# INFSCI2125 Final Project: Gnutella P2P Network Analysis

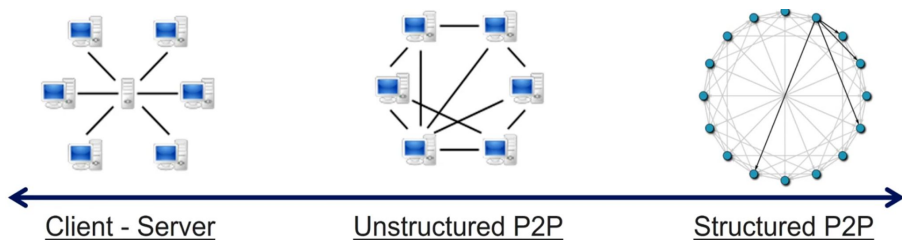Yang

# Gnutella Network – Unstructured P2P
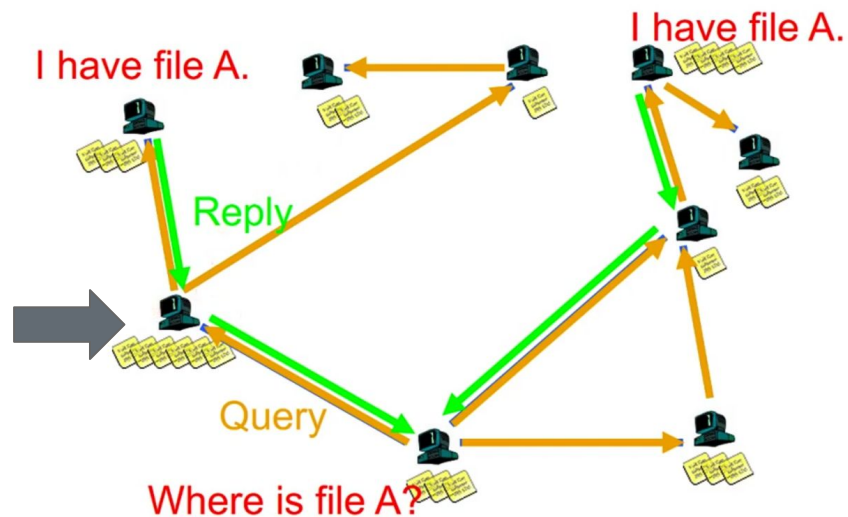


Technical goals:

- Self-organization(vs. centralized )
- Availability (vs. failure, censorship)
- Load balancing (vs. freeloading, locality)
- Anonymity (vs. monitoring)

Challenges:

- Node churn
- Control overheads
- Malicious nodes

Client - Server        Unstructured P2P        Structured P2P

# Gnutella Network – More Details



Mechanism - Query Flooding:

- Join the network: ping-pong neighbours
- Publish: no need
- Search: flood query
- Fetch : direct download from peer

TTL-limited search (usually 7 hops)

- Query dies after
- Only works well for common objects

Pros:
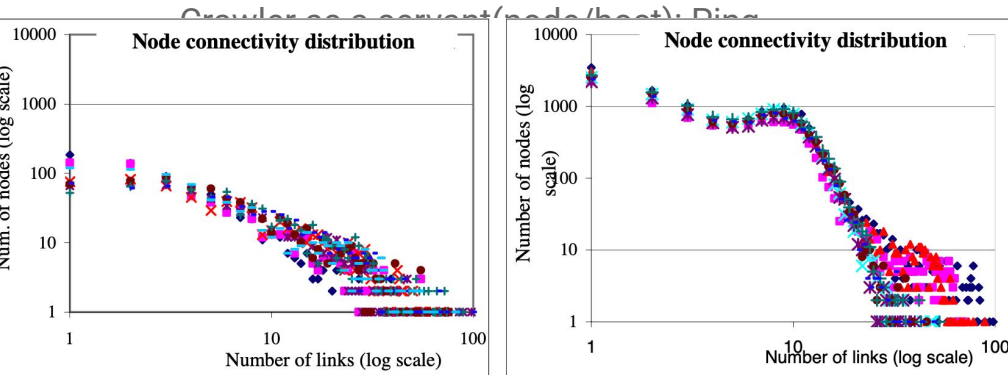
- Fully decentralized
- Search cost distributed

Cons:

- Search Scope is O(N)
- Search time can't be determined
- Nodes often leave, network unstable

# Literature Review – Mapping the Gnutella Network
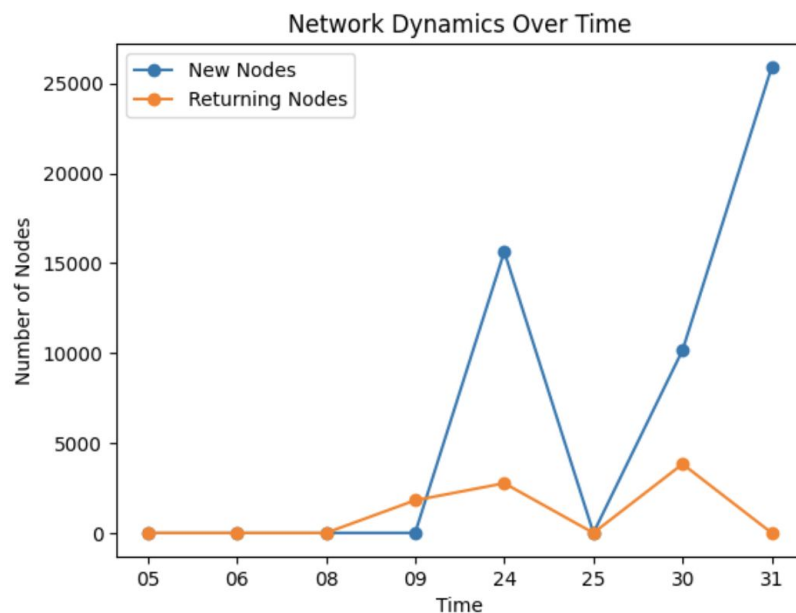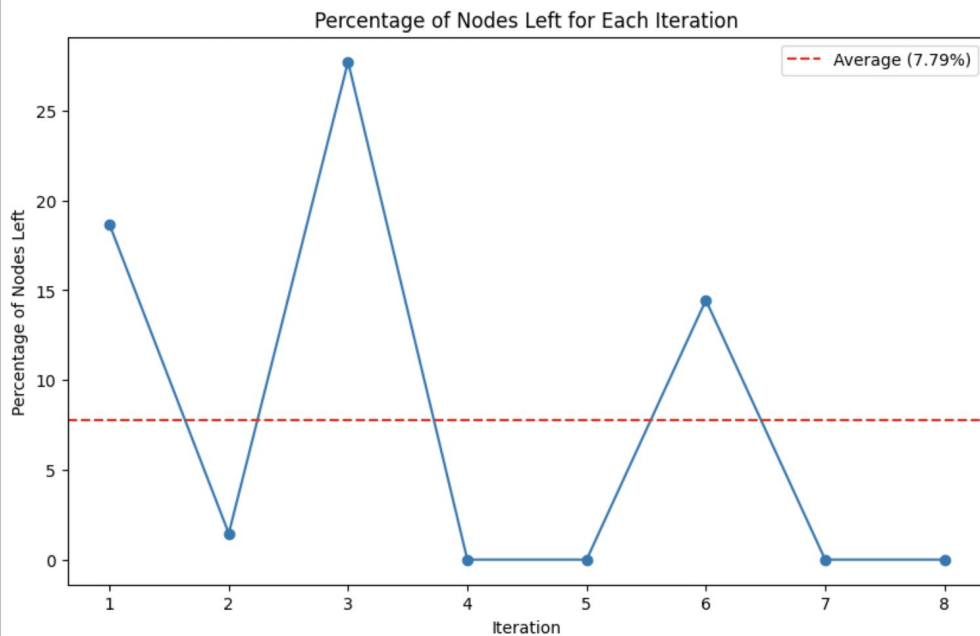
Data

Crawler as a servent(node/host); Ping



because the node has set limited number of TCP connections, or the node has left the network.

Findings:

1. Dynamic network: 40% of the nodes leave the network in less than 4 hours, only 25% of the nodes are alive for more than 24 hours.
2. Massive overhead traffic in the early stage with more than 55%(Nov 2000). Dropped to 8% later with newer implementations(June 2001).
3. Power-law distribution before Nov 2000; too few nodes with low connectivity to form a pure power-law network afterwards
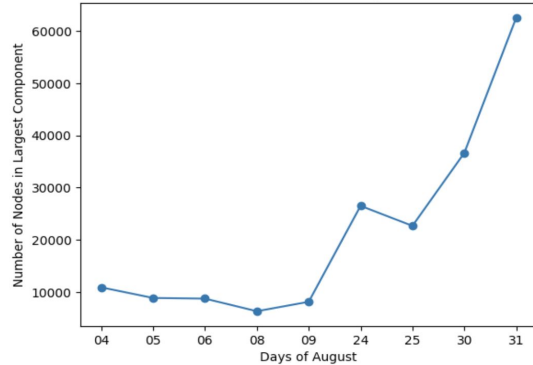
# Churn Rate



Percentage of Nodes Left for Each Iteration
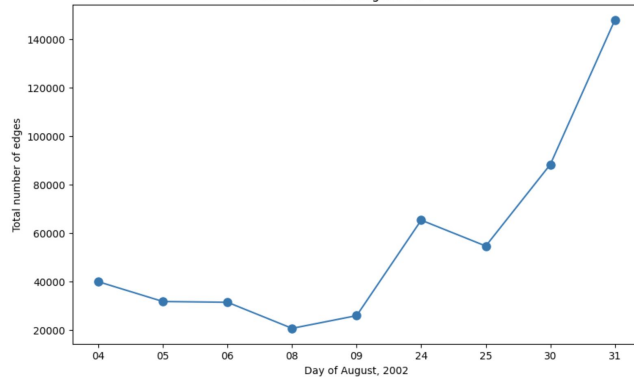
Network Dynamics Over Time

Days: ['04', '05', '06', '08', '09', '24', '25', '30', '31']

# Giant Component

Number of Nodes in Largest Component vs. Time



Total number of edges over time



Percentages of nodes in the largest connected component for each graph:

Graph 0: 100.00%

Graph 1: 99.95%

Graph 2: 100.00%

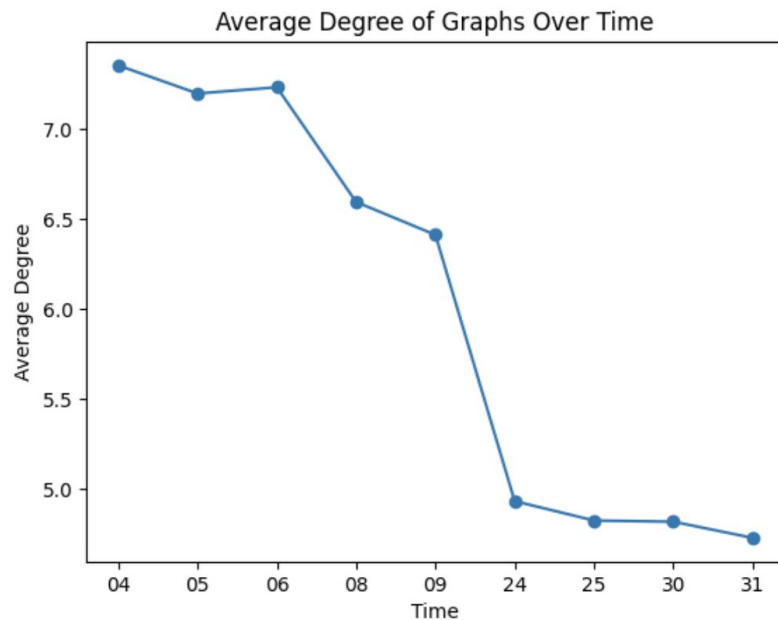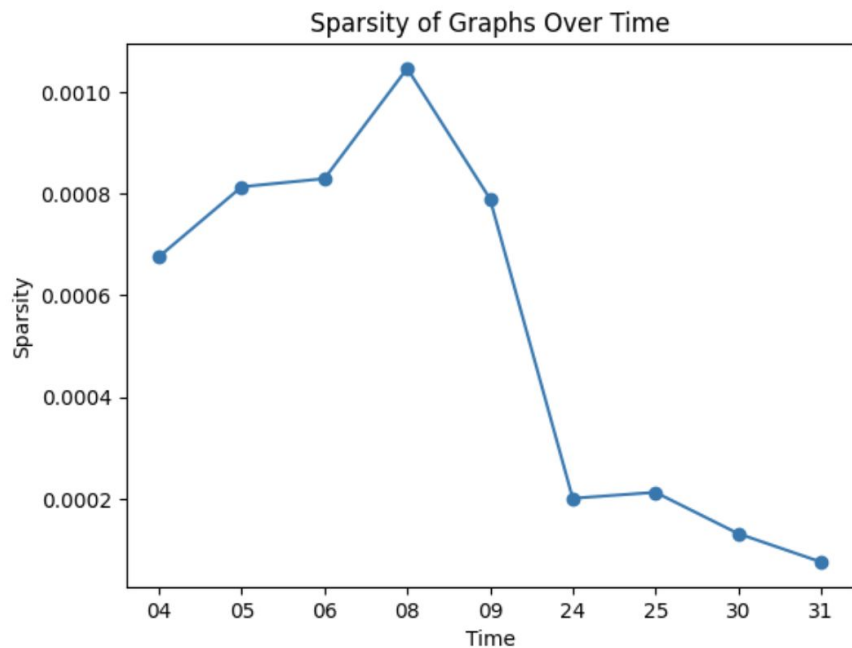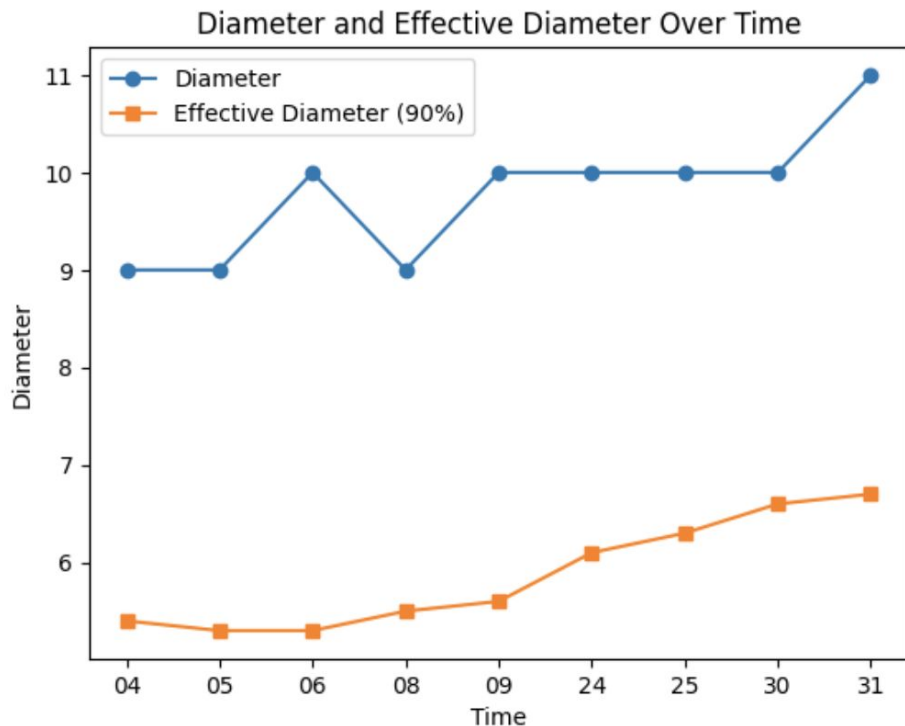Graph 3: 99.97%

Graph 4: 99.88%

Graph 5: 99.92%

Graph 6: 99.89%

Graph 7: 99.90%
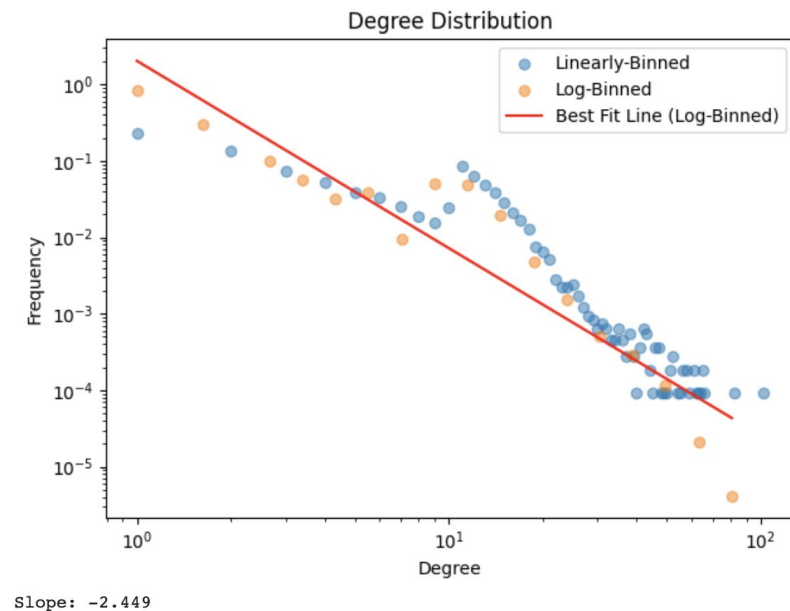
Graph 8: 99.96%

# Sparsity
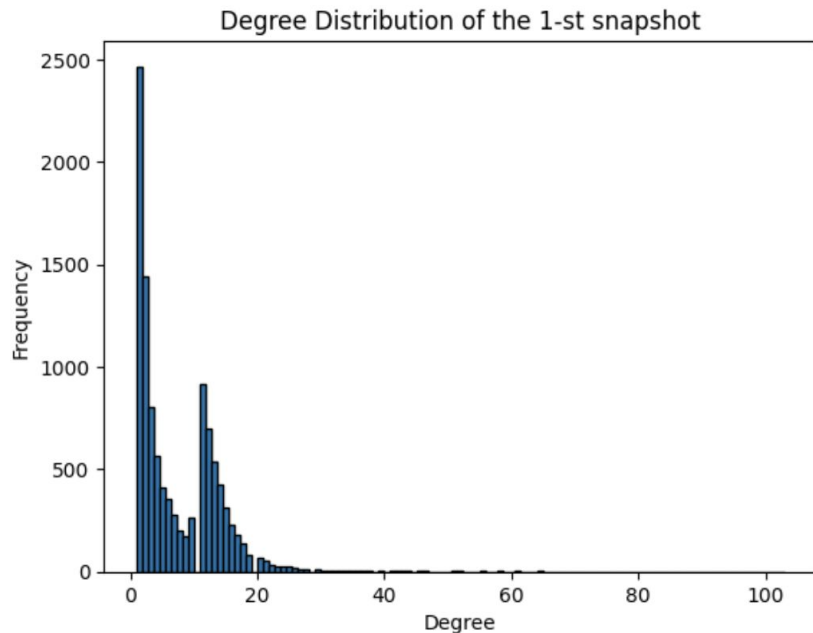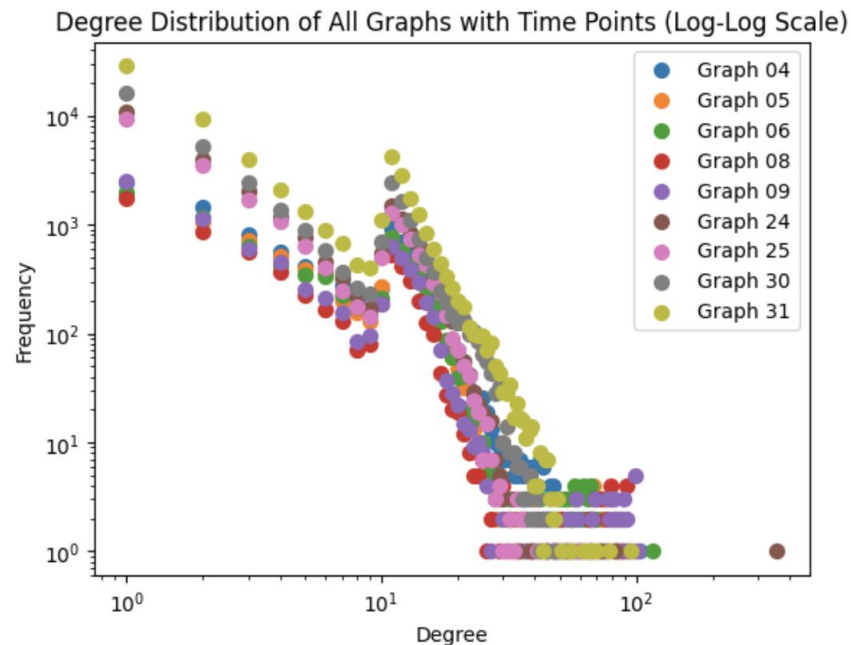
# Diameter



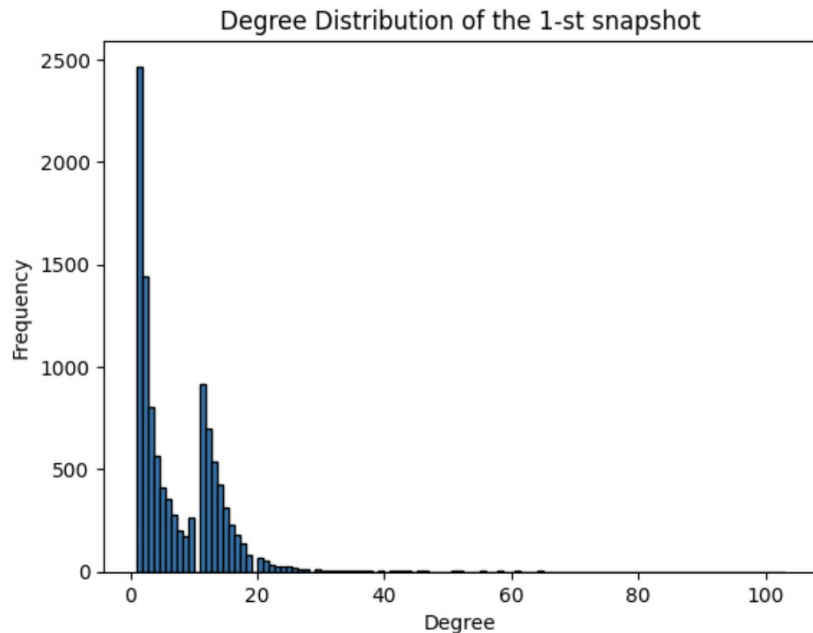Diameter and Effective Diameter Over Time

- Network expanding faster than densification
- No effects on query due to TTL 7 hops limitation.

# Scale–Free: Power–law Distribution

# Scale–Free: Power–law Distribution



Degree Distribution of the 1-st snapshot

Degree Distribution of All Graphs with Time Points (Log-Log Scale)

# Scale–Free: Preferential Attachment

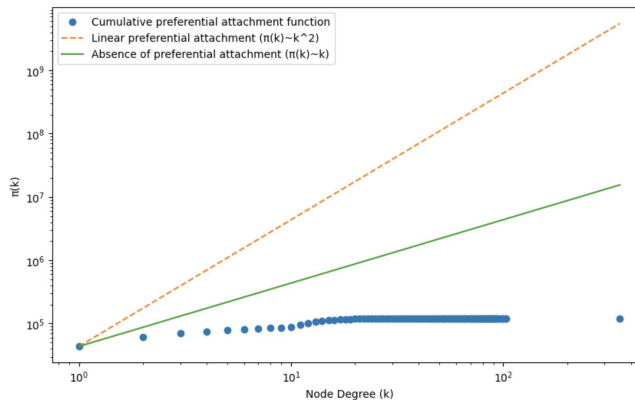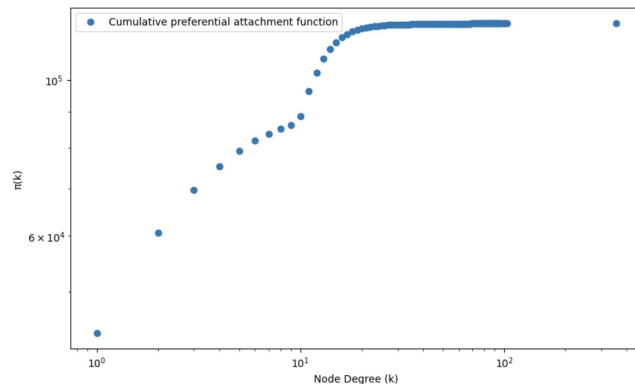**Days:** `['04', '05', '06', '08', '09', '24', '25',`
`'30', '31']`

**Hypothesis 1**
The likelihood to connect to a node depends on that node's degree $k$. This is in contrast with the random network model, for which $\Pi(k)$ is independent of $k$.

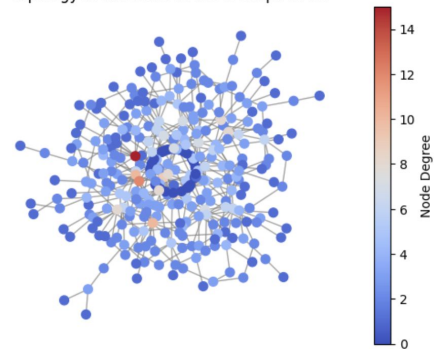**Hypothesis 2**
The functional form of $\Pi(k)$ is linear in $k$.

We detect k–dependance, thus it in line with H1.
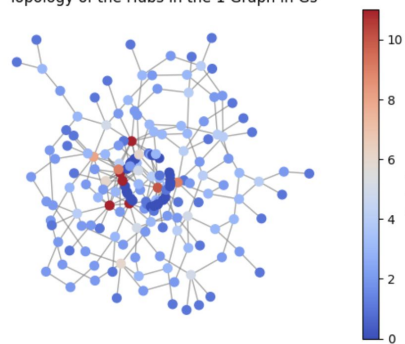Grows slower than ($\pi(k)\sim k2$), thus sublinear.
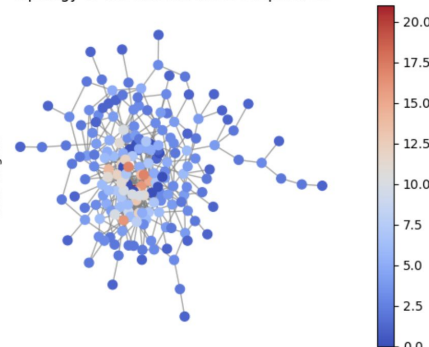
# Hubs – Threshold 20
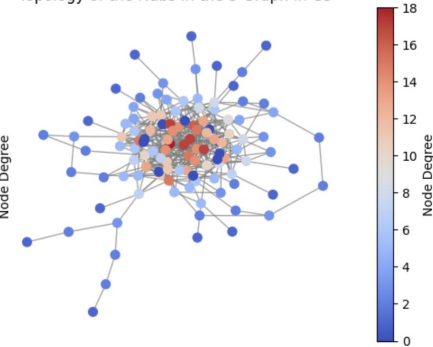


Topology of the Hubs in the 0 Graph in Gs

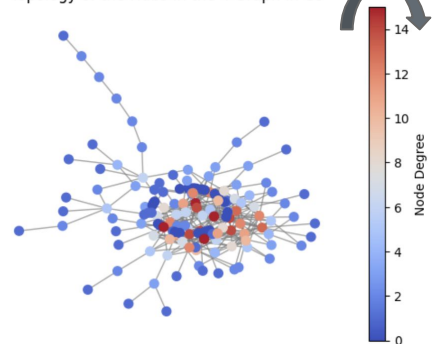Topology of the Hubs in the 1 Graph in Gs

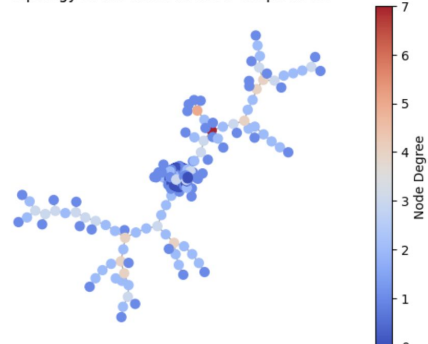Topology of the Hubs in the 2 Graph in Gs

Topology of the Hubs in the 3 Graph in Gs

Topology of the Hubs in the 4 Graph in Gs

Topology of the Hubs in the 5 Graph in Gs

Topology of the Hubs in the 6 Graph in Gs

Topology of the Hubs in the 7 Graph in Gs

Days: ['04', '05', '06', '08', '09', '24', '25', '30', '31']

# Hubs – Threshold 50



Topology of the Hubs in the 0 Graph in Gs

Topology of the Hubs in the 1 Graph in Gs

Topology of the Hubs in the 2 Graph in Gs

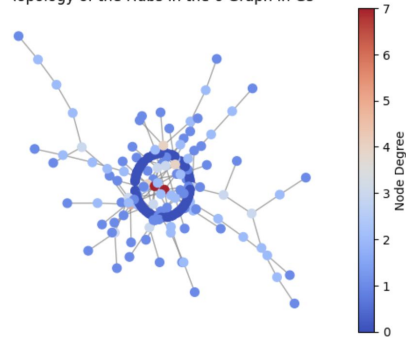Topology of the Hubs in the 3 Graph in Gs

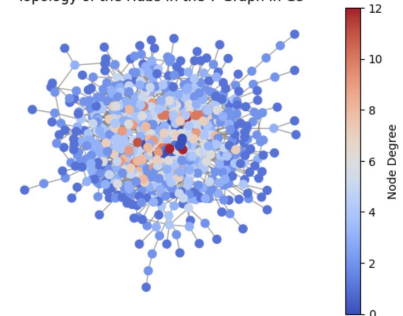Topology of the Hubs in the 4 Graph in Gs

Topology of the Hubs in the 5 Graph in Gs

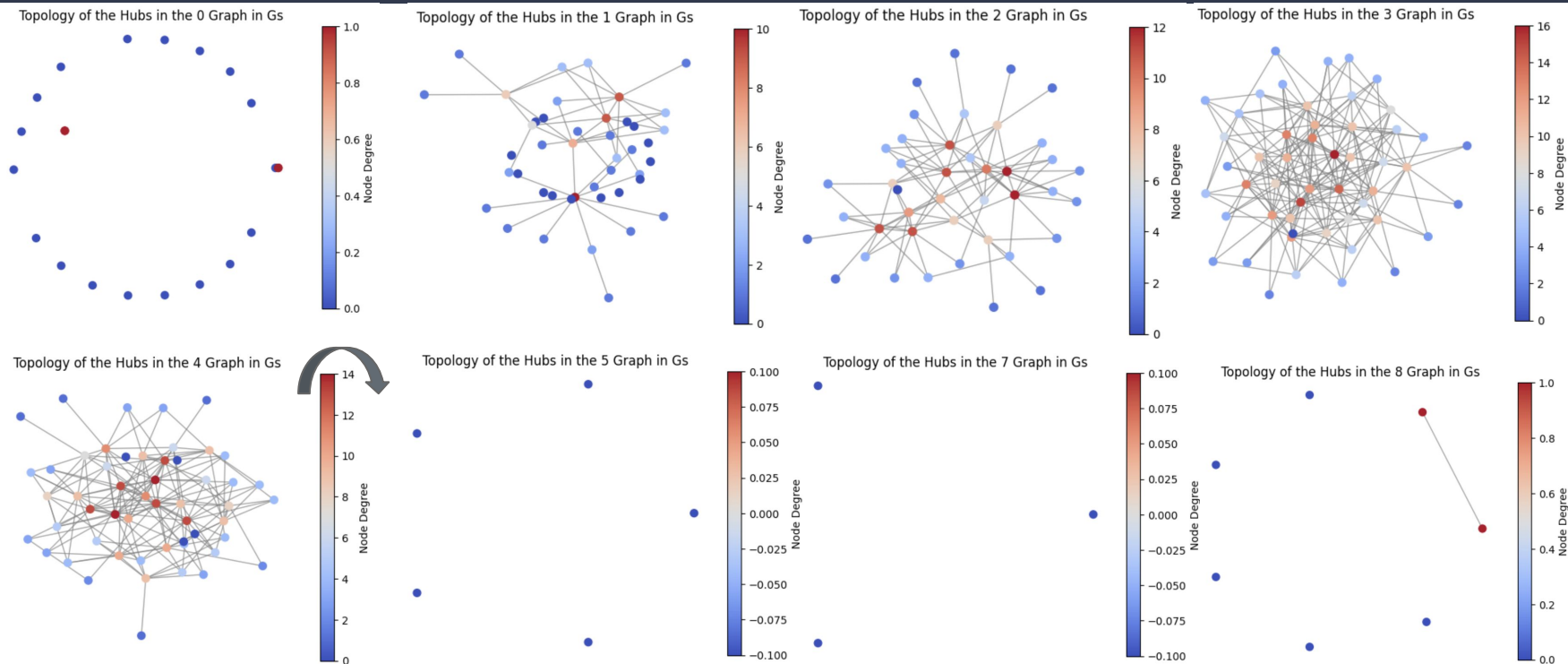Topology of the Hubs in the 7 Graph in Gs

Topology of the Hubs in the 8 Graph in Gs

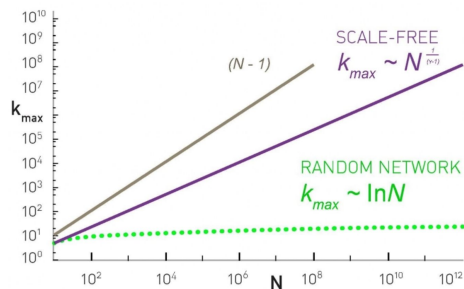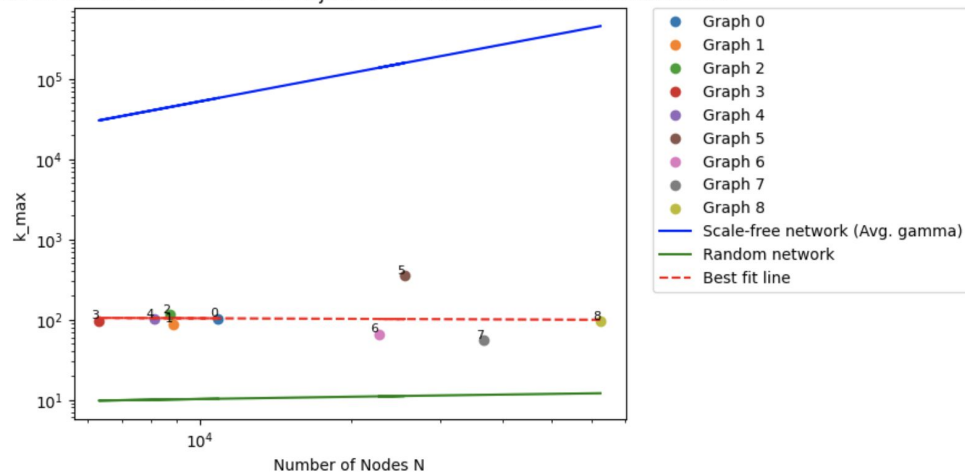Days: ['04', '05', '06', '08', '09', '24', '25', '30', '31']

# Hubs – k_max vs. N



k_max vs. Number of Nodes N with Adjusted Scale-free and Random Network Lines

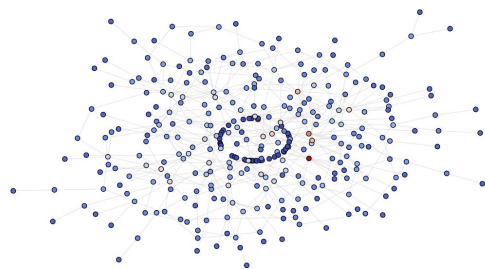gamma = [1.668, 1.674, 1.673, 1.754, 1.780, 1.986, 1.994, 2.036, 2.063]

Average gamma = 1.848

**Anomalous Regime** ($\gamma$ = 2): or $\gamma$ = 2 the degree of the biggest hub grows linearly with the system size, i.e. kmax ~ N. This forces the network into a hub and spoke configuration in which all nodes are close to each other because they all connect to the same central hub. In this regime the average path length does not depend on N.
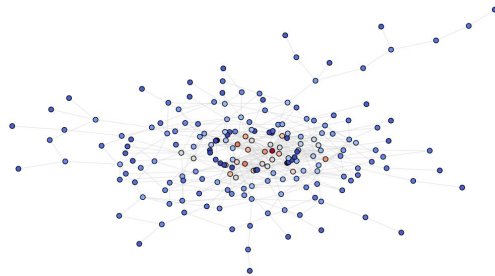
# Future Work

- More consistent dataset in time-series
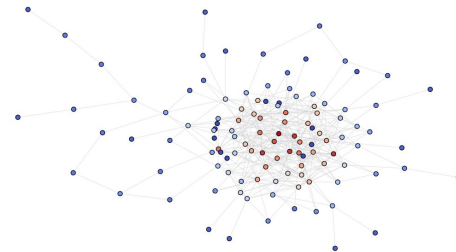- Community detection
- SIR simulation on Gnutella network
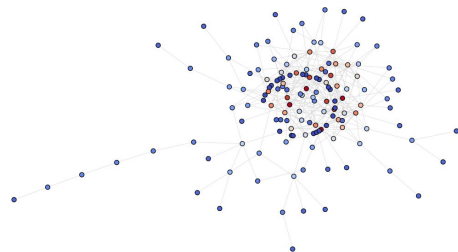
Topology of the Hubs in the 0-th Graph in Gs

Topology of the Hubs in the 2-th Graph in Gs

Topology of the Hubs in the 3-th Graph in Gs
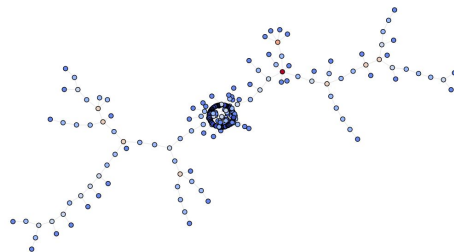
Topology of the Hubs in the 4-th Graph in Gs

change

Topology of the Hubs in the 5-th Graph in Gs

Topology of the Hubs in the 7-th Graph in Gs

Node Degree