The Caltech Graph has n = 2000 nodes and m = 124130 edges.
It was generated using the following selection policy:

Rules:
 - Only visit pages within the Caltech domain.
 - Visit each page only once.
 - Ignore pages that are not HTML (ignoring multimedia/data files as well).
 - Ignore the parameters in dynamic URLs to visit each dynamic page only once.
 - Removes nodes that fail to be crawled due to interruptions/errors.
 - Use breadth-first search algorithm to prioritize between urls.
Pros:
 - Finds pages closer to the starting URL.
 - Keeps diameter of the graph small via BFS.
 - Saves time by visiting each page only once.
Cons:
 - Potentially skips unique content of dynamic web pages.
 - Misses any links within multimedia/data files.
 - Misses any links that are not in the HTML source code.
 - Misses any new or changed links that are added after each page is visited.

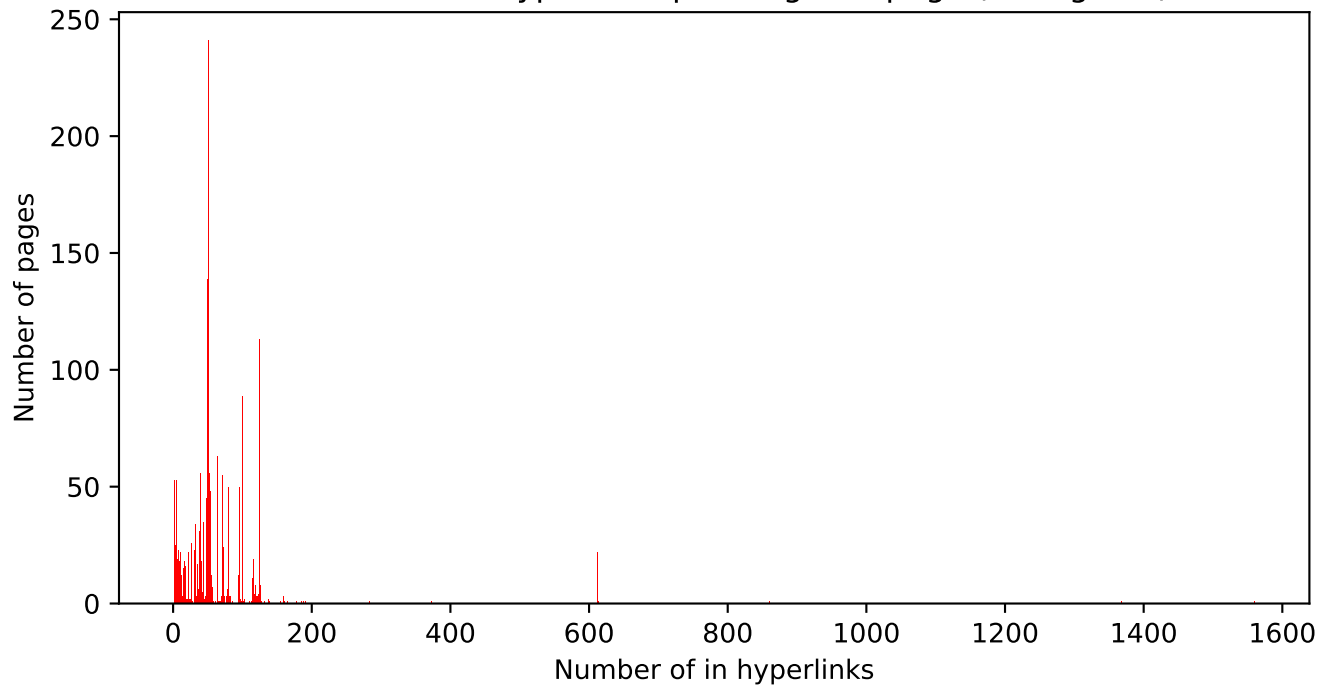Clustering and Diameter analysis (treating the graph as undirected)

Global clustering coefficient: 0.4857047474903301

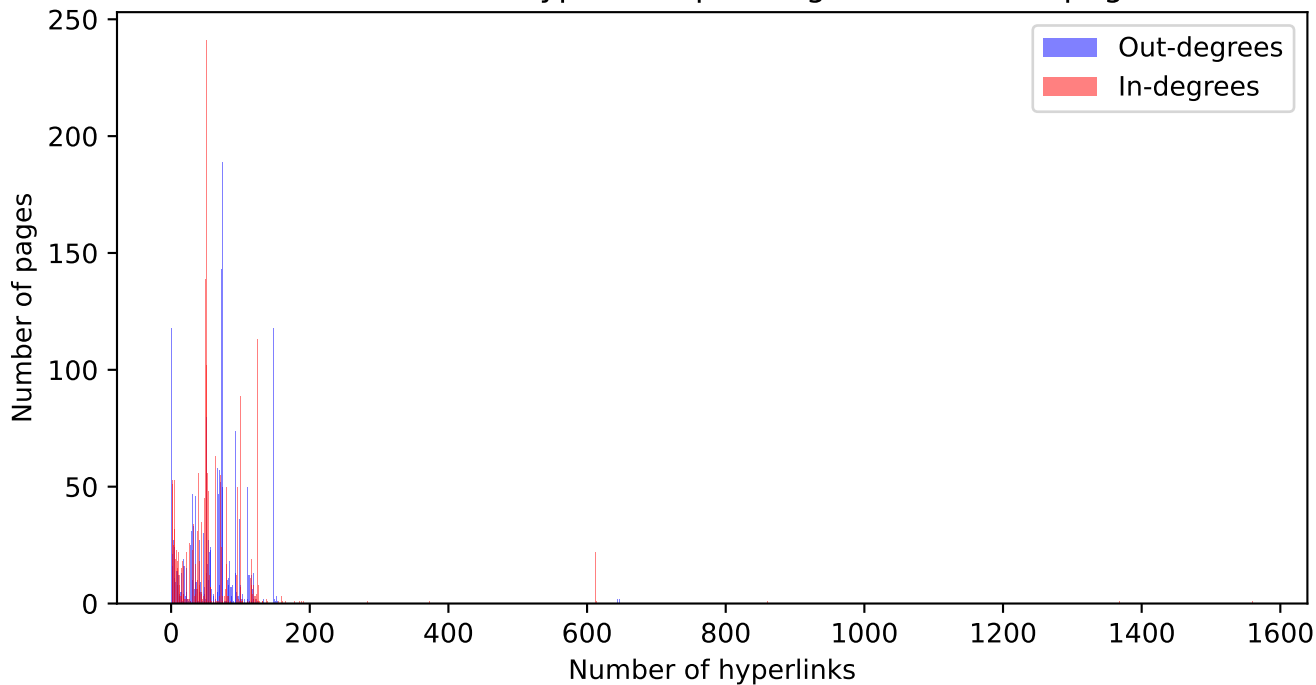Average clustering coefficient: 0.898791351073389

Maximum diameter: 4

Average diameter: 2.257902451225613

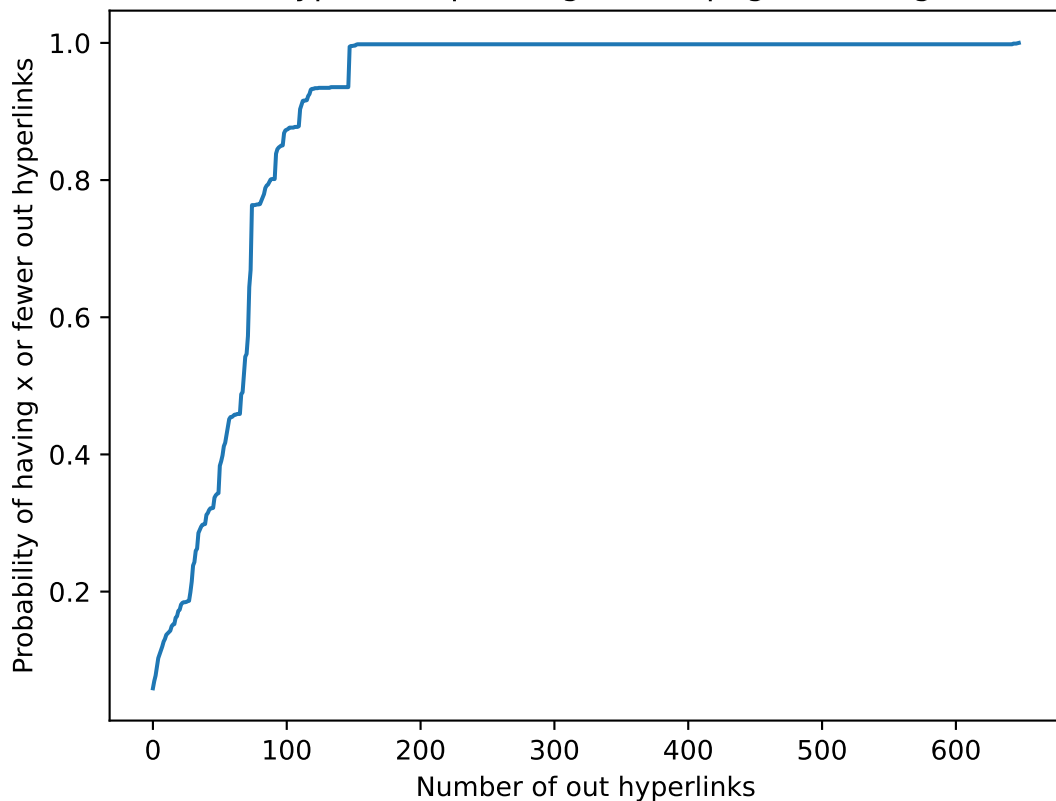Distribution of hyperlinks pointing from a page (out-degrees)

Distribution of hyperlinks pointing to a page (in-degrees)

Distribution of hyperlinks pointing to and from a page

CDF of hyperlinks pointing from a page (out-degrees)

CDF of hyperlinks pointing to a page (in-degrees)