

WUSTL Networks: Finding Student Communities

By Min Choi, Zac Christensen, Mason Hall





Representing the Data





Student Responses

We gathered initial data by 1. **asking** students at dining halls around campus about where they just came from, and 2. sending out a **survey** asking for the same information

Potential issues with this method:

- All data was collected using fixed sources and dining halls as targets
- This was on purpose to verify our hypothesis that dining halls on campus are central nodes of student flows on campus.
- We combined halls into areas, and didn't ask about study spaces without dining
- In reality, student flow is not always through dining halls, and often happens like:
 - Study area -> study area
 - Study area -> dining area
 - Dining area -> study area
 - Dining area -> dining area



Creating files

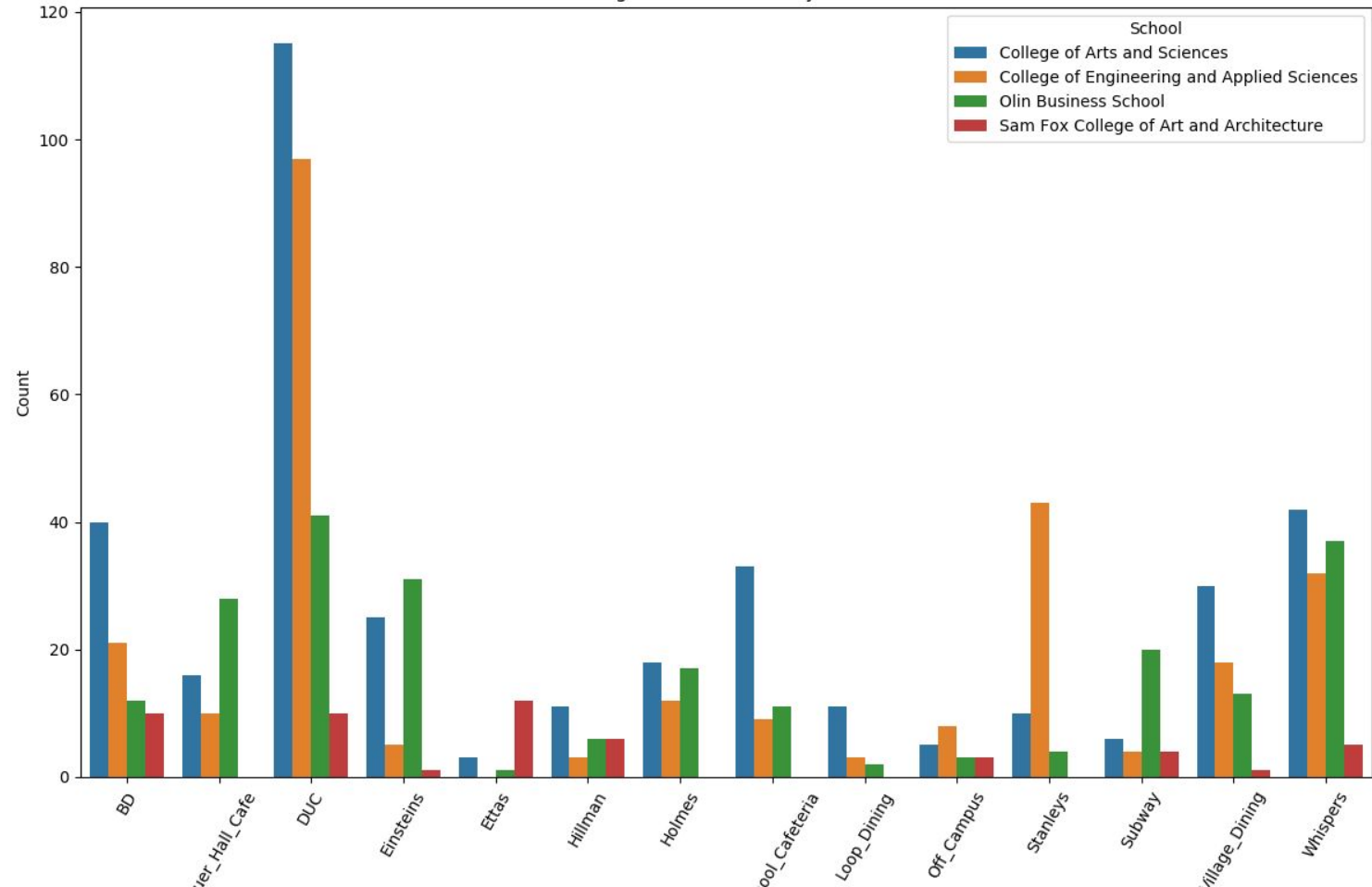
Pandas and CSVs were used to manipulate and store all of the data that was found

Allowed for several things:

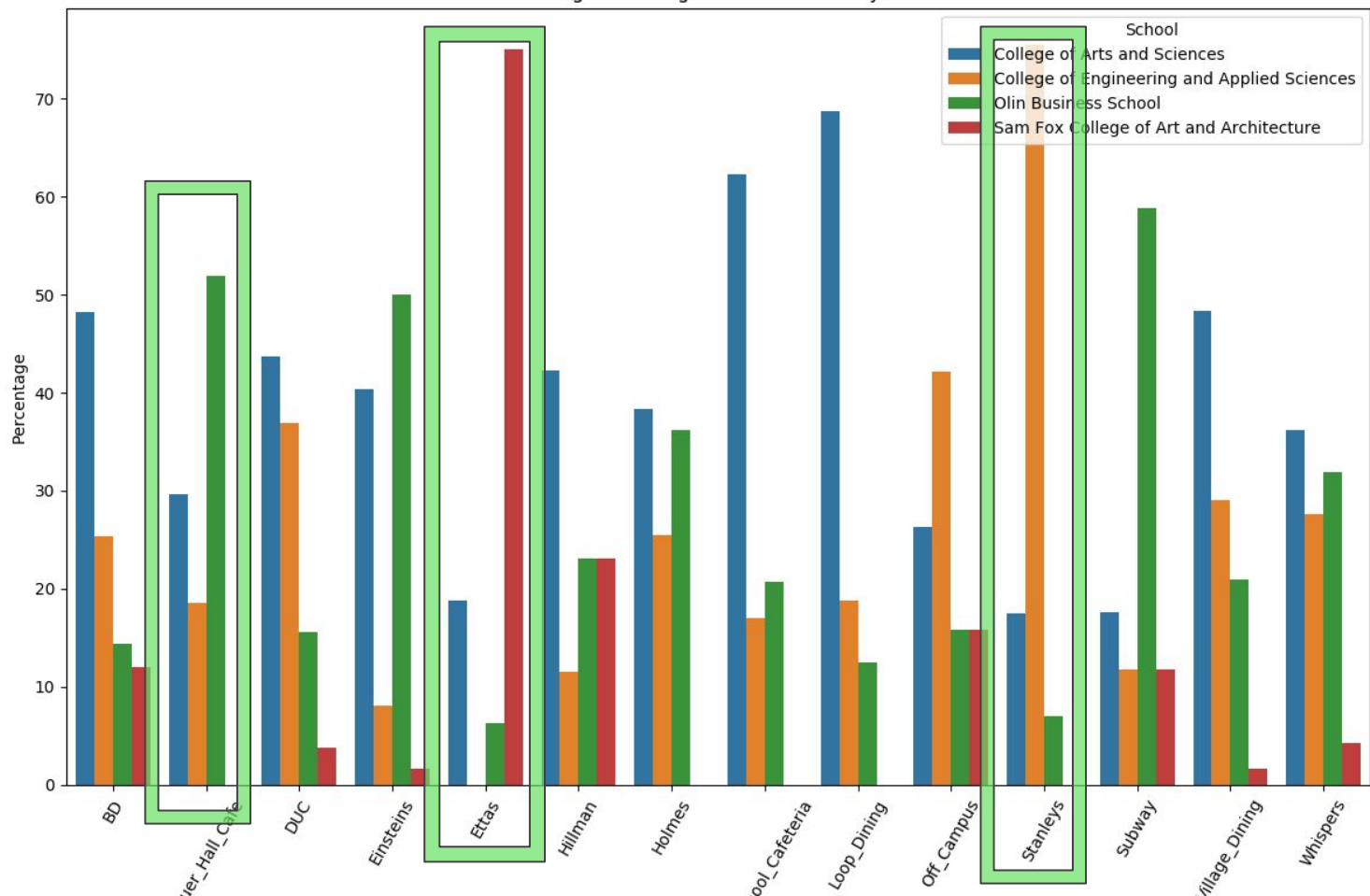
- Easy manipulation of groups/data points
- Removal of empty nodes/responses
- Easy transfer into graphing/network packages
- Familiarity through Python

NetworkX was used due to the vast range of analysis possible

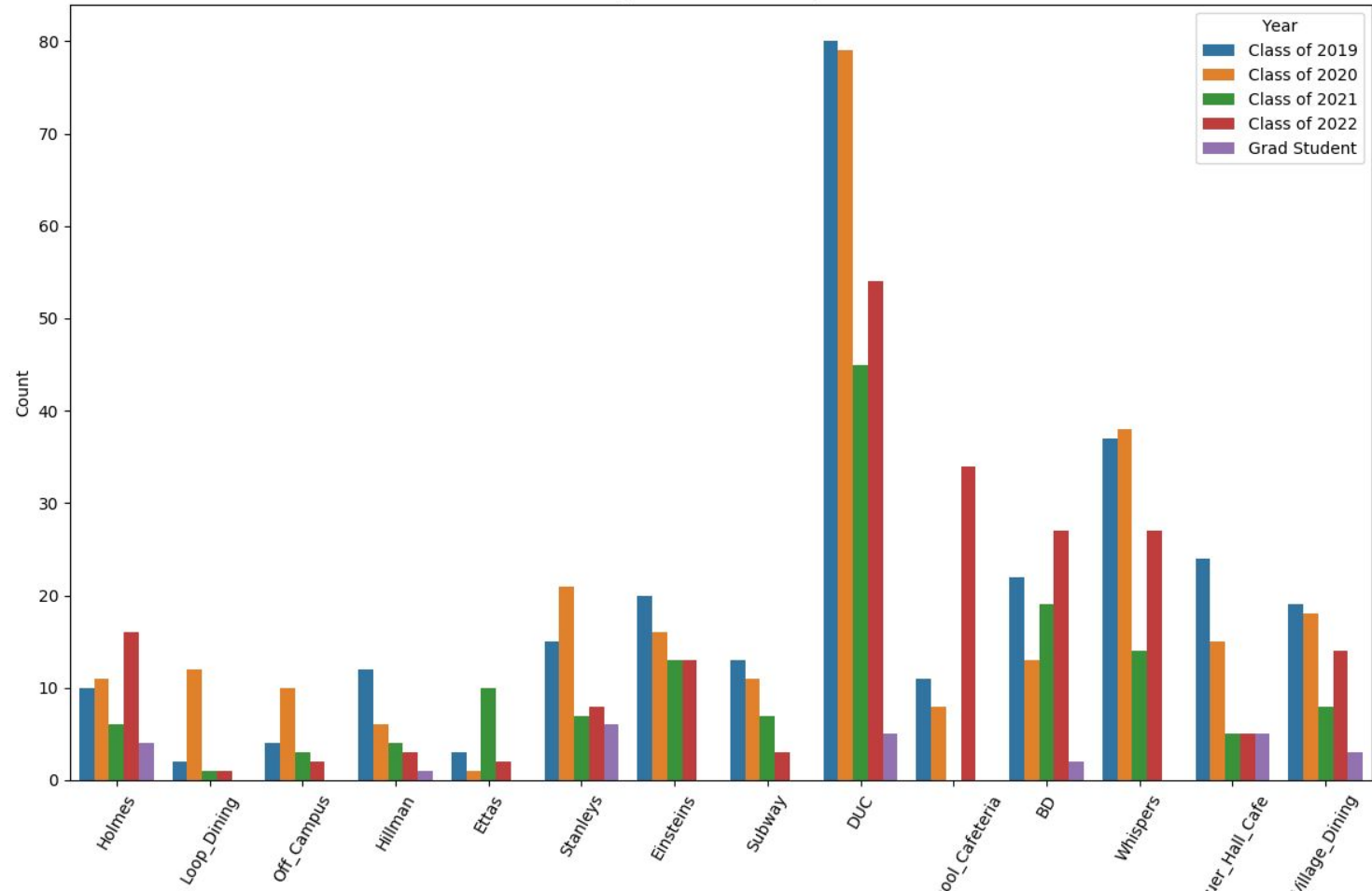
Dining Hall Attendance by School



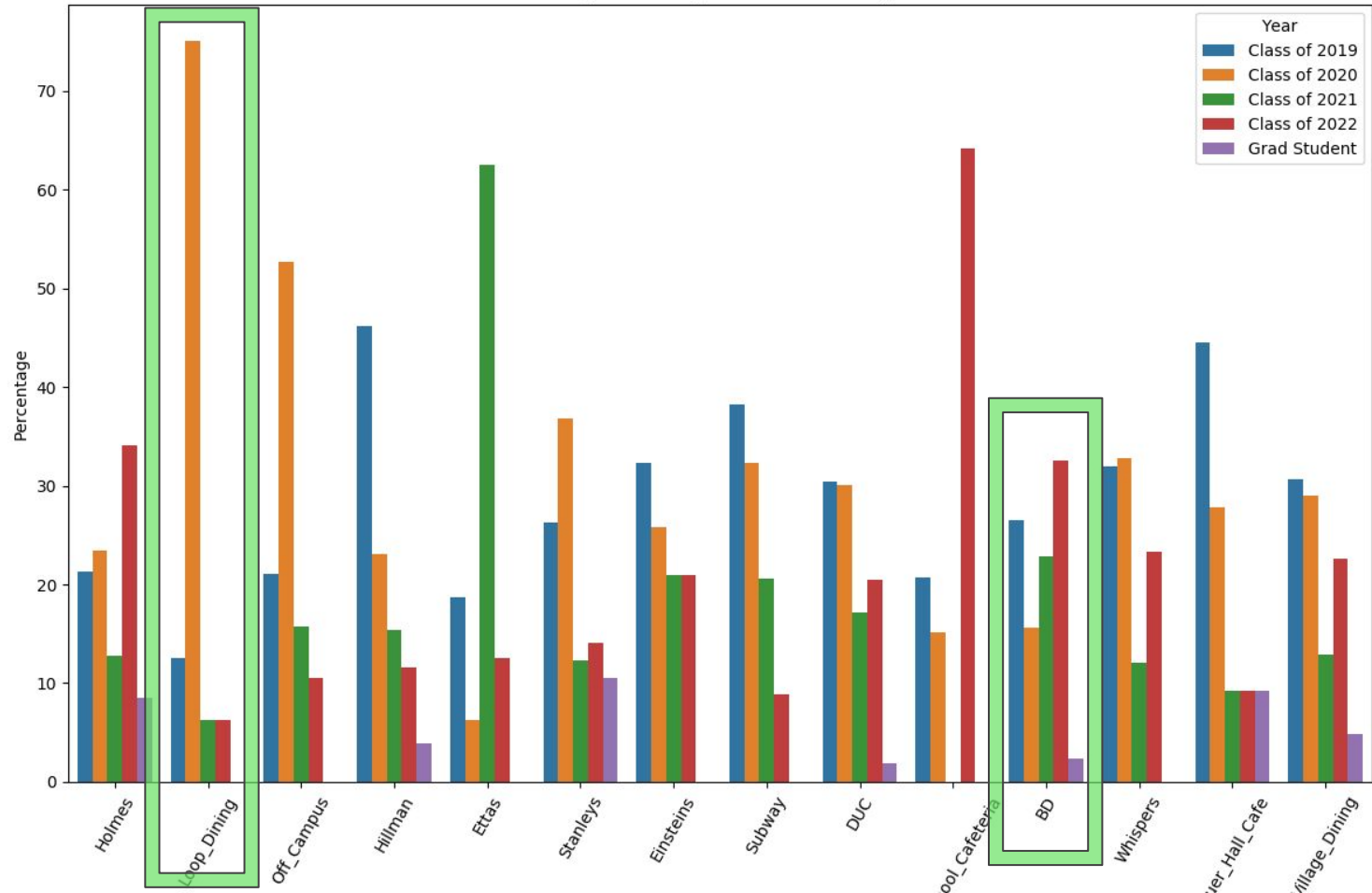
Percentage of Dining Hall Attendance by School

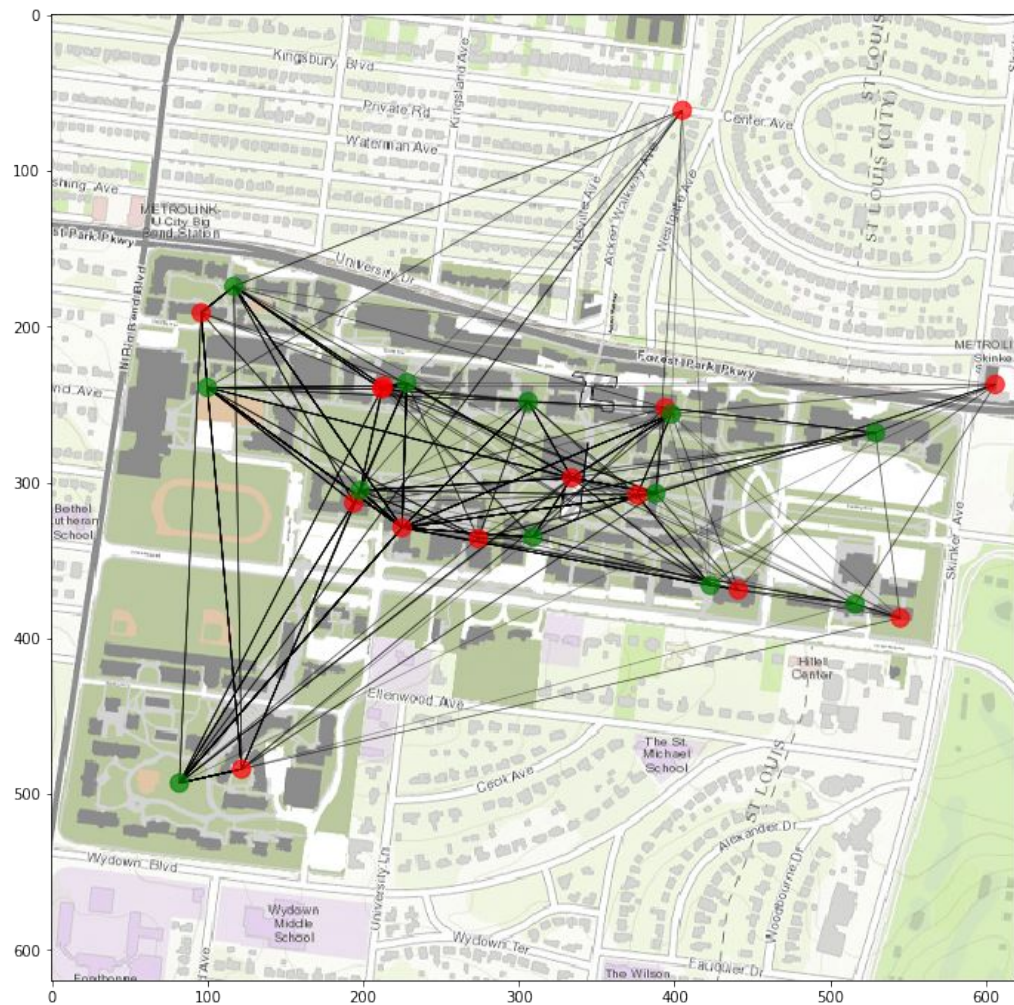


Dining Hall Attendance by Year



Percentage of Dining Hall Attendance by Year

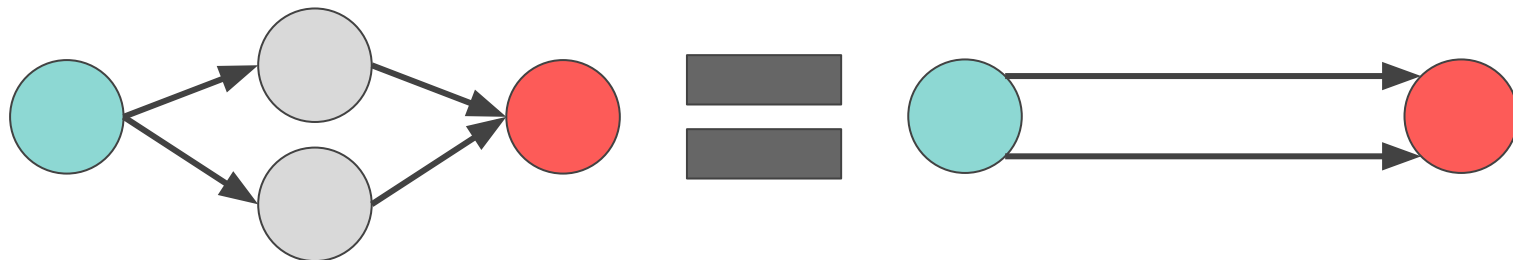






Avoiding Edge Weights

Edge Weights are unnecessary





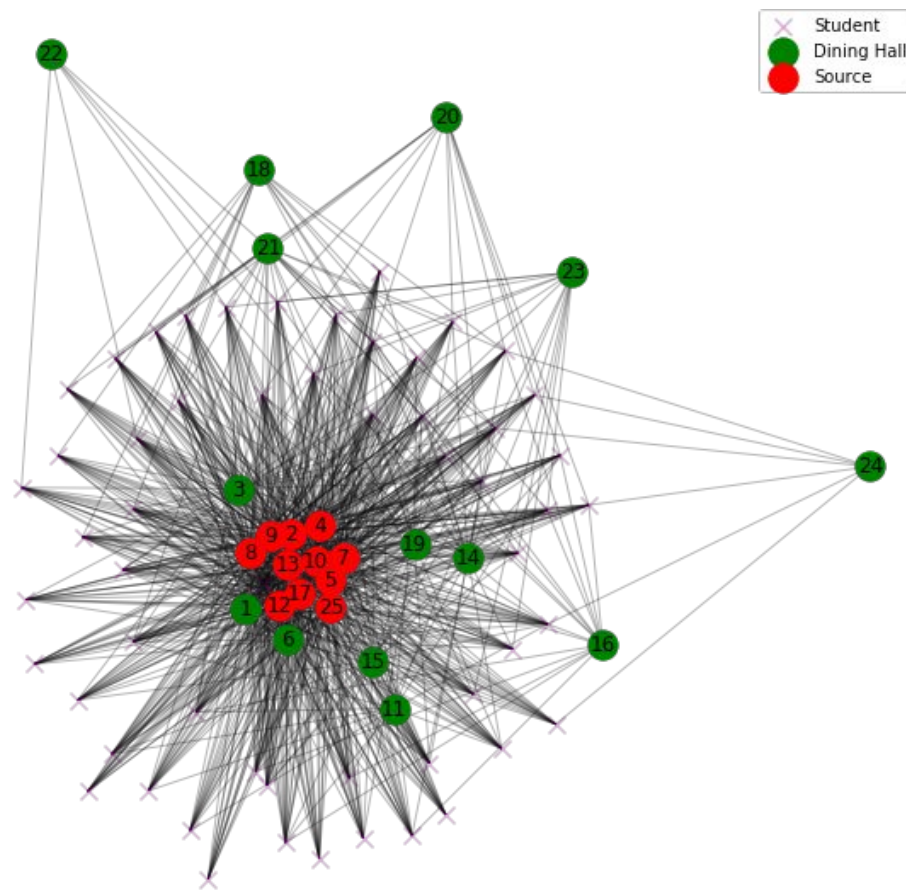
Graph Format

Three different formats utilized:

1. Pathways as nodes
2. Students as nodes
3. Students as edges (multigraph)



Student Nodes

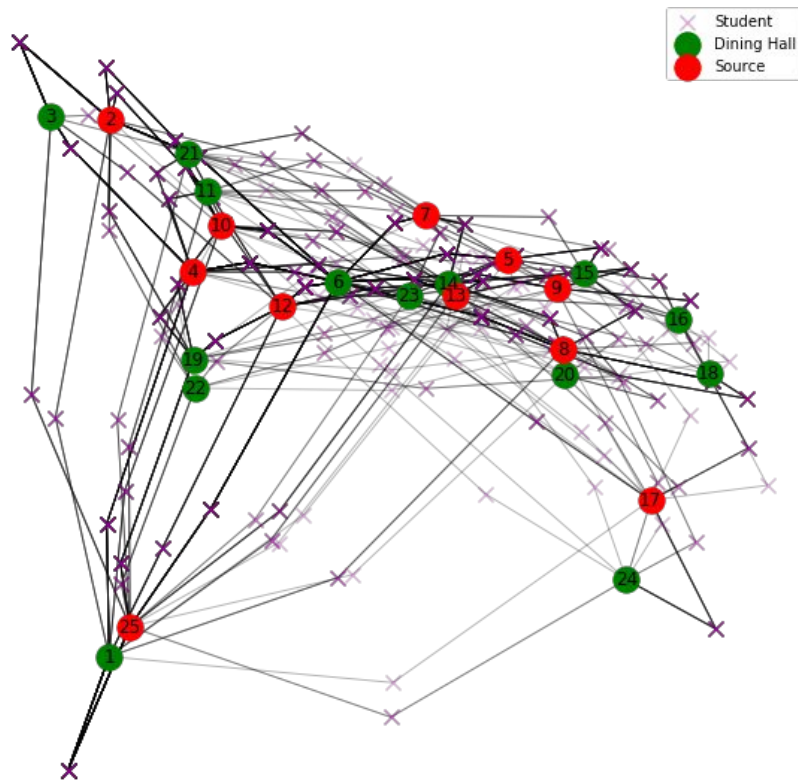




Pathway Nodes

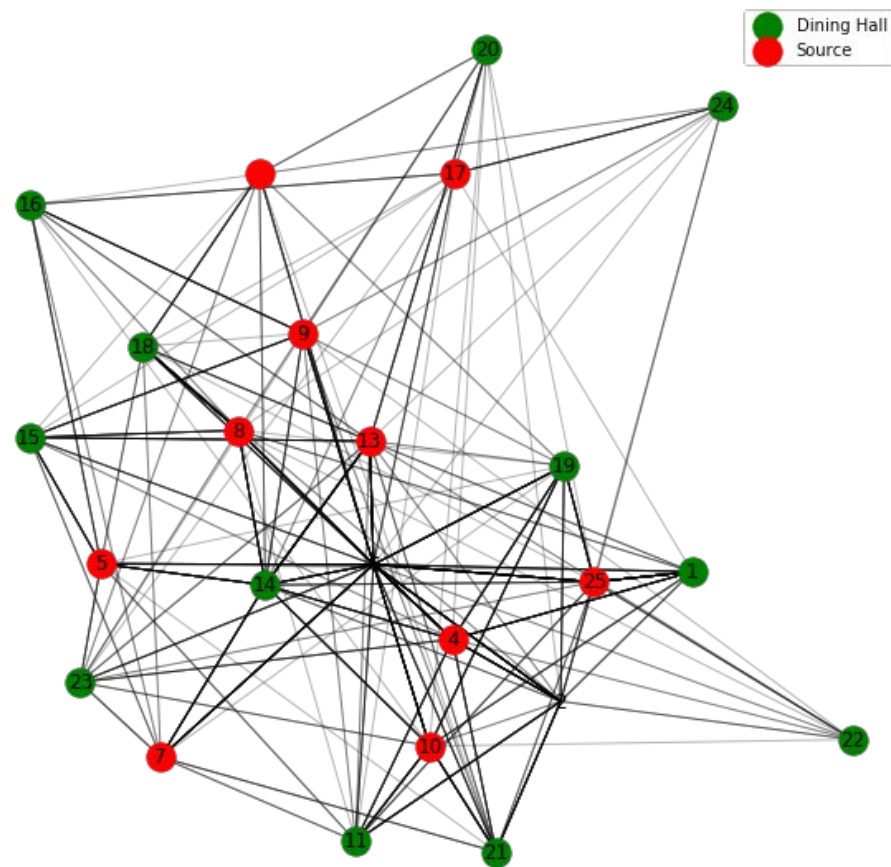
Technically students are nodes, but multiple nodes per student

Each edge has a node, to avoid being a multigraph

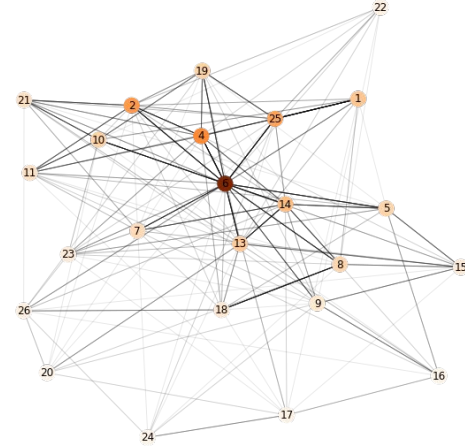
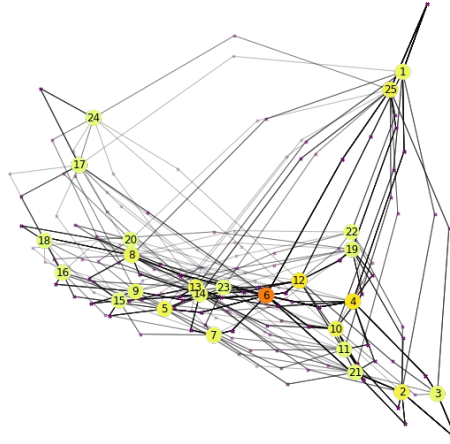
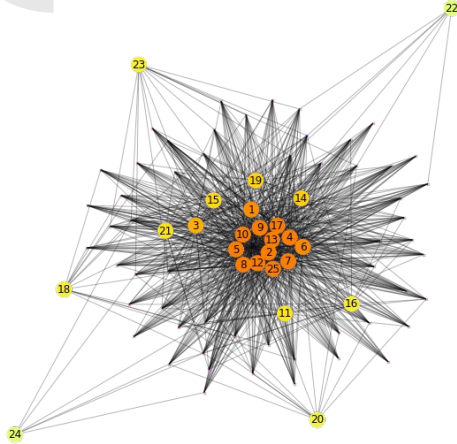




Multigraph

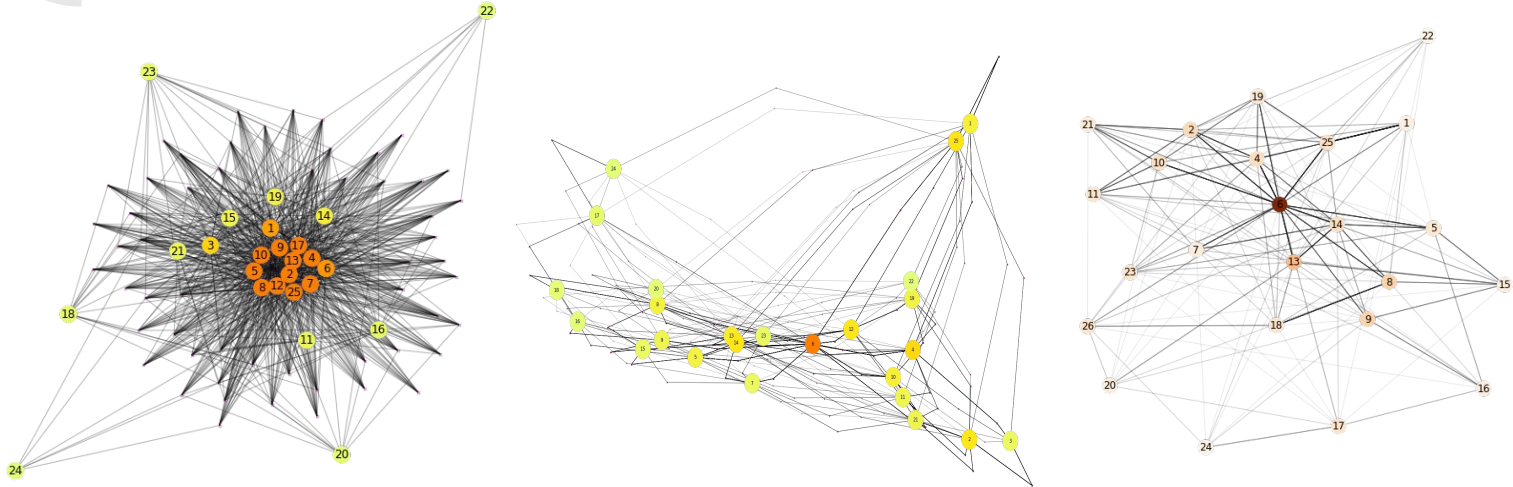


Eigenvector Centrality



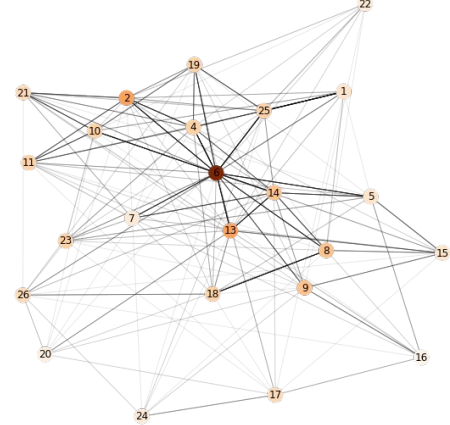
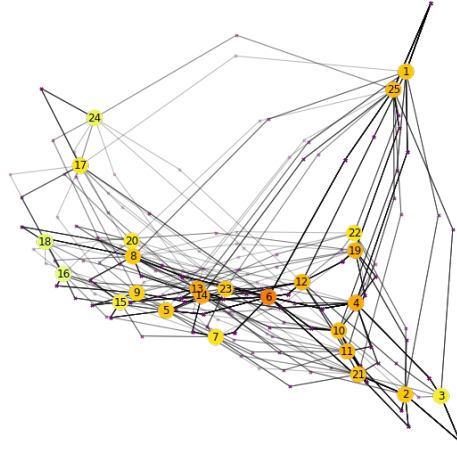
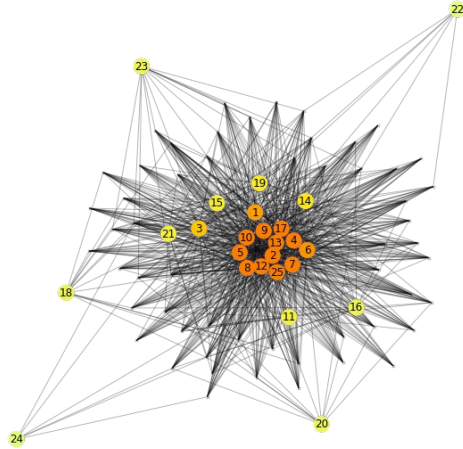
Student Nodes	Pathway Nodes	Multigraph
South 40: 0.182	DUC: 0.635	DUC: 0.650
Brookings: 0.182	AC & Frat Row: 0.162	AC & Frat Row: 0.337
Brown: 0.182	DUC Area: 0.151	Village: 0.302

Betweenness Centrality



Student Nodes	Pathway Nodes	Multigraph
South 40: 0.0415	DUC: 0.390	DUC: 0.115
Village: 0.0415	AC & Frat Row: 0.135	Brookings: 0.042
DUC: 0.0415	South 40: 0.114	Engineering: 0.032

Closeness Centrality



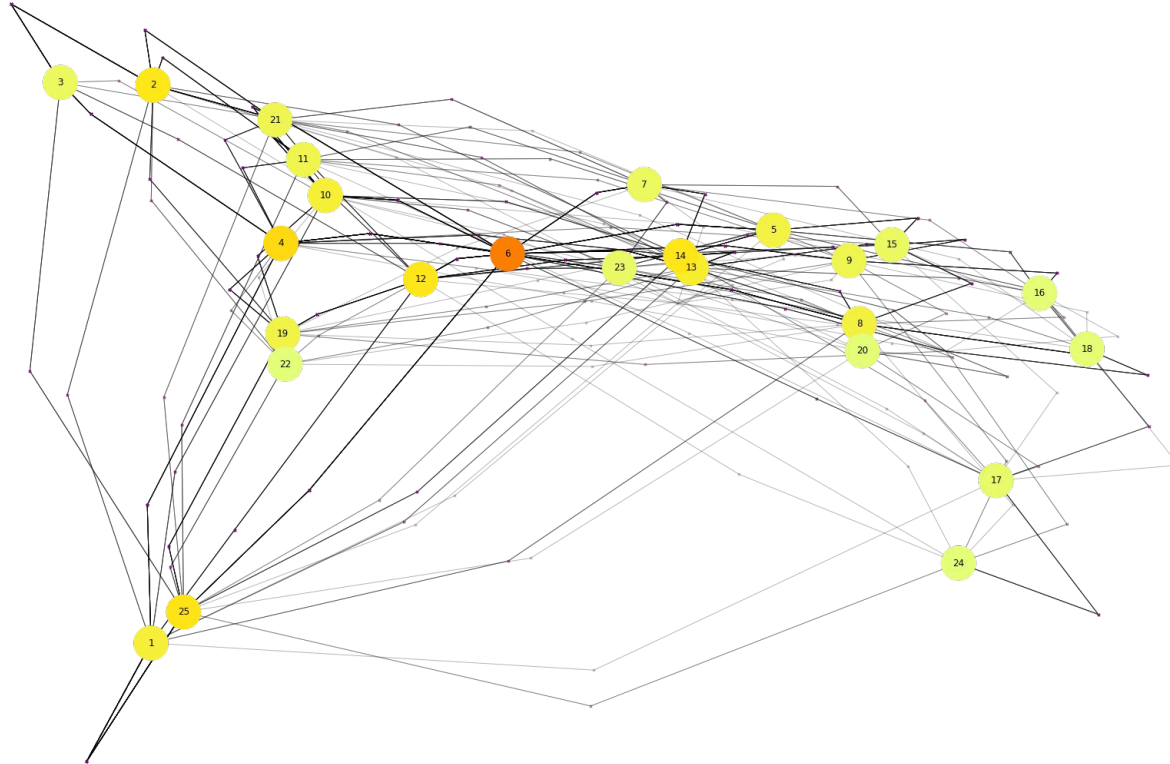
Student Nodes	Pathway Nodes	Multigraph
South 40: 0.762	DUC: 0.408	DUC: 0.92
Village: 0.762	Whispers: 0.357	Village: 0.719
DUC: 0.762	Brookings: 0.354	Brookings: 0.719



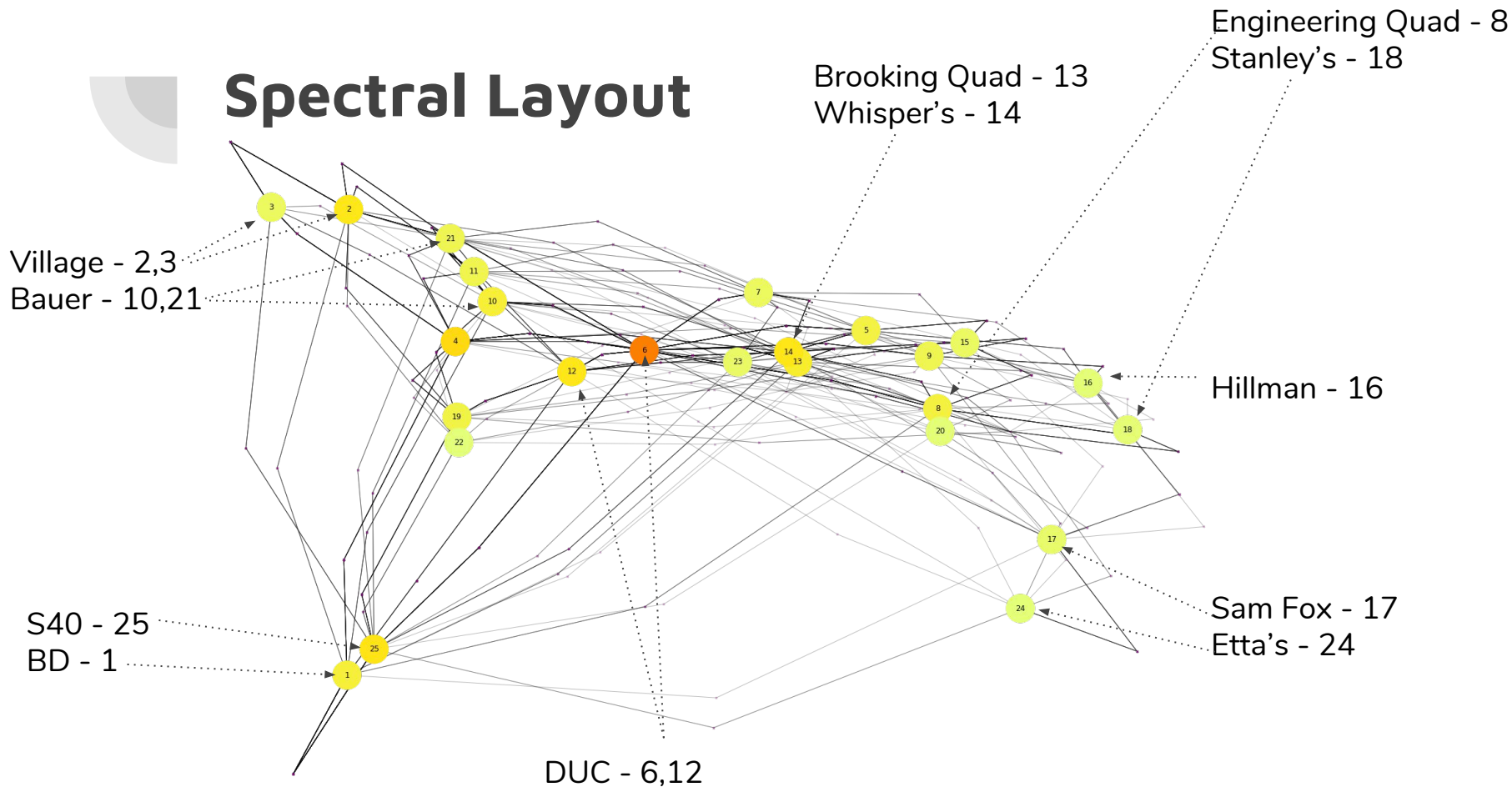
On the Student Nodes Graph

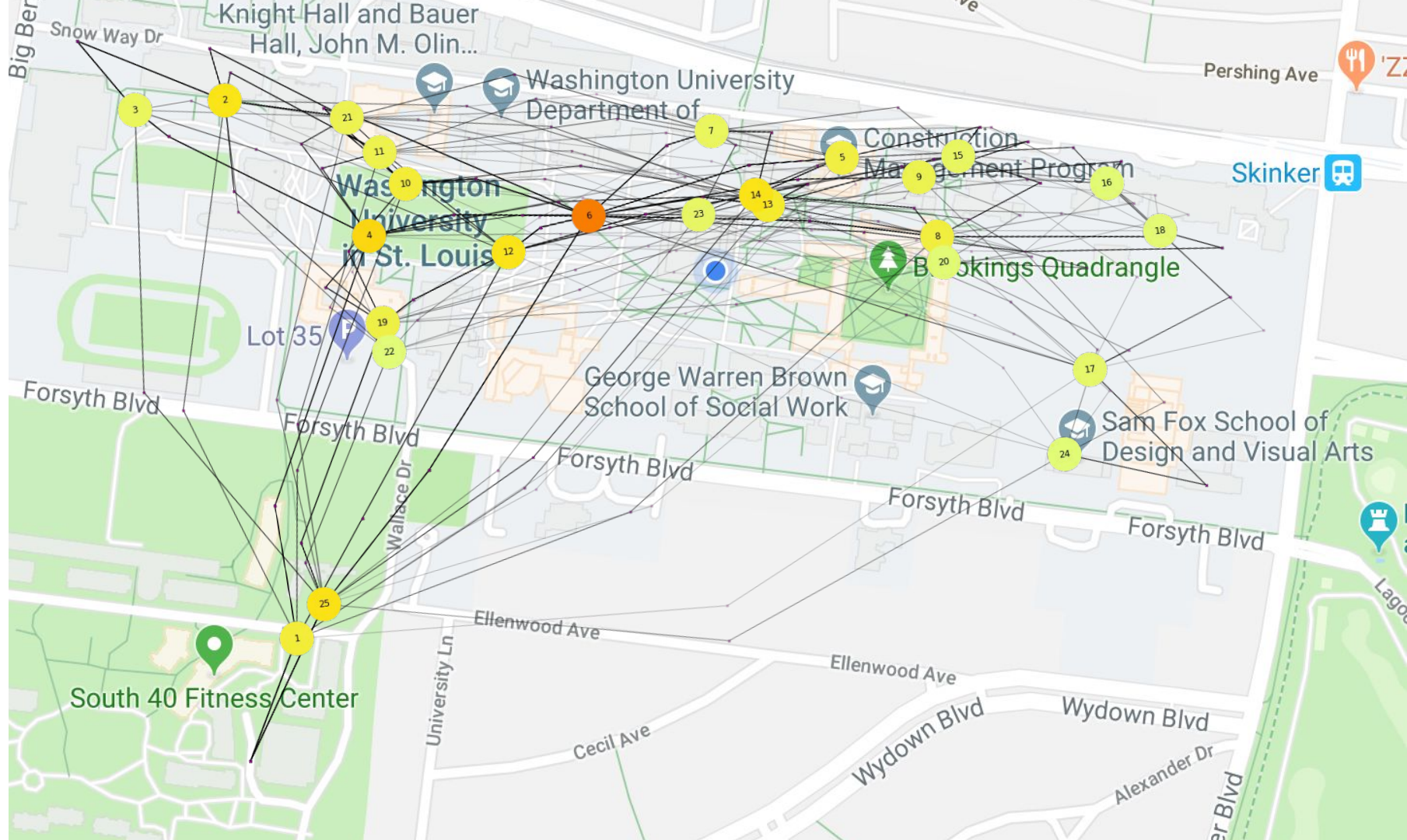
The DUC was always the Dining Hall with the highest centrality, regardless of calculation method

Pathways - Spectral Layout

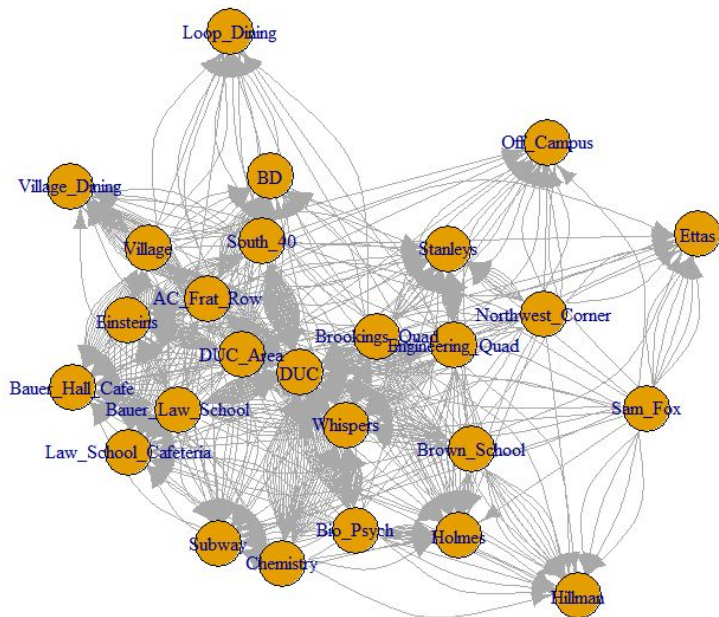


Spectral Layout





Community detection using R igraph library



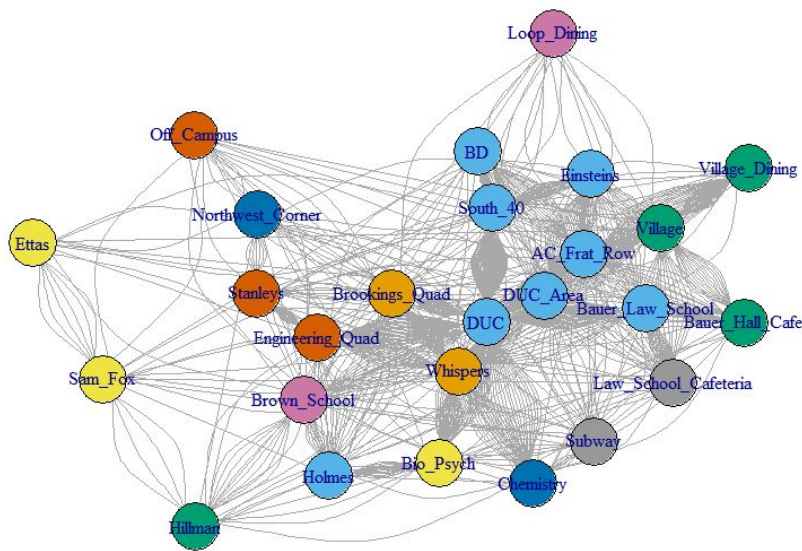
The igraph library in R gives some handy community detection and plotting tools

We have adopted 2 ways that we learned in class

- Betweenness-based clustering
- Modularity Maximization

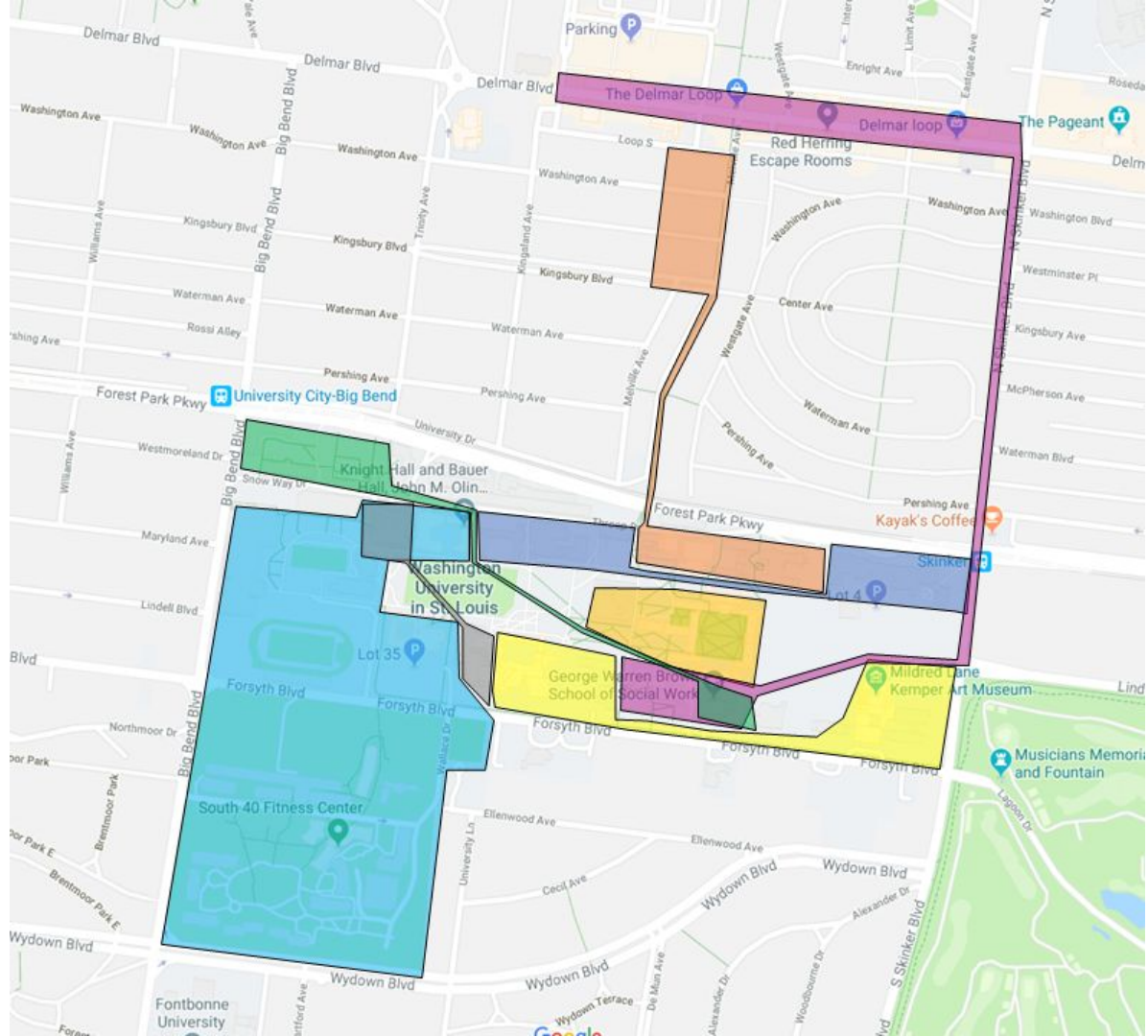
Graph layout used: Fruchterman-Reingold-Grid

Community detection using R igraph library - Betweenness-based clustering

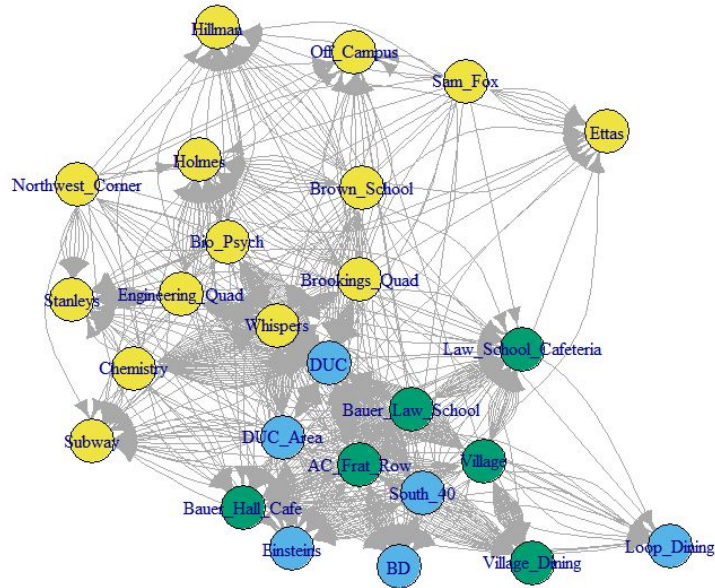


Interesting findings:

- Bridges exist among different segments of buildings on campus, which can be mostly explained by the schools that students are in
- Two anomalies
 - Law school and Subway are in a cluster
 - Hillman and Village / Bauer Hall Cafe are in a cluster
 - Maybe due to the sampling bias - our sample population is majorly undergraduates, so might give a biased measure on graduate school buildings

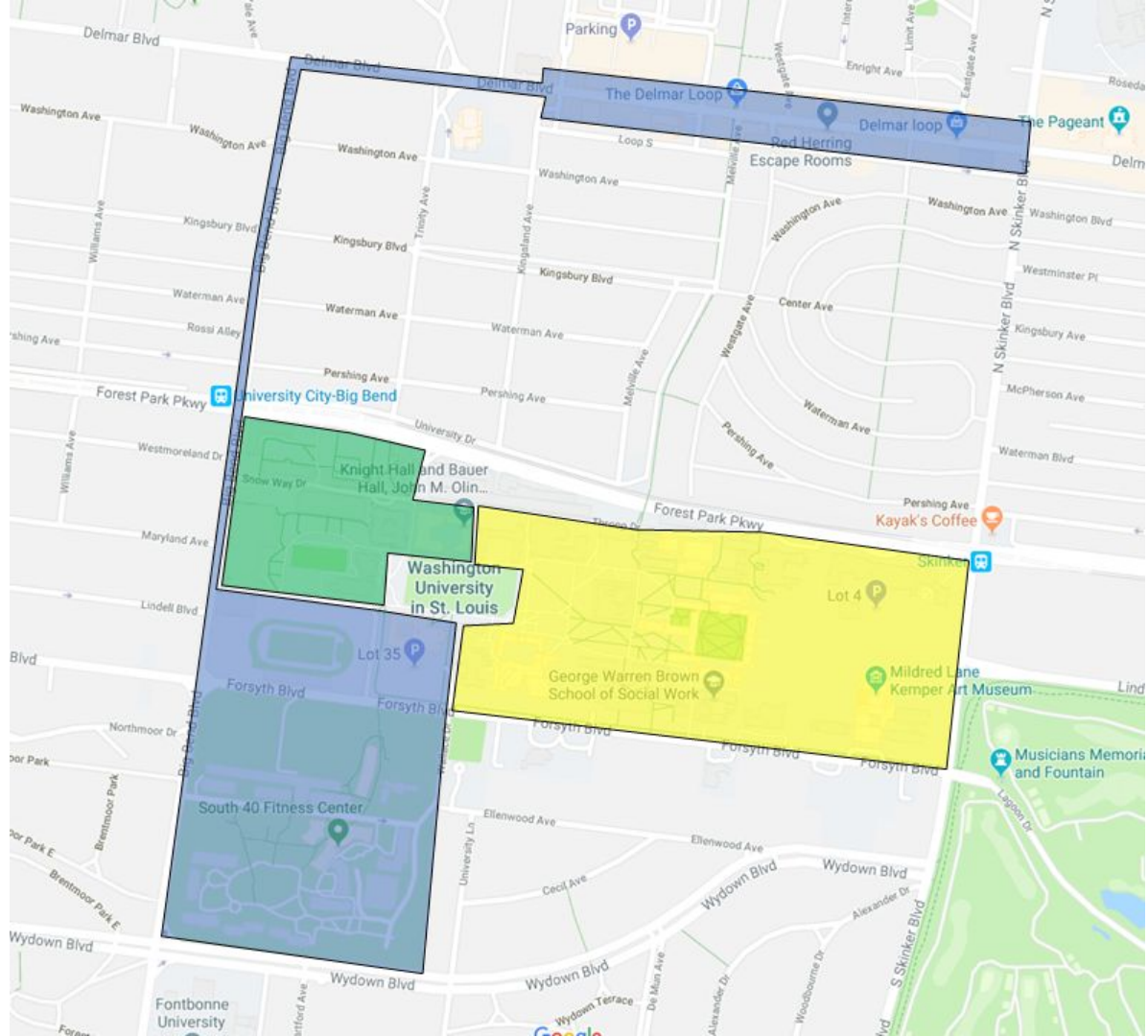


Community detection using R igraph library - Modularity Maximization



Interesting findings:

- Shows the communities based on dense connection to each other (more student movements among buildings) but sparse connection with nodes in other communities
- Basically makes sense with actual distances among buildings
- If we were to have more data, we could have gotten better resolution in terms of communities (micro-communities could have been detected)





Conclusion

There is a definite student flow on campus through dining areas

Whisper's is less important than we theorized, the DUC is core to student flow



Moving forward

If we were to have more data, we could have gotten better resolution on community detection.

More options for sources and targets

Collecting over a longer period of time