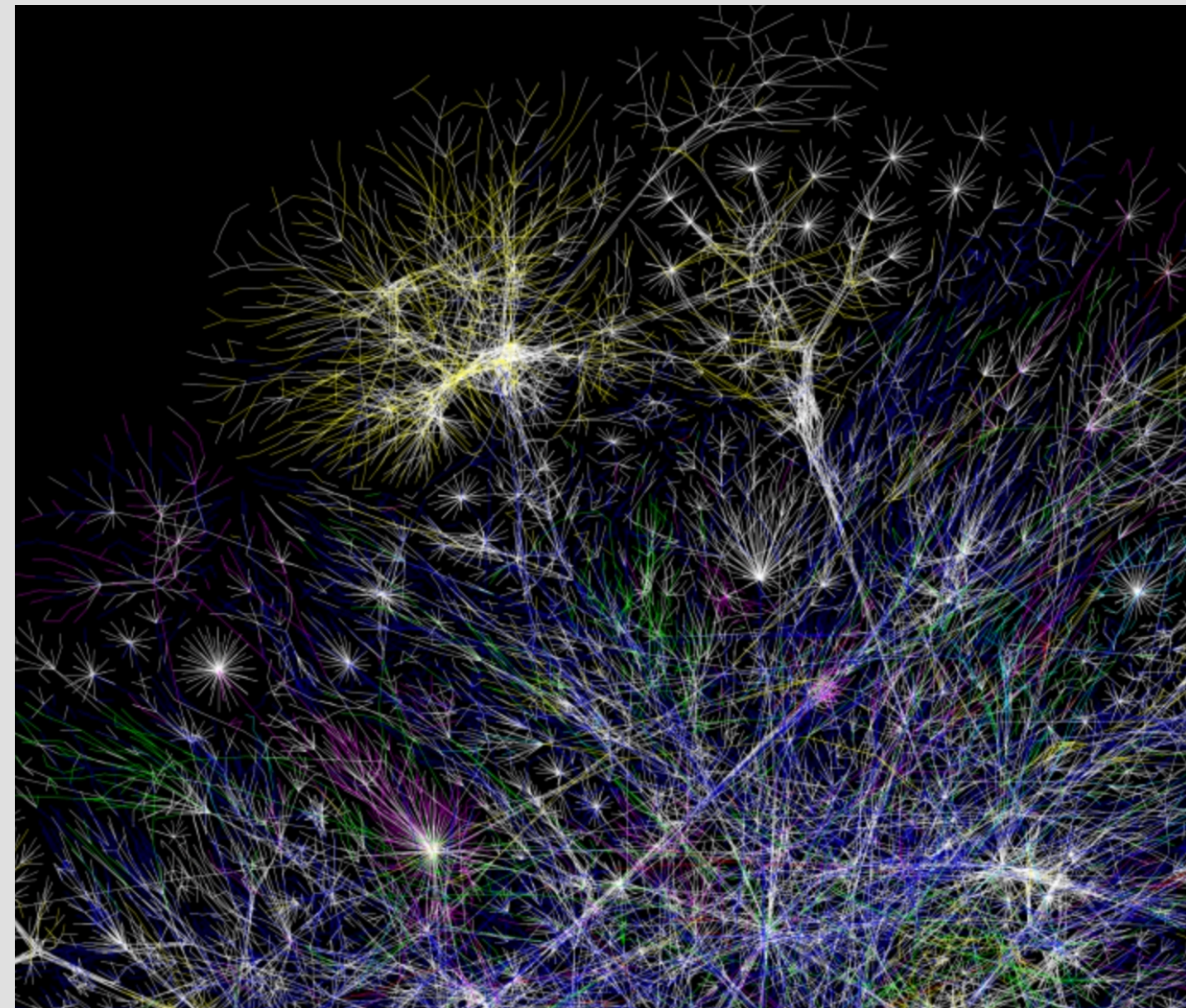# Combinatorial $\ell_p$-norm Correlation Clustering

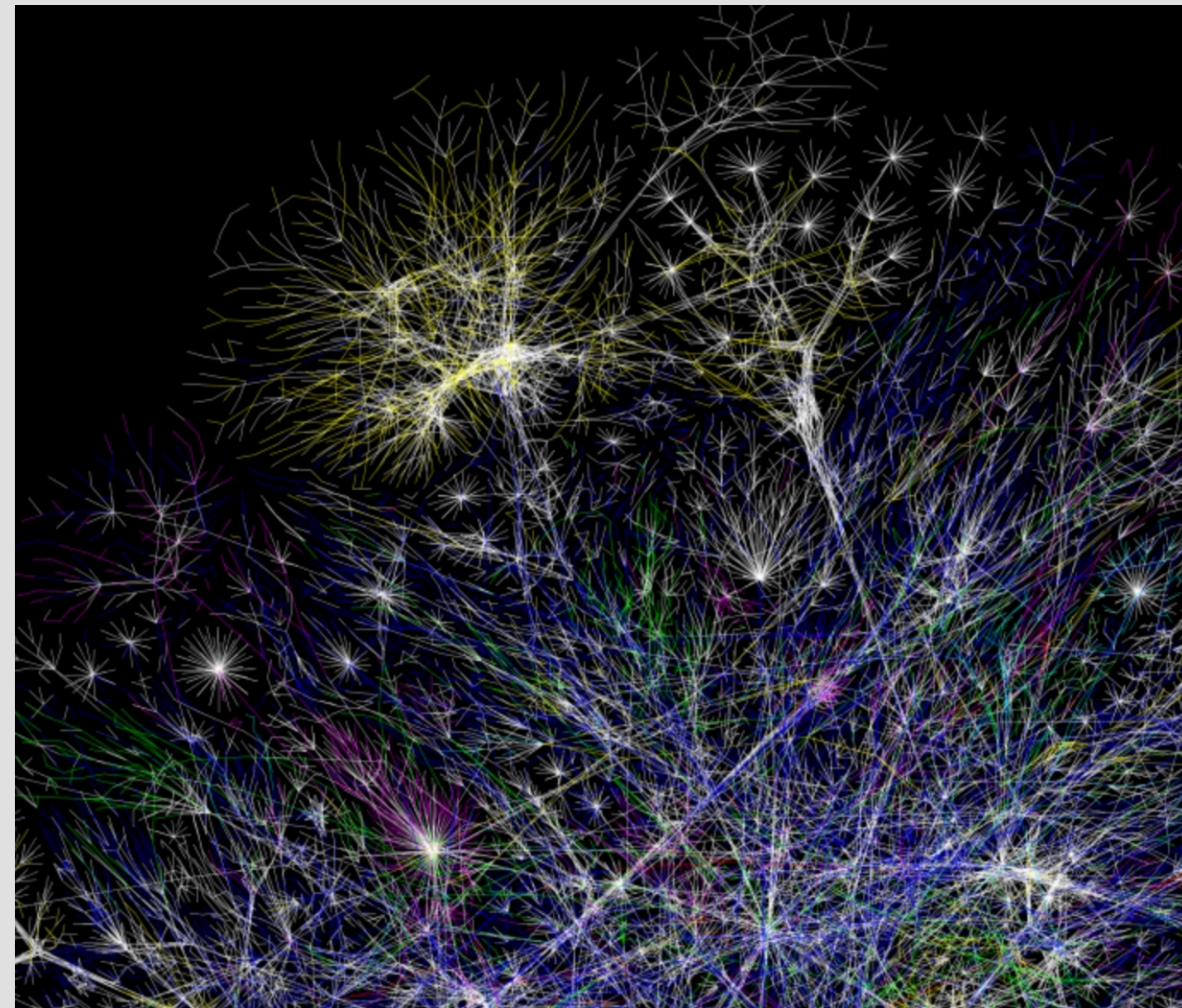**Sami Davies (UC Berkeley/ Simons)**, Benjamin Moseley (CMU), Heather Newman (CMU)

# Community detection



✦ Creating large-scale maps with meta nodes

✦ Understanding *community* vs *aggregate* features

✦ Identifying topological/ spectral properties

# Community detection

## Many different models



✦ Hierarchical clustering — Best for data with underlying heirarchy

✦ Minimum cut clustering — Fixed # of clusters

✦ Girvan Newman algorithm — Runtime $O(m^2n)$

✦ Modularity maximization — Doesn't find small clusters

✦

# Community detection

## Many different models

✦ Hierarchical clustering — Best for data with underlying heirarchy

✦ Minimum cut clustering — Fixed # of clusters

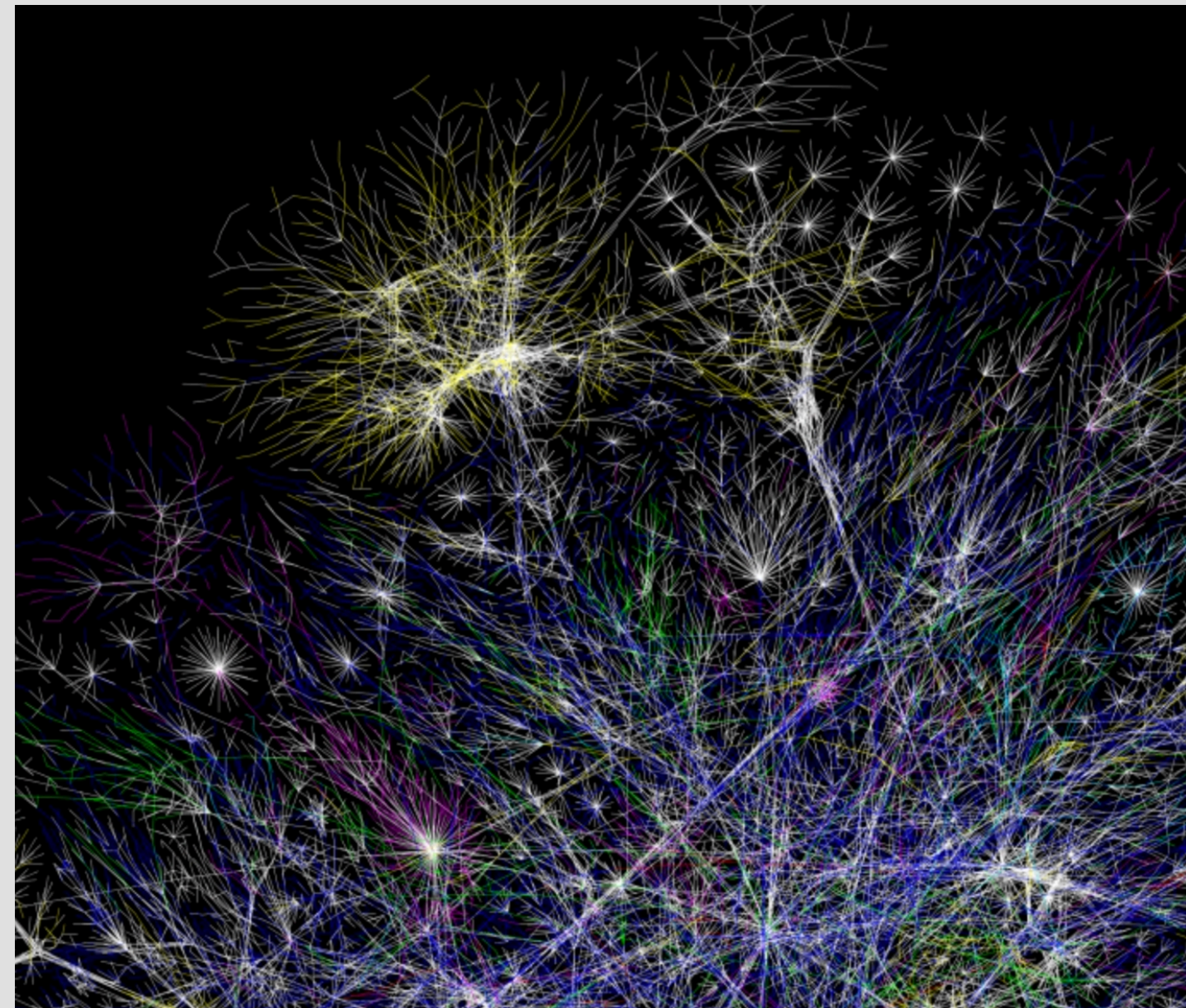✦ Girvan Newman algorithm — Runtime $O(m^2 n)$
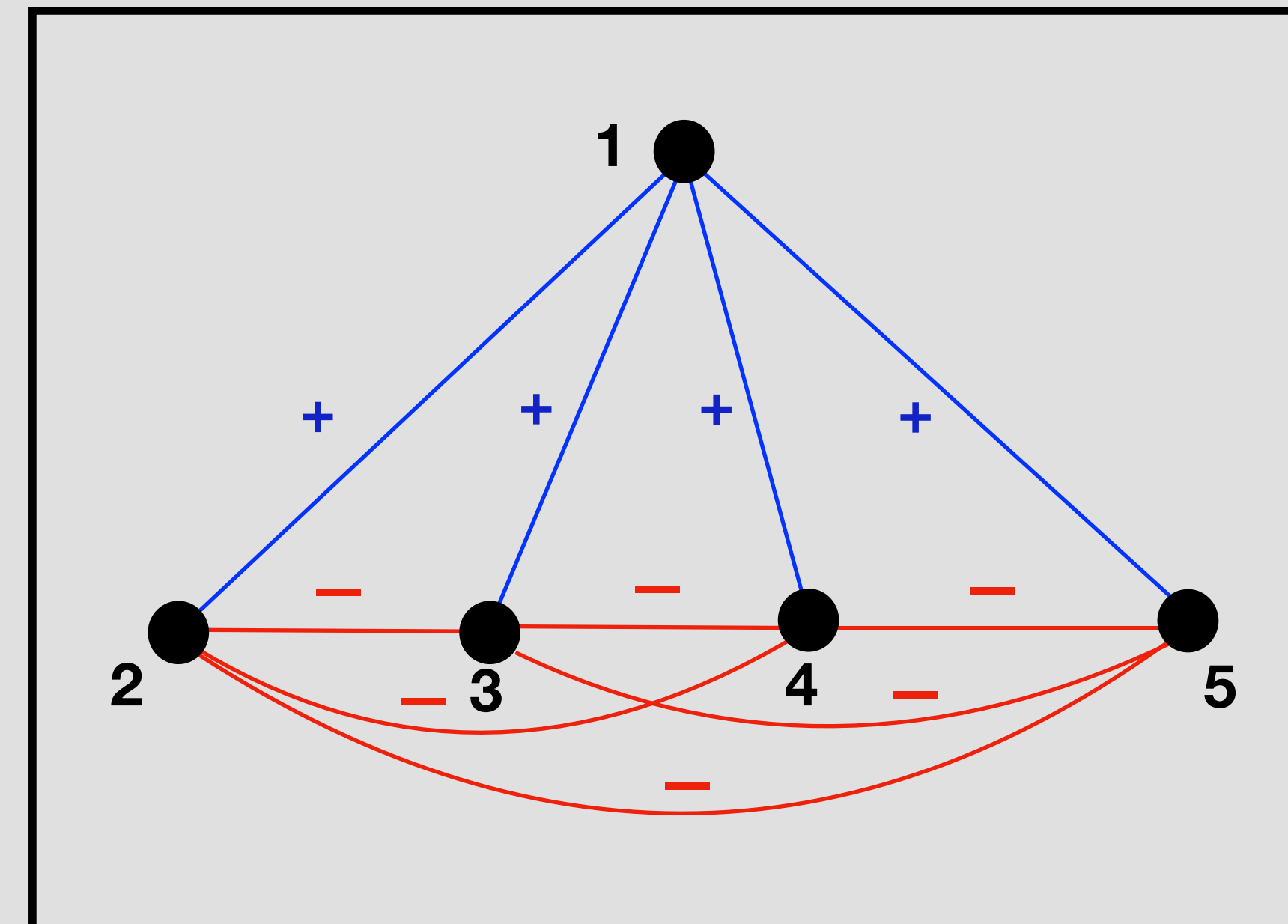
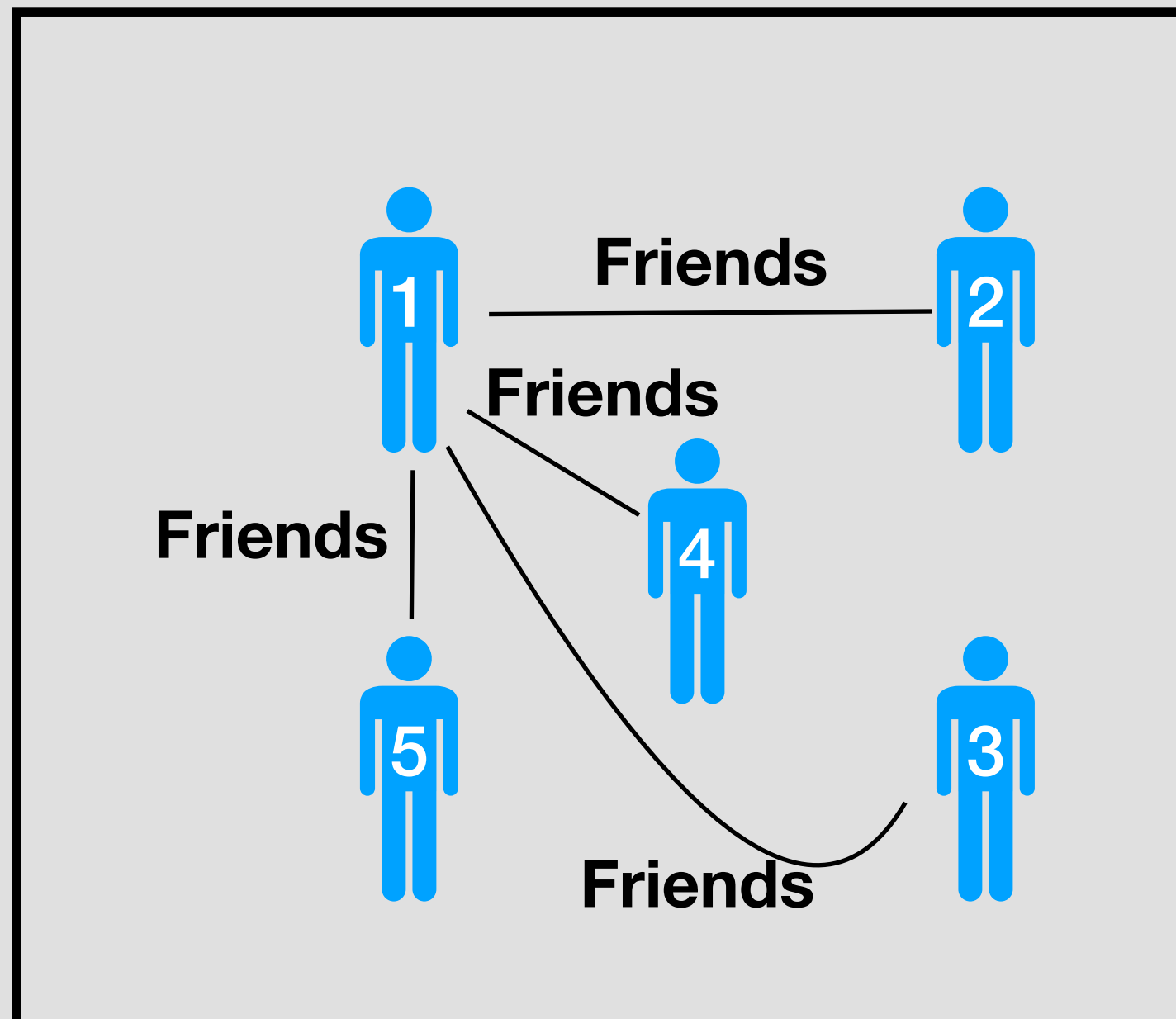✦ Modularity maximization — Doesn't find small clusters

✦ **Correlation clustering**

# Correlation clustering

**Model**:

- ‣ Cluster *similar* nodes together, separate *dissimilar* nodes
- ‣ No pre-fixed # of clusters, complete unweighted graph

# Correlation clustering

**Model**:

‣ Cluster *similar* nodes together, separate *dissimilar* nodes

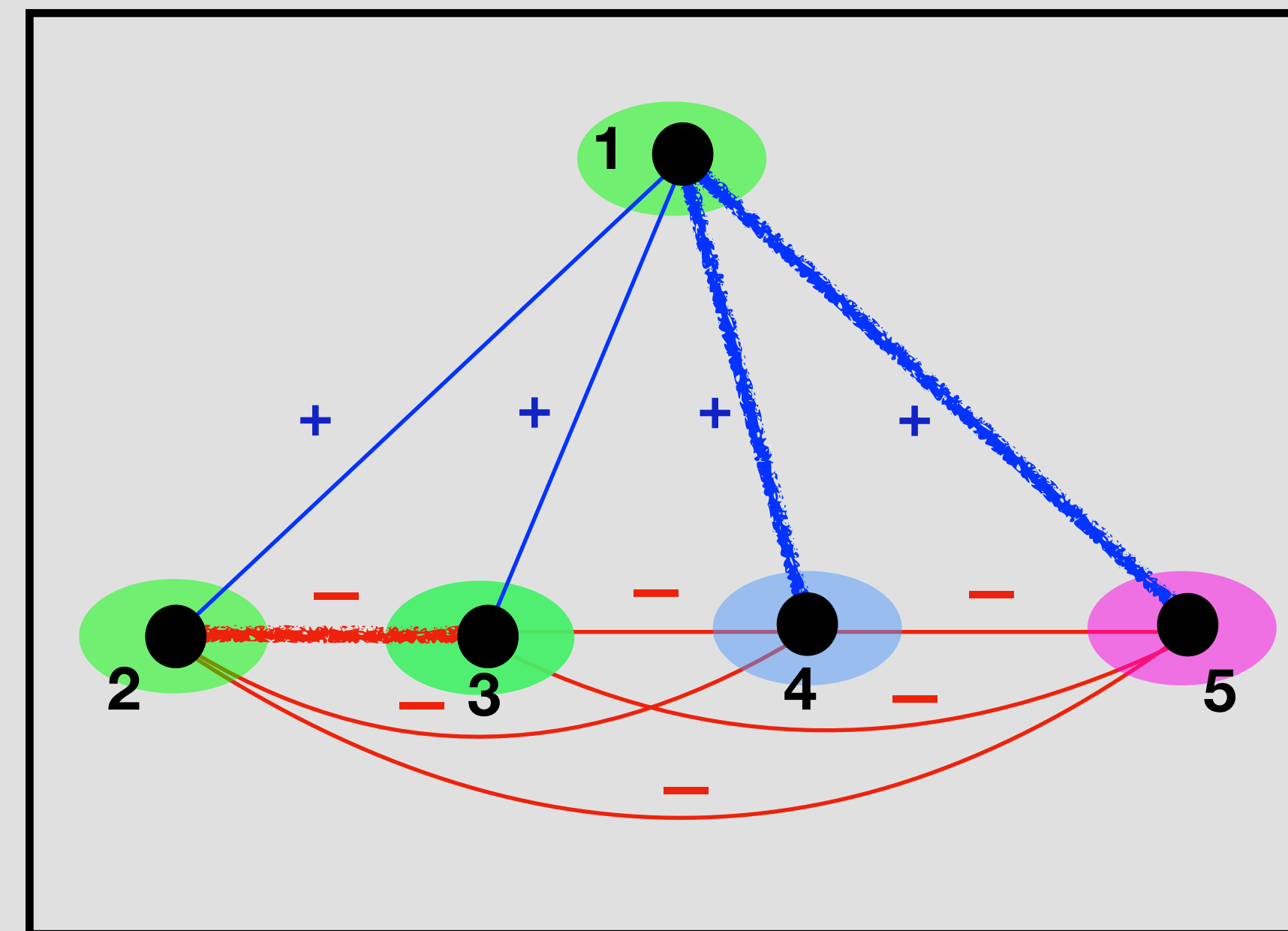‣ No pre-fixed # of clusters, complete unweighted graph

‣ Edge *(u,v)* in **disagreement** w.r.t $C$ if

✦ (+) with $u, v$ different clusters or

✦ (−) with $u, v$ same cluster

Original objective for Correlation Clustering = minimize # of edges in disagreement

# $\ell_p$ Correlation clustering

**Model**:
- ‣ Cluster **similar** nodes together, separate **dissimilar** nodes
- ‣ No pre-fixed # of clusters, complete unw...
- ‣ Edge *(u,v)* in **disagreement** w.r.t *C* if
  - ✦ (+) with *u, v* different clusters or
  - ✦ (−) with *u, v* same cluster
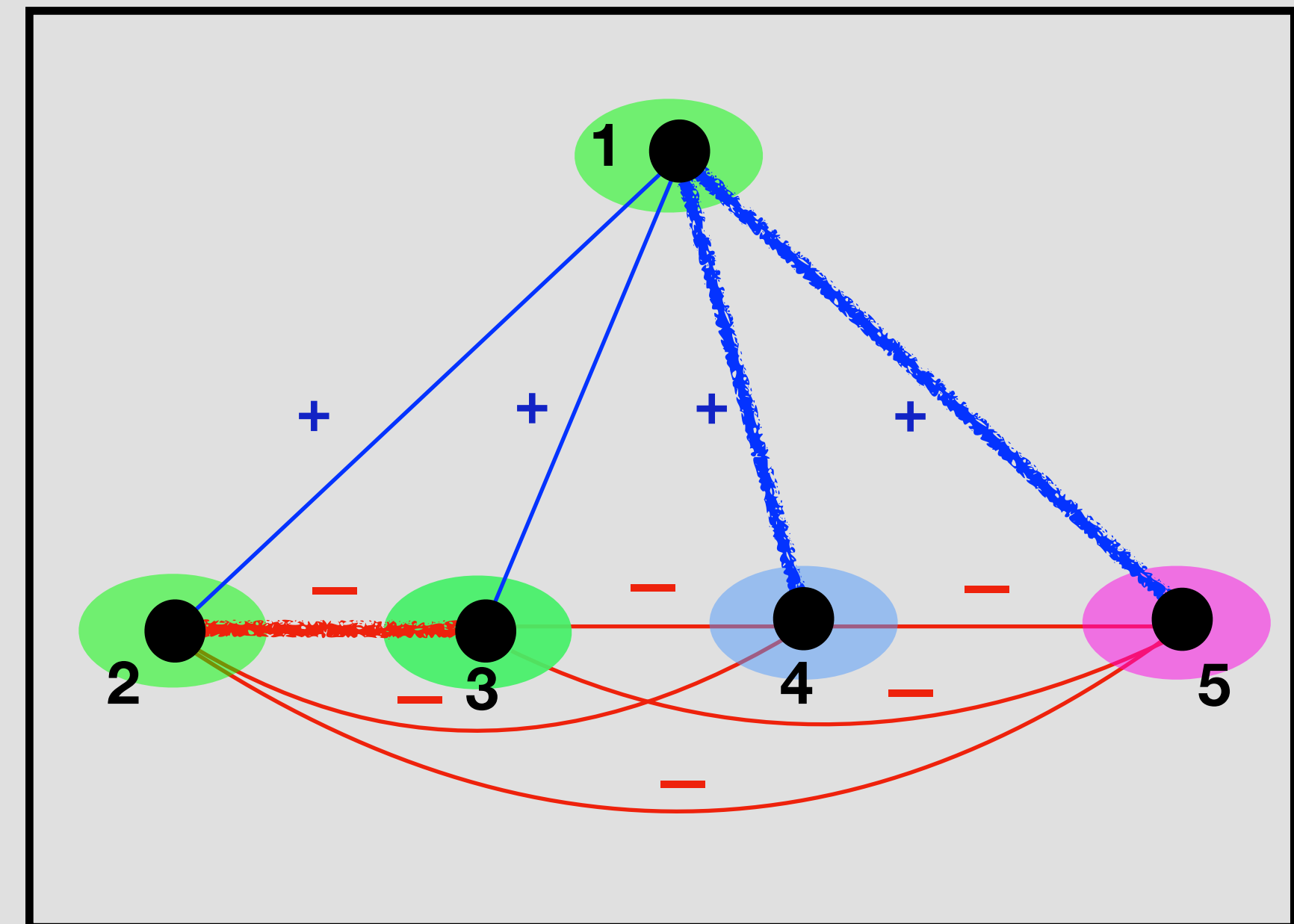
- ‣ $y_i^C$ = # disagreements w.r.t. C incident to $v_i$
- ‣ Goal: find $\mathbf{argmin}_C ||y^C||_p$

*p≥1*

$\ell_1$ = original cc
$\ell_\infty$ = min max norm

$p$ small = global obj   $\leftrightarrow$   $p$ large = local/fair  obj

# Previous work

**For $\ell_1$-norm (original) objective:**

▸ Introduced by [Bansal, Blum, Chawla '04]

▸ Linear time Pivot algorithm gives 3-apx

[Ailon, Charikar, Newman JACM08] [Chierichetti, Dalvi, Kumar KDD14]

▸ APX-hard

[Charikar, Guruswami, Wirth JCSS05]

▸ Many other active threads of research!

[Ahmadi, Khuller, Saha IPCO19] [Veldt ICML22] [Cohen-Addad, Lee, Li, Newman FOCS23]

# Previous work

**For $\ell_1$-norm (original) objective:**

‣ Introduced by [Bansal, Blum, Chawla '04]

‣ Linear time Pivot algorithm gives 3-apx

 [Ailon, Charikar, Newman JACM08] [Chierichetti, Dalvi, Kumar KDD14]

‣ APX-hard

 [Charikar, Guruswami, Wirth JCSS05]

‣ Many other active threads of research!

 [Ahmadi, Khuller, Saha IPCO19] [Veldt ICML22] [Cohen-Addad, Lee, Li, Newman FOCS23]

**For general $\ell_p$-norm objectives:**

‣ 5-approximation algorithm; NP-hard

 [Puleo, Milenkovic ICML16], [Charikar, Gupta, Schwartz IPCO17], [Kalhan, Makarychev, Zhou ICML19]

‣ Techniques round solution to a convex program

# Previous work

**Pivot algorithm**

- Randomly choose a *pivot (*unclustered vertex)
- Make new cluster with pivot and all its unclustered positive neighbors
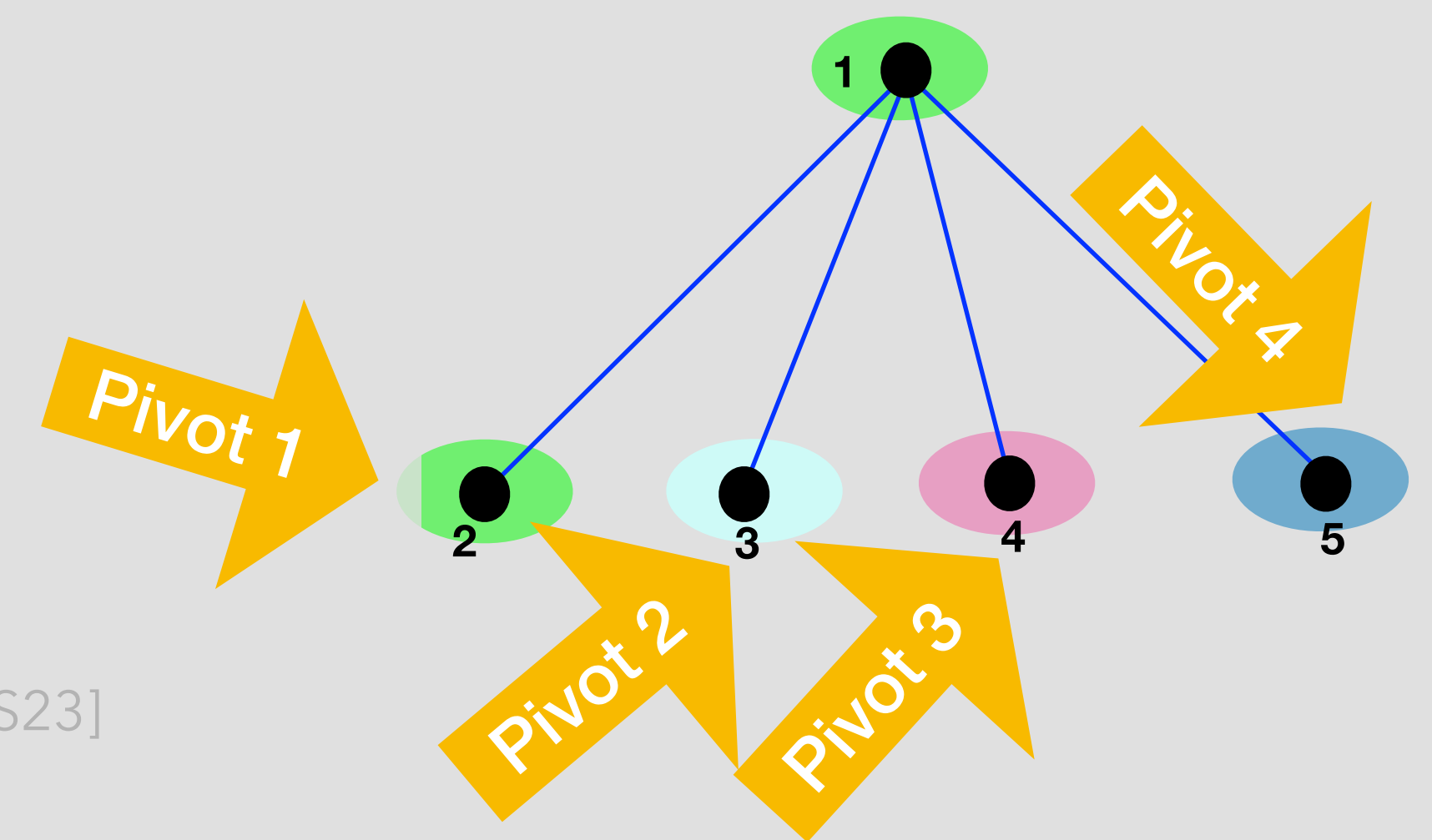
**For $\ell_1$-norm (original) objective:**

▸ Introduced by [Bansal, Blum, Chawla '04]

▸ Linear time Pivot algorithm gives 3-apx

[Ailon, Charikar, Newman JACM08] [Chierichetti, Dalvi, Kumar KDD14]

▸ APX-hard

[Charikar, Guruswami, Wirth JCSS05]

▸ Many other active threads of research!

[Ahmadi, Khuller, Saha IPCO19] [Veldt ICML22] [Cohen-Addad, Lee, Li, Newman FOCS23]

**For general $\ell_p$-norm objectives:**

▸ 5-approximation algorithm; NP-hard

[Puleo, Milenkovic ICML16], [Charikar, Gupta, Schwartz IPCO17], [Kalhan, Makarychev, Zhou ICML19]

▸ Techniques round solution to a convex program

# Trouble with the convex program

$\ell_p$-norm correlation clustering algs solve a convex program

Solving **metric constrained** LPs on large networks is slow!

Not very amenable to online settings

Solution specific to **one fixed** $\ell_p$-norm

Work on solving CC LPs fast only scales to graphs with few thousand vertices!

[Ruggles et al. '20], [Sonthalia & Gilbert '20], [Veldt '22]

**All-norms objective** = simultaneously optimize all $\ell_p$-norms
Introduced by [Azar, Epstein, Richter, Woeginger '04]

**Universal algorithms** produce a solution good for many objs
In, e.g., Steiner tree, TSP, clustering

# Trouble with the convex program

$\ell_p$-norm correlation clustering algs solve a convex program

Solving **metric constrained** LPs on large networks is slow!

Work on solving CC LPs fast only scales to graphs with few thousand vertices!

[Ruggles et al. '20], [Sonthalia & Gilbert '20], [Veldt '22]

Not very amenable to online settings

Solution specific to **one fixed** $\ell_p$-norm

**All-norms objective** = simultaneously optimize all $\ell_p$-norms
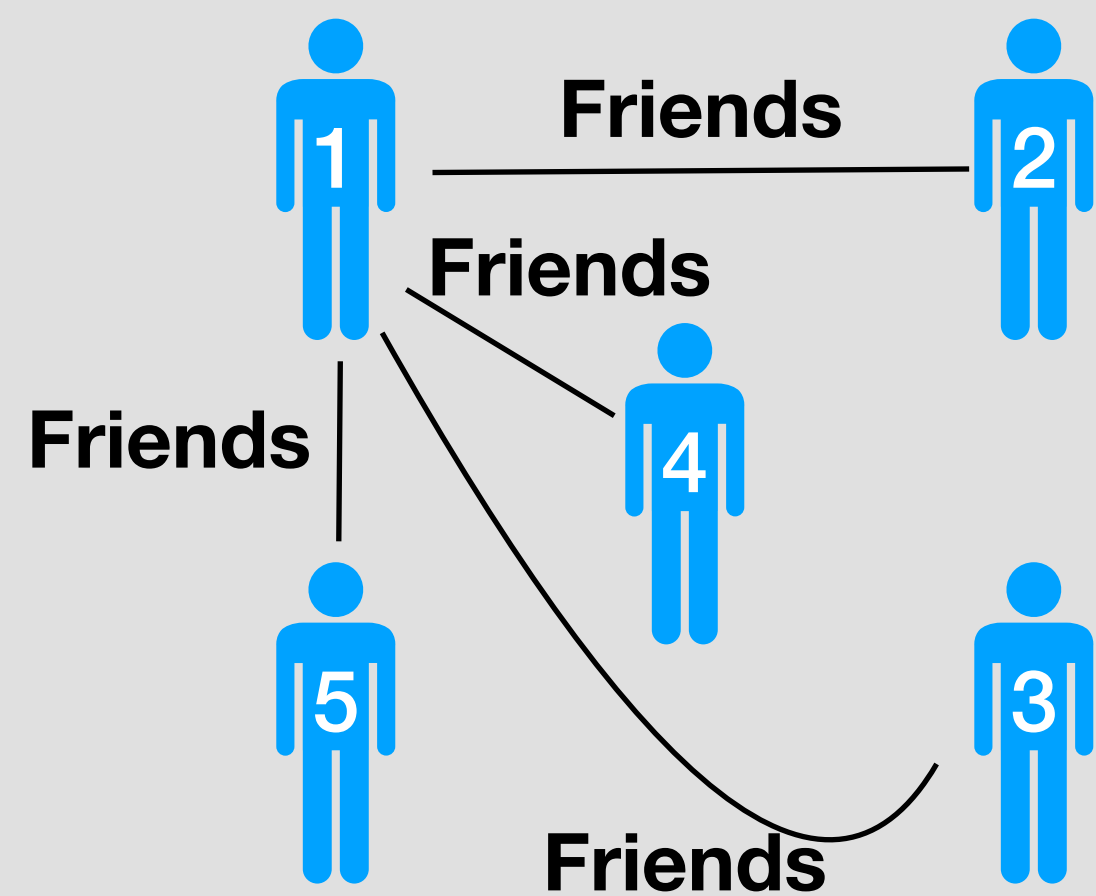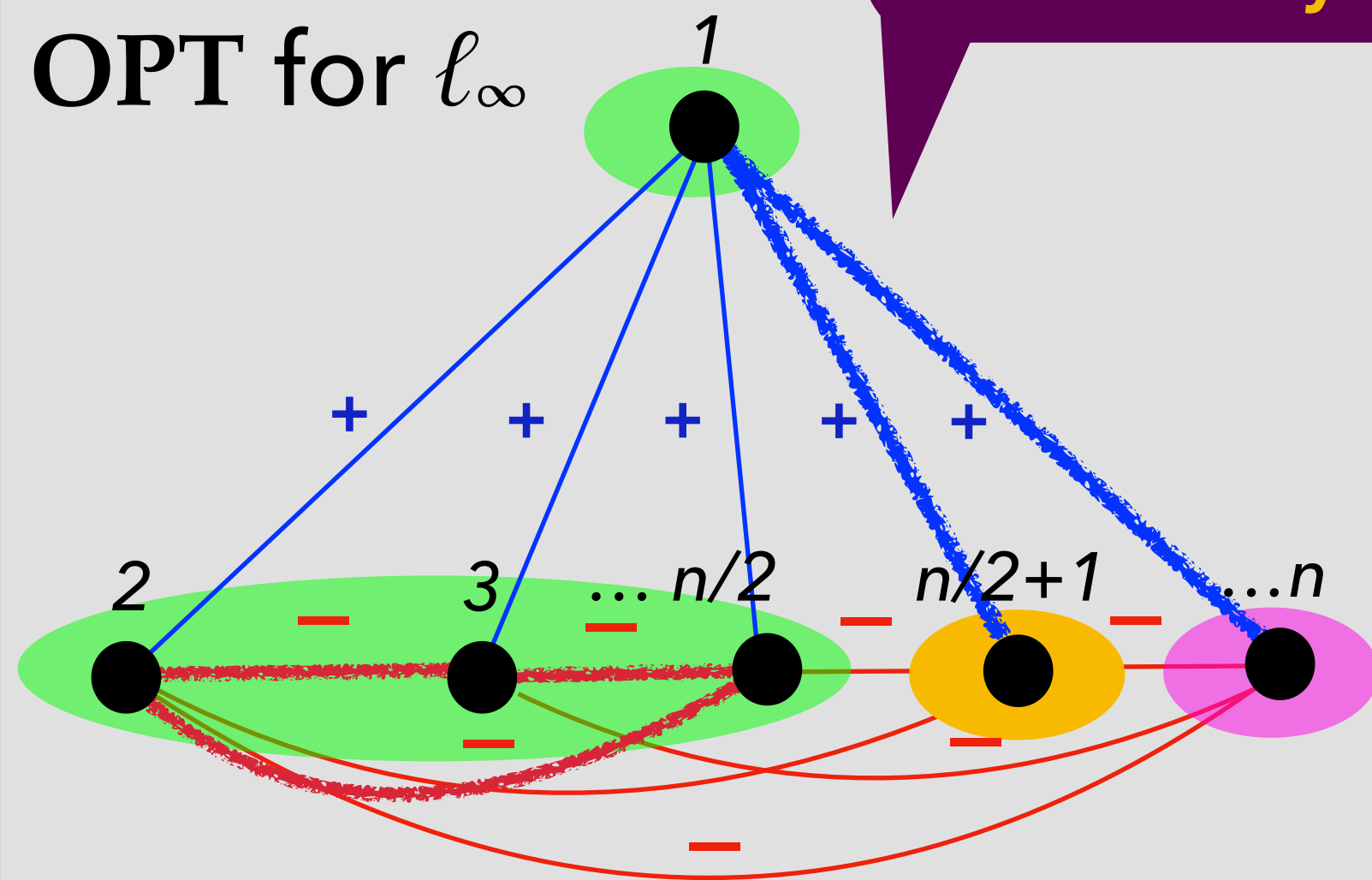Introduced by [Azar, Epstein, Richter, Woeginger '04]

**Universal algorithms** produce a solution good for many objs
In, e.g., Steiner tree, TSP, clustering

# Trouble with the convex program

**OPT for one $\ell_p$-norm can be really bad for others!**

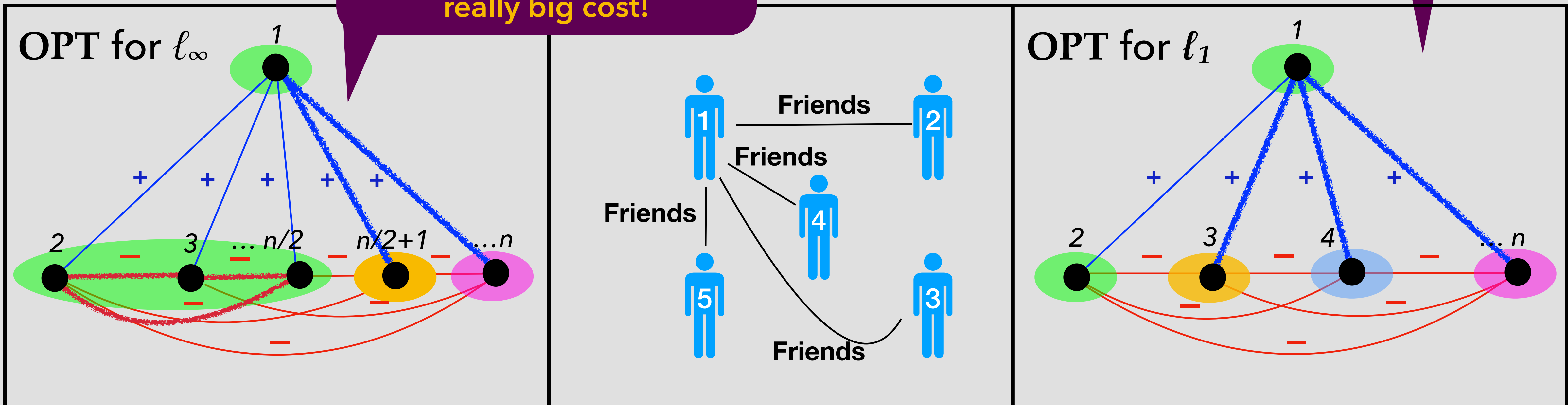Cost for $\ell_1$ norm is $\theta(n^2)$, **really big cost!**

OPT for $\ell_\infty$

1

$+$ $+$ $+$ $+$ $+$

2    3    ... n/2    n/2+1    ...n

$-$  $-$  $-$  $-$  $-$

$-$  $-$

$-$

1 —Friends— 2

Friends

Friends

5    4    3

Friends

# Trouble with the convex program

$\ell_p$-norm correlation clustering algs solve a convex program

Solving *metric constrained* LPs on large networks is slow!

Not very amenable to online/ streaming settings

Solution specific to *one fixed $\ell_p$-norm*

# Our combinatorial approach

"Fast Combinatorial Algorithms for Min Max Correlation Clustering"
*ICML23*

"One Partition Approximating All $\ell_p$-norm Objectives in Correlation Clustering"
*In sub*

Initial constant was 40
Heidrich, Irmai, Andres built off us, improve to 4!

(1) Develop faster *O(1)*-apx alg for min max objective;
⤷ near-linear time on networks with small positive degree

(2) Find *simultaneously O(1)*-apx clustering for all $\ell_p$-norm objs

(3) Algorithms in the online setting

Not possible for *k*-center & *k*-median
[Alamdari & Shmoys WAOA17]

*In progress*

16

# Today

✦ **Introduction** (the model, prior work, our results) 🔥

✦ The correlation metric (constructing a "guess" for the fraction solution, an inherent asymmetry) 🔥 🔥

✦ Proof sketch for the $\ell_\infty$-norm 🔥 🔥 🔥

✦ Adjusting the correlation metric (regular graphs are easy, dealing with negative edges ) 🔥 🔥

✦ Conclusions (mainly vibes) 🔥

# Today

# Previous techniques

## Integer convex program

Not (generally) practical to solve

$x_{uv} = 0$ then $u, v$ same cluster
$x_{uv} = 1$ then $u, v$ different clusters

$$\min ||y||_p$$

$$y(u) = \sum_{v \in N_u^+} x_{uv} + \sum_{v \in N_u^-} (1 - x_{uv}) \qquad \forall u \in V$$

$$x_{uv} \leq x_{vw} + x_{uw} \qquad \forall u, v, w \in V$$

$$x_{uv} \in \mathbb{Z}_{\geq 0} \qquad \forall u, v \in V$$

## Convex program relaxation

Can be solved efficiently

$$\min ||y||_p$$

$$y(u) = \sum_{v \in N_u^+} x_{uv} + \sum_{v \in N_u^-} (1 - x_{uv}) \qquad \forall u \in V$$

$$x_{uv} \leq x_{vw} + x_{uw} \qquad \forall u, v, w \in V$$

$$0 \leq x_{uv} \leq 1 \qquad \forall u, v \in V$$

## Integer convex program

Not (generally) practical to solve

$x_{uv} = 0$ then $u, v$ same cluster

$x_{uv} = 1$ then $u, v$ d

$y(u) = \sum_{v \in N_u^+} x_{uv} + \sum_{v \in N_u^-}$

$x_{uv} \leq x_{vw} + x_{uw}$ $\quad \forall u, v, w \in V$

$x_{uv} \in \mathbb{Z}_{\geq 0}$ $\quad \forall u, v \in V$

**Step 1: Solve convex program**
**Step 2: "Round" fractional solution to integral one**

## Convex program relaxation

Can be solved efficiently

$$\min ||y||_p$$

$$y(u) = \sum_{v \in N_u^+} x_{uv} + \sum_{v \in N_u^-} (1 - x_{uv}) \qquad \forall u \in V$$

$$x_{uv} \leq x_{vw} + x_{uw} \qquad \forall u, v, w \in V$$

$$0 \leq x_{uv} \leq 1 \qquad \forall u, v \in V$$

# Previous techniques

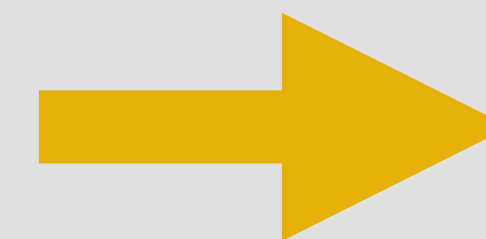Input: semi-metric $x$ on $V$

Let $r = 1/5$
While there is some unclustered vertex
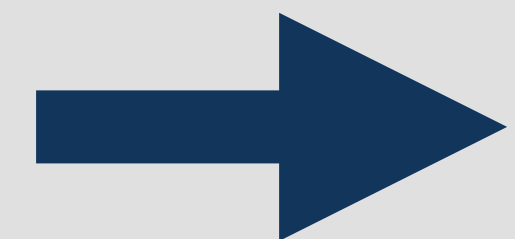    Find "densest" cluster with center $c^*$ and radius $r$
    Create cluster $C$ around $c^*$ with radius $2r$
Return clusters

Convex program for $\ell_p$ correlation clustering

$$\min ||y||_p$$

$$y(u) = \sum_{v \in N_u^+} x_{uv} + \sum_{v \in N_u^-} (1 - x_{uv}) \qquad \forall u \in V$$

$$x_{uv} \le x_{vw} + x_{uw} \qquad\qquad \forall u, v, w \in V$$

$$0 \le x_{uv} \le 1 \qquad\qquad \forall u, v \in V$$

Constraints induce a semi-metric space

**LP solution**

Rounding algorithm by Kalhan, Makarychev, Zhou

**Clustering**

21

# Correlation metric

**Correlation metric for $\ell_\infty$:**

(1) satisfies triangle inequality

(2) has $\displaystyle\sum_{v\in N_u^+} d_{uv} + \sum_{v\in N_u^-} (1 - d_{uv}) \leq O(1) \cdot \max_{w\in V} y(w)$

$$\sum_{u\in V} \left( \sum_{v\in N_u^+} d_{uv} + \sum_{v\in N_u^-} (1 - d_{uv}) \right)^p \leq O(1) \cdot \sum_{w\in V} y(w)^p$$

$$\min ||y||_p$$

$$y(u) = \sum_{v\in N_u^+} x_{uv} + \sum_{v\in N_u^-} (1 - x_{uv}) \qquad \forall u \in V$$

$$x_{uv} \leq x_{vw} + x_{uw} \qquad \forall u, v, w \in V$$

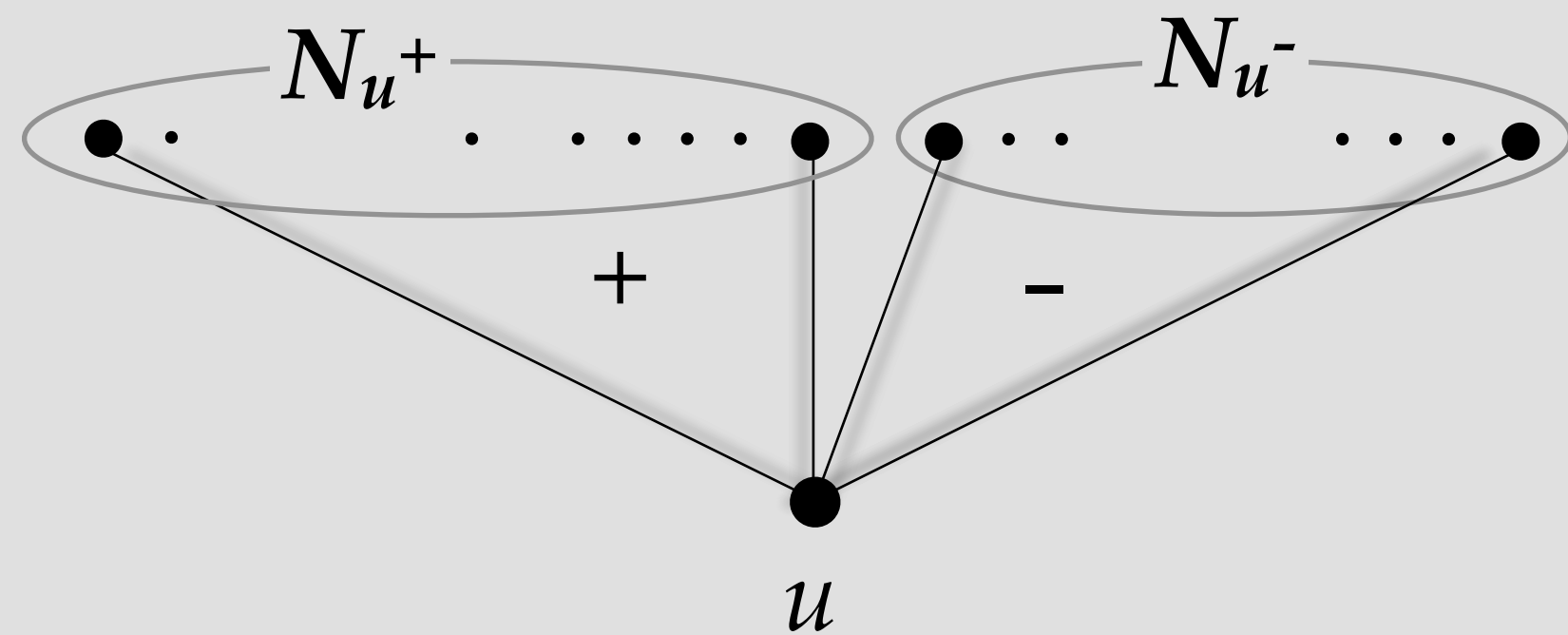$$0 \leq x_{uv} \leq 1 \qquad \forall u, v \in V$$

Input correlation metric $d_{uv}$, an apx for $x_{uv}$

**LP solution**

Rounding algorithm by Kalhan, Makarychev, Zhou

**Clustering**

# Correlation metric

▸ $N_u{}^+ = (+)$ neighbors of $u$, $N_u{}^- = (-)$ neighbors of $u$
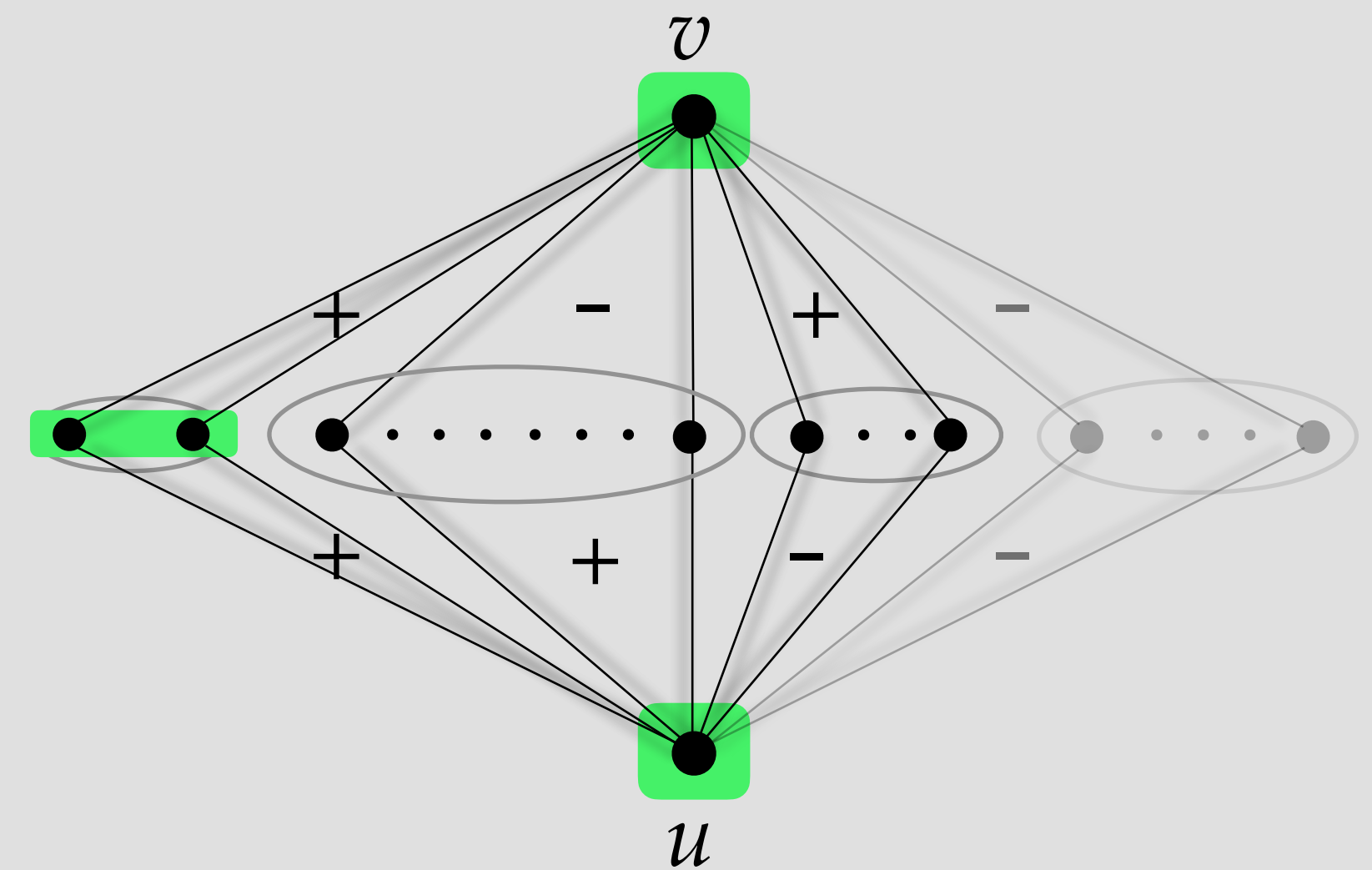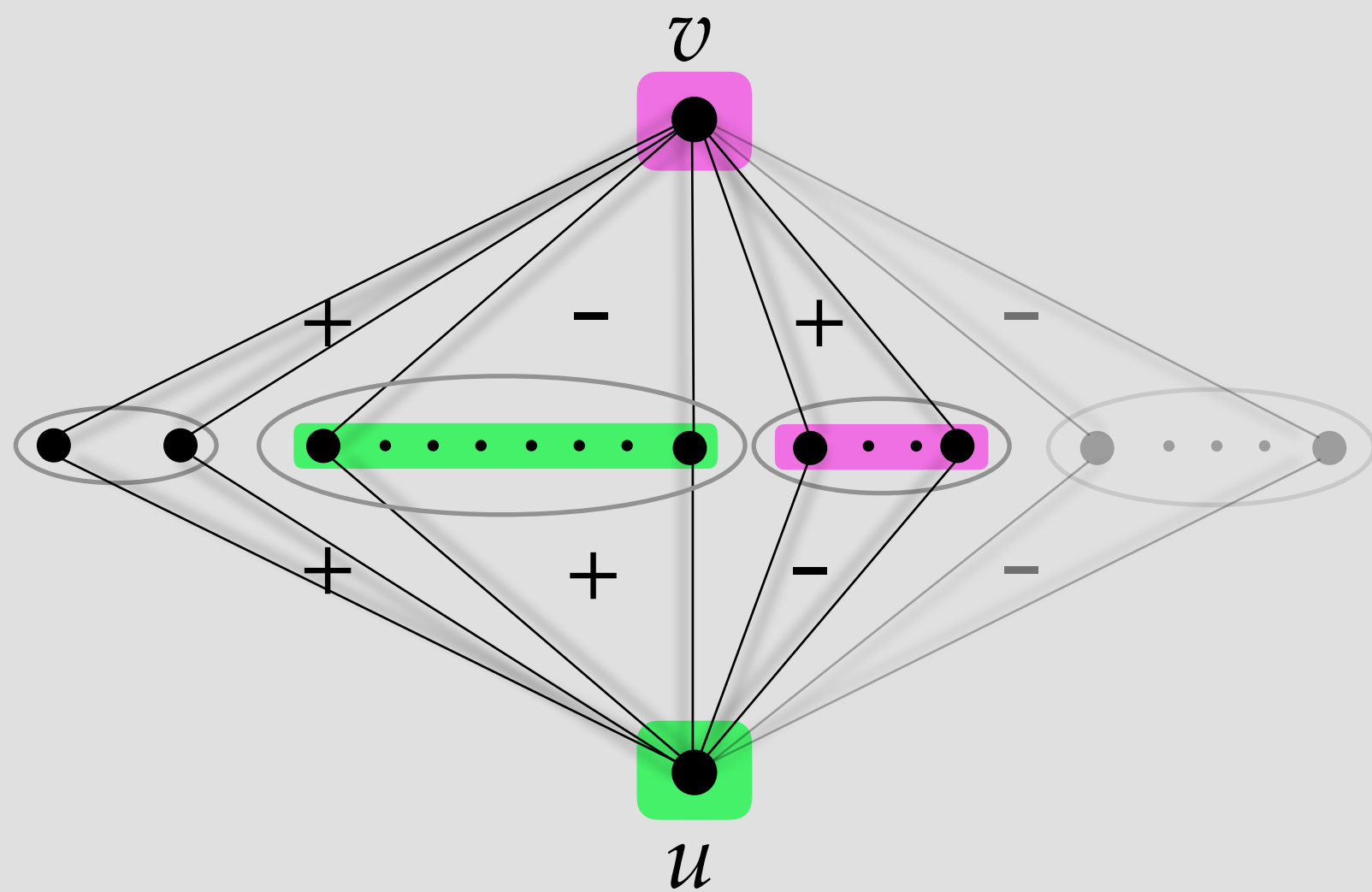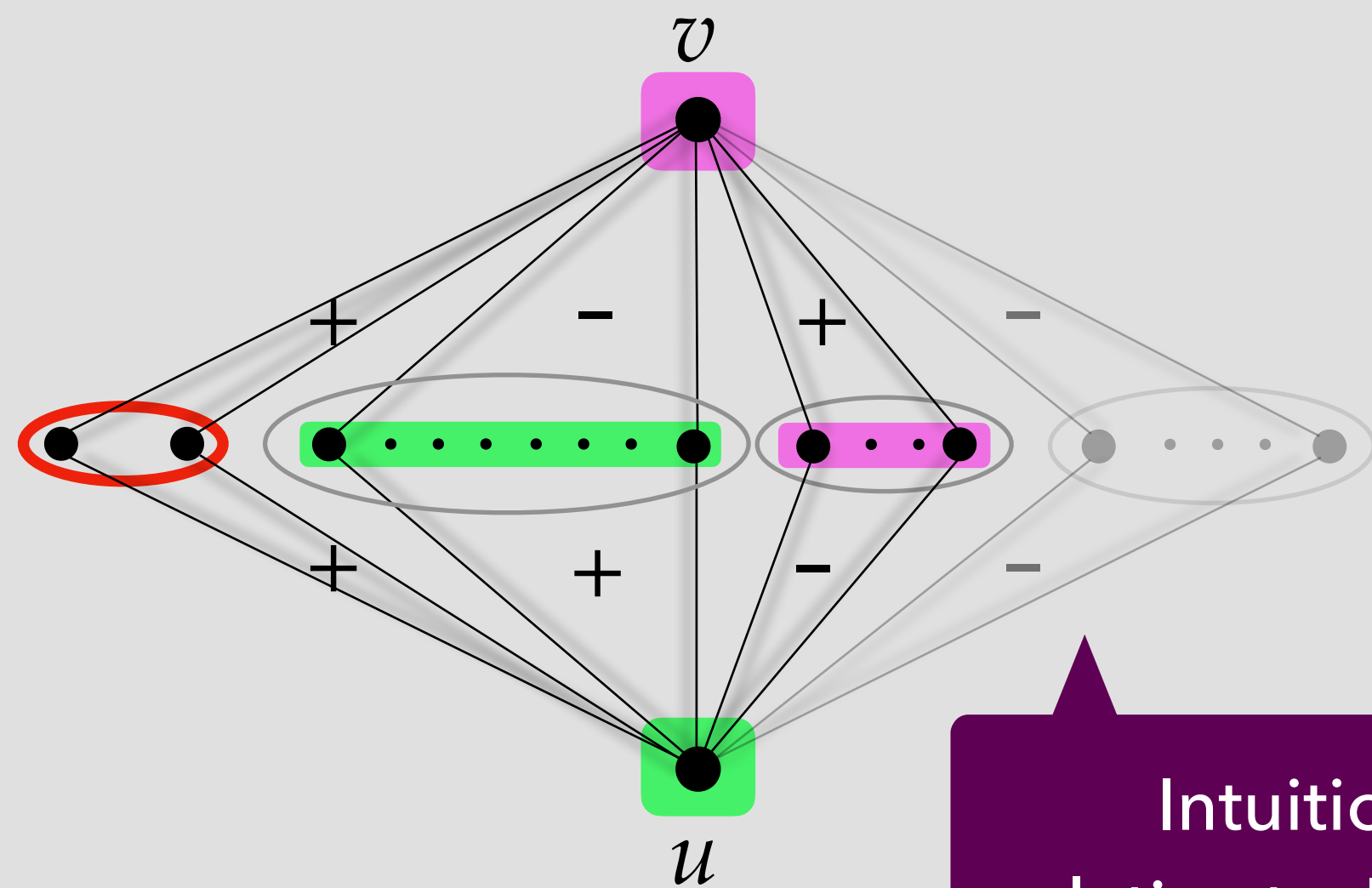
# Correlation metric

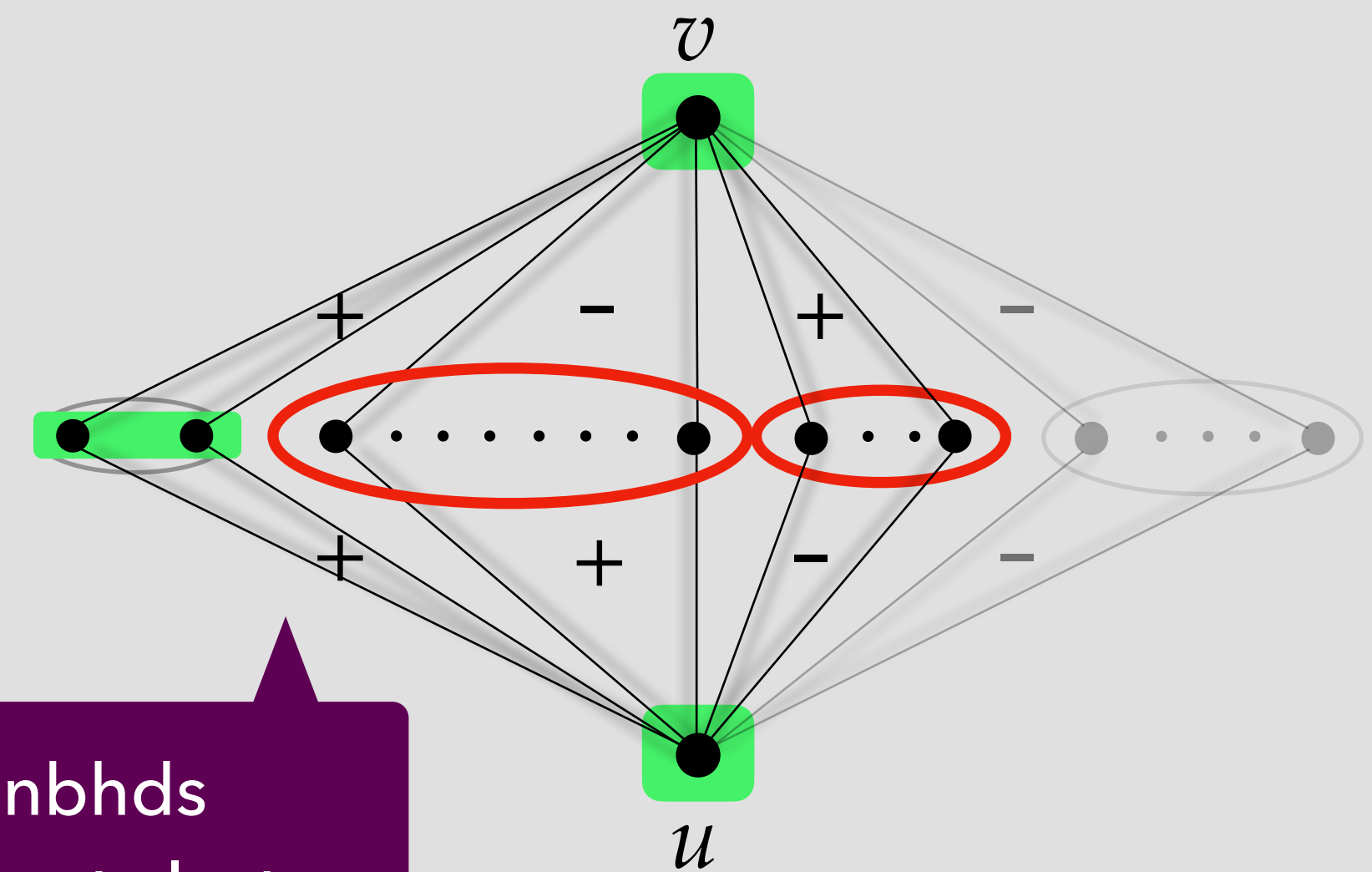‣ $N_u{}^+ = (+)$ neighbors of $u$, $N_u{}^- = (-)$ neighbors of $u$

# Correlation metric

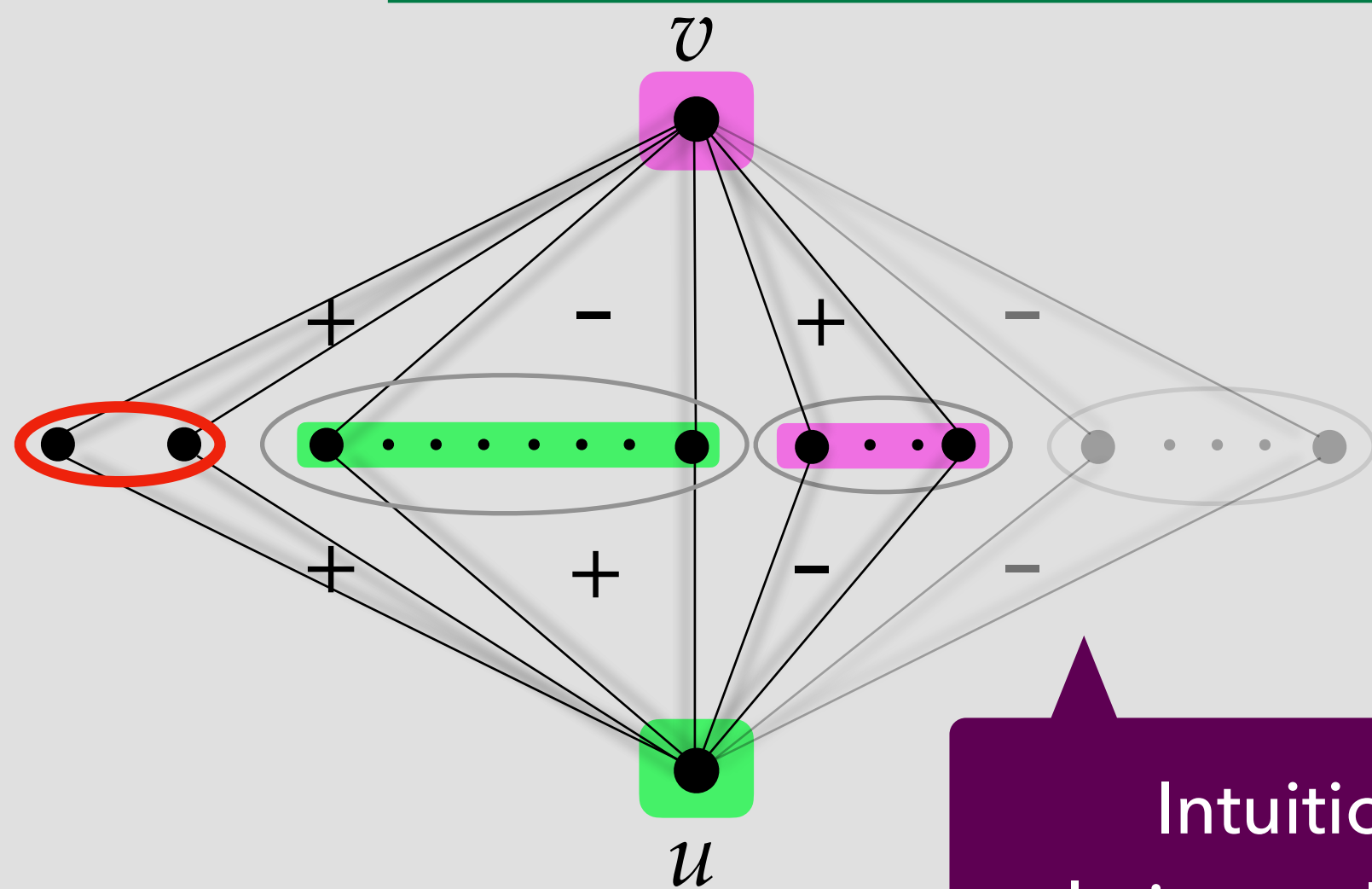- $N_u^+ = (+)$ neighbors of $u$, $N_u^- = (-)$ neighbors of $u$



Intuition: if $u$ and $v$ have large mixed nbhds relative to $|N_u^+ \cup N_v^+|$, want them in different clusters

# Correlation metric

Very coarse approximation for probability
Pivot separates $u,v$

Correlation metric $= d_{uv} = 1 - \dfrac{|N_u^+ \cap N_v^+|}{|N_u^+ \cup N_v^+|} = \dfrac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$

$v$

$+$  $-$  $+$  $-$

$+$  $+$  $-$  $-$

$u$

**Pivot algorithm**
- Randomly choose a *pivot (*unclustered vertex)
- New cluster with pivot + all its unclustered + neighbors

1

Pivot 1

2  3  4  5

Intuition: if $u$ and $v$ have large mixed nbhds
relative to $|N_{u^+} \cup N_{v^+}|$, want them in different clusters

# Correlation metric

Correlation metric fast to compute, time $O(n^\omega)$.

↳ Even faster on sparse graphs $O(n \cdot \Delta^2 \cdot \log n)$
↳ Further sped up on any graph with sampling procedure

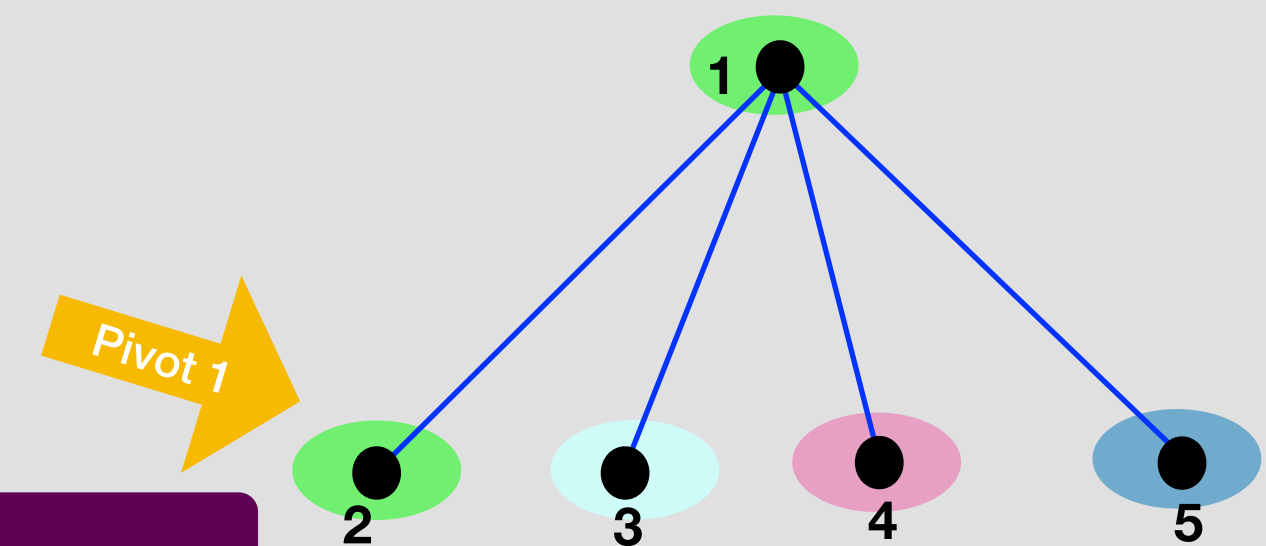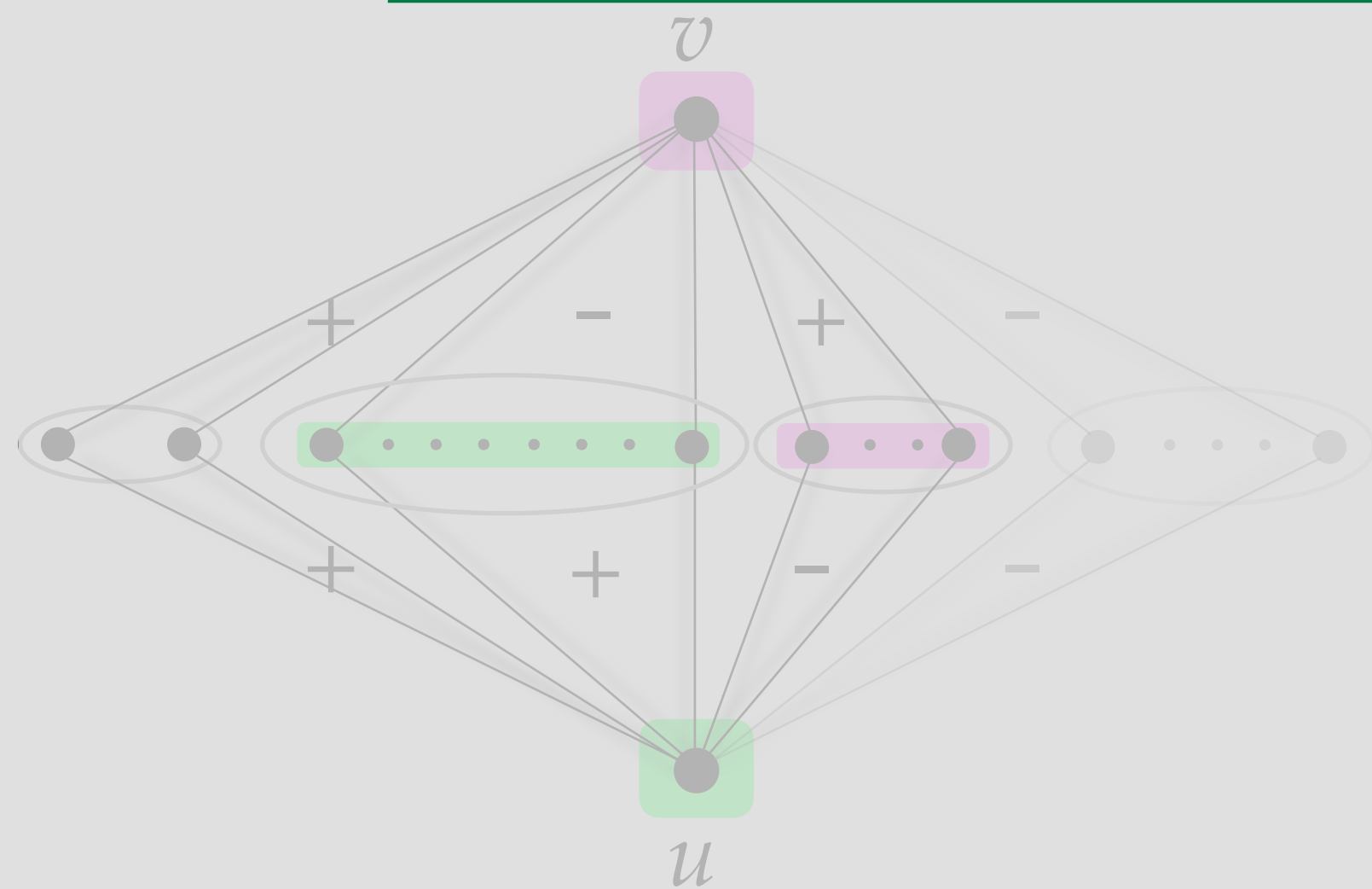Works as is for $\ell_\infty$ norm objective

$$\text{Correlation metric} = d_{uv} = 1 - \frac{|N_u^+ \cap N_v^+|}{|N_u^+ \cup N_v^+|} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

$v$

$u$

Correlation metric

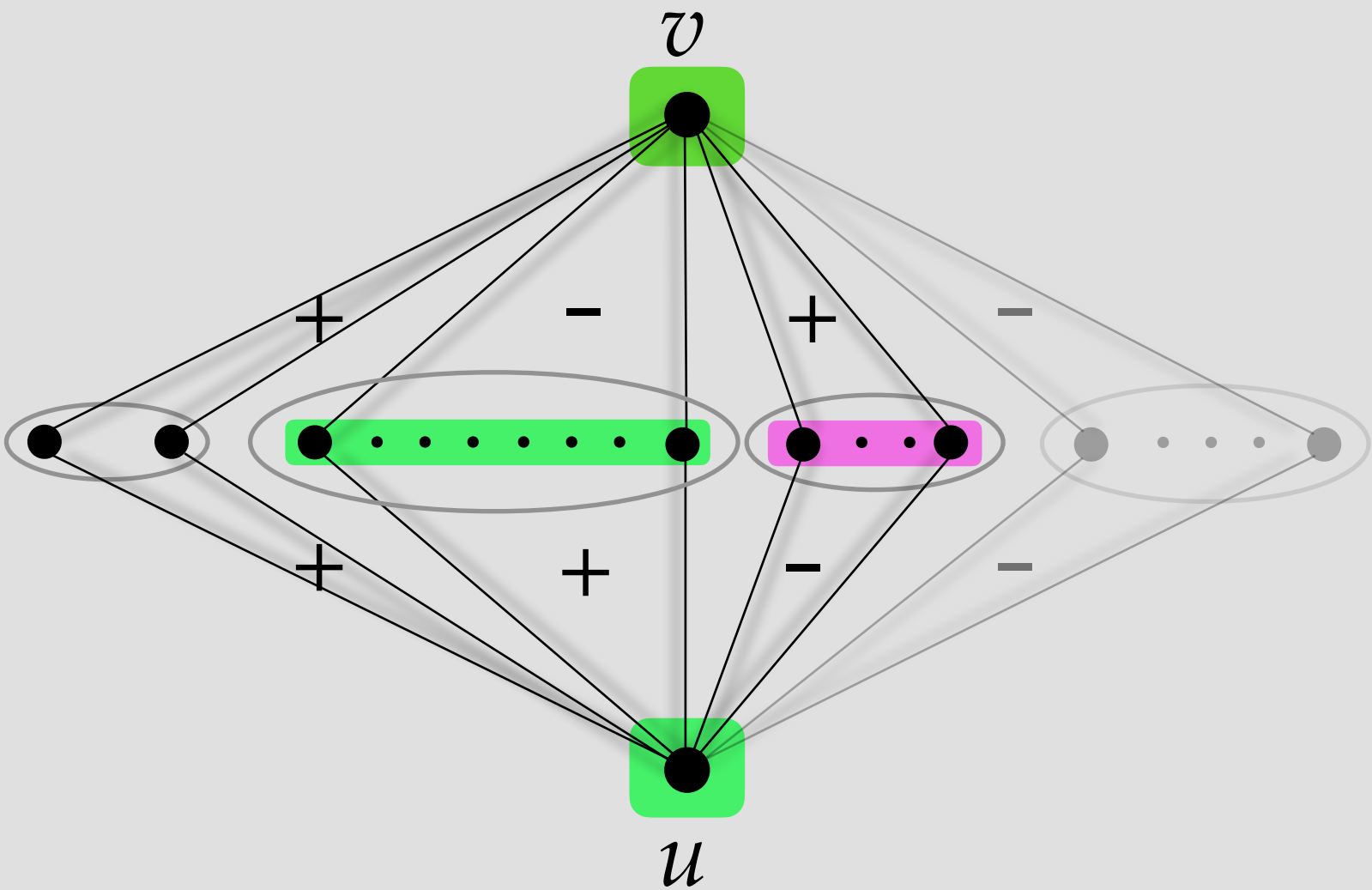Rounding algorithm by Kalhan, Makarychev, Zhou

Clustering

# Today

# Correlation metric for $\ell_\infty$

$$d_{uv} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

**Want to show for $\ell_\infty$:** $\displaystyle\sum_{v \in N_u^+} d_{uv} + \sum_{v \in N_u^-} (1 - d_{uv}) \leq O(1) \cdot \max_{w \in V} y(w)$

Easy to bound positive edges for $\ell_\infty$ objective!

$$\sum_{v \in N_u^+} d_{uv} \leq \sum_{v \in N_u^+ \cap C(u)} \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|} + \sum_{v \in N_u^+ \cap \overline{C(u)}} 1$$

$$d_{uv} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

**Want to show for $\ell_\infty$:** $\displaystyle\sum_{v \in N_u^+} d_{uv} + \sum_{v \in N_u^-} (1 - d_{uv}) \leq O(1) \cdot \max_{w \in V} y(w)$

Recall $y(u)$ = # disagreements incident to $u$

Easy to bound positive edges for $\ell_\infty$ objective!

$$\sum_{v \in N_u^+} d_{uv} \leq \sum_{v \in N_u^+ \cap C(u)} \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|} + \sum_{v \in N_u^+ \cap \overline{C(u)}} 1$$
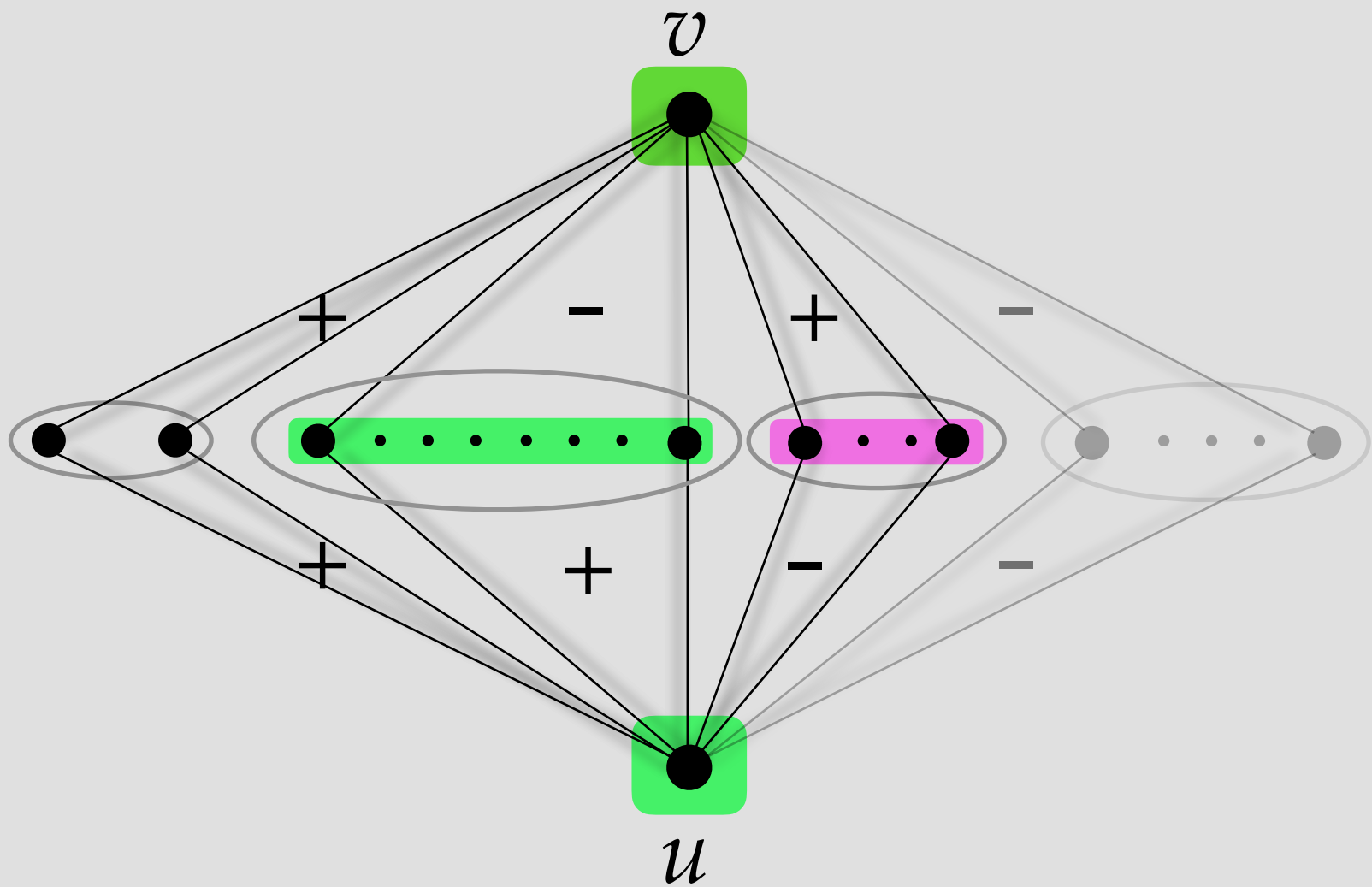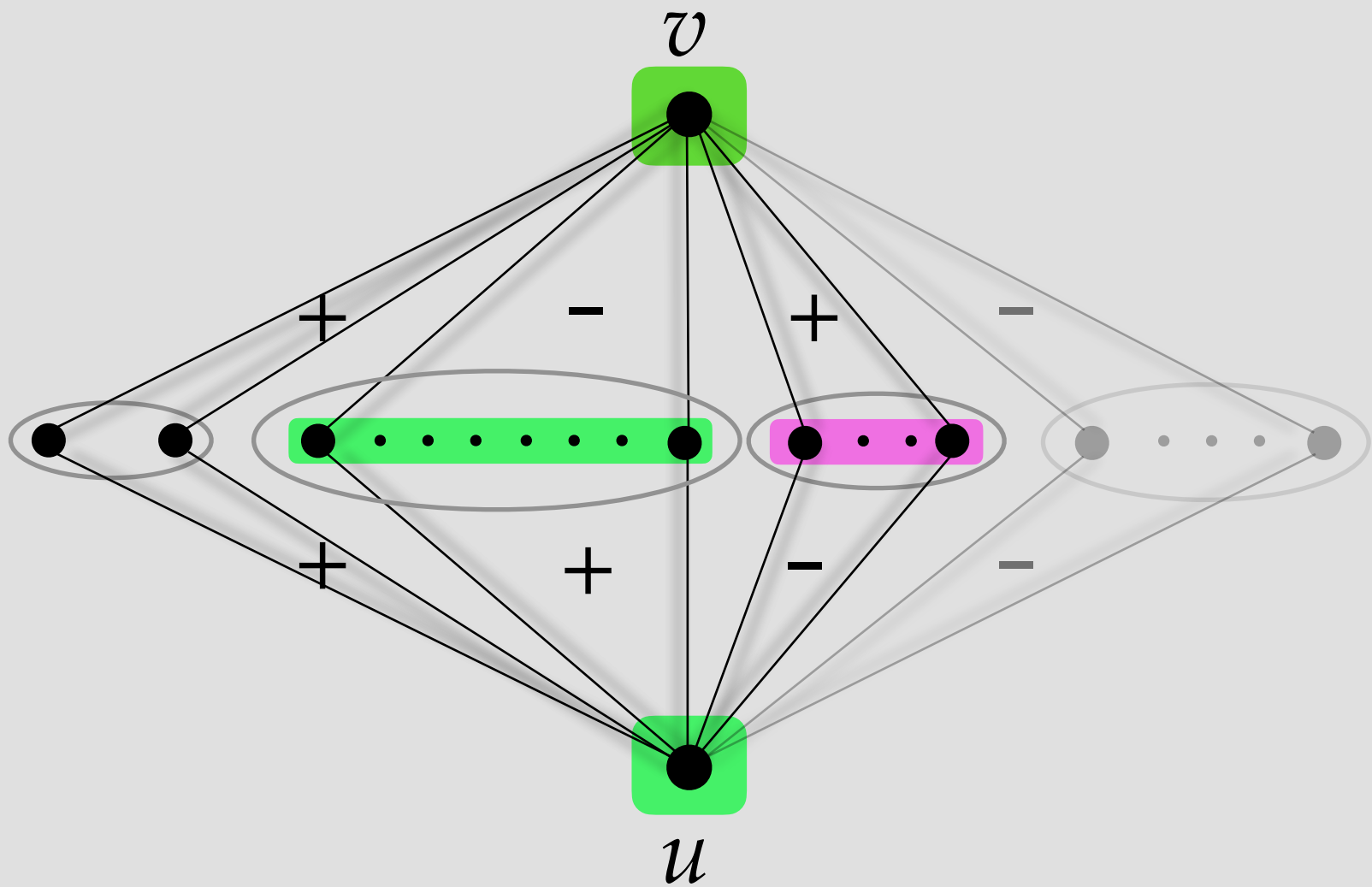
$$\leq \frac{1}{|N_u^+|} \sum_{v \in N_u^+ \cap C(u)} (|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|) + y(u)$$

# Correlation metric for $\ell_\infty$

$$d_{uv} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

**Want to show for $\ell_\infty$:** $\displaystyle\sum_{v \in N_u^+} d_{uv} + \sum_{v \in N_u^-} (1 - d_{uv}) \leq O(1) \cdot \max_{w \in V} y(w)$

Recall $y(u)$ = # disagreements incident to $u$

Easy to bound positive edges for $\ell_\infty$ objective!

$$\sum_{v \in N_u^+} d_{uv} \leq \sum_{v \in N_u^+ \cap C(u)} \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|} + \sum_{v \in N_u^+ \cap \overline{C(u)}} 1$$

$$\leq \frac{1}{|N_u^+|} \sum_{v \in N_u^+ \cap C(u)} (|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|) + y(u)$$

$$\leq \frac{1}{|N_u^+|} \sum_{v \in N_u^+ \cap C(u)} (y(u) + y(v)) + y(u)$$

$$\leq 2y(u) + \max_z y(z) \leq 3 \cdot \mathrm{OPT} \,.$$



31

# Correlation metric for $\ell_\infty$

$$d_{uv} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

**Want to show for $\ell_\infty$:** $\displaystyle\sum_{v \in N_u^+} d_{uv} + \sum_{v \in N_u^-} (1 - d_{uv}) \leq O(1) \cdot \max_{w \in V} y(w)$

Easy to bound positive edges for $\ell_\infty$ objective!

$$\sum_{v \in N_u^+} d_{uv} \leq 3 \cdot \text{OPT}.$$

# Correlation metric for $\ell_\infty$

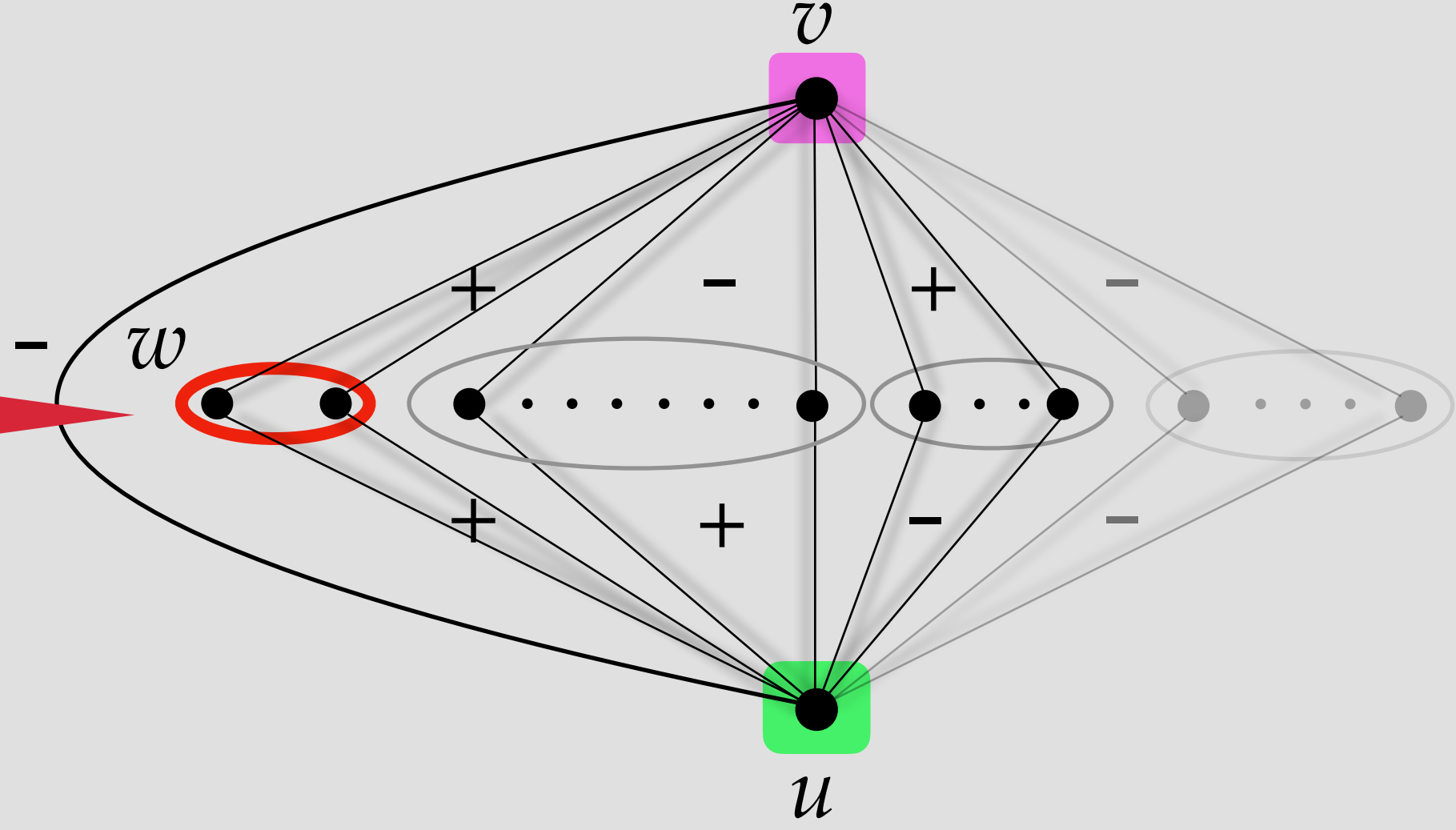$$d_{uv} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

**Want to show for $\ell_\infty$:** $\displaystyle\sum_{v \in N_u^+} d_{uv} + \sum_{v \in N_u^-} (1 - d_{uv}) \leq O(1) \cdot \max_{w \in V} y(w)$

## Bound on negative edges

$$\sum_{v \in N_u^-} (1 - d_{uv}) = \sum_{v \in N_u^- \cap C(u)} (1 - d_{uv}) + \sum_{v \in N_u^- \cap \overline{C(u)}} (1 - d_{uv})$$

$$= y(u) + \sum_{v \in N_u^- \cap \overline{C(u)}} \frac{|N_u^+ \cap N_v^+|}{n - |N_u^- \cap N_v^-|}$$

Easy to bound positive edges for $\ell_\infty$ objective!

$$\sum_{v \in N_u^+} d_{uv} \leq 3 \cdot \mathrm{OPT}.$$

Every $w$ in $|N_{u^+} \cap N_{v^+}|$ incident to an edge in disagreement, charge to carefully chosen $v^*(w)$ in $C(w)$

# Today

✦ ~~Introduction~~ ~~(the model, prior work, our results)~~ 🔥

✦ ~~The correlation metric~~ ~~(constructing a "guess" for the fraction solution, an inherent asymmetry)~~ 🔥🔥

✦ ~~Proof sketch for the $\ell_\infty$-norm~~ 🔥🔥🔥

✦ Adjusting the correlation metric (regular graphs are easy, dealing with negative edges ) 🔥🔥

✦ Conclusions (mainly vibes) 🔥

34

# Adjusted correlation metric

$$d_{uv} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

→ Simultaneous approximation for $\ell_1$- and $\ell_\infty$-norm objectives

For ***regular graphs***, correlation metric $O(1)$–apxs $\ell_1$-norm

✦ Proof via dual fitting!

✦ Problem is when graph is far from regular

Must ***adjust*** correlation metric for non-regular graphs for general $\ell_p$-norms

$d_{uv} = 2/3$ for all $u,v$ in $\{2,\dots,n\}$, so fractional cost w.r.t $d$ is $\theta(n^2)$

# Adjusted correlation metric

$$\text{Correlation metric} = d_{uv} = 1 - \frac{|N_u^+ \cap N_v^+|}{|N_u^+ \cup N_v^+|} = \frac{|N_u^+ \cap N_v^-| + |N_u^- \cap N_v^+|}{|N_u^+ \cup N_v^+|}$$

For *all* $\ell_p$ norms

- If negative edge *(u,v)* has $d_{uv}$ *>0.7,* update $d_{uv} \leftarrow 1$

- For *u* with $|N_u^- \cap \{v : d_{uv} \leq 0.7\}| \geq \frac{10}{3}\Delta_u$, update $d_{uv} \leftarrow 1$

**Adjusted correlation metric**

Rounding algorithm by Kalhan, Makarychev, Zhou

**Clustering**

36

# Today

# Summary

$\ell_p$-norm correlation clustering algs solve a convex program

Solving *metric constrained* LPs on large networks is slow!

Not very amenable to online settings

Solution specific to *one fixed* $\ell_p$-norm

# Summary

$\ell_p$-norm correlation clustering algs solve a convex program

Solving *metric constrained* LPs on large networks is slow!

Not very amenable to online settings

Solution specific to *one fixed* $\ell_p$-norm

Combinatorial techniques can resolve these issues

# Summary

Result 1: *O(1)*-apx alg with run-time *O(min{ n·Δ²·log n , nω})*. Near-linear for sparse graphs.

Result 2: ∃ an alg producing a clustering that is *O(1)*-apx for all $\ell_p$-norms, simultaneously.

Result 3: (In progress, probably true) Given a random ε–fraction of the network, ∃ a *semi-online algorithm* that for any $\ell_p$-norm objective produces a *O(log n)*-competitive algorithm.

Correlation clustering has interesting combinatorial structure that can be exploited

# What's next?

▸ 🧑‍💻**In progress:** Extend to a semi-online setting

↪ Factor depends on $p$. For $p=\infty$, the algorithm is $\theta(\log n)$-competitive

▸ 🔥**Hot conjecture:** Exists a combinatorial alg simultaneously 4-approximating all $\ell_p$-norms running in $O(n^\omega)$ time

▸ 😶‍🌫️ **Broader Qs:**

1. Combinatorial algorithms by designing "approximate LP solution"

2. Further study on the all-norms objective

# Thank you!

samidavies@berkeley.edu