

---

# Scaling Overlapping Clustering



**Kyle Kloster**  
(now at NCSU)

**Merrielle Spain**  
**Stephen Kelley**





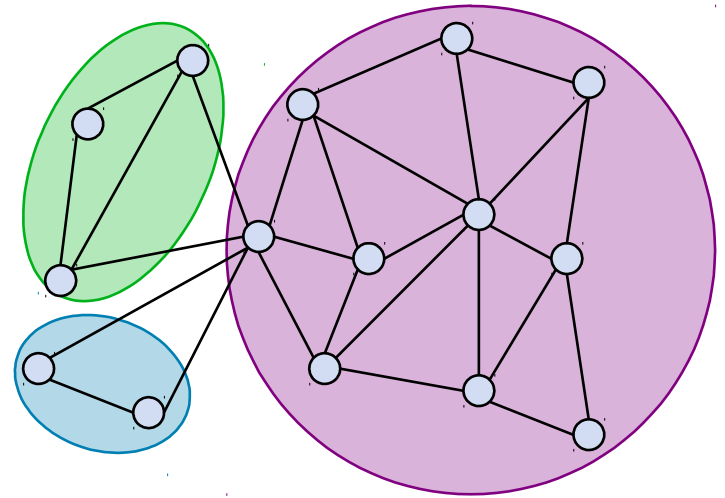
# Scaling graph analytics

## *Broader Problem*

- Many global graph analytics scale poorly to large graphs

## *Case study*

- Speed up overlapping clustering:
  - Link Clustering
  - Baselines: Infomap, BigCLAM





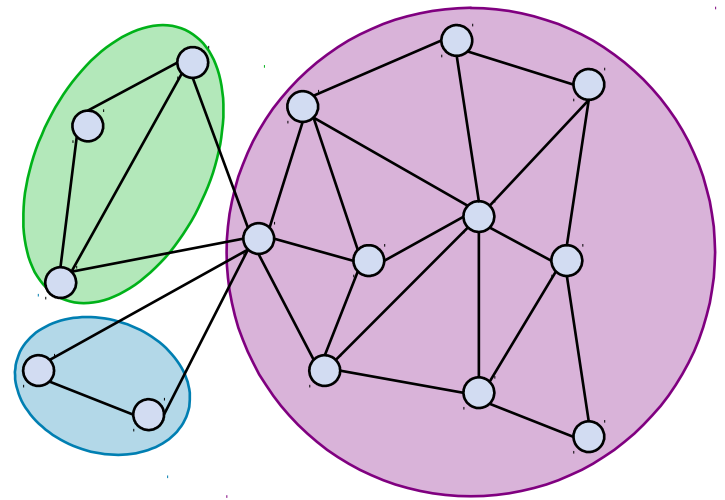
# Scaling graph analytics

## *Broader Problem*

- Many global graph analytics scale poorly to large graphs

## *Case study*

- Speed up overlapping clustering:
  - Link Clustering
  - Baselines: Infomap, BigCLAM



## *Our problem*

- Even state-of-the-art algorithms scale poorly

## *This talk*

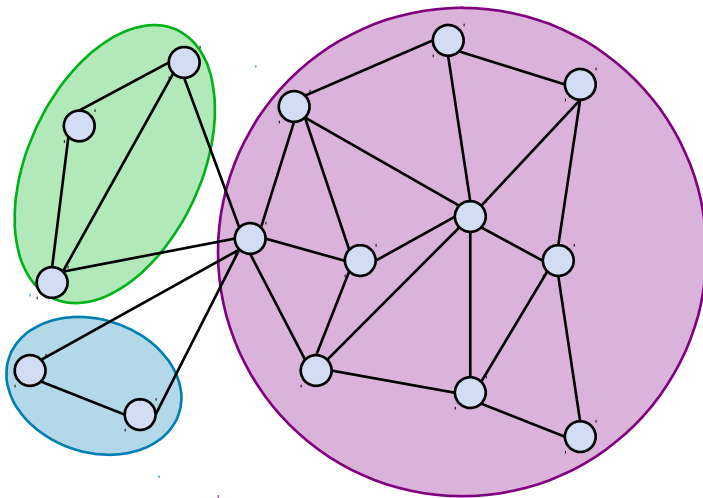
- Tight complexity analysis -> identify, alleviate bottlenecks



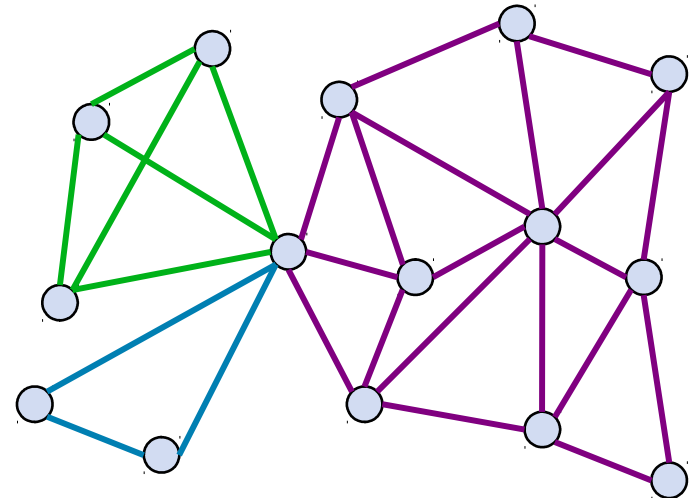
# Our focus: Link Clustering

- Detect communities by partitioning the *edges*
- A node belongs to all communities that its edges do!

## Node clustering

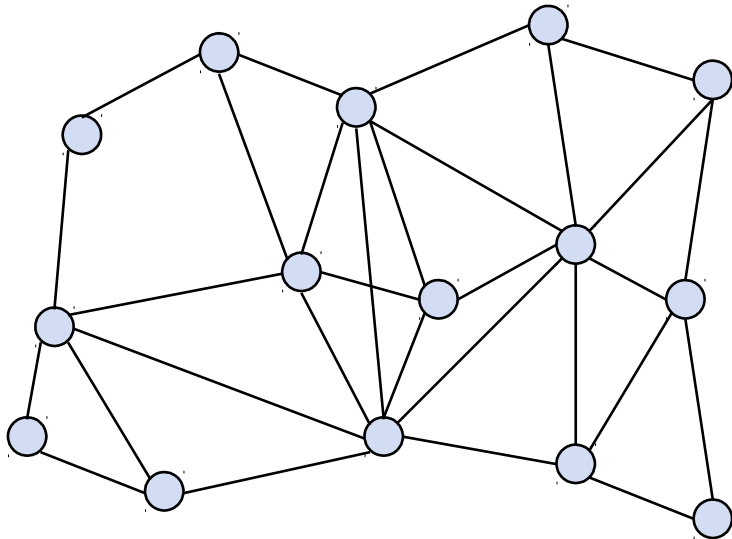


## Link clustering



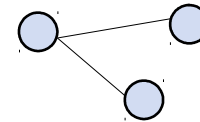


# Which property describes scaling?



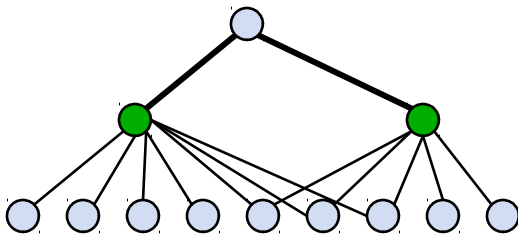
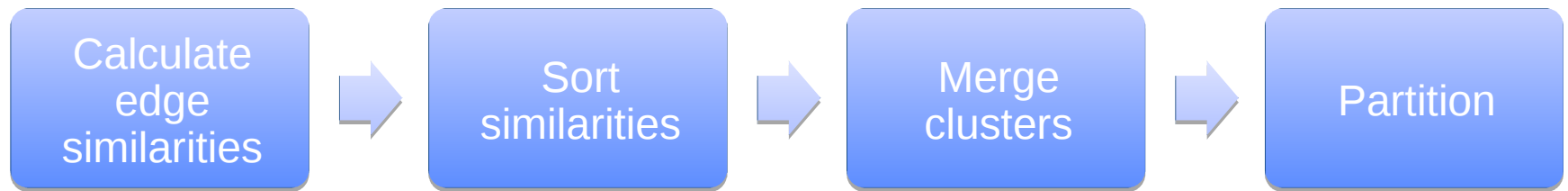
## Features:

- Edges
- Maximum node degree
- Nodes
- Wedges





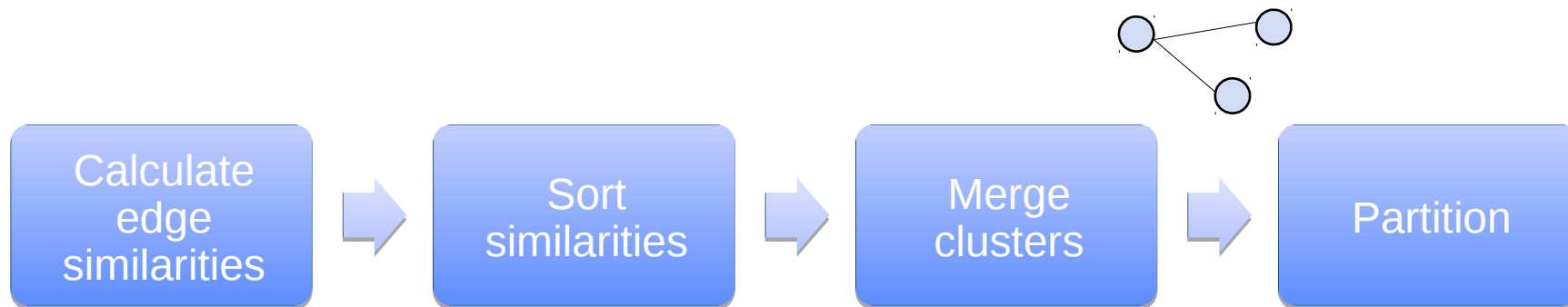
# Link Clustering - Overview





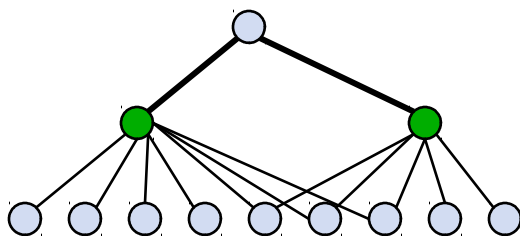
# Link Clustering

Compute Jaccard Similarity of any 2 edges that share a node.



**Jaccard Similarity:**

$$\frac{|\text{neighborhood intersection}|}{|\text{neighborhood union}|}$$

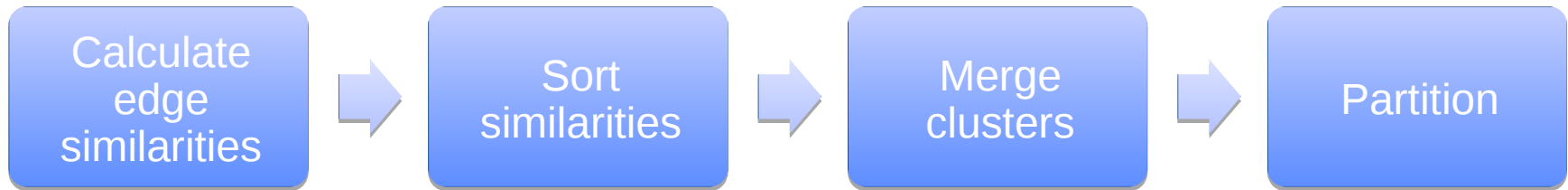




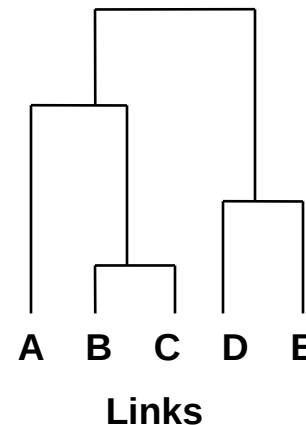
# Link Clustering

Sort all pairs of edges (wedges) by their Jaccard Similarity.

Combine clusters with highly similar edges.



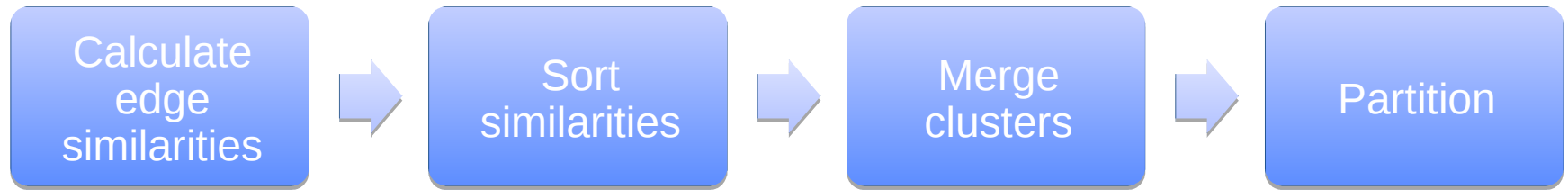
B,C  
D,E  
A,C  
A,B  
B,D  
C,E  
A,D  
A,E  
B,E  
C,D







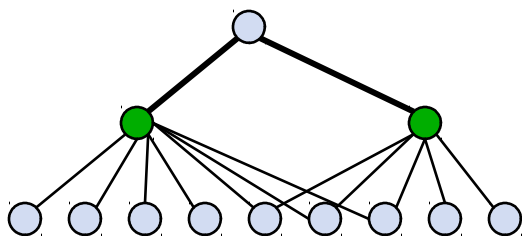
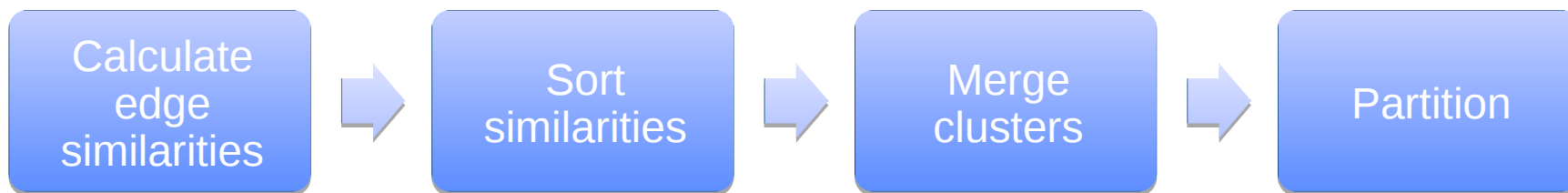
# Link Clustering: tight complexity



$$\mathcal{O}(\text{wedges}) \quad \mathcal{O}(\text{wedges} \times \log_2(\text{wedges}))$$



# Link Clustering



Number of wedges at a node =  
(degree choose 2)  $\sim$  degree<sup>2</sup>

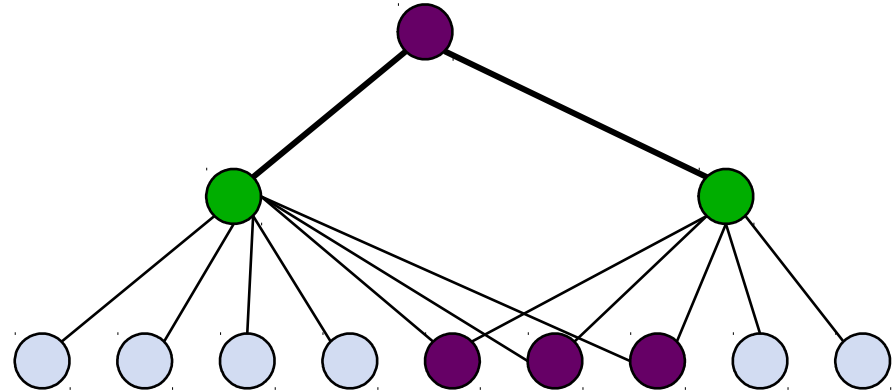
Large degree nodes (hubs) bad...



# Tight complexity identifies bottlenecks

- *Sorting* = dominant subroutine
- Runtime scales with wedges ...
  - Hubs cause  $\sim O(n^2)$  wedges!

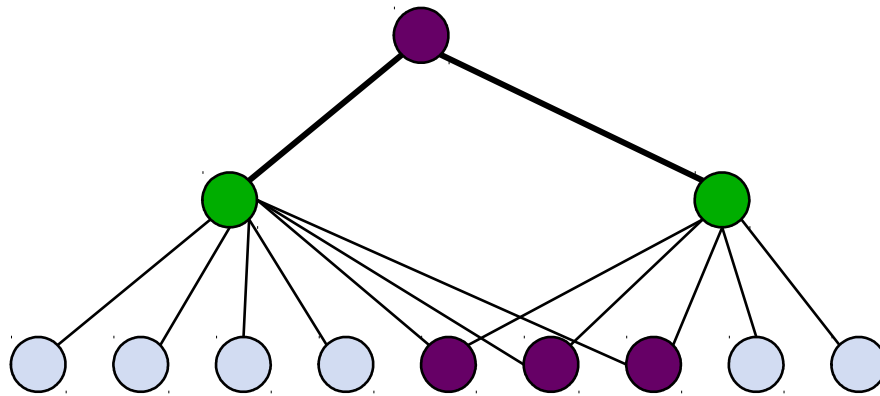
**Bottlenecks: hubs, sorting**





# Similarity scores are structured

## Jaccard Similarity:



- is a fraction: denominator, numerator are integers in [ 1 , 2\*max degree ]
- We show this implies small number of distinct scores...
- Enables linear-time sort!
- “Shelf” sort: bucket/counting

$$\mathcal{O}(\text{wedges} \times \log_2(\text{wedges})) \longrightarrow \mathcal{O}(\text{wedges})$$



# Shelf Sort

$$4 \times \textit{degree}^2$$



...



2



1



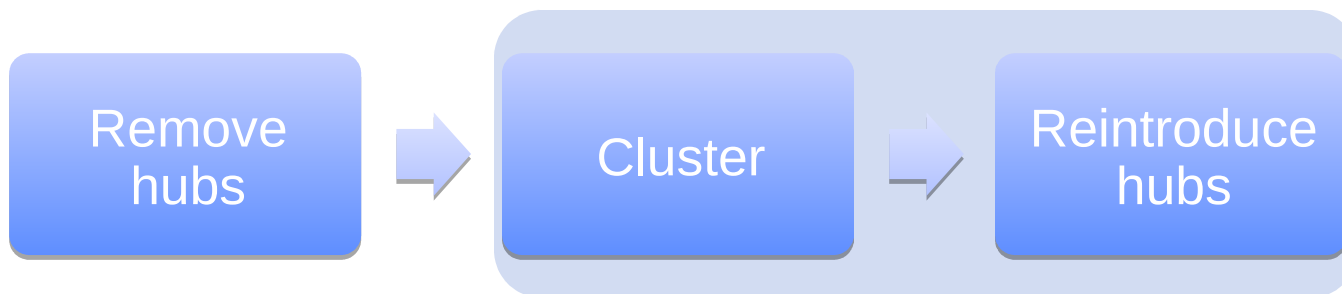
# Handling Hub Bottlenecks

- **Claim: large degree nodes (hubs) can be removed**
  - often meaningless (Spammers)
  - often uninformative (Twitter: Bieber, Obama; bio examples)
  - if important, worth processing individually
- **Hubs drive up runtime**
  - Hub  $\rightarrow \sim n^2$  wedges
- **Hubs can obscure smaller-scale community structure**





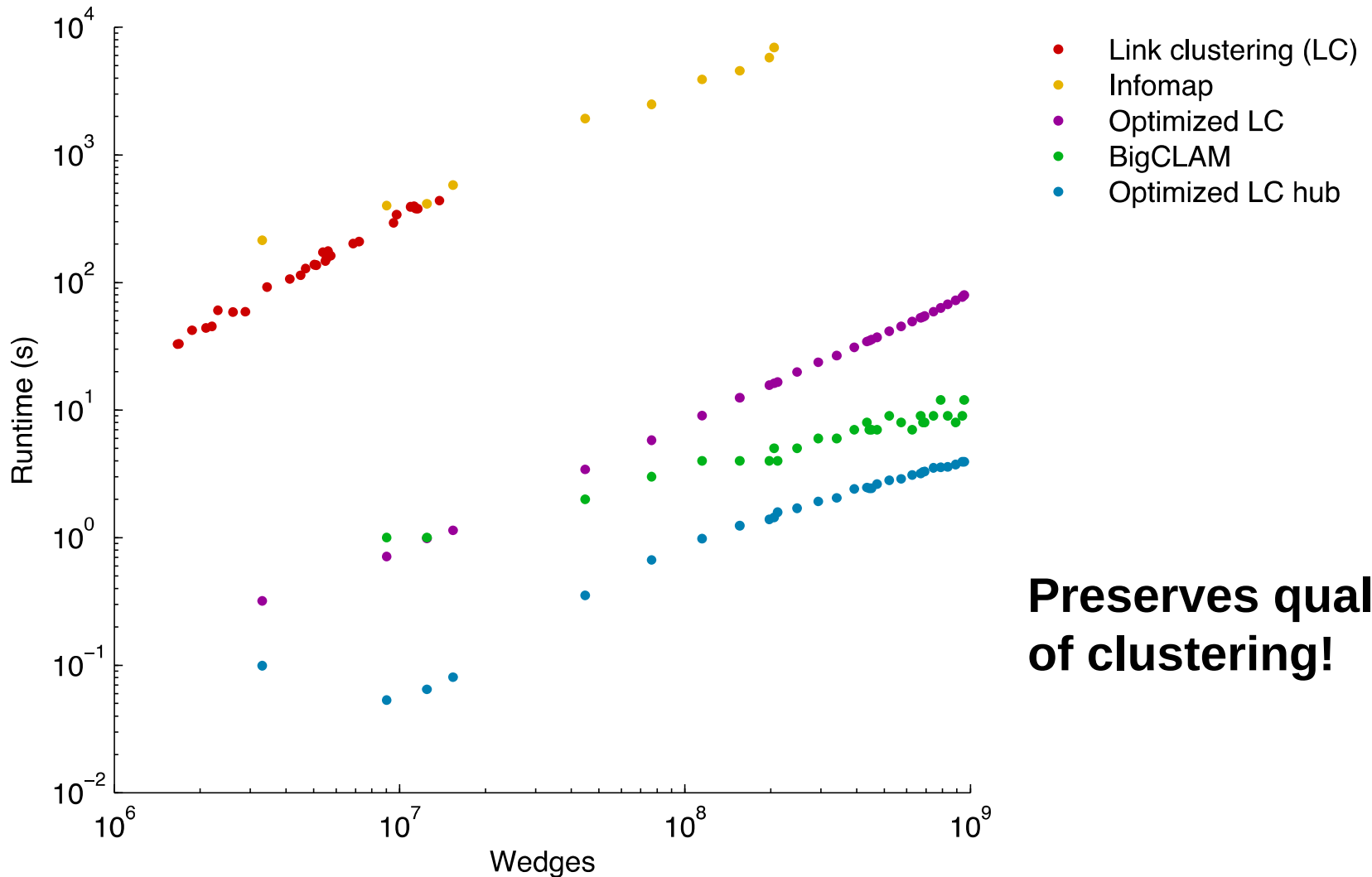
# Hub Removal



- **Add hubs only to clusters they strongly connect to**



# ~10,000 X faster



**Preserves quality  
of clustering!**





# Conclusions

## Scaling graph analytics:

- **Tighter complexity analysis can reveal bottlenecks**
  - Runtime =  $O(\text{sorting wedges})$
  - Dominant subroutines, problematic graph features
- **Exploiting bottleneck analysis**
  - Remove bottleneck graph feature; handle separately?
  - Graphs are discrete  $\rightarrow$  structure! (e.g. sorting = linear)
- **Context can guide handling bottleneck**
  - are your bottlenecks disposable? (spam, Justin Bieber)
  - ... worth processing separately?