pkl', 'wb') as f:
45     pickle.dump(G, f)
46 print("Graph saved to caltech_web_graph_np.pkl")
47
48 # Plot the graph
49 plt.figure(figsize=(10, 10))
50 nx.draw(G,
51     pos=nx.spring_layout(G, k=0.15, iterations=20),
52     node_size=10,
53     width=0.1,
54     arrowsize=5,
55     with_labels=False,
56     node_color='green',
57     edge_color='gray',
58     alpha=0.6)
59 plt.title("Caltech Web Crawl (Sequential)")
60 plt.axis('off')
61 plt.savefig('caltech_web_graph_np.png', dpi=500)

```

Listing 3: Sequential Crawler Implementation

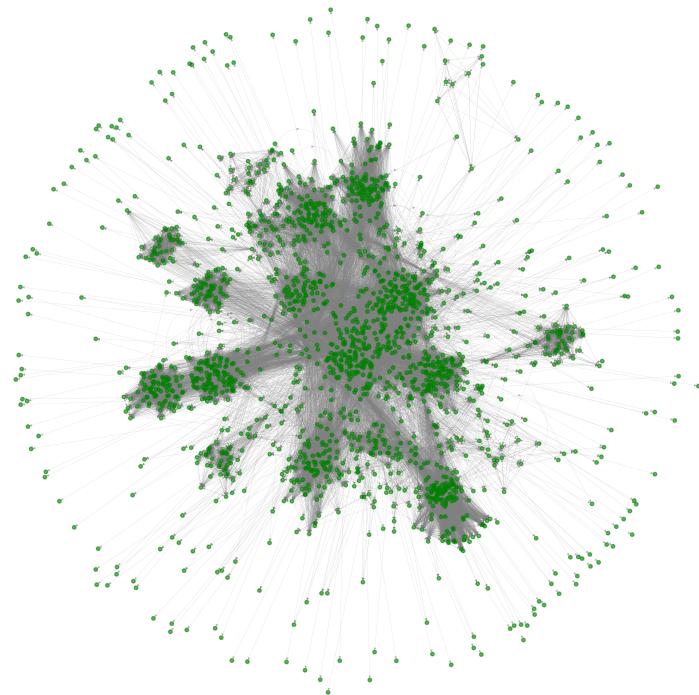


Figure 6: Visualization of the Caltech web graph (Sequential).

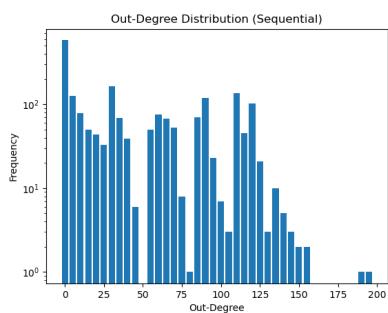


Figure 7: Sequential Out-Degree Histogram

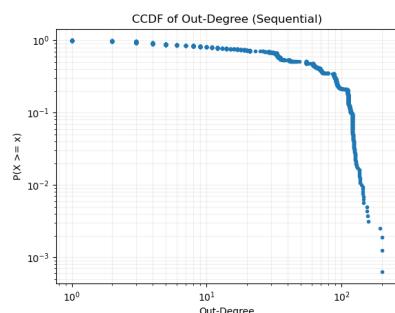


Figure 8: Sequential Out-Degree CCDF (Log-Log)

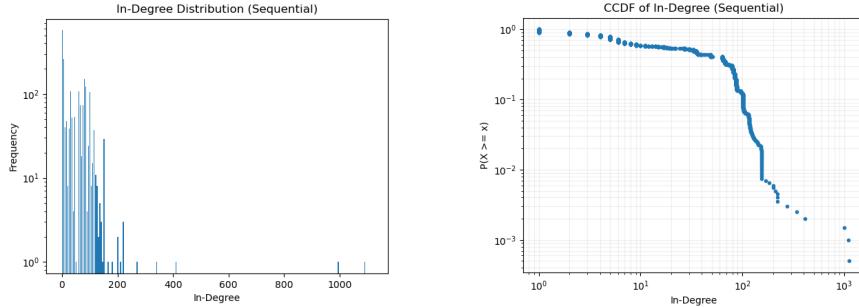


Figure 9: Sequential In-Degree Histogram

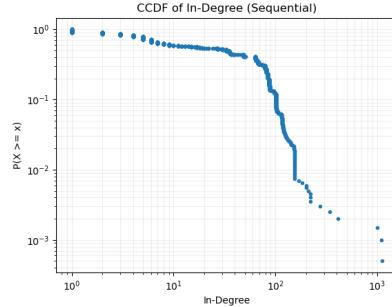


Figure 10: Sequential In-Degree CCDF (Log-Log)

The metrics for the sequential graph are:

- **Average Clustering Coefficient:** 0.7526
- **Overall Clustering Coefficient:** 0.6106
- **Maximal Diameter:** 4
- **Average Diameter:** 2.5954

**Conclusion:** The results from the sequential crawler are very similar to the parallel version. The graph structure, degree distributions, and clustering coefficients show no significant differences, suggesting that the parallel implementation correctly captures the graph properties and the observed distributions are intrinsic to the crawled subgraph.

## 1.2(a) Flight to the Queen Bee

The following paths were discovered on music-map.com. In both instances, the first path found was also the shortest.

- **Nikolai Rimsky-Korsakov → Beyoncé**  
Path: Nikolai Rimsky-Korsakov → Mozart → Distant Fires Burning → Chanyeol → BTS → Arctic Monkeys → Ed Sheeran → Adele → Beyoncé
- **Lyle Mays → Taylor Swift**  
Path: Lyle Mays → John Taylor → Mick Fleetwood → U2 → The Beatles → One Direction → Olivia Rodrigo → Taylor Swift

## 1.2(b) Know your Professors

The following paths were found using the “Co-authorship Path” tool on csauthors.net.

1. **Eric Mazumdar → Rudolf E. Kálmán** (Distance: 5)  
Source: <https://csauthors.net/distance/eric-mazumdar/rudolf-e-kalman>

- Eric Mazumdar co-authored **10 papers** with Adam Wierman
- Adam Wierman co-authored **8 papers** with Linqi Guo
- Linqi Guo co-authored **4 papers** with Guangyi Shi
- Guangyi Shi co-authored **1 paper** with Kazuo Toraichi
- Kazuo Toraichi co-authored **2 papers** with Rudolf E. Kálmán

2. **Georgia Gkioxari → Paul Erdős** (Distance: 3)

Source: <https://csauthors.net/distance/georgia-gkioxari/paul-erdos>

- Georgia Gkioxari co-authored **12 papers** with Jitendra Malik
- Jitendra Malik co-authored **2 papers** with Fan Chung Graham
- Fan Chung Graham co-authored **7 papers** with Paul Erdős

## 1.2(d) Citations

I attempted to find paths for this section for over 30 minutes using multiple tools, including *Inciteful* and large language models like *Gemini*.

- **Inciteful:** This tool was problematic because it did not strictly follow the directed edge requirement (A cites B). Instead, it often connected papers if they simply cited the same source, resulting in invalid paths for this specific task. It found a distance of 3, but this was based on undirectional logic.
- **Gemini / LLMs:** Gemini proposed paths such as *Zipf (1936)* → *Simon (1955)* → *Price (1965)* → *Granovetter (1973)* and *Kalman (1960)* → *Bryson & Ho (1969)* → *Rumelhart et al. (1986)* → *LeCun et al. (1998)*. However, verifying these specific citation links manually was extremely difficult and time-consuming, as it required accessing the full text of each paper to confirm the bibliography entries.

Due to these difficulties and the unreliability of automated tools for strict directed citation verification, I was unable to definitively confirm a valid path.