

Abstract

We perform an analysis on GitSED, or the GitHub Socially Enhanced Dataset. In particular, we study the network topology of collaboration networks of repositories predominantly written in specific programming languages via classical network analysis techniques. Additionally, we evaluate the dataset in the context of recommendation systems, exploring its potential for tasks such as developer and repository recommendation. Our results suggest that while modeling Github collaboration as a social network enables insights to be derived using traditional community detection and Graph Neural Network techniques, re-expressions of these relationships may provide better performance depending on the task.

Github Repository: <https://github.com/dankeg/GithubGraphML>

Introduction

In recent years, GitHub has emerged as one of the most popular platforms for collaborative software development, hosting millions of repositories across a wide range of domains. Beyond version control, GitHub also functions as a complex socio-technical ecosystem shaped by collaborative interactions among developers. This rich interaction network, all of which is accessible via the GitHub GraphQL API, presents a unique opportunity for studying collaboration patterns and enabling intelligent tools such as recommendation systems.

One notable study by the Federal University of Minas Gerais utilized this opportunity to create the GitHub Socially Enhanced Dataset (GitSED) [1]. This dataset provides a comprehensive view of developer interactions, which we assume to be primarily defined as pull requests, along with metadata on the associated developers and repositories. GitSED also includes metrics that estimate the strength of social coding collaborations.

In this paper, we leverage GitSED to conduct our own analysis of the GitHub collaboration network. Specifically, we investigate the structural properties of these networks across different programming languages and assess the dataset's suitability for powering recommendation systems focused on developer and repository suggestions.

Overall, our findings indicate that modeling GitHub solely as a simple social network may overlook important structural nuances and contextual factors, which may be better captured by alternative representations such as bipartite networks.

Approach

Our approach consists of the following steps. First, we analyze the programming language-specific collaboration networks provided by GitSED, as well as the aggregated collaboration network, using classical network analysis metrics, community detection techniques, and visualization. Then, we evaluate the suitability of the dataset for recommendation tasks by naively training a graph neural network (GNN) on the combined collaboration network to perform node classification and edge prediction. Finally, out of pure curiosity, we repeat the edge prediction task on a bipartite augmentation of the combined collaboration network where both developers and repositories are represented as nodes.

Dataset

We base our analysis on the GitHub Socially Enhanced Dataset (GitSED), a curated dataset capturing developer interactions on GitHub up through June 2019. This dataset encompasses data from over 32 million developers and 8.5 million repositories, all of which are primarily written in Assembly, JavaScript, Pascal, Python, Ruby, and Visual Basic. In addition to repository and developer metadata, GitSED also includes social collaboration metrics that quantify the strength and nature of developer interactions. For the purposes of this study, however, we focus on the collaboration networks themselves and do not engage with these social metrics. We also do not analyze the Ruby collaboration network due to its disproportionate size to the other collaboration network with respect to its popularity.

In this dataset, A collaboration social network is modeled as a weighted undirected graph $G(V, E^w)$, where V represents developers and an edge $e^{ij} \in E^w$ exists between developers i and j if they have both contributed to the same repository and engaged in some form of collaborative interaction. The edge weight w reflects the estimated strength of collaboration, based on shared interactions such as pull requests.

The specifics of the data contained within the dataset can be described in the relational schema provided by the original GitSED paper (Figure 1). It is important to note that the number_commits metric in the schema does not represent the total number of commits made to a repository during an interaction, but rather the total number of lines modified in the interaction.

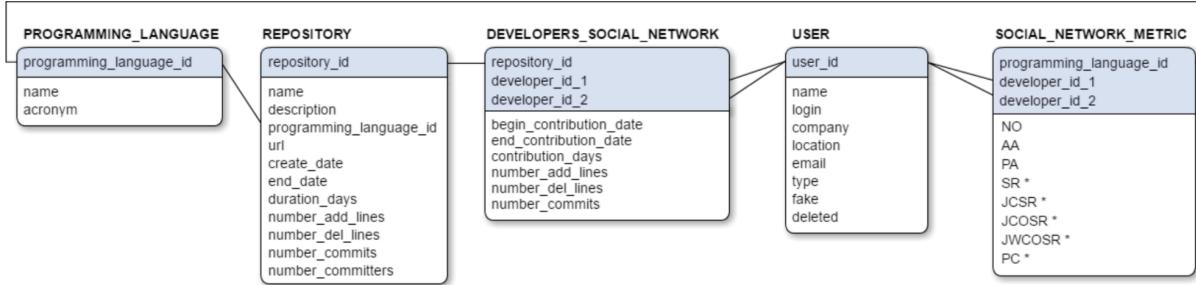


Figure 1: Dataset schema: table names, identifiers (primary keys), attributes and foreign keys (lines connecting tables). *Metric also combined with tieness (T) and resource allocation (RA), i.e., the table has two extra columns for each combination.

Classical Metrics

To better understand the structure of the collaboration networks derived from GitSED, we compute a set of well-established network metrics as follows.

Number of Vertices and Edges. These refer to the number of developers n and interactions m in the graph. This metric is useful in getting a grasp on the scale and density of the network.

$$n = |V|, m = |E|$$

Average Degree. The degree k of a developer v represents the number of other developers the developer have collaborated with, while the average degree \bar{k} is simply the average across all developers $v \in V$. A larger \bar{k} indicates a more interconnected or dense collaboration network.

$$\bar{k} = \frac{2|E|}{|V|}$$

Standard Deviation of Degree. The standard deviation of degree σ measures the variation of degree across all developers. A low standard deviation indicates a more uniform network, while a high one suggests the presence of hubs.

$$\sigma = \sqrt{\frac{1}{n} \sum_{v \in V} (k_v - \bar{k})^2}$$

Global Clustering Coefficient. The global clustering coefficient C evaluates the overall tendency of nodes to form closed triangles. A higher value indicates a stronger presence of community-like structure across the entire network.

$$3 \times \frac{\text{number of triangles}}{\text{number of connected triplets}}$$

Average Local Clustering Coefficient. The local clustering coefficient of a developer quantifies the likelihood that the developer's collaborators have also collaborated with one another. This metric has a similar meaning to the global clustering coefficient, both indicative of the presence of cohesive communities within the network. Let N_v be the neighborhood of node v .

$$\bar{c} = \sum_{v \in V} \frac{|\{(i, j)\}|}{k_v(k_v - 1)} : i, j \in N_v, (j, k) \in E$$

Average Shortest Path Length. The Average Shortest Path Length measures how closely connected developers are within a network, reflecting the efficiency of information flow between them. The all pairwise shortest path distances used to calculate this metric are found using Johnson's algorithm. Let $d(i, j)$ be the shortest path from developer i and j .

$$l = \frac{1}{n(n-1)} \sum_{i \neq j \in V} d(i, j)$$

Maximum Betweenness Centrality. The betweenness centrality quantifies how often a developer appears on the shortest paths between pairs of other developers, serving as a measure of their importance in facilitating communication within the network. Computing the maximum betweenness centrality reveals whether the network has a highly central contributor who acts as a key bridge. Let σ_{st} be the number of shortest paths from s to t , and $\sigma_{st}(v)$ be the number of shortest paths from s to t that pass through v .

$$C_B^{max}(v) = max \left(\left\{ \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}, \forall v \in V \right\} \right)$$

Community Detection

In addition to classical network metrics, we performed community detection to uncover any deeper structural patterns or hidden properties within the GitHub collaboration networks. For this, we used Stochastic Block Models (SBMs), which works by assuming that each node, or developer, belongs to some latent community, and then fitting the model on the maximum likelihood estimation on the probability of an edge forming between any two nodes

based on their community memberships. For a better explanation SBMs consult the graph-tool documentation [2], the library we used for community detection and analysis.

Community Metrics

Once communities were detected using the Stochastic Block Model (SBM), we computed two metrics to quantify and validate some of the structural characteristics and intuition about the GitHub collaboration network.

Average Community Size. This metric reflects the mean number of developers assigned to each detected community. It provides a high-level view of how collaboration is distributed across the network. Smaller values suggest localized, tightly-knit groups of contributors, while larger sizes may indicate broader, more collaborative clusters. This helps us gauge the overall granularity of the inferred community structure. Let K be the number of communities detected and B_k be the k^{th} community.

$$S_{avg} = \frac{1}{K} \sum_{k=1}^K |B_k|$$

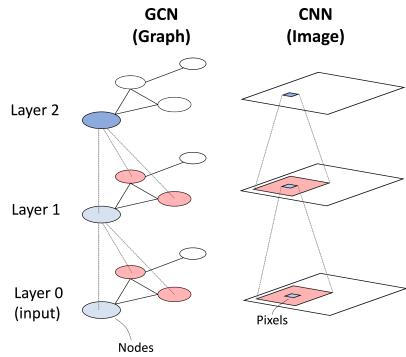
F1 Score: Communities vs. Repositories. To evaluate the alignment between the detected communities and the actual organization of collaboration on GitHub, we compared each developer’s community assignment with their associated repositories. We treated repository membership as ground truth labels and computed the F1 score between communities and repositories. A high F1 score suggests that the detected communities correspond well to real-world collaborative structures, while lower scores may reveal more diffuse or cross-cutting collaboration patterns not fully captured by repository boundaries. Let R be the set of repositories, T_r be the set of developers that contribute to R , and $F1_{kr}$ be the F1 score between the community k and repository r .

$$\begin{aligned} Precision_{kr} &= \frac{|T_r \cap B_k|}{|B_k|}, \quad Recall_{kr} = \frac{|T_r \cap B_k|}{|T_r|} \\ F1_{kr} &= 2 \frac{Precision_{kr} \cdot Recall_{kr}}{Precision_{kr} + Recall_{kr}} \end{aligned}$$

GCN

Fully feed-forward graph neural networks, which function similarly to their non-graph counterparts, suffer from poor efficacy on graph-based tasks. For one, the inflexible structure results in nodes being treated independently, with minimal avenues for aggregating information about graph structure. This poor feature capture as it relates to graph structure results in poor generalization: nodes may appear to be “identical” due to insufficient information capture regarding its neighborhood.

To address these limitations, we leveraged Graph Convolutional Networks [4]. These layers convolve over the structure of the graph, aggregating and integrating information not only from the node provided as input, but also from its neighbors. Each layer of convolution enables nodes an additional hop away to be integrated. Similar to the more familiar task of Convolutional Neural Networks on images, the size of the convolution presents a tradeoff: selecting a kernel the size of an image effectively results in a feed-forward network, and likewise implementing too many GCN layers results in the entire graph being aggregated. GCNs can otherwise be constructed identically to feed-forward neural networks, with node-specific features being extracted such as the number of commits, and the number of contribution days.



GNN

Node Classification

To further investigate the relationship between graph structure and the behavior of Github users, we constructed a node classification problem. Given the primary link strength features of the graph, commits and contribution days, combined with a bias term, we seek to predict the programming languages used by individual users (represented by nodes).

While this problem has limited scope in application upon real-world Github social networks, as in order for links to exist between nodes the particular repositories a user contributes to and thus the languages they use is trivial to derive, this task seeks to solidify the intuition developed from the community detection analysis: in that users of different languages have significant differences in contribution behavior. A model that accurately classifies nodes supports this hypothesis, while poor performance suggests a lack of such a link, and necessitates further investigation.

Given the large number of combinations of programming languages, with a total of 63 unique combinations, we opted to model this as a multi-label classification: with the model predicting the usage of each language individually. Thus, PyTorch's BCEWithLogitsLoss was leveraged during training. Two convolutional layers will be used, with a hidden dimension of size 16 and ReLU for activation between the layers. Thus, per some node that we seek to classify, the model will aggregate information regarding nodes at most 2 hops away. Training was performed using a random split of nodes into a train and test set, using a randomly generated mask on an adjacency matrix derived from the graph-tools representation.

Generic Link Prediction

Given the representation of Github collaboration as a social network, creating a recommendation engine is a useful task. Leveraging GCNs, we can construct this as a link prediction problem: given two nodes, does a link exist between them? If so, users should both be recommended to work together on a particular repository. If not, it's unlikely that a collaboration between the users will be effective. Thus, given the primary link strength features of the graph, commits and contribution days, combined with a bias term, we seek to predict whether two users share a link, and thus whether they should be recommended to collaborate. This is framed as a simple binary classification task, and thus also leverages BCEWithLogitsLoss.

An initial attempt at this task resulted in instability and low performance, as discussed in our presentation, leading to investigating Bipartite Link Prediction as an alternative. Further investigation revealed several causes for these issues, primarily linked to the representation of the graph, in terms of training, test, and the encoding of the graph by the model, but also limitations in the previous model construction. As discussed in our analysis, addressing these representational challenges led to significant improvements in model performance in both the generic and bipartite approaches.

Starting with model design, the number of hidden dimensions was increased from 16 to 64 to improve feature extraction. Additionally, a drop-out layer was added to improve generalization by addressing overfitting. Improvements were also made to the process for generating the train-test split, with the Pytorch method RandomLinkSplit leveraged over a naive randomly generated mask. This ensures that graph structure is kept wholly intact, enabling training on a more accurate representation, rather than ignoring the existence of nodes not in the split.

A key difference between our introductory node classification task and link prediction is the training of embeddings, which involves an encoder-decoder architecture: enabling a richer representation of graph behavior and structure. The original representation, effectively a single tensor layer which internally utilized a simple dot-product based similarity calculation, lacked the sufficient complexity to effectively represent the graph. As a result, this has been replaced with two feed-forward multi-layer perceptrons, one for encoding and one for decoding, with the goal of being able to develop more nuanced and effective representations.

Bipartite Link Prediction

Initial investigations into Generic Link Prediction resulted in poor performance, with the model displaying high instability during training, and exceptionally low performance during evaluation: with a less than 50% accuracy on eval despite a greater than 50% accuracy on test. Thus, a bipartite representation of the network was investigated for link prediction in order to determine whether a different representation would result in greater effectiveness.

Conversion of the original social network to a bipartite representation is a straightforward process. In the original network, each edge between nodes represents a repository which both users contribute to. This edge is broken out into a separate node, with two new edges constructed from each of the user nodes. This process is repeated for all edges, creating new nodes for a specific repository if one does not already exist.

Due to the structural limitations of Bipartite networks, our pre-existing model was modified to leverage SAGEConv layers, as opposed to generic GraphConv layers [3]. Due to an aggregation strategy which in part leverages the mean of features, rather than the traditionally used sum, it is better able to generalize to atypical graph structures including but not limited to bipartite networks. Model design and training is otherwise unmodified from the improved generic approach discussed above.

Results

Classical & Community Analysis

Classical Metrics																
	Graph							Largest Component								
Metrics	Number vertices	Number edges	Average degree	Degree standard deviation	Global clustering coefficient	Average clustering coefficient	Average shortest path length	Maximum betweenness centrality	Number vertices	Number edges	Average degree	Degree standard deviation	Global clustering coefficient	Average clustering coefficient	Average shortest path length	Maximum betweenness centrality
Assembly	9,353	22,044	4.71378	0.13845	0.62004	0.38552	∞	0.00206	537	6,480	24.13408	1.83893	0.61566	0.05788	3.94714	0.62884
JavaScript	1,152,467	5,905,002	10.24759	0.06313	0.27539	0.66617	N/A	N/A	499,588	5,218,099	20.88961	0.14419	0.27501	0.50935	N/A	N/A
Pascal	4,241	12,008	5.66282	0.20567	0.77347	0.41798	∞	0.01127	719	3387	9.42142	0.56440	0.76611	0.12523	7.60052	0.39362
Perl	41,141	280,291	13.62587	0.19875	0.49981	0.66179	∞	0.03482	18,316	224,752	24.54160	0.40455	0.47106	0.45767	5.31946	0.17571
Python	711,232	9,252,164	26.01729	0.20780	0.45395	0.64705	N/A	N/A	328,631	8,485,915	51.64403	0.44115	0.44955	0.50847	N/A	N/A
Visual Basic	7,419	12,327	3.32309	0.10336	0.39889	0.35823	∞	0.00033	179	2,081	23.25140	1.86942	0.54590	0.04761	2.97207	0.56890
Combined	1,796,692	15,483,836	17.23594	0.09484	0.41031	0.67392	N/A	N/A	874,309	14,534,055	33.24695	0.19340	0.41021	0.54185	N/A	N/A

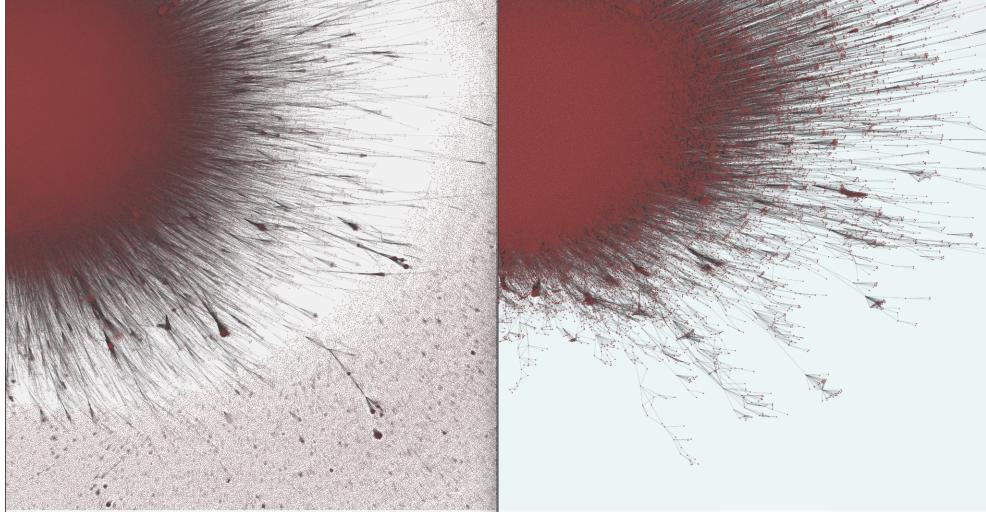
After running our analyses, we computed all the metrics outlined in the approach section of this paper, as well as other visualizations, for each programming language-specific collaboration network provided by the GitSED dataset plus the combined network. From these metrics and visualizations, there are several interesting observations and insights could be taken from the data.

It is quite obvious, both from intuition and from the visualizations, that each language-specific collaboration network exhibits distinctive structural properties. For instance, the Assembly collaboration network stands out with a particularly stands out with a particularly unusual characteristic—its average local clustering coefficient (0.38552) is nearly half of the global clustering coefficient (0.62004), a discrepancy not observed in any of the other networks. This suggests that while some global structure or community patterns may exist, local neighborhoods of developers are relatively sparse or fragmented. This pattern is likely a reflection of the collaborative nature of Assembly itself: as a low-level language, it is notoriously difficult to work with in teams, which may result in many small, disconnected components where collaboration is limited or isolated.

In contrast, the JavaScript collaboration network presents a very different structure. Compared to networks that are on a similar scale, such as Python, the JavaScript collaboration network has a much more uniformly distributed and interconnected topology, shown by its relatively low global clustering coefficient (0.27539) and a visual comparison between the JavaScript and Python network. This likely reflects JavaScript’s widespread, if not only use, in web development where open-source contributions and teamwork are more common, leading to greater global connectivity.

Community Metrics

Metrics	Average Community Size	Maximum F1 Score
Assembly	222.69048	0.50720
JavaScript	1293.45342	0.24396
Pascal	121.17143	0.50730
Perl	252.39877	0.44514
Python	896.88777	0.25692
Visual Basic	190.23077	0.41598
Combined	1181.25707	0.24273

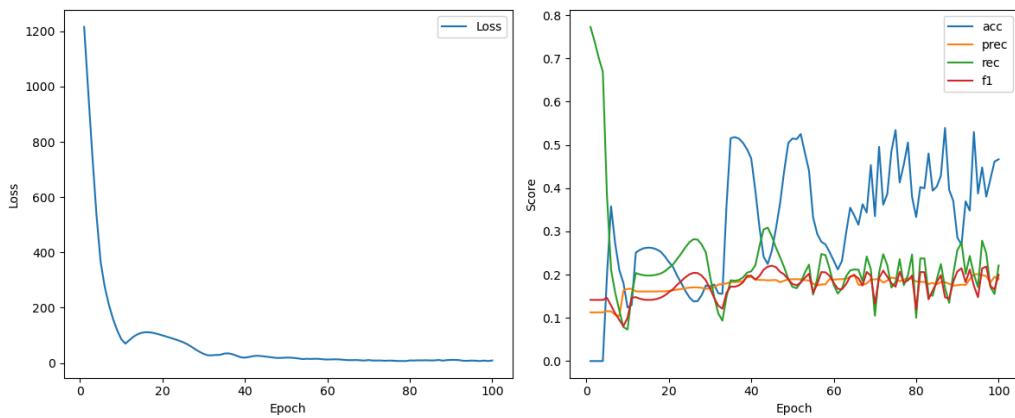


Another clear pattern that can be observed from the metrics is that these collaboration networks share several properties with stereotypical social networks. In particular, they all exhibit characteristically high clustering coefficients and contain large connected components relative to their overall density. Interestingly, however, the maximum betweenness centrality appears to be negatively correlated with the scale of the graph, suggesting that these networks do not follow a classic hub-and-spoke structure. In larger collaboration networks, influence and connectivity seem to be more evenly distributed rather than concentrated in a few central developers.

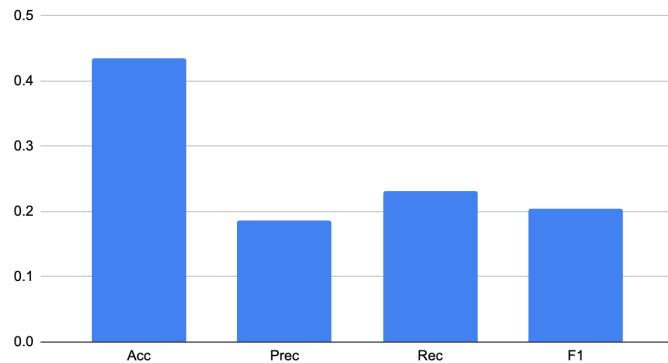
Finally, through our community detection analysis (metrics in top right corner), we observed that community-repository alignment is high within smaller communities, but this correlation diminishes significantly as community size increases. This suggests that while small groups of developers tend to collaborate within the same repositories, larger communities may span multiple projects, reflecting more complex or diffuse collaboration patterns.

Node Classification

Training was performed on a M1 Max Macbook Pro, with an average training time of 10 minutes. Random restarts had a significant impact on model performance, and thus to provide a baseline the performance of the lowest performance model was analyzed. The results are summarized in the following plots:



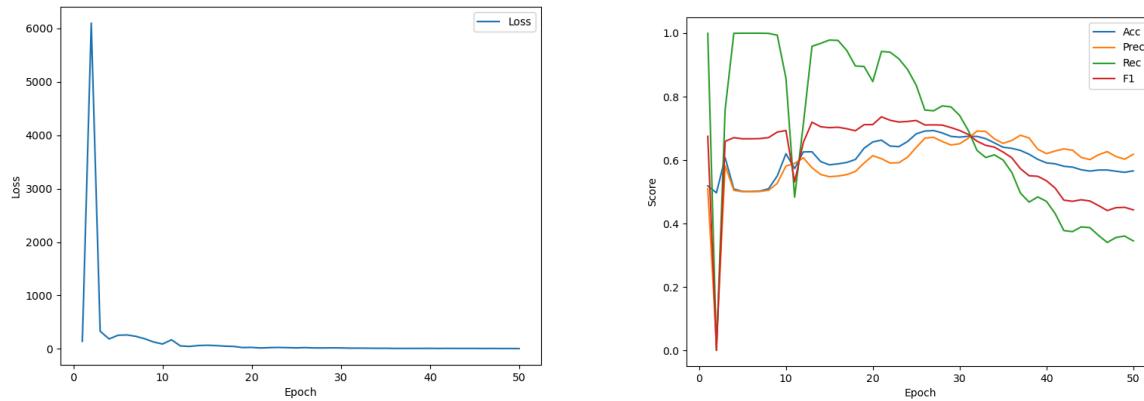
Node Classification Evaluation Metrics



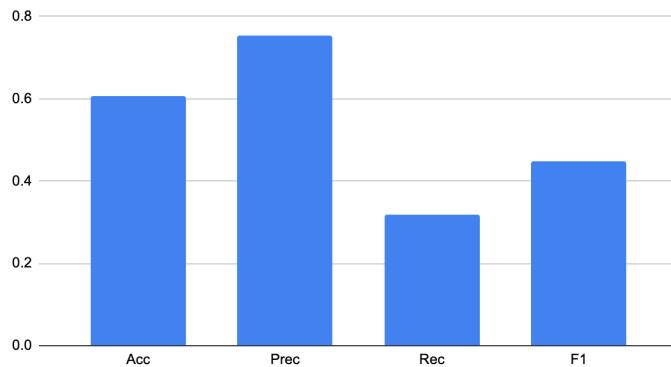
Overall, considering the number of combinations of labels that each node can be classified as, the model performed quite well, with an accuracy of 0.4351. In comparison, the model with the highest performance had an accuracy of 0.55, which is a meaningful but not overly large gap. Both training metrics and loss demonstrate a degree of oscillation, which does not meaningfully reduce using a smaller learning rate. However, more performant models demonstrated this behavior to a smaller degree, indicating that it may be representative of poor feature extraction. This is supported by the atypical loss curve, which has a large plummet followed by marginal improvements rather than a smooth convergence.

Generic Link Prediction

Training was performed on a M1 Max Macbook Pro, with an average training time of 50 minutes. Random restarts once again had a noticeable impact on model performance, and thus to provide a baseline the performance of the lowest performance model was analyzed. The results are summarized in the following plots:



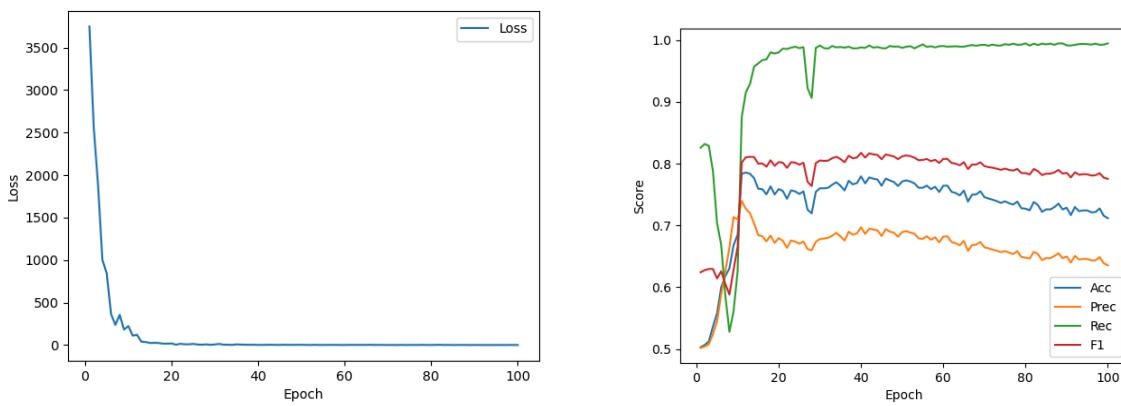
Generic Link Pred. Evaluation Metrics

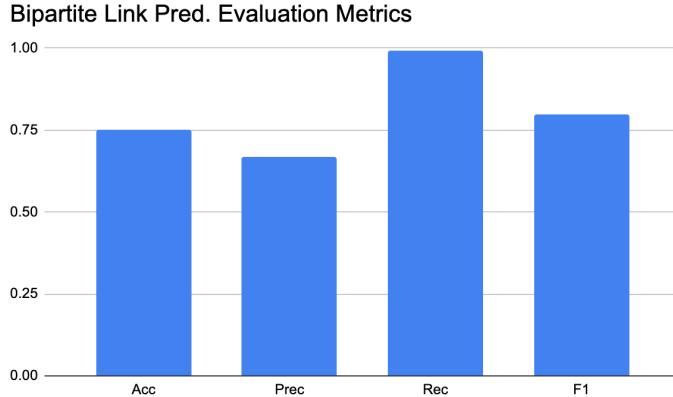


Analyzing the results, we see a clear dramatic improvement from the previous performance on this task, with the most performant model having an accuracy of 0.73. Additionally, examining the training vs test metrics, we see a minimal difference in performance. This indicates that the model now generalizes significantly better, in comparison to the nearly or worse than random predictions seen prior to revisions. Similar to the node classification model, the worst performing model has strange patterns in the loss, with both the loss and training metrics containing oscillations. For instance, the first few iterations go from low loss to an extremely high value which is reflected in a dip in the training metrics, before continuing to decrease: removing this outlier the loss improvement would appear as a very shallow curve. Additionally, regardless of these oscillations, precision remains significantly higher than the other metrics, which may suggest imbalances in the generation of positive and negative edges for prediction.

Bipartite Link Prediction

Training was performed on a M1 Max Macbook Pro, with an average training time of 50 minutes. Given that Bipartite Link Prediction uses the same link splitting technique as the generic version, random restarts continue to have a noticeable impact on model performance. The results are summarized in the following plots:





The most performant bipartite link prediction model had an accuracy of 0.83, again demonstrating a moderate range of effectiveness as previously seen. Comparing the performance of this approach to the revised generic approach and the bipartite approach pre-revision, it is clear that both the model and representational corrections, in addition to the bipartite representation, have a noticeable positive impact on efficacy. Besides having good training metrics which generalize well to the test set, this approach is considerably more stable, with the magnitude of oscillations in loss and training metrics significantly reduced. A potential concern is the high recall, which may also be due to imbalances in the divisions of positive and negative edges.

Conclusion

Based on this work, there are several avenues for future exploration. The variability of GCN performance suggests that random restarts and the random test-train data split seems to have a significant impact on model performance and generalizability. This behavior persisted for RandomLinkSplit despite its associated performance benefits, indicating that a more intentional approach is warranted. To leverage this finding to improve model performance, we can investigate alternative methods of splitting: such as splitting by cluster rather than individual nodes. Such an approach may enable the network to better interpret the different morphologies which may naturally exist, whereas the random split currently being used would disrupt this behavior.

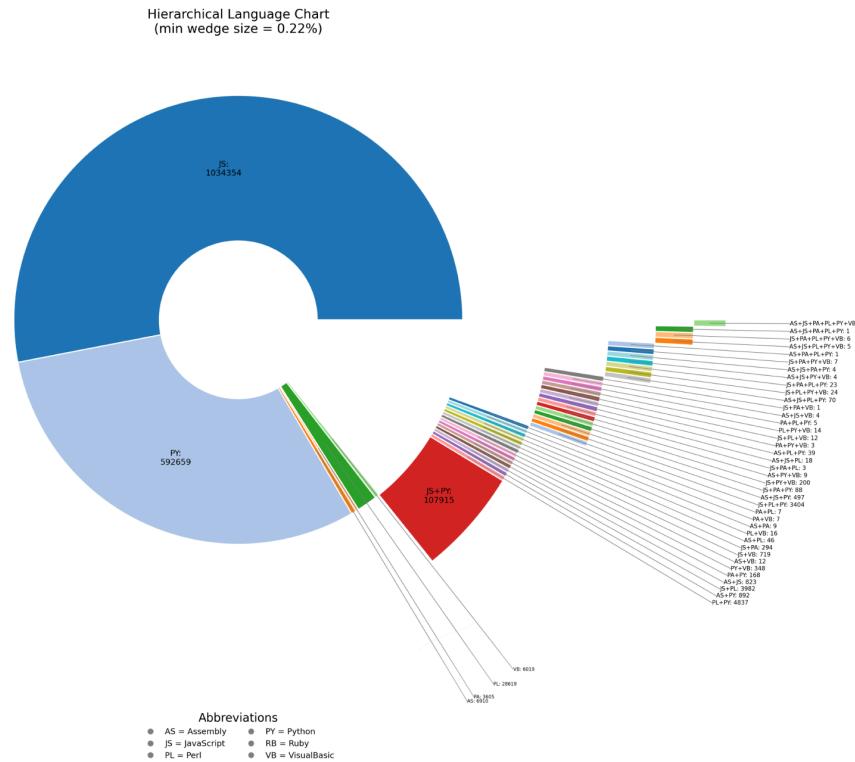
Additionally, our analysis identifies that for GNN related tasks, a typical social network representation may not be the most appropriate. However, initial community detection investigations helped illuminate crucial differences in the behavior of different collaborator networks, such as the highly connected nature of JavaScript compared to Python, regardless of their similar sizes. Thus, it may be worthwhile to further explore community detection using more complicated algorithms, as well as analyses especially suited for social networks such as cascading behavior.

References

- [1] Natércia A. Batista, Michele A. Brandão, Gabriela B. Alves, Ana Paula Couto da Silva, and Mirella M. Moro. 2017. Collaboration strength metrics and analyses on GitHub. In Proceedings of the International Conference on Web Intelligence (WI '17). Association for Computing Machinery, New York, NY, USA, 170–178.
<https://doi.org/10.1145/3106426.3106480>.
- [2] Peixoto, T. P. (2025). *Inferring modular network structure*. In *graph-tool 2.97 documentation*. Retrieved April 25, 2025, from <https://graph-tool.skewed.de/static/doc/demos/inference.html>.
- [3] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In Advances in Neural Information Processing Systems 30 (NeurIPS '17). Curran Associates, Inc., Red Hook, NY, USA, 1024–1034. <https://doi.org/10.48550/arXiv.1706.02216>
- [4] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In Proceedings of the International Conference on Learning Representations (ICLR '17).
<https://doi.org/10.48550/arXiv.1609.02907>

Appendix

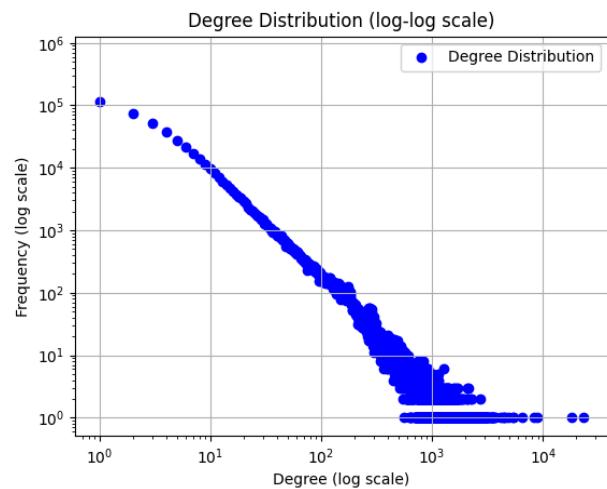
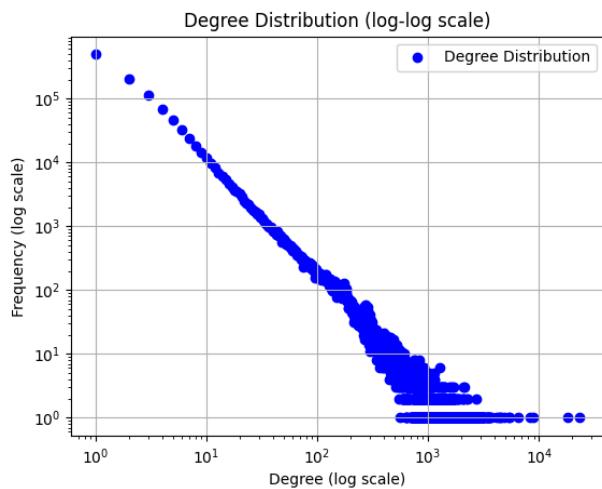
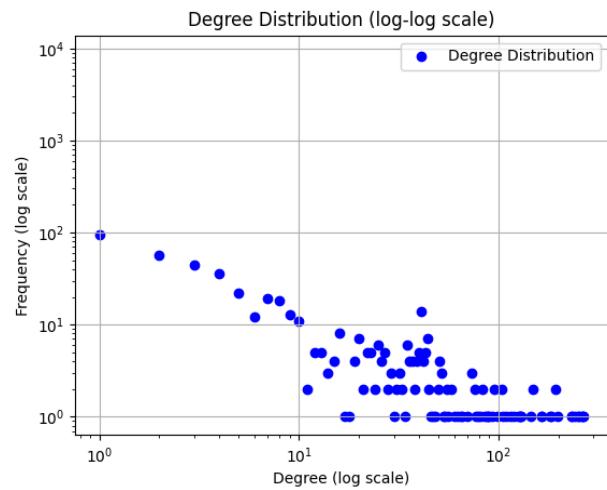
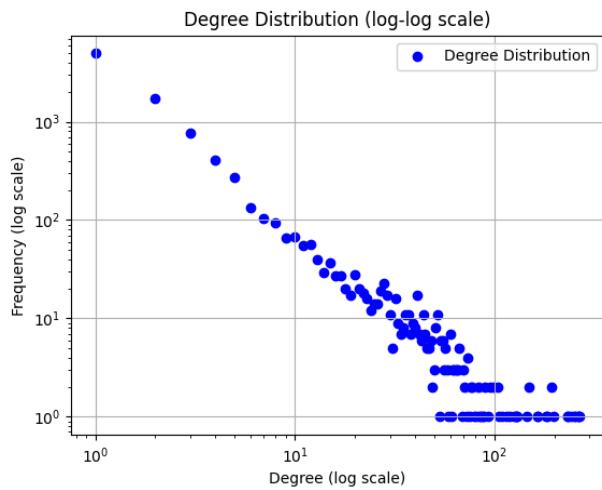
Developer Repository Language Contribution Distribution



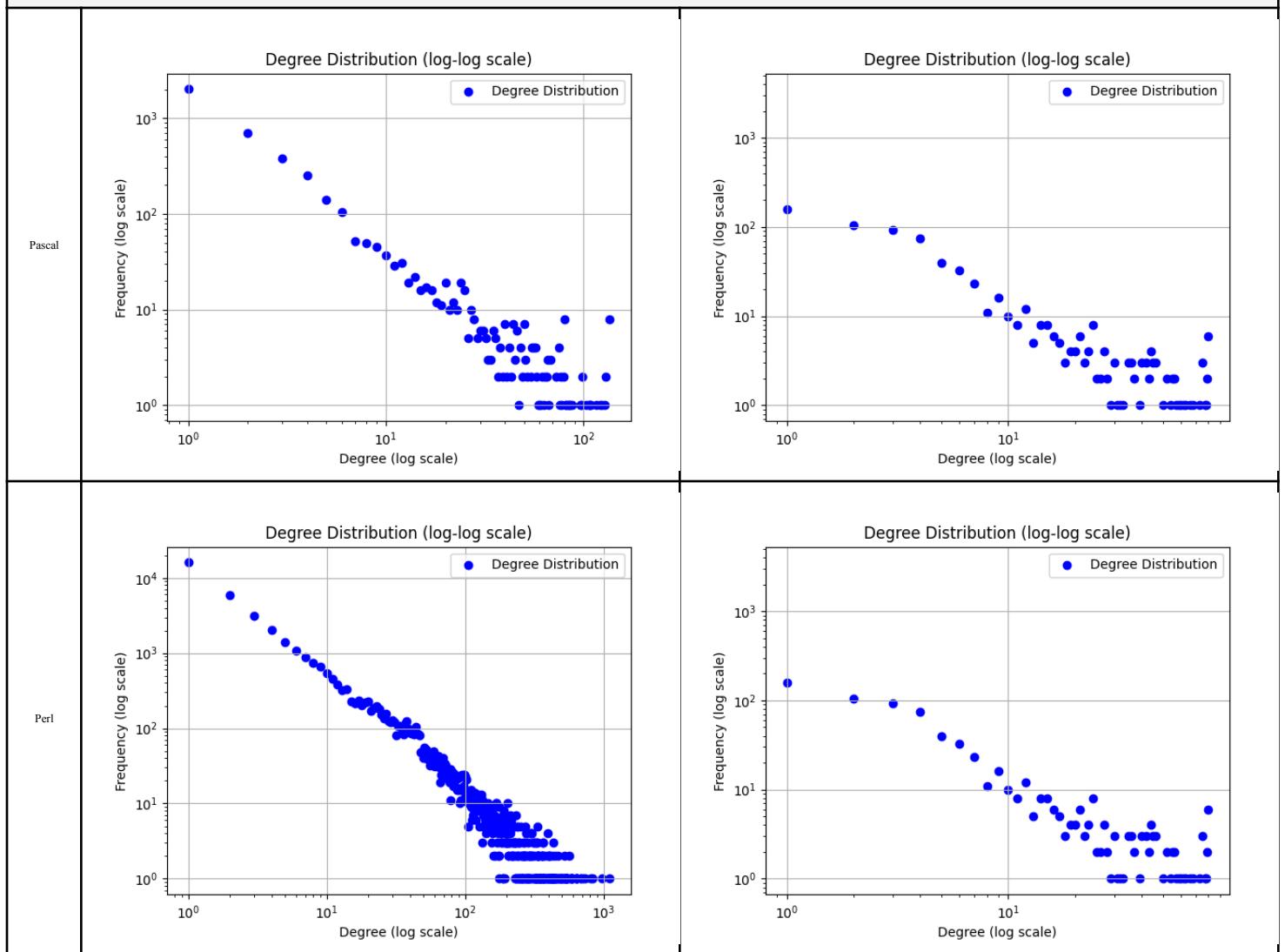
Classical Metrics

	Graph							Largest Component								
Metrics	Number vertices	Number edges	Average degree	Degree standard deviation	Global clustering coefficient	Average clustering coefficient	Average shortest path length	Maximum betweenness centrality	Number vertices	Number edges	Average degree	Degree standard deviation	Global clustering coefficient	Average clustering coefficient	Average shortest path length	Maximum betweenness centrality
Assembly	9,353	22,044	4.71378	0.13845	0.62004	0.38552	∞	0.00206	537	6,480	24.13408	1.83893	0.61566	0.05788	3.94714	0.62884
JavaScript	1,152,467	5,905,002	10.24759	0.06313	0.27539	0.66617	N/A	N/A	499,588	5,218,099	20.88961	0.14419	0.27501	0.50935	N/A	N/A
Pascal	4,241	12,008	5.66282	0.20567	0.77347	0.41798	∞	0.01127	719	3387	9.42142	0.56440	0.76611	0.12523	7.60052	0.39362
Perl	41,141	280,291	13.62587	0.19875	0.49981	0.66179	∞	0.03482	18,316	224,752	24.54160	0.40455	0.47106	0.45767	5.31946	0.17571
Python	711,232	9,252,164	26.01729	0.20780	0.45395	0.64705	N/A	N/A	328,631	8,485,915	51.64403	0.44115	0.44955	0.50847	N/A	N/A
Visual Basic	7,419	12,327	3.32309	0.10336	0.39889	0.35823	∞	0.00033	179	2,081	23.25140	1.86942	0.54590	0.04761	2.97207	0.56890
Combined	1,796,692	15,483,836	17.23594	0.09484	0.41031	0.67392	N/A	N/A	874,309	14,534,055	33.24695	0.19340	0.41021	0.54185	N/A	N/A

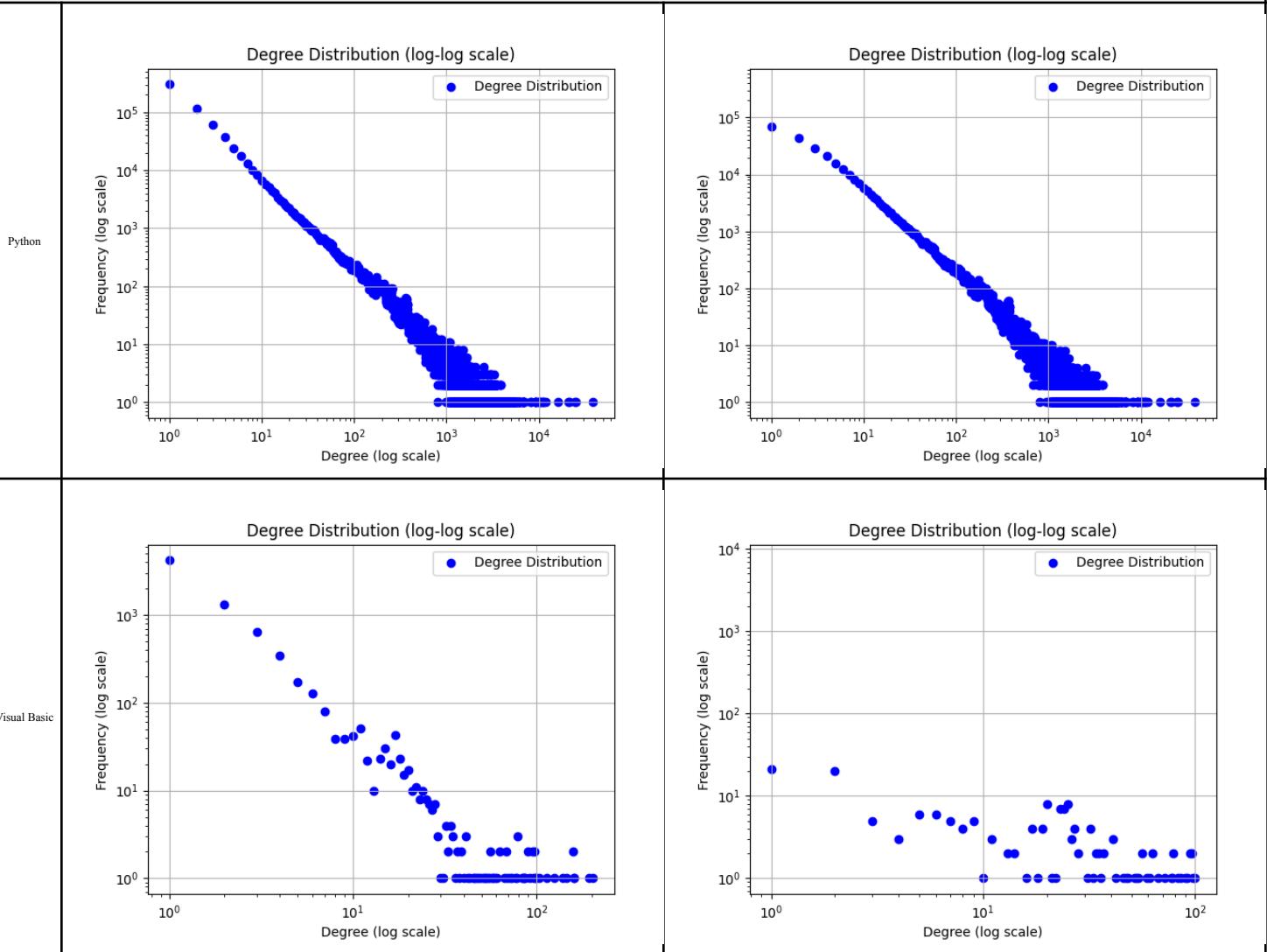
Degree Distribution Plots



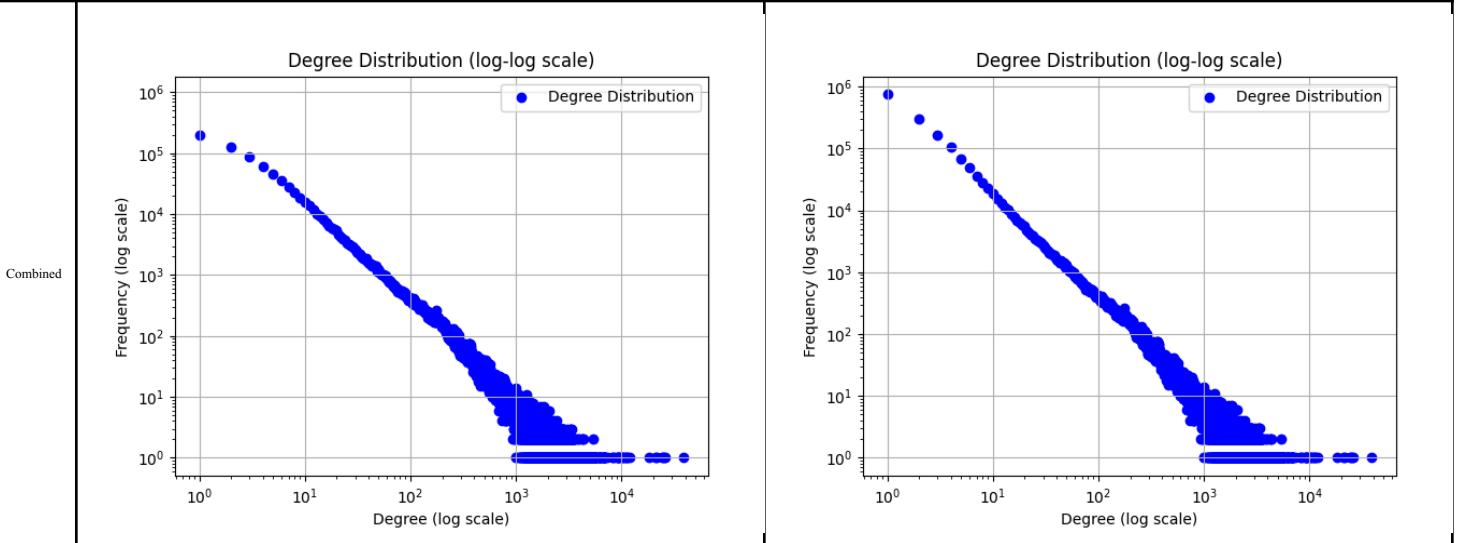
Classical Metrics



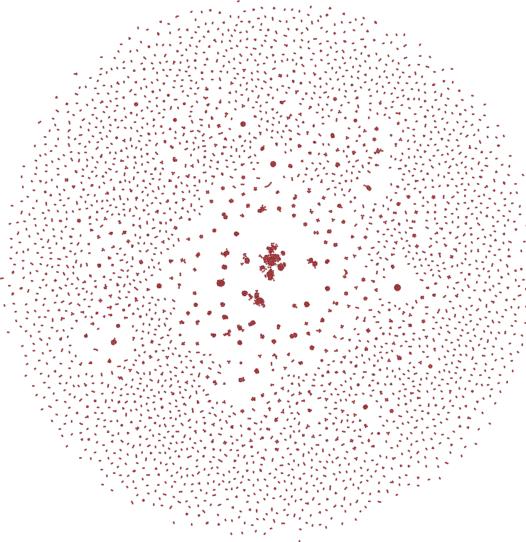
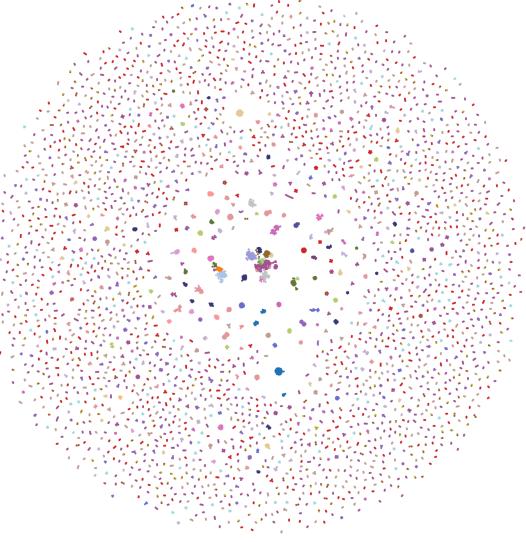
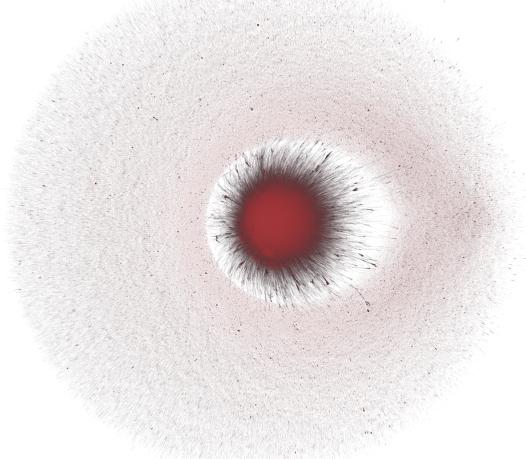
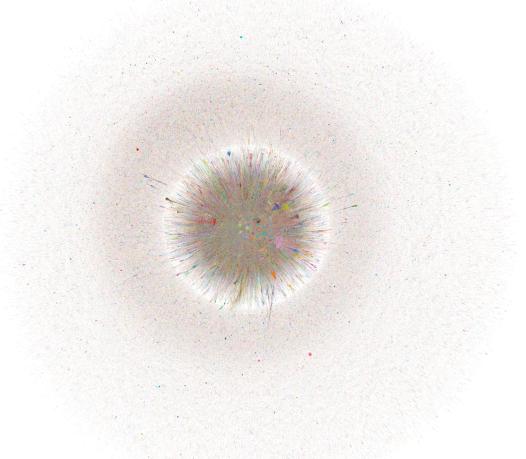
Classical Metrics



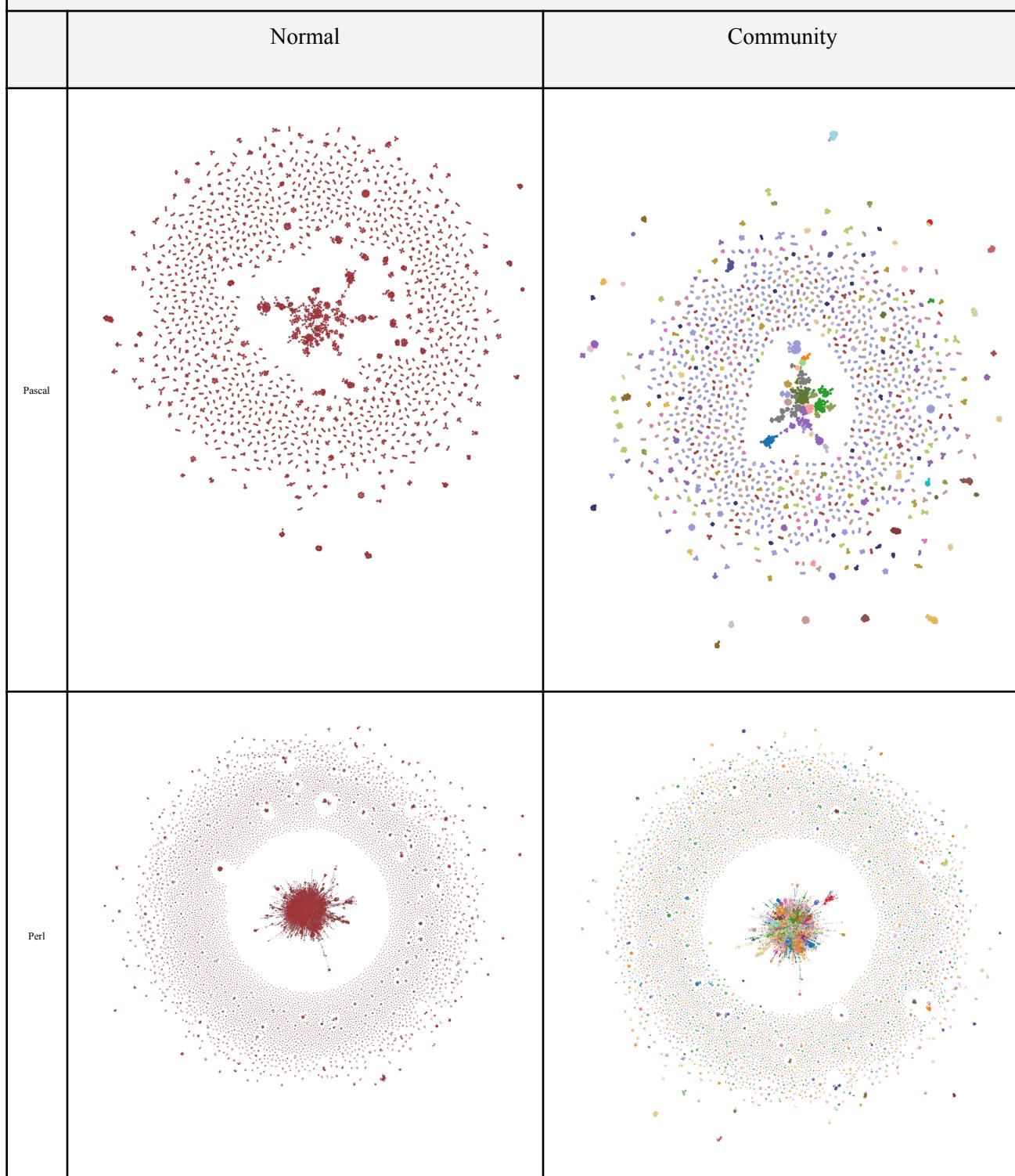
Classical Metrics



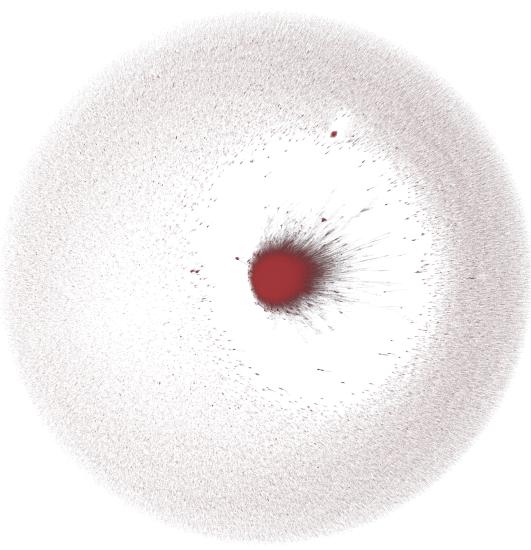
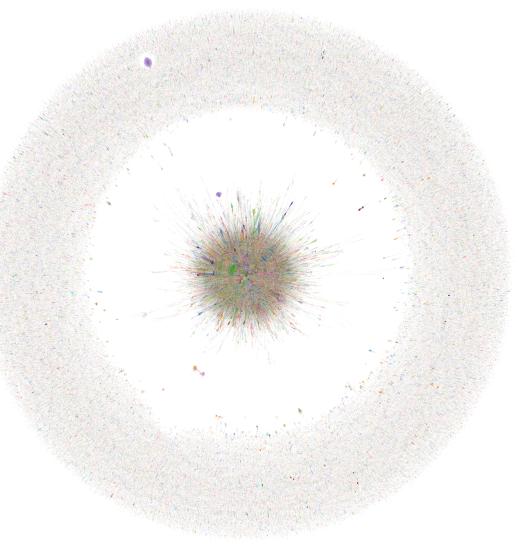
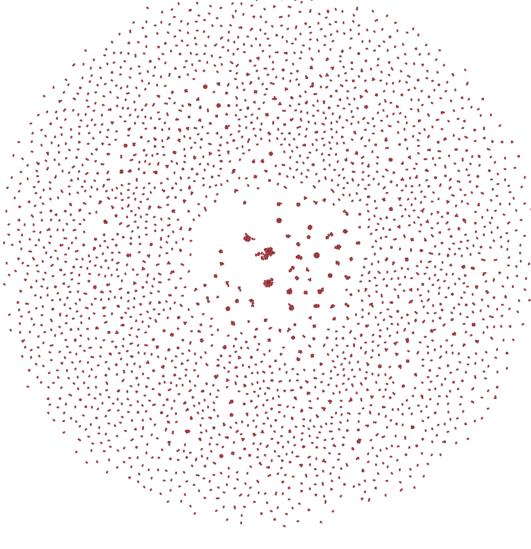
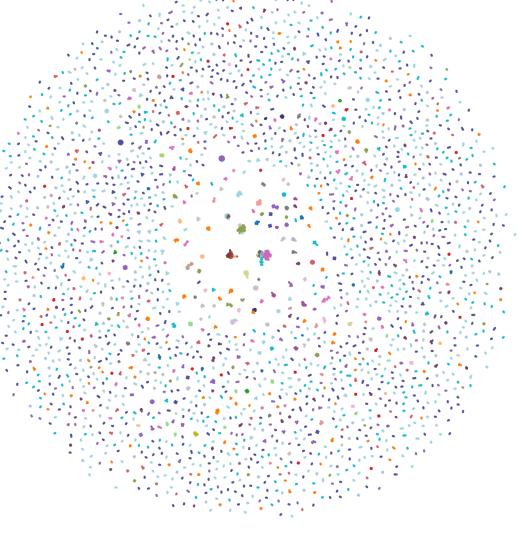
Network Visualizations

	Normal	Community
Assembly		
JavaScript		

Network Visualizations



Network Visualizations

	Normal	Community
Python		
Visual Basic		

Network Visualizations

