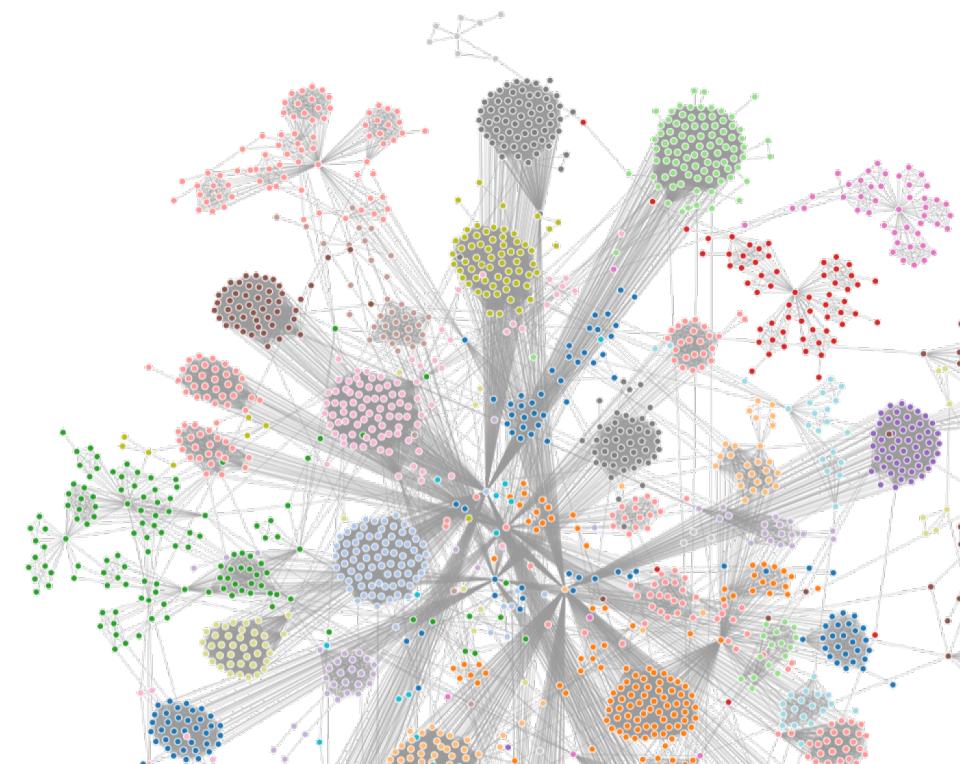
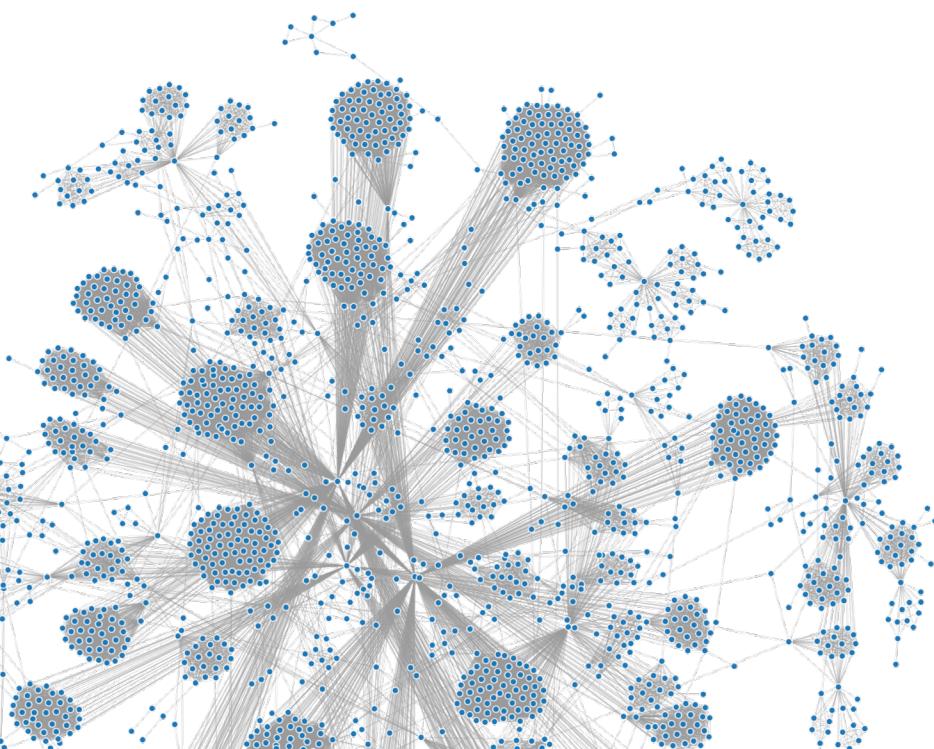
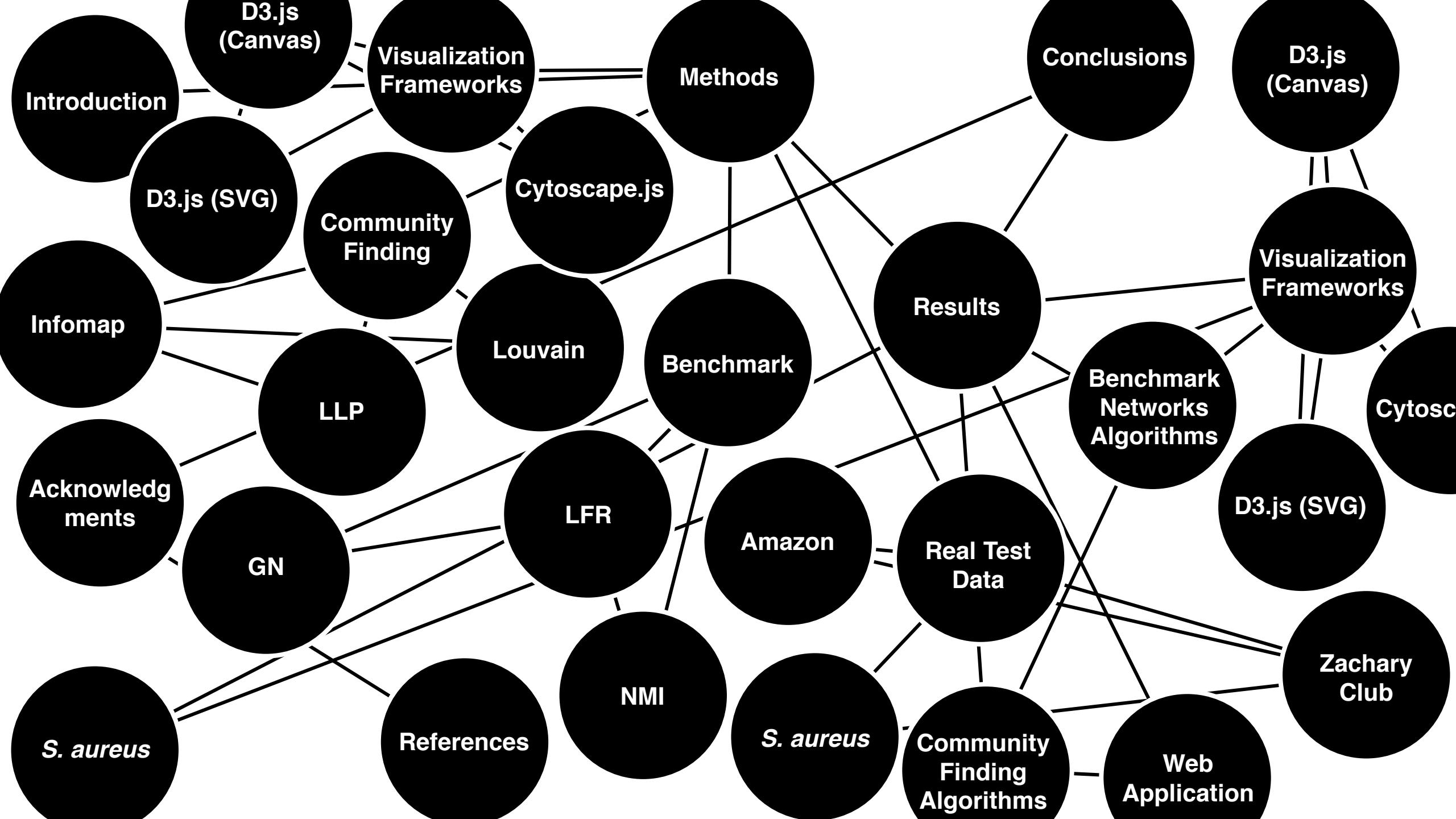


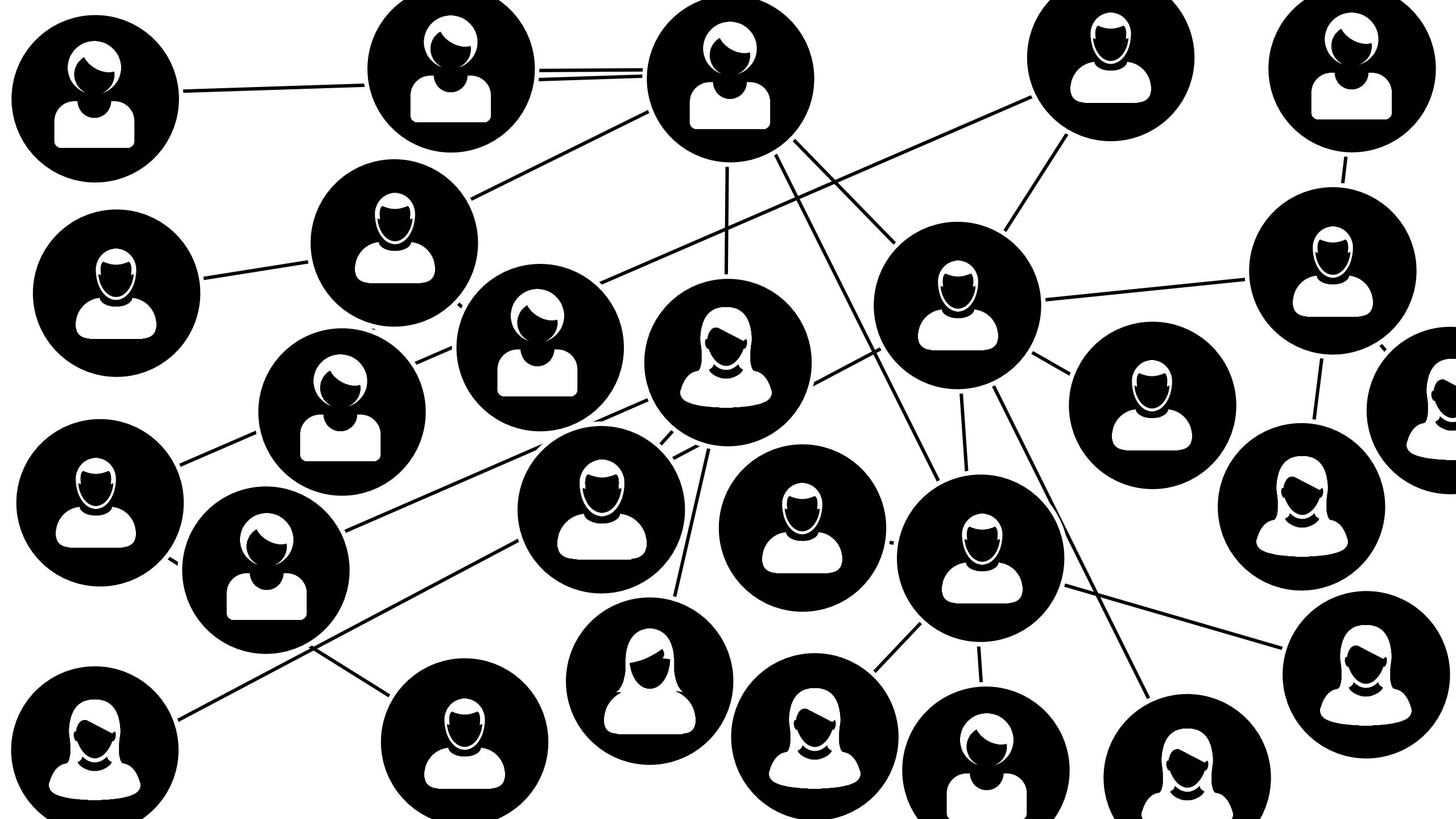


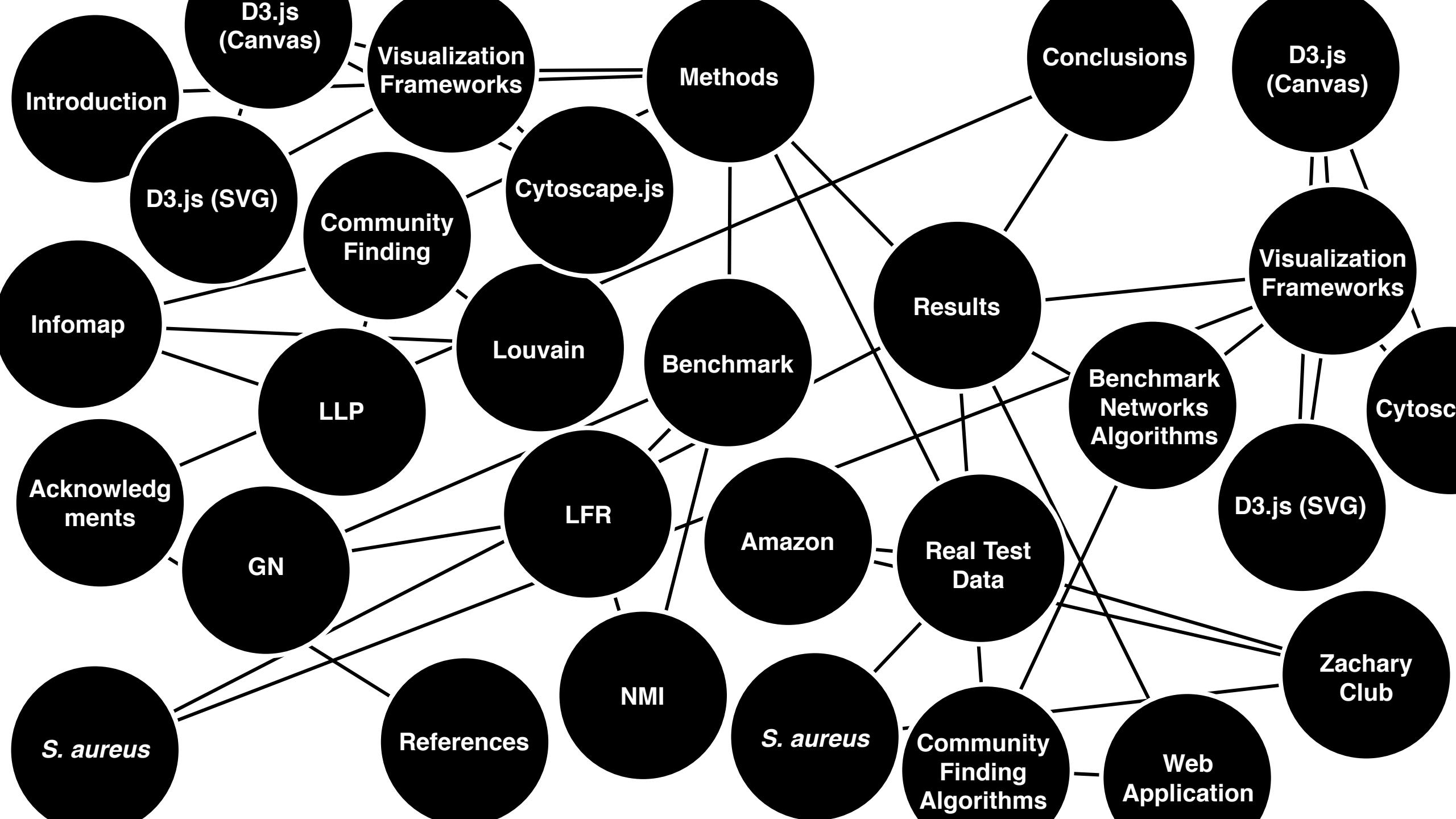
# Community Finding with Applications on Phylogenetic Networks

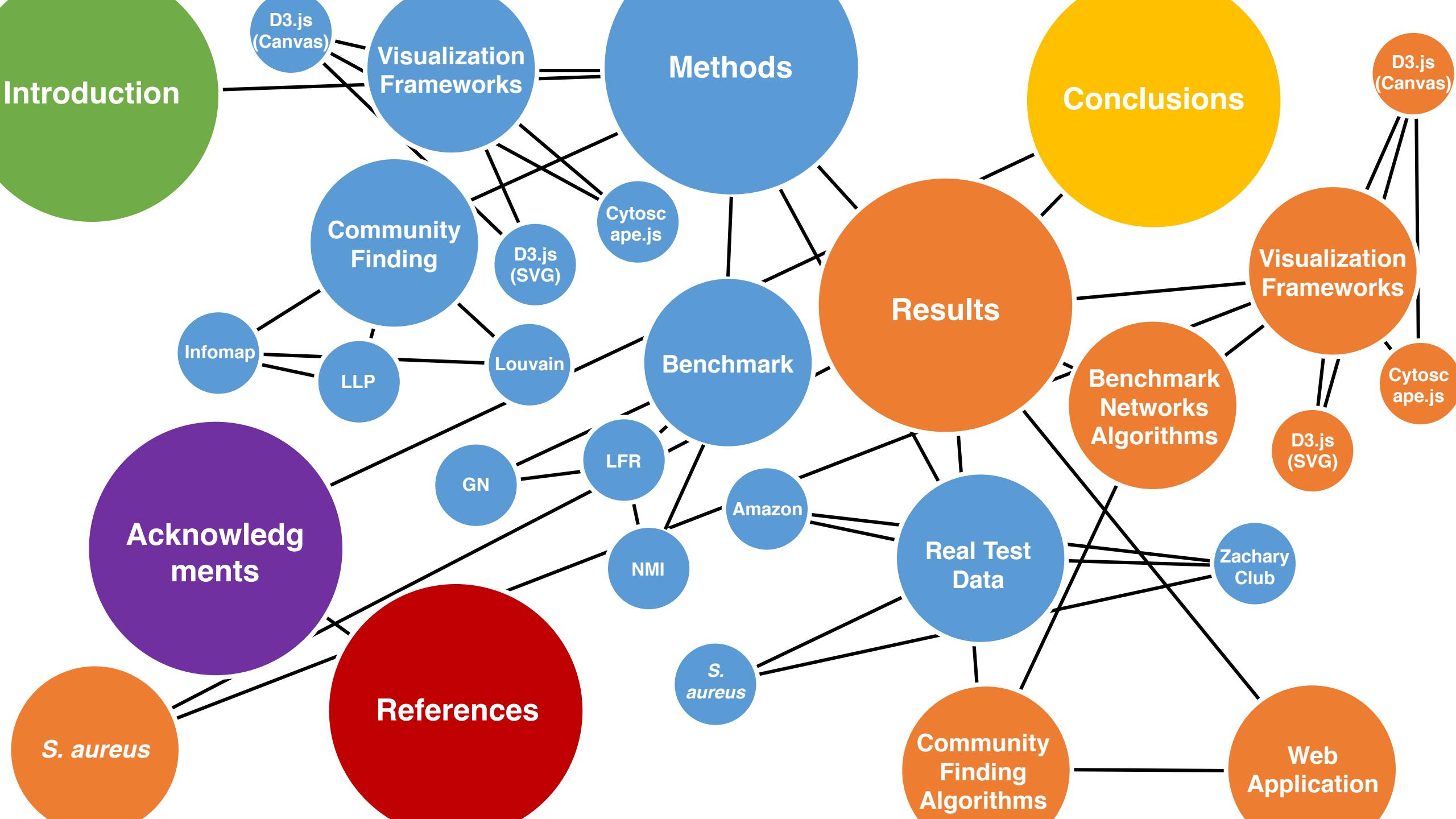
Luís Rita











# Introduction

*“The world’s most valuable resource is no longer oil, but data.”*

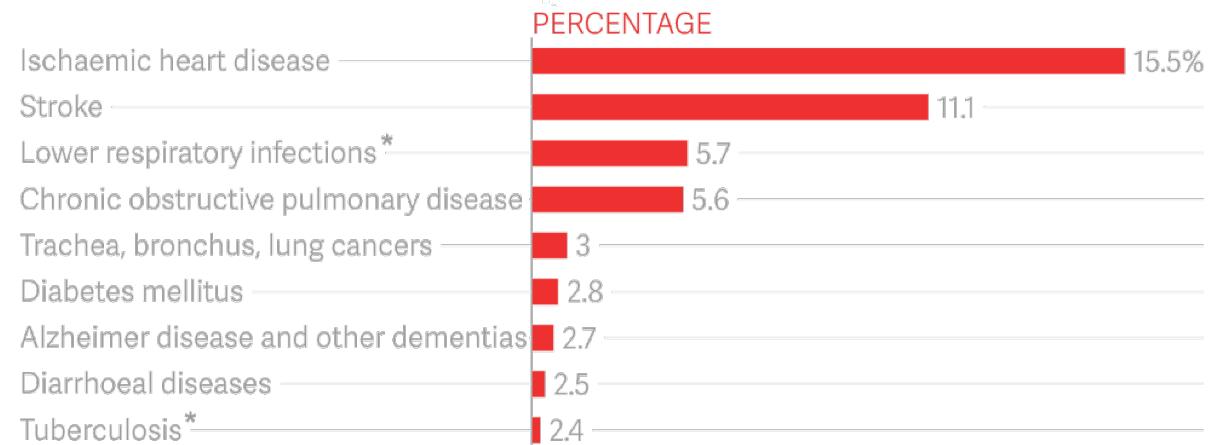
The Economist

*“The rapid advances in NGS technologies and capabilities has proven and further promises to be a game changer in diagnostic microbiology (...)"*

BDQ Journal (Elsevier)

*“Ten threats to global health in 2019 – Antimicrobial Resistance”*

World Health Organization



# Methods

## Community Finding

### Louvain

```

1: Compute  $Mod_{new} = Mod(M)$ 
2: repeat
3:    $Mod = Mod_{new}$ 
4:   Randomize the order of vertices.
5:   for all  $i \in V^k$  do
6:      $best_{com} = M_i^k$ ;  $best_{increase} = 0$ 
7:     for all  $M' \in C^k$  do
8:        $M_i^k = M_i^k \setminus \{i\}$ 
9:        $\Sigma_{tot}^{M_i^k} = \sum_{\alpha \in M_i^k} k_\alpha - k_i$ ;  $\Sigma_{tot}^{M'_i} = \sum_{\alpha \in M'_i} k_\alpha + k_i$ 
10:       $\Sigma_{in} = \Sigma_{tot}^{M_i^k} - \sum k_{i,j}, (i, j) \in E^k, i \in C_i^k \text{ and } j \notin C_i^k$ 
11:       $k_{i,in} = \sum_{\alpha \in M'_i} k_{i,\alpha}$ 
12:      If  $\delta Mod_{M_i^k \rightarrow M'_i} > best_{increase}$  then
13:         $best_{increase} = \delta Mod_{M_i^k \rightarrow M'_i}$ 
...

```

$$M = \frac{1}{2m} \sum_{i,j} \left( A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

$$\Delta M m = k_{i,in} - \frac{\Sigma_{tot} k_i}{2m}$$

### Infomap

```

1: Compute  $L_{new} = L(M)$ 
2: repeat
3:    $L = L_{new}$ 
4:   Randomize the order of vertices.
5:   for all  $u \in V^k$  do
6:     if  $M'_u = argmin(\delta L_{M_u^k \rightarrow M'_u^k}) < 0$  then
7:        $M_u^k = M_u^k \setminus \{u\}$ ;  $M'_u^k = M'_u^k \cup \{u\}$ 
8:        $p^{M_u^k} = \sum_{\alpha \in M_u^k} p_\alpha - p_u$ ;  $p^{M'_u^k} = \sum_{\alpha \in M'_u^k} p_\alpha + p_u$ 
9:       update  $q^{M_u^k}$ ; update  $q^{M'_u^k}$ 
10:      end if
11:    end for
12:    Compute  $L_{new} = L(M)$ 
13:    until No vertex movement or  $L - L_{new} < \theta$ 
14:    Compute  $L_{new} = L(M)$ 
15:    if  $L - L_{new} < \theta$  then
...

```

$$L(M) = \left( \sum_{m \in M} q_m \right) \log \left( \sum_{m \in M} q_m \right) - 2 \sum_{m \in M} q_m \log(q_m)$$

$$- \sum_{\alpha \in V} p_\alpha \log(p_\alpha) + \sum_{m \in M} \left( q_m + \sum_{\alpha \in m} p_\alpha \right) \log \left( q_m + \sum_{\alpha \in m} p_\alpha \right)$$

### Layered Label Propagation

```

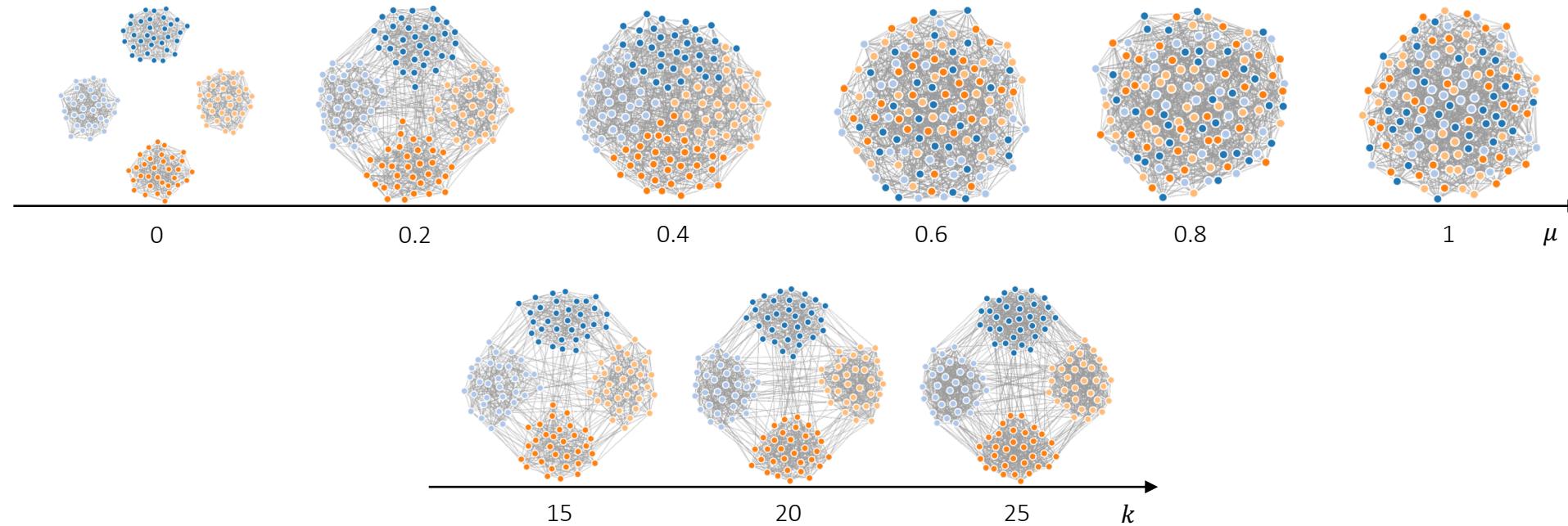
1:  $k = 0$ 
2: for all  $i \in V^0$  do
3:    $M_i^0 = \{i\}$ 
4: end for
5: repeat
6:   Randomize the order of vertices.
7:   for all  $u \in V^k$  do
8:     for all  $i \in labels(u)$  do
9:       if  $k_i(u) - \gamma(v_i(u) - k_i(u)) > k_{i-1}(u) - \gamma(v_{i-1}(u) - k_{i-1}(u))$  then
10:          $M_u^k = M_u^k \setminus \{u\}$ ;  $M'_u^k = M'_u^k \cup \{u\}$ 
11:       end if
12:     end for
13:   end for
14:    $k = k + 1$ 
15: until No vertex movement or  $k < max_{iteration}$ 

```

$$k_i - \gamma(v_i - k_i)$$

# Methods

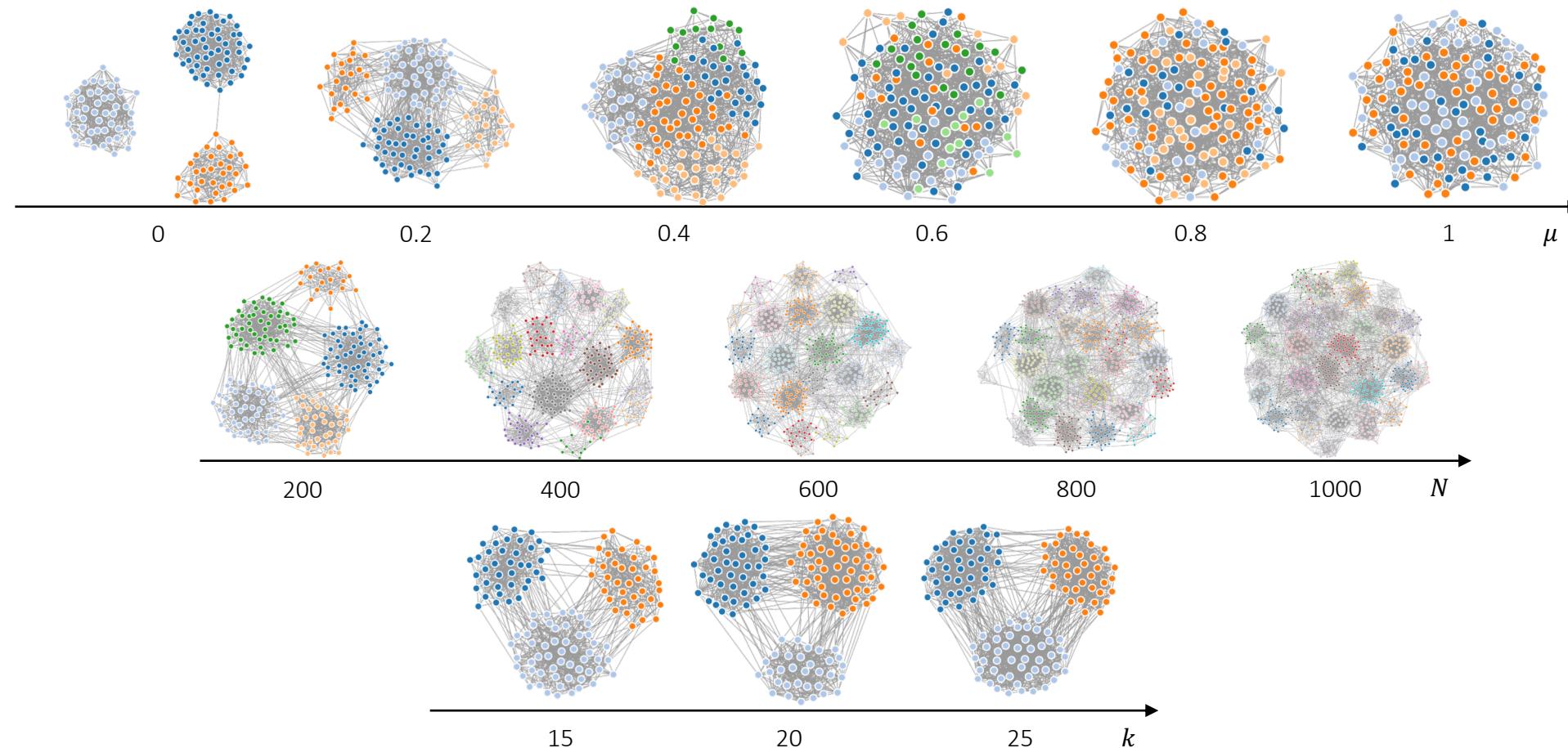
## Benchmark Networks



Girvan - Newman

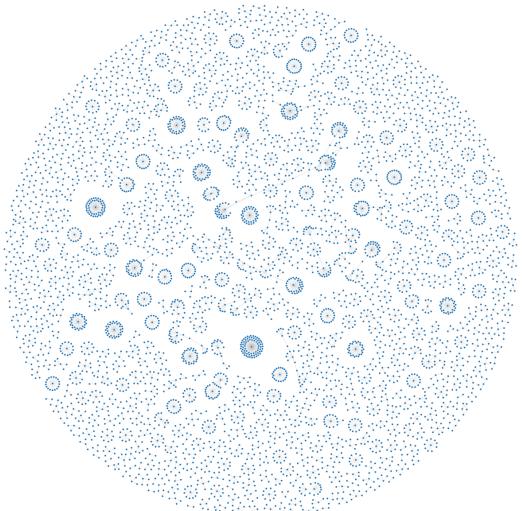
# Methods

## Benchmark Networks

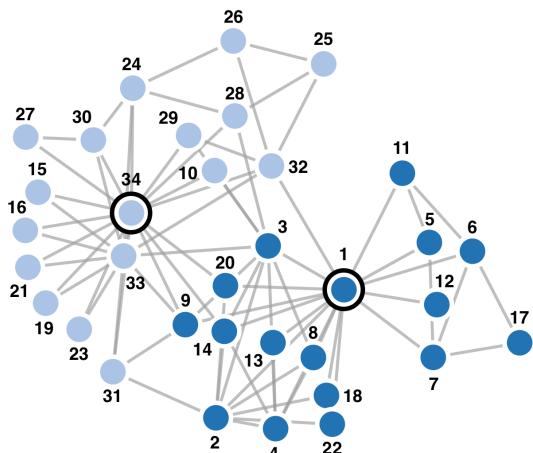


# Methods

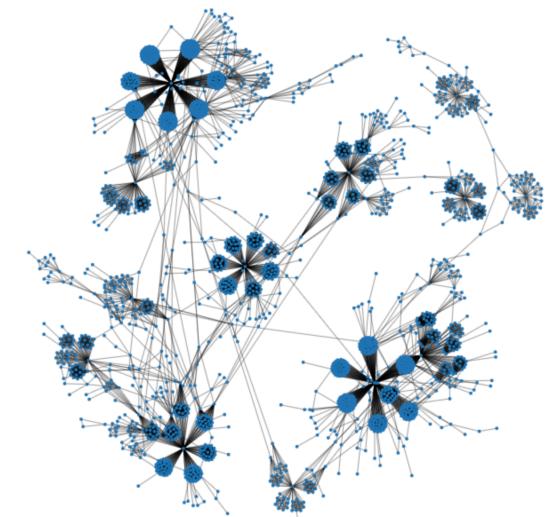
## Real Test Data



Amazon



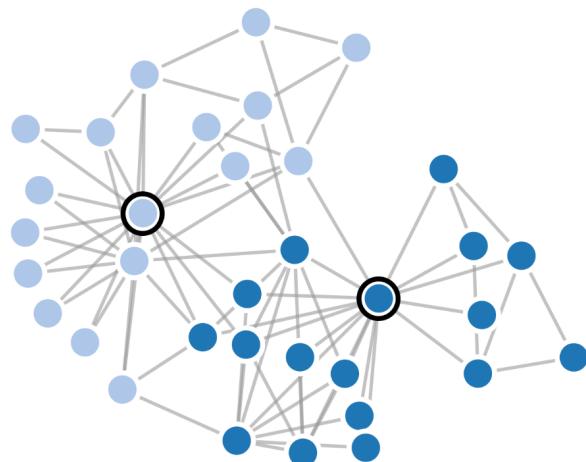
Zachary's Karate Club



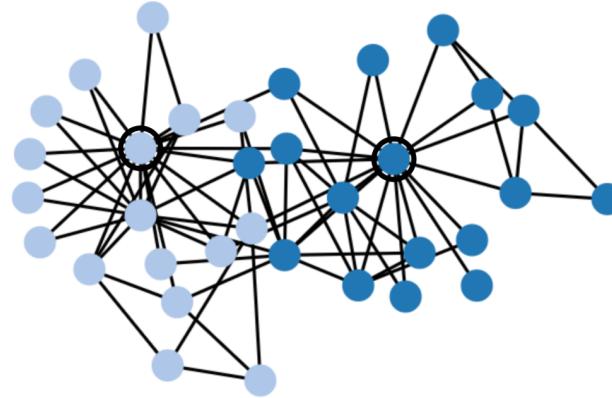
*Staphylococcus aureus*  
MLST SLV Network

# Methods

## Visualization Frameworks



D3.js & SVG



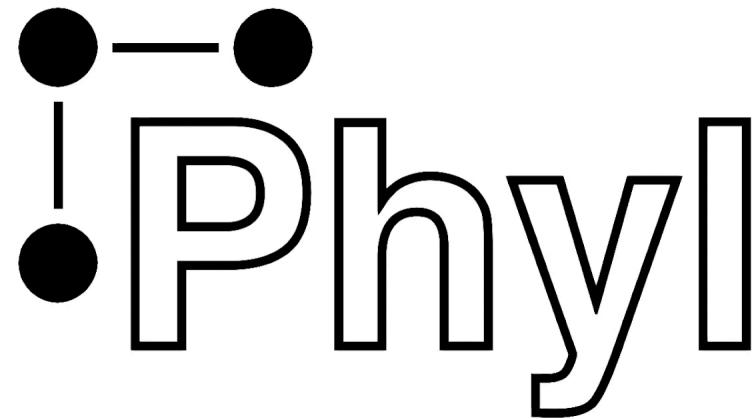
D3.js & Canvas



Cytoscape.js

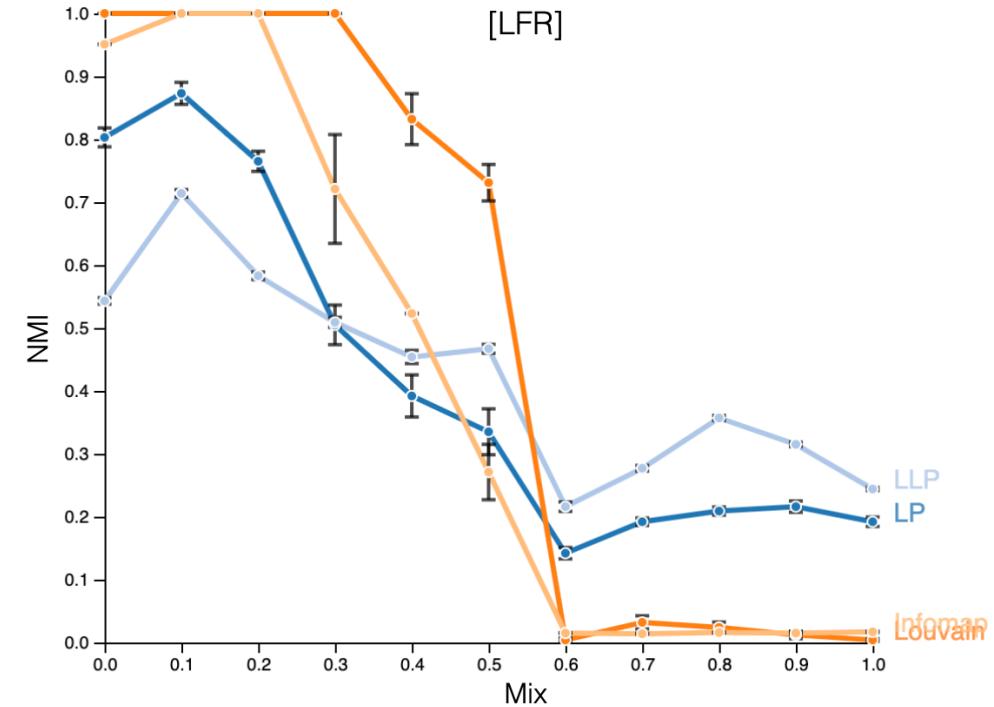
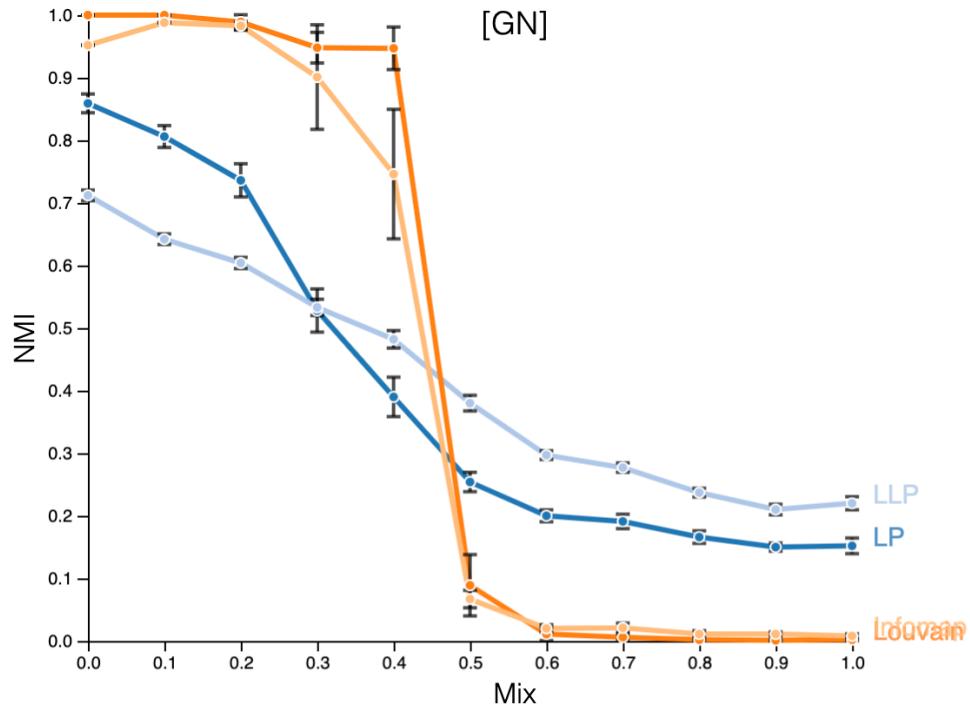
# Results

Web Application



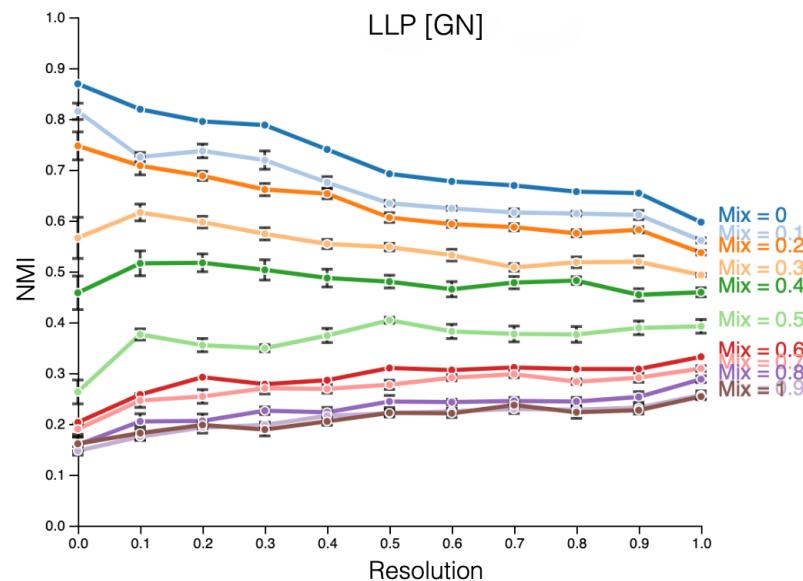
# Results

## Community Finding Algorithms

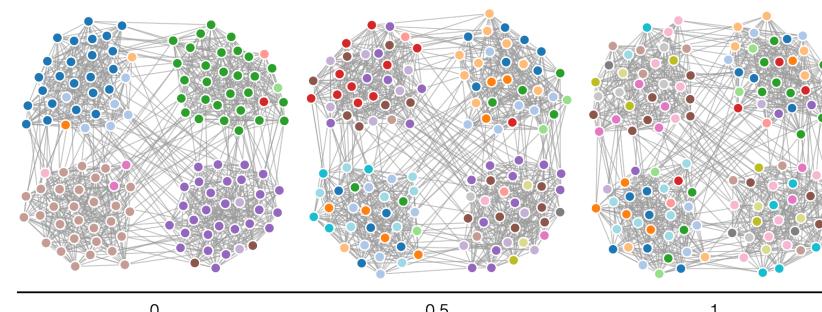
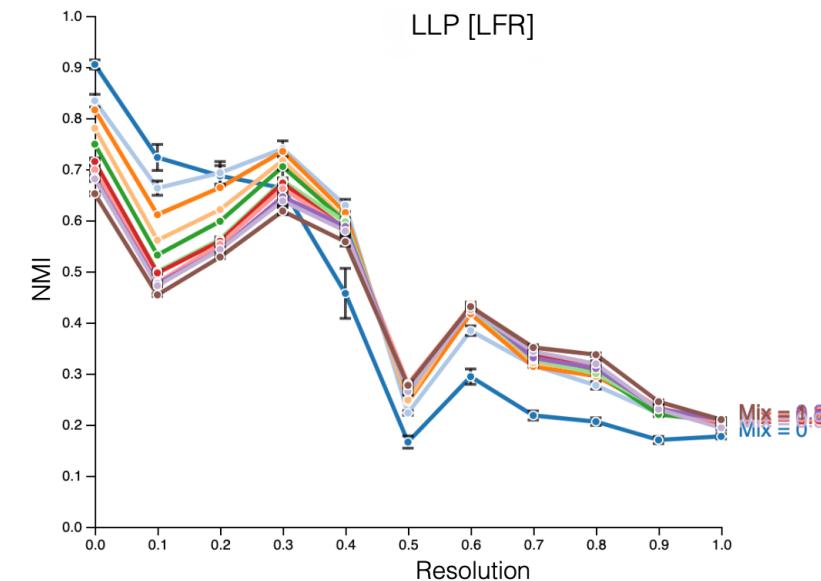


# Results

## Community Finding Algorithms

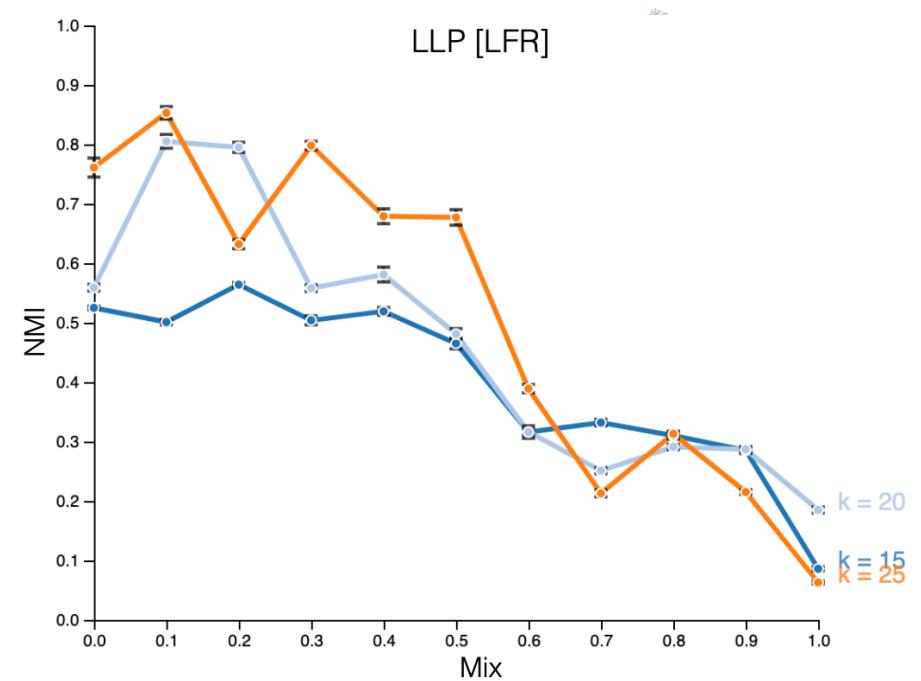
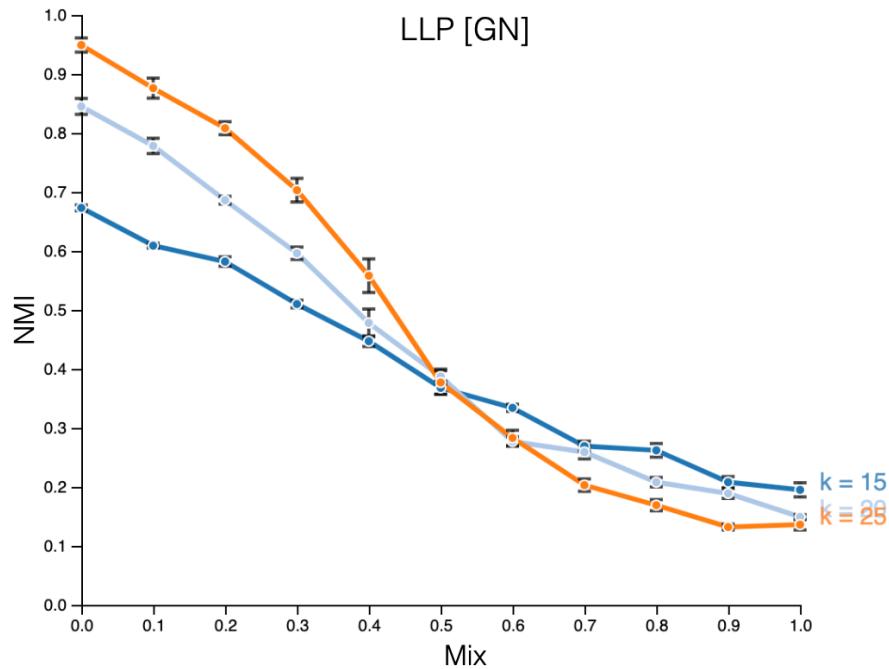


$k_i - \gamma(v_i - k_i)$



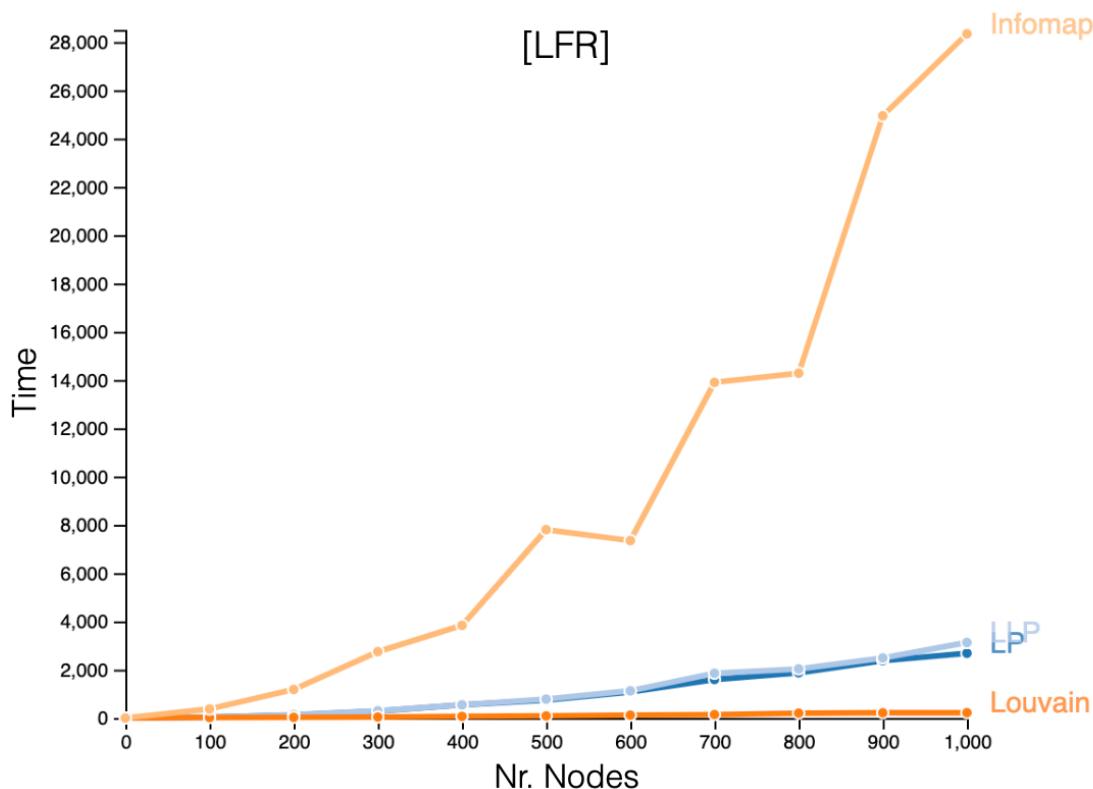
# Results

## Community Finding Algorithms



# Results

## Community Finding Algorithms



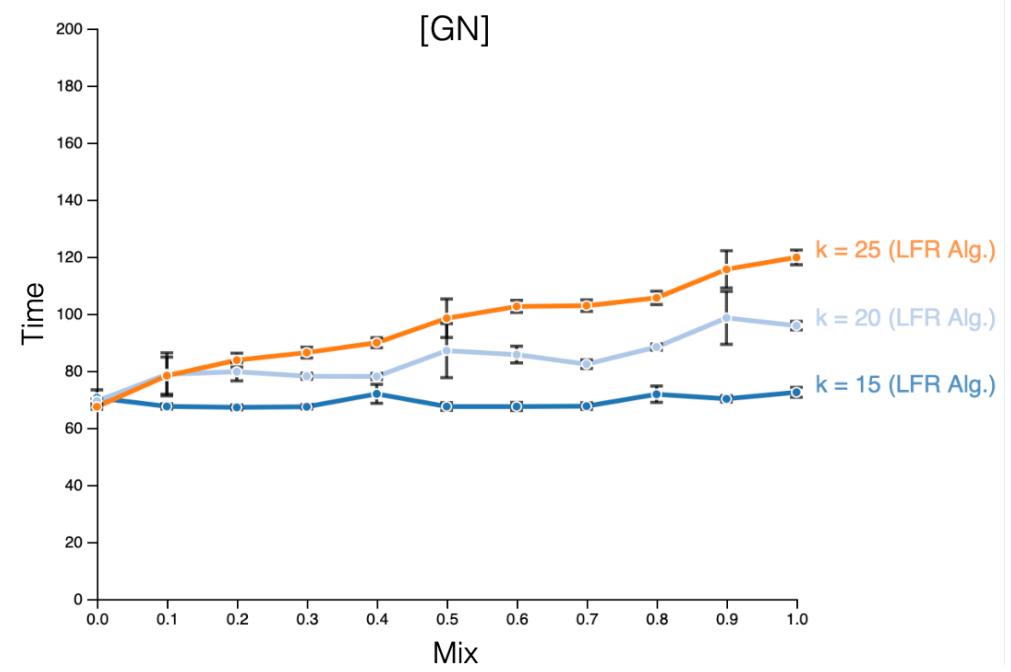
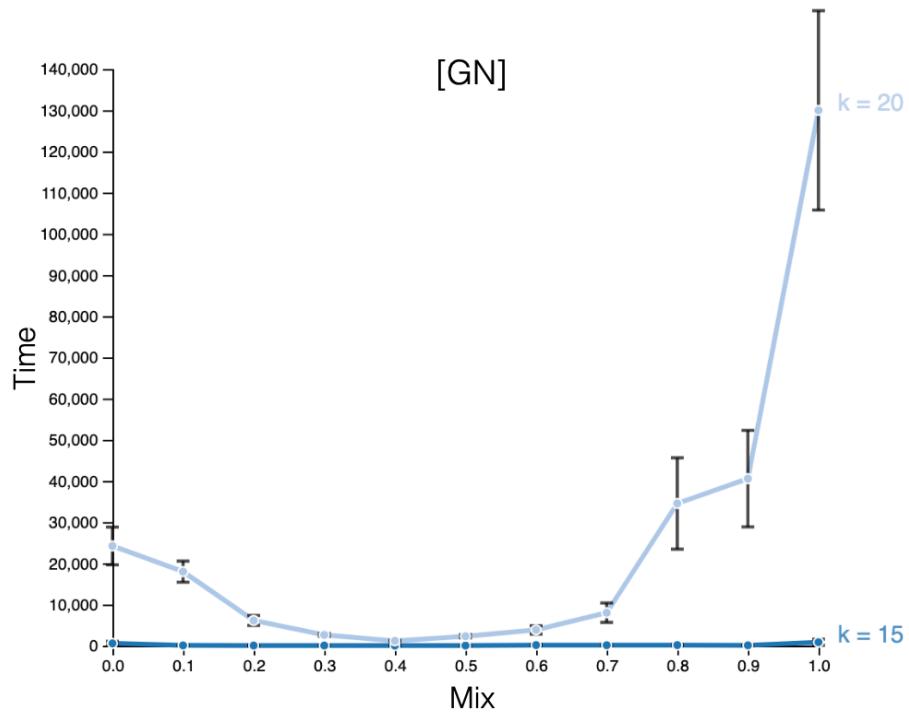
Algorithm	Time Complexity
Louvain	$O(L)$
Infomap	$O(N \log N)$
LP	$O(L)$
LLP	$O(L)$

$$\Delta Mm = k_{i,in} - \frac{\sum_{tot} k_i}{2m}$$
$$= \left( \sum_{m \in M} q_m \right) \log \left( \sum_{m \in M} q_m \right) - 2 \sum_{m \in M} q_m \log(q_m)$$
$$- \sum_{\alpha \in V} p_\alpha \log(p_\alpha) + \sum_{m \in M} \left( q_m + \sum_{\alpha \in m} p_\alpha \right) \log \left( q_m + \sum_{\alpha \in m} p_\alpha \right)$$

$$k_i - \gamma(v_i - k_i)$$

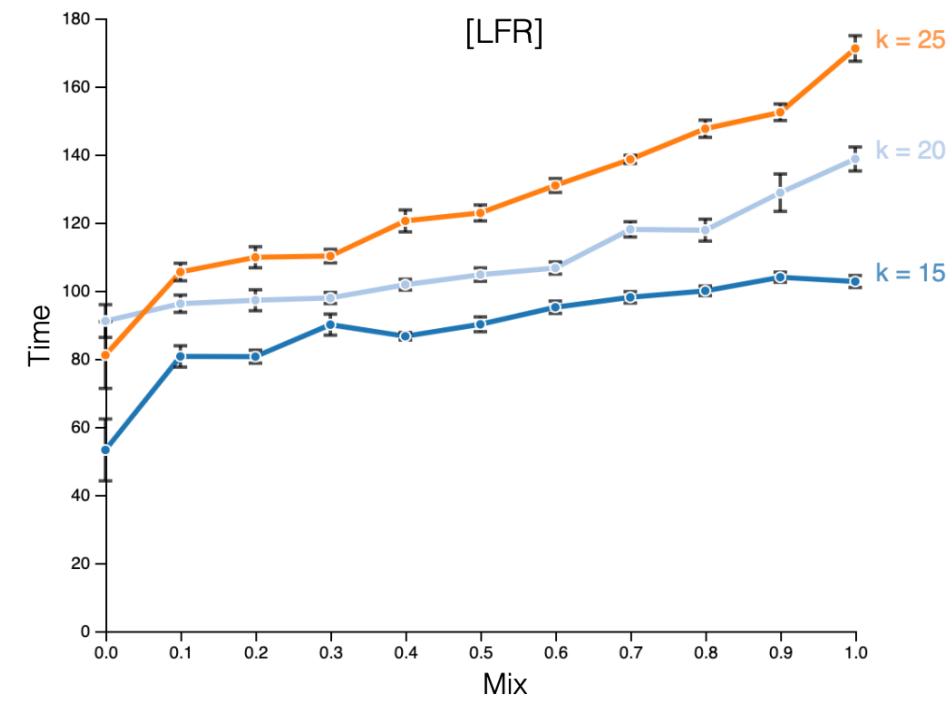
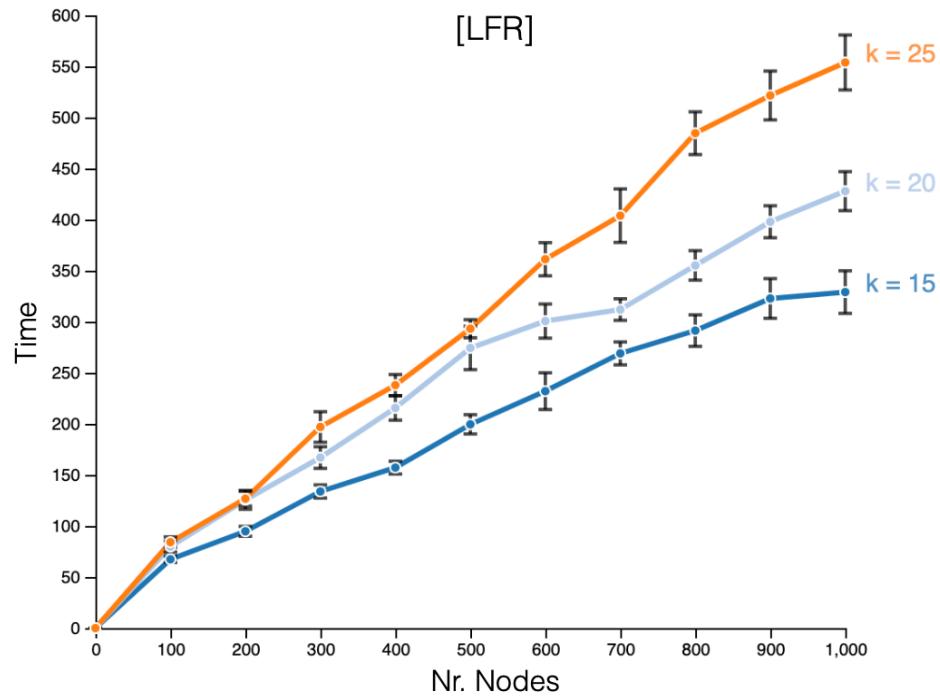
# Results

## Benchmark Networks Algorithms



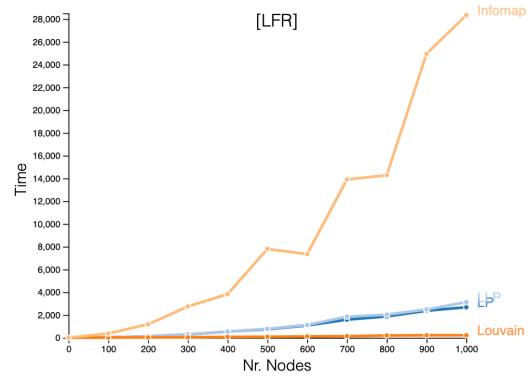
# Results

## Benchmark Networks Algorithms

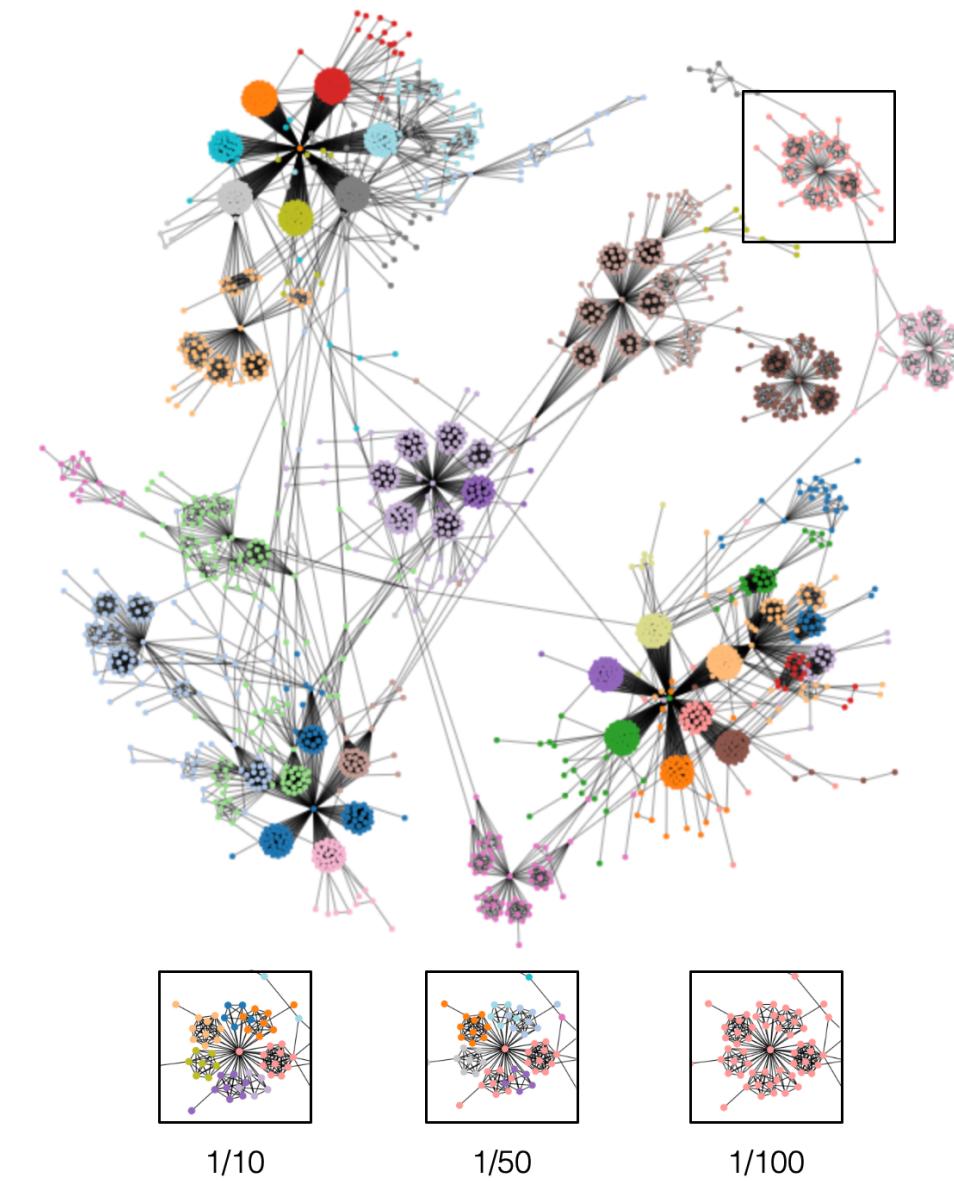
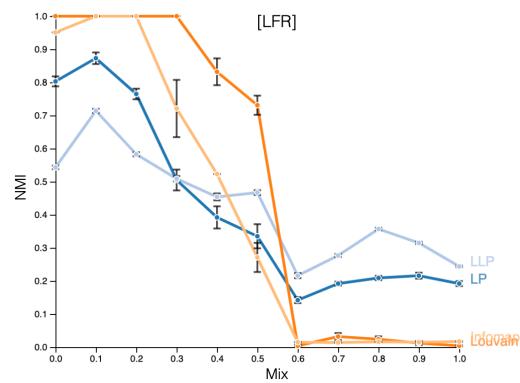
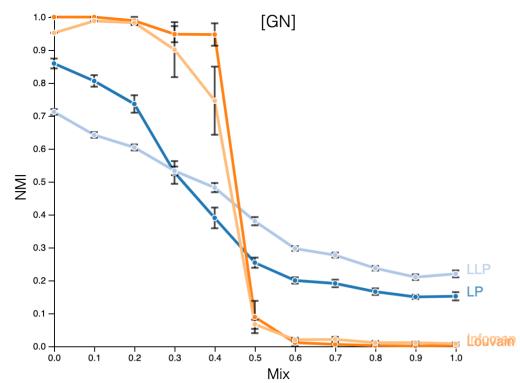


# Results

## *Staphylococcus aureus*



Louvain & D3.js (Canvas)  
★★★★★

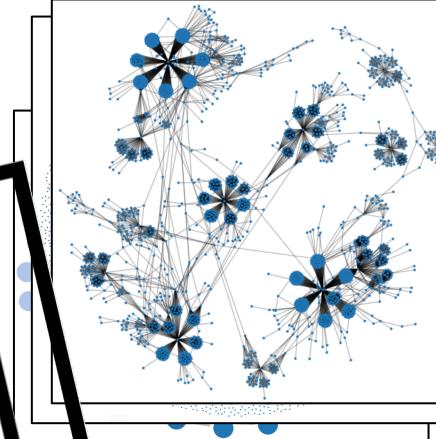
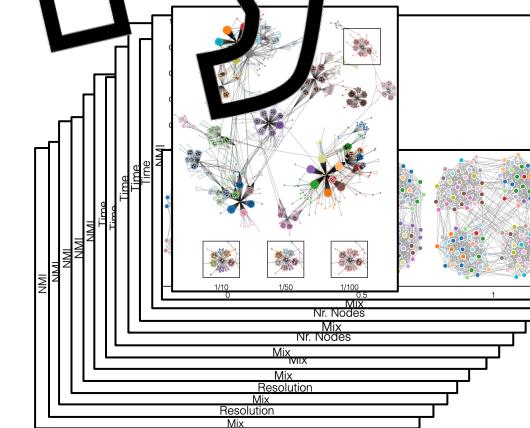
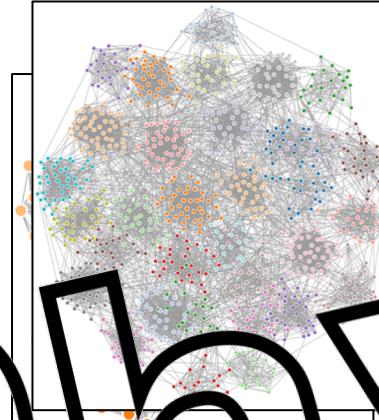
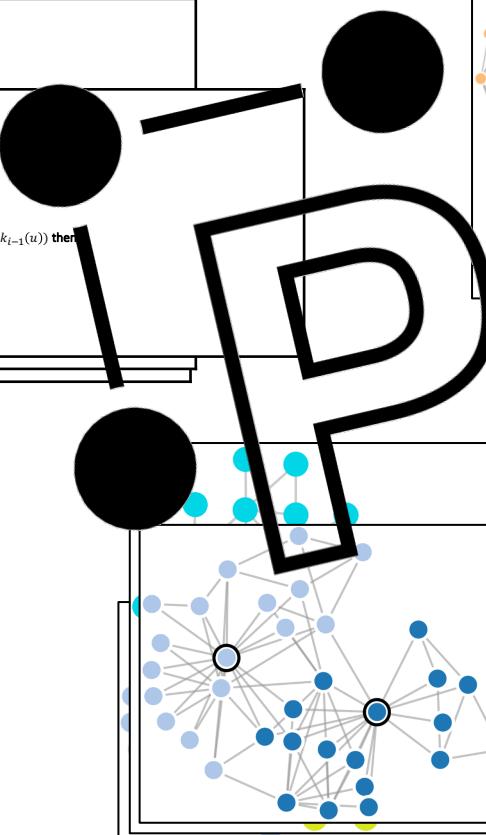


# Conclusions

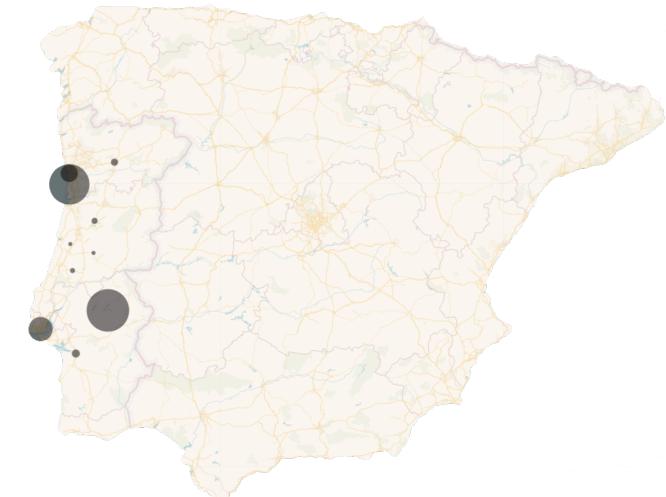
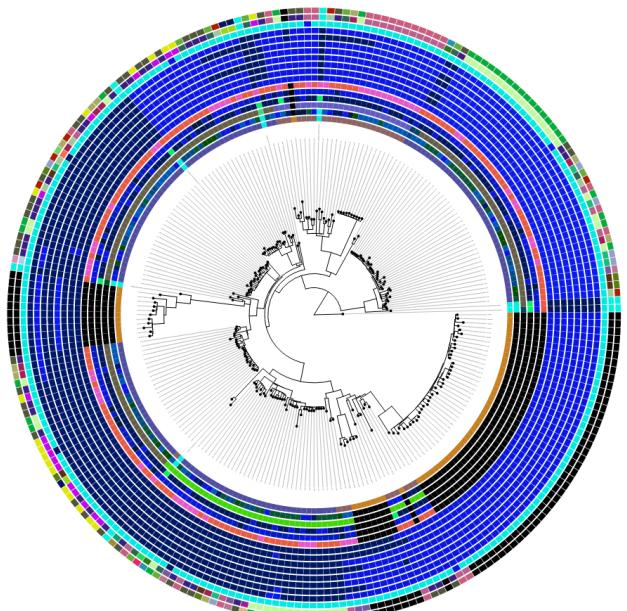
```

1: Compute  $Mod_{new} = Mod(M)$ 
2: repeat
3:    $Mod = Mod_{new}$ 
4:   Randomize the order of vertices.
5:   for all  $i \in V^k$  do
6:      $best_{com} = M_i^k$ ;  $best_{increase} = 0$ 
7:     for all  $M' \in C^k$  do
8:        $M_i^k = M_i^k \setminus \{i\}$ 
9:        $\Sigma_i^k = M_i^k \cup \Sigma_i^{k-1}$ 
10:       $L_i^k = L(M_i^k)$ 
11:      if  $L_i^k > best_{increase}$  then
12:         $best_{increase} = L_i^k$ 
13:         $best_{com} = M_i^k$ 
14:      end if
15:    end for
16:     $k = k + 1$ 
17:  until No vertex movement or  $k < max_{iteration}$ 
18:   $V^{k+1} \leftarrow C^k$ ;  $E^{k+1} \leftarrow e(C_{lk}, C_R)$ ;  $G^{k+1} = (V^{k+1}, E^{k+1})$ ;  $k = k + 1$ 
19:   $V \leftarrow C^k$ ;  $L \leftarrow L_i^k$ ;  $G = (V, L)$ ;  $k = k + 1$ 

```



# Future Work



# Acknowledgments

João  
Carriço

01/18

Vítor  
Borges



11/17

Alexandre  
Francisco

10/18

# References

- [Slide 6] "The world's most valuable resource is no longer oil, but data," *The Economist*, 6 May 2017. [Online]. Available: <https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>. [Accessed 28 May 2019];
- [Slide 6] "Ten threats to global health in 2019," World Health Organization, [Online]. Available: <https://www.who.int/emergencies/ten-threats-to-global-health-in-2019>. [Accessed 4 June 2019];
- [Slide 6] "Top 10 Leading Causes of Death Globally," [Online]. Available: <https://www.theatlas.com/charts/HkLaDreuW>. [Accessed 12 May 2019];
- [Slide 6] Y. Motro and J. Moran-Gilad, "Next-generation sequencing applications in clinical bacteriology," *Biomolecular Detection and Quantification*, vol. 14, p. 1–6, 2017;
- [Slide 7] V. D. Blondel, J.-L. Guillaume, R. Lambiotte and E. Lefebvre, "Fast unfolding of communities in large networks," *J. Stat. Mech.* (2008) P10008, p. 12, 2008;
- [Slide 7] M. Rosvall, D. Axelsson and C. T. Bergstrom, "The map equation," *The European Physical Journal Special Topics*, vol. 178, no. 1, pp. 13-23, 2009;
- [Slide 7] P. Boldi, M. Rosa, M. Santini and S. Vigna, "Layered label propagation: a multiresolution coordinate-free ordering for compressing social networks," in *WWW '11 Proceedings of the 20th international conference on World wide web*, 2011;
- [Slide 8] M. Girvan and M. E. J. Newman, "Community structure in social and biological networks," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 99, no. 12, pp. 7821-7826, 2002;
- [Slide 9] A. Lancichinetti, S. Fortunato and F. Radicchi, "Benchmark graphs for testing community detection algorithms," *Physical review. E, Statistical, nonlinear, and soft matter physics.*, vol. 78, no. 4, 2008;
- [Slide 9] A. Lancichinetti, S. Fortunato and J. Kertesz, "Detecting the overlapping and hierarchical community structure of complex networks," *New Journal of Physics*, vol. 11, 2009;
- [Slide 10] J. Yang and J. Leskovec, "Defining and Evaluating Network Communities based on Ground-truth," in *Proceedings of 2012 IEEE International Conference on Data Mining (ICDM)*, 2012;
- [Slide 11] W. Zachary, "An Information Flow Model for Conflict and Fission in Small Groups," *Journal of anthropological research*, vol. 33, 1976;
- [Slide 11] "Staphylococcus aureus MLST Databases," PubMLST, 5 June 2019. [Online]. Available: <https://pubmlst.org/saureus/>. [Accessed 5 June 2019];
- [Slide 14] A. Lancichinetti and S. Fortunato, "Benchmarks for testing community detection algorithms on directed and weighted graphs with overlapping communities," *Physical review. E, Statistical, nonlinear, and soft matter physics.*, vol. 80, 2009;