

Robust Online Correlation Clustering

Rudy Zhou

Carnegie Mellon University

Silvio Lattanzi, Benjamin Moseley, Sergei Vassilvitskii, Yuyan Wang, Rudy Zhou

Robust Online Correlation Clustering


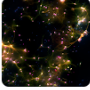
Neural Information Processing Systems (NeurIPS) 2021.

Motivation: Google News

Brain Cells in a Dish Learn How to Play Pong - IGN
IGN · 9 hours ago

- Lab-grown brain cells play Pong computer game
Al Jazeera English · 17 hours ago



[View Full Coverage](#)





Signs of Water on Mars Might Actually Be an Indication of Something Else
SciTechDaily · 8 hours ago

- First Martian life likely broke the planet with climate change, made themselves extinct
Livescience.com · 16 hours ago

[View Full Coverage](#)





Weather forces delay for NASA astronauts returning from space station on SpaceX capsule
CNN · 10 hours ago




NASA DART Mission Successfully Smashes Asteroid Into New Path
The New York Times · 2 days ago

- Op-Ed: Good news for a change — NASA proves there's a defense against killer asteroids
Los Angeles Times · 19 hours ago · Opinion

[View Full Coverage](#)



Space 'fingerprint' created by stars, NASA James Webb Telescope finds



⋮
?

Motivation: Google News

Brain Cells in a Dish Learn How to Play Pong - IGN
IGN · 9 hours ago

- Lab-grown brain cells play Pong computer game
Al Jazeera English · 17 hours ago

[View Full Coverage](#)

Signs of Water on Mars Might Actually Be an Indication of Something Else
SciTechDaily · 8 hours ago

- First Martian life likely broke the planet with climate change, made themselves extinct
Livescience.com · 16 hours ago

[View Full Coverage](#)

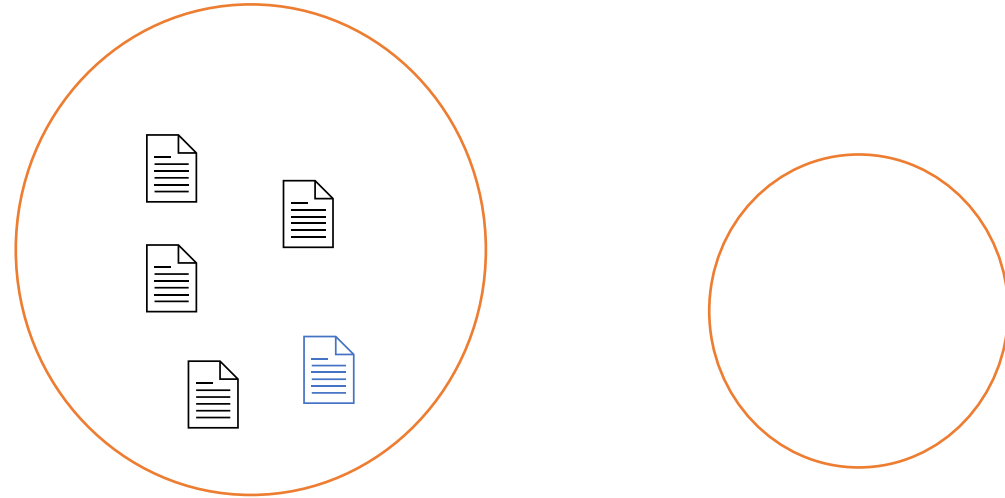
Weather forces delay for NASA astronauts returning from space station on SpaceX capsule
CNN · 10 hours ago

NASA DART Mission Successfully Smashes Asteroid Into New Path
The New York Times · 2 days ago

- Op-Ed: Good news for a change — NASA proves there's a defense against killer asteroids
Los Angeles Times · 19 hours ago · Opinion

[View Full Coverage](#)

Space 'fingerprint' created by stars, NASA James Webb Telescope finds



⋮
?

Motivation: Google News

Brain Cells in a Dish Learn How to Play Pong - IGN
IGN · 9 hours ago

- Lab-grown brain cells play Pong computer game
Al Jazeera English · 17 hours ago

[View Full Coverage](#)

Signs of Water on Mars Might Actually Be an Indication of Something Else
SciTechDaily · 8 hours ago

- First Martian life likely broke the planet with climate change, made themselves extinct
Livescience.com · 16 hours ago

[View Full Coverage](#)

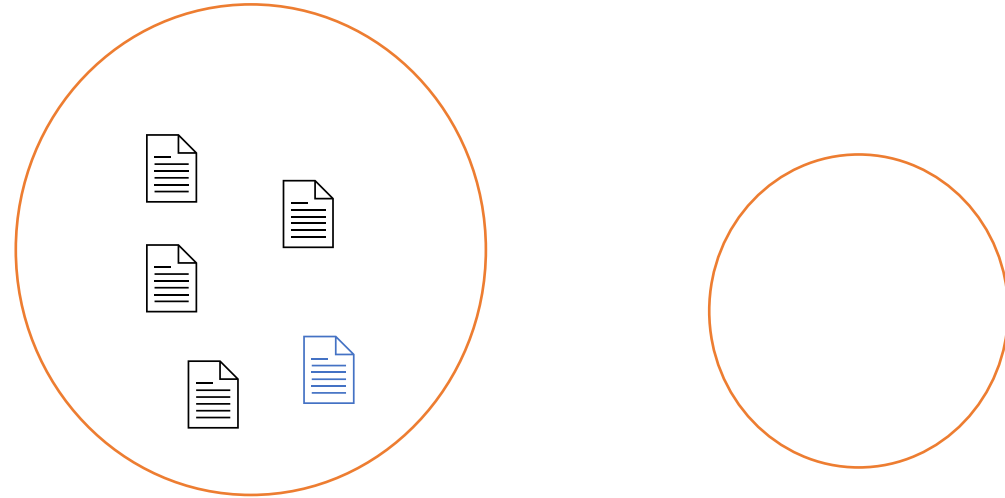
Weather forces delay for NASA astronauts returning from space station on SpaceX capsule
CNN · 10 hours ago

NASA DART Mission Successfully Smashes Asteroid Into New Path
The New York Times · 2 days ago

- Op-Ed: Good news for a change — NASA proves there's a defense against killer asteroids
Los Angeles Times · 19 hours ago · Opinion

[View Full Coverage](#)

Space 'fingerprint' created by stars, NASA James Webb Telescope finds



Motivation: Google News

Brain Cells in a Dish Learn How to Play Pong - IGN
IGN · 9 hours ago

- **Lab-grown brain cells play Pong computer game**
Al Jazeera English · 17 hours ago

[View Full Coverage](#)

Signs of Water on Mars Might Actually Be an Indication of Something Else
SciTechDaily · 8 hours ago

- **First Martian life likely broke the planet with climate change, made themselves extinct**
Livescience.com · 16 hours ago

[View Full Coverage](#)

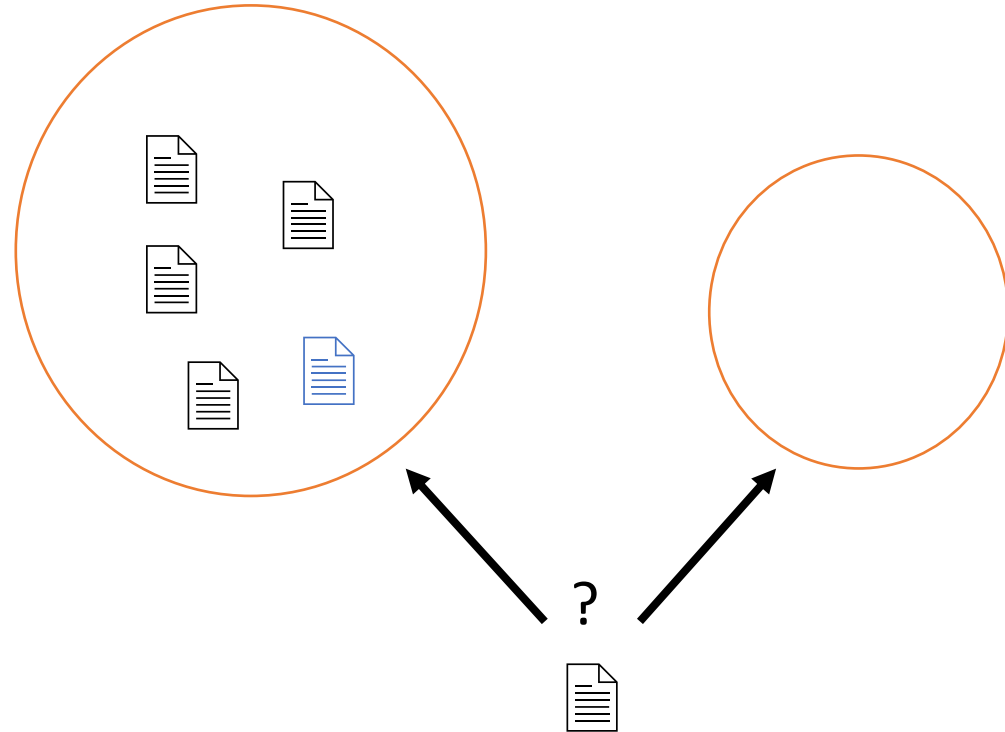
Weather forces delay for NASA astronauts returning from space station on SpaceX capsule
CNN · 10 hours ago

NASA DART Mission Successfully Smashes Asteroid Into New Path
The New York Times · 2 days ago

- **Op-Ed: Good news for a change — NASA proves there's a defense against killer asteroids**
Los Angeles Times · 19 hours ago · Opinion

[View Full Coverage](#)

Space 'fingerprint' created by stars, NASA James Webb Telescope finds



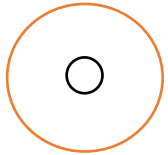
Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



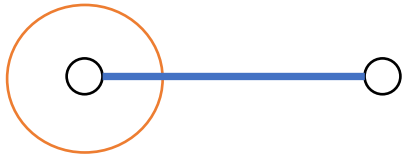
Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



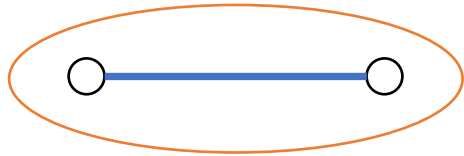
Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



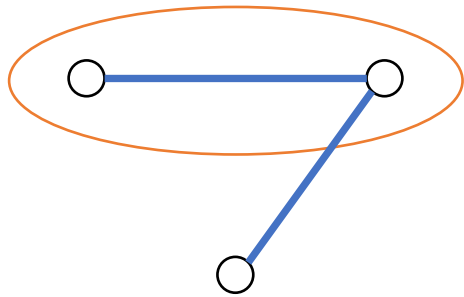
Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



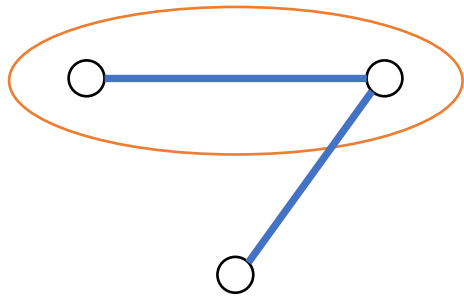
Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Online Correlation Clustering

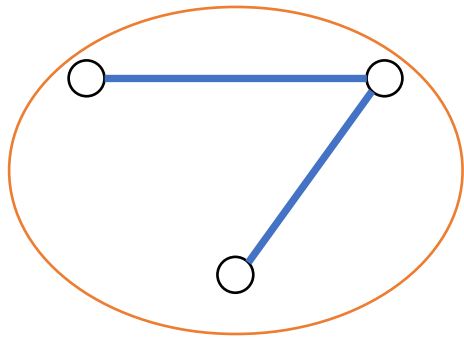
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

Online Correlation Clustering

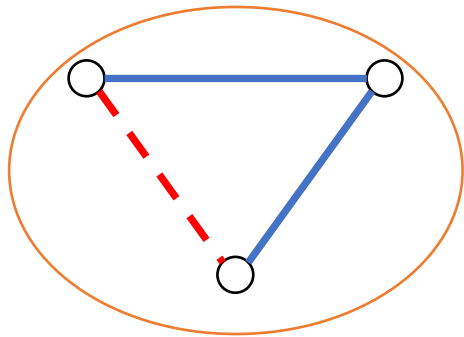
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize #(*disagreements*) = #(edges across) + #(non – edges within)

Online Correlation Clustering

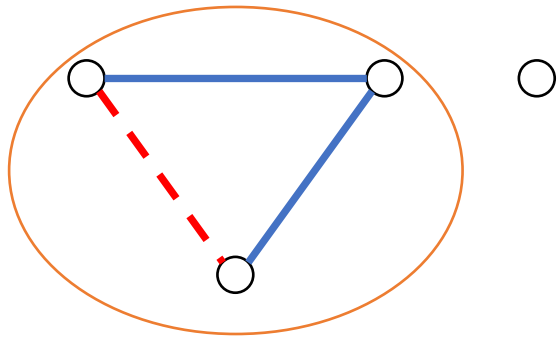
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

Online Correlation Clustering

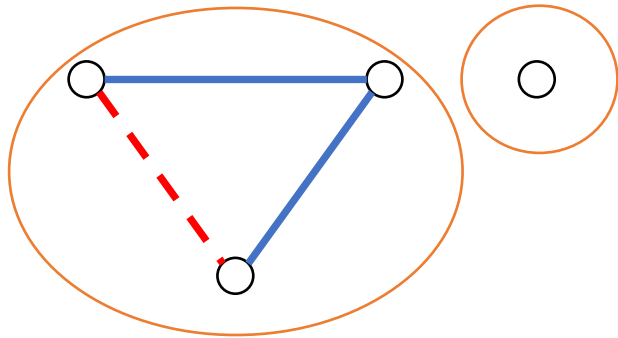
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

Online Correlation Clustering

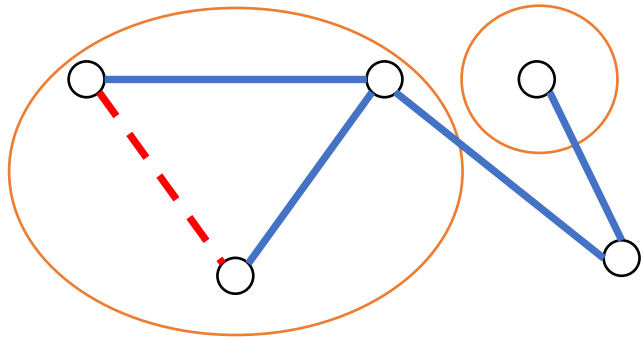
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

Online Correlation Clustering

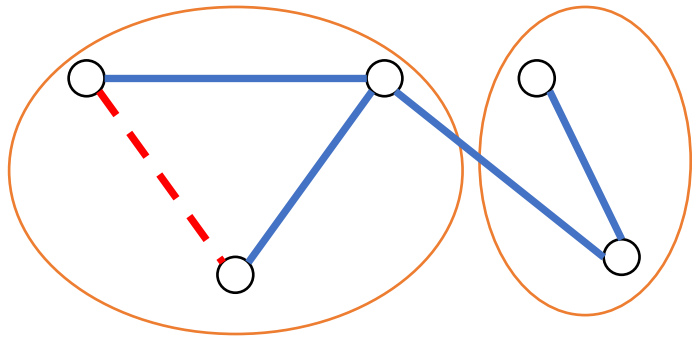
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

Online Correlation Clustering

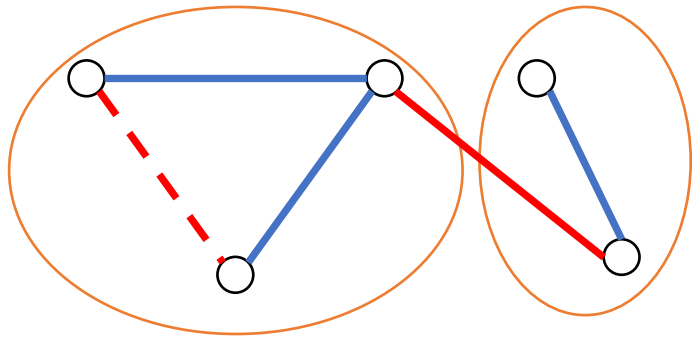
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

Online Correlation Clustering

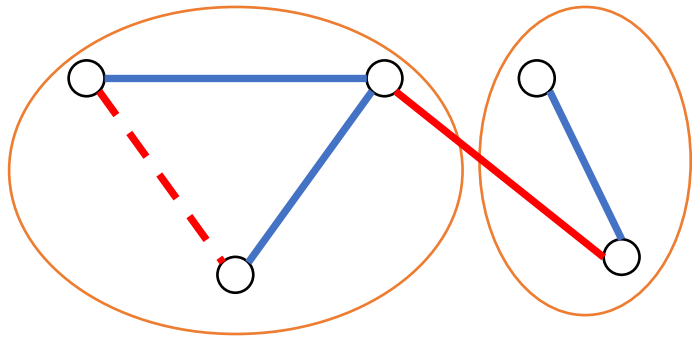
- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



Minimize #(*disagreements*) = #(edges across) + #(non – edges within)

Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster



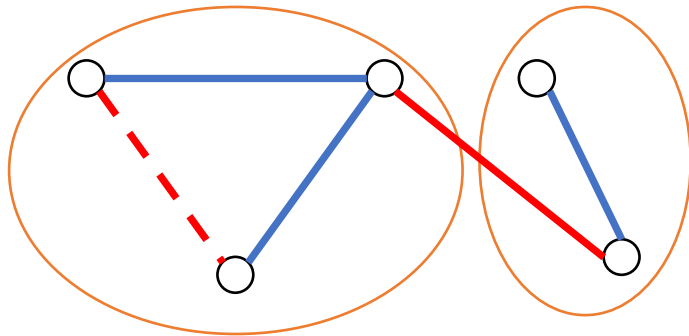
Minimize $\#(\textit{disagreements}) = \#(\textit{edges across}) + \#(\textit{non - edges within})$

...compared to optimal offline clustering that knows entire graph

Online Correlation Clustering

- Vertices arrive online; reveal edges to previous arrivals
- Assign vertex to existing cluster or make new singleton cluster

Competitive Ratio: An algorithm is ***c*-competitive** if for any input graph and arrival order:
 $\#(\text{disagreements by ALG}) \leq c \cdot \#(\text{disagreements by OPT})$

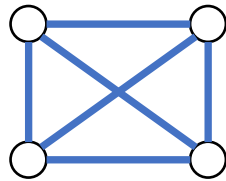


Minimize $\#(\text{disagreements}) = \#(\text{edges across}) + \#(\text{non-edges within})$

...compared to optimal offline clustering that knows entire graph

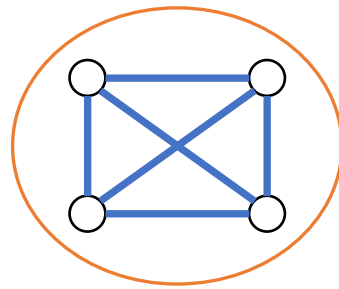
Prior Work

- Every online algorithm is $\Omega(\# \text{ vertices})$ -competitive



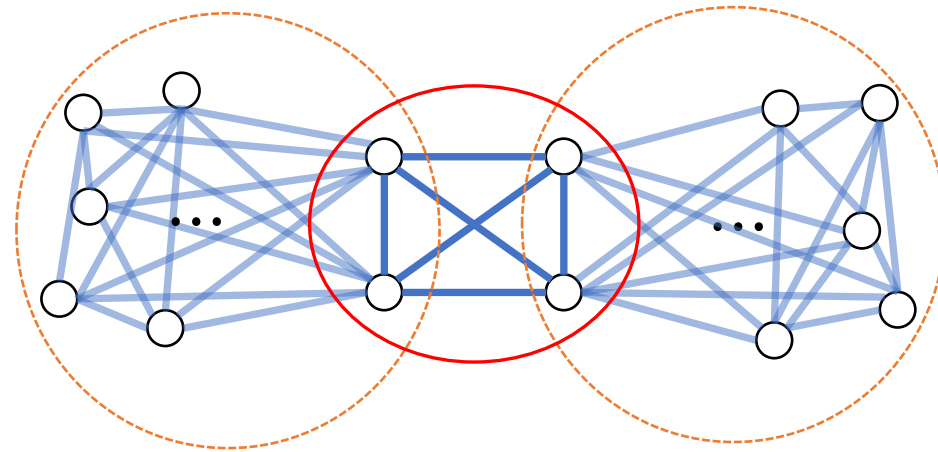
Prior Work

- Every online algorithm is $\Omega(\# \text{ vertices})$ -competitive



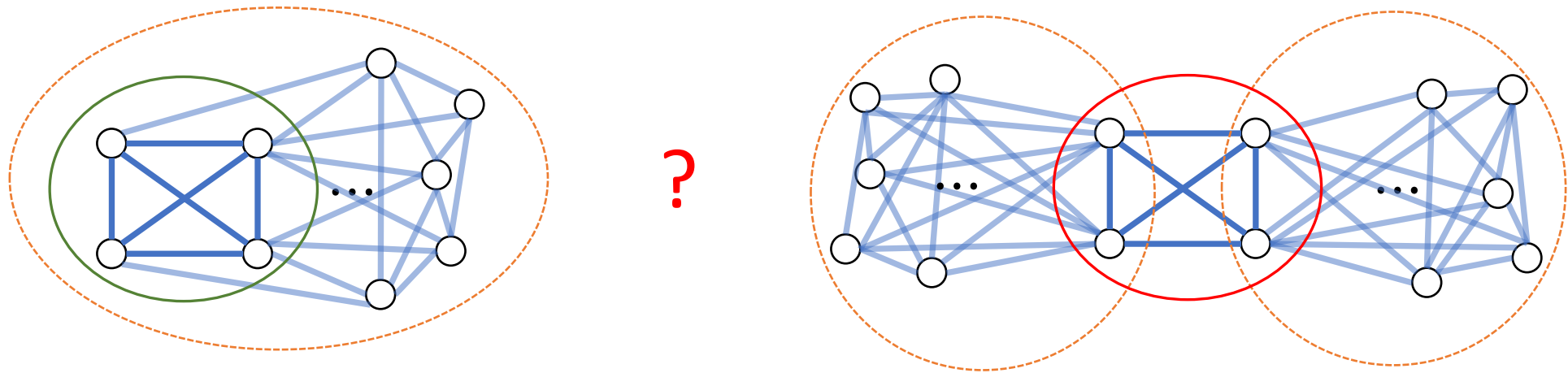
Prior Work

- Every online algorithm is $\Omega(\# \text{ vertices})$ -competitive



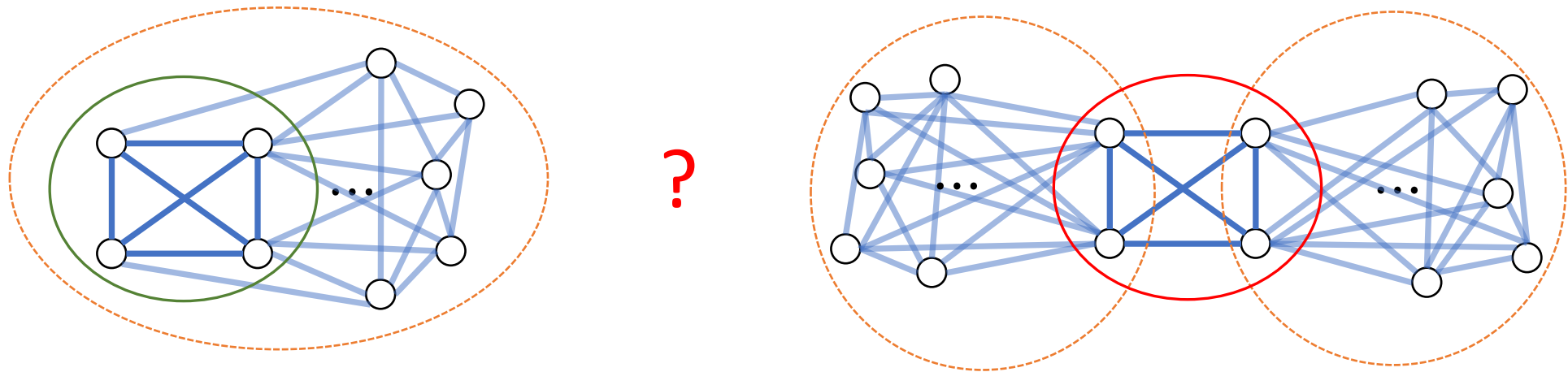
Prior Work

- Every online algorithm is $\Omega(\# \text{ vertices})$ -competitive



Prior Work

- Every online algorithm is $\Omega(\# \text{ vertices})$ -competitive



How to overcome online lower bound?

How to overcome online lower bound?

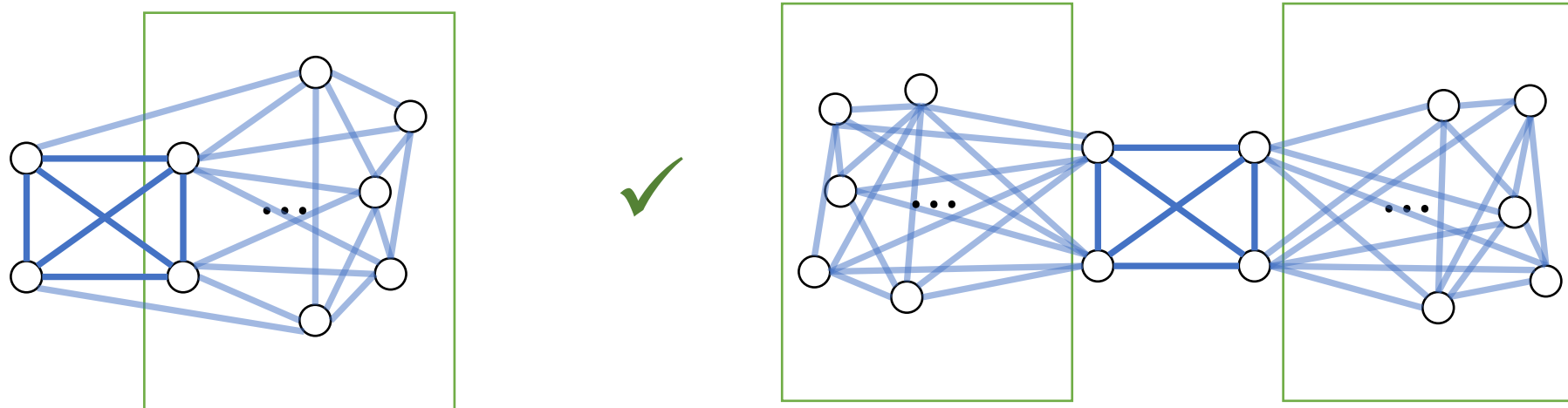
- Augment algorithm with **historical data**
- Introduce **new online model**
 - Historical data should be (partially) related to online arrivals
 - No other assumptions on online arrival order

Semi-Online Model

- Two phases: Offline and Online
- **Offline Phase:** corrupted random subgraph of ϵ -fraction of vertices revealed offline
 - Adversary chooses α -fraction of vertices
 - $(\epsilon - \alpha)$ -fraction of remaining vertices are randomly chosen
- **Online Phase:** remaining vertices arrive online

Semi-Online Model

- Two phases: Offline and Online
- **Offline Phase:** corrupted random subgraph of ϵ -fraction of vertices revealed offline
 - Adversary chooses α -fraction of vertices
 - $(\epsilon - \alpha)$ -fraction of remaining vertices are randomly chosen
- **Online Phase:** remaining vertices arrive online



Our Contribution

- Introduce **semi-online model** for sequential decision-making problems

Main Theorem: We design an algorithm for semi-online correlation clustering that is $O\left(\frac{1}{\epsilon - \alpha}\right)$ – competitive*.

* assuming $\alpha \leq \frac{\epsilon}{2}$

Our Contribution

- Introduce **semi-online model** for sequential decision-making problems

Main Theorem: We design an algorithm for semi-online correlation clustering that is $O\left(\frac{1}{\epsilon - \alpha}\right)$ – competitive*.

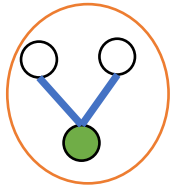
- ... and any semi-online algorithm must be $\Omega\left(\frac{1}{\epsilon - \alpha}\right)$ – competitive*
- ... and the theoretical results are predictive of practice

* assuming $\alpha \leq \frac{\epsilon}{2}$

Algorithm

- **Pivot Algorithm**

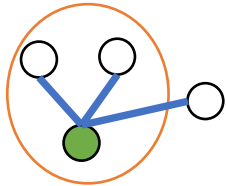
- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



Algorithm

- **Pivot Algorithm**

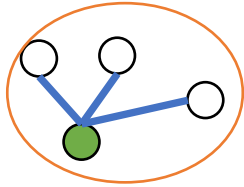
- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



Algorithm

- **Pivot Algorithm**

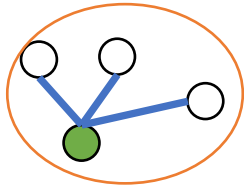
- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



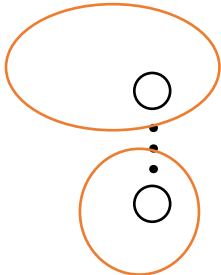
Algorithm

- **Pivot Algorithm**

- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



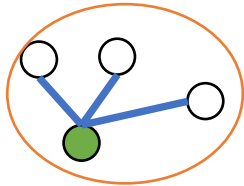
- Else make v a Pivot, and v starts its own cluster



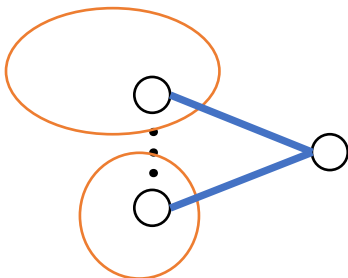
Algorithm

- **Pivot Algorithm**

- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



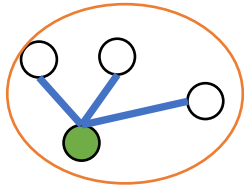
- Else make v a Pivot, and v starts its own cluster



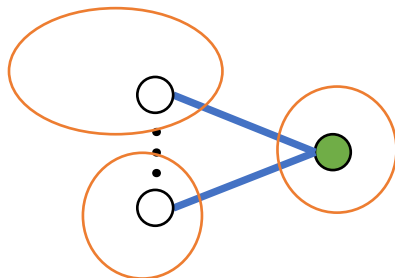
Algorithm

- **Pivot Algorithm**

- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



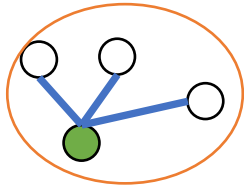
- Else make v a Pivot, and v starts its own cluster



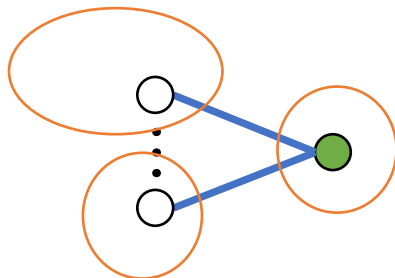
Algorithm

- **Pivot Algorithm**

- Maintain collection of vertices called **Pivots**
- Consider vertices in some order
- If v has an edge to a previous Pivot, then v joins the first such Pivot's cluster



- Else make v a Pivot, and v starts its own cluster

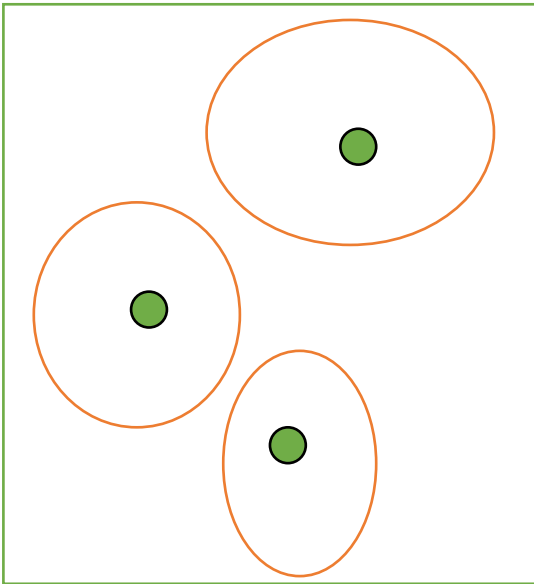


Our Algorithm: Run Pivot in **random order** in offline phase; then continue in **arrival order** in online phase

Analysis Overview

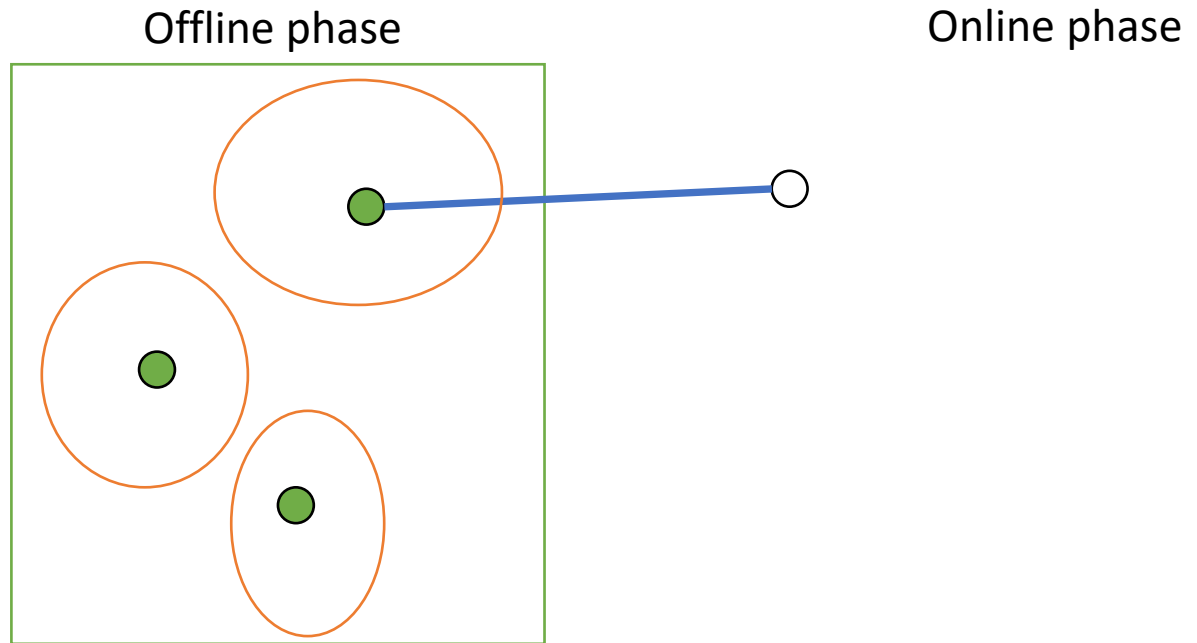
- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals

Offline phase



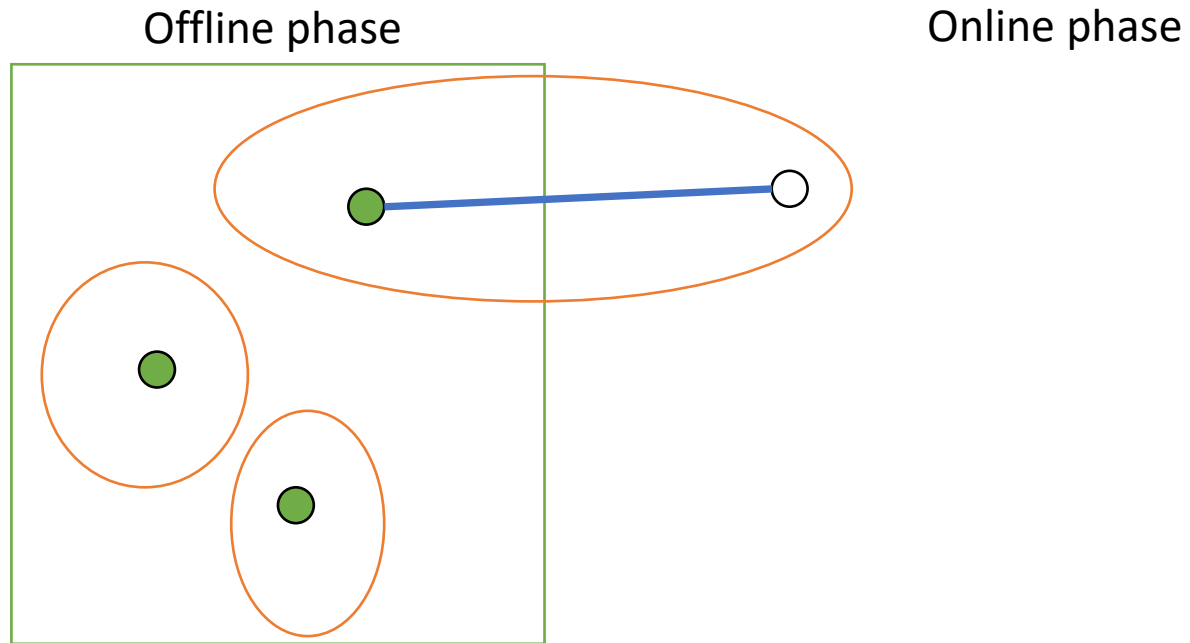
Analysis Overview

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals



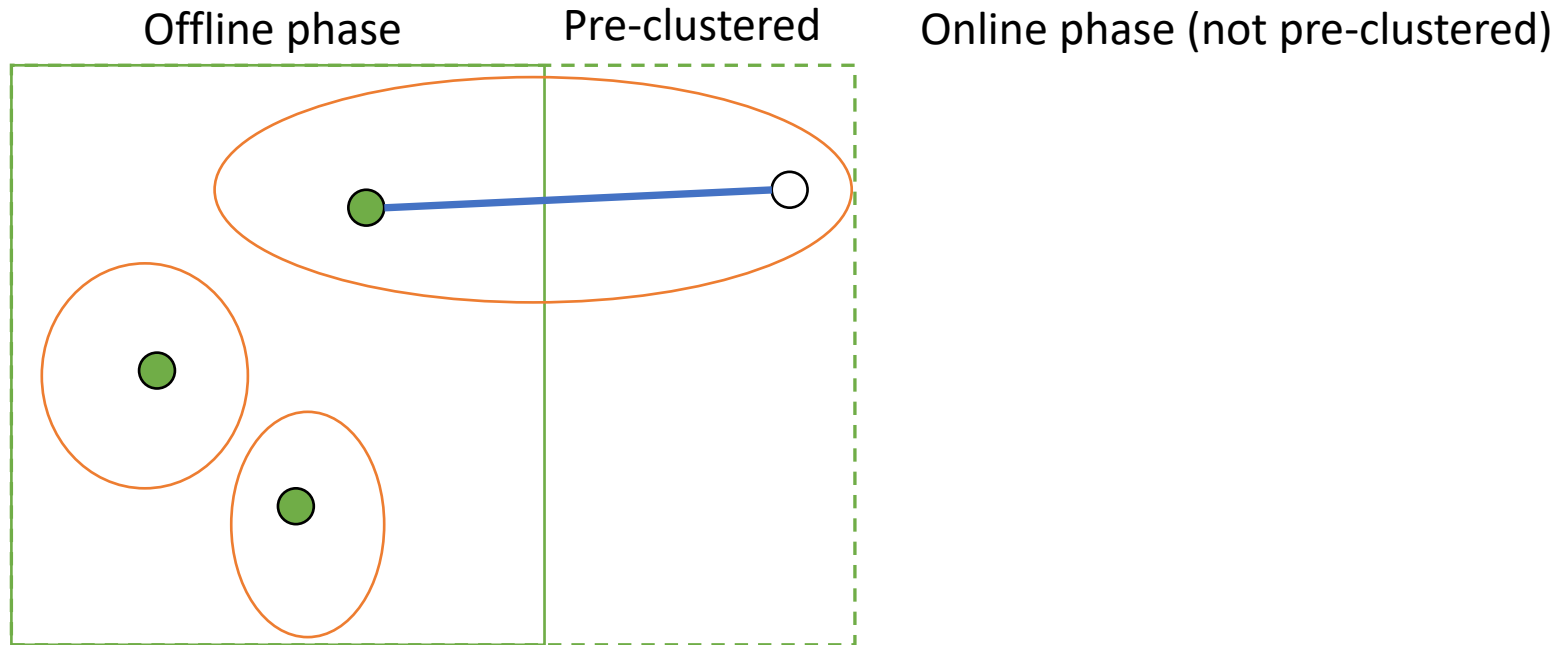
Analysis Overview

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals



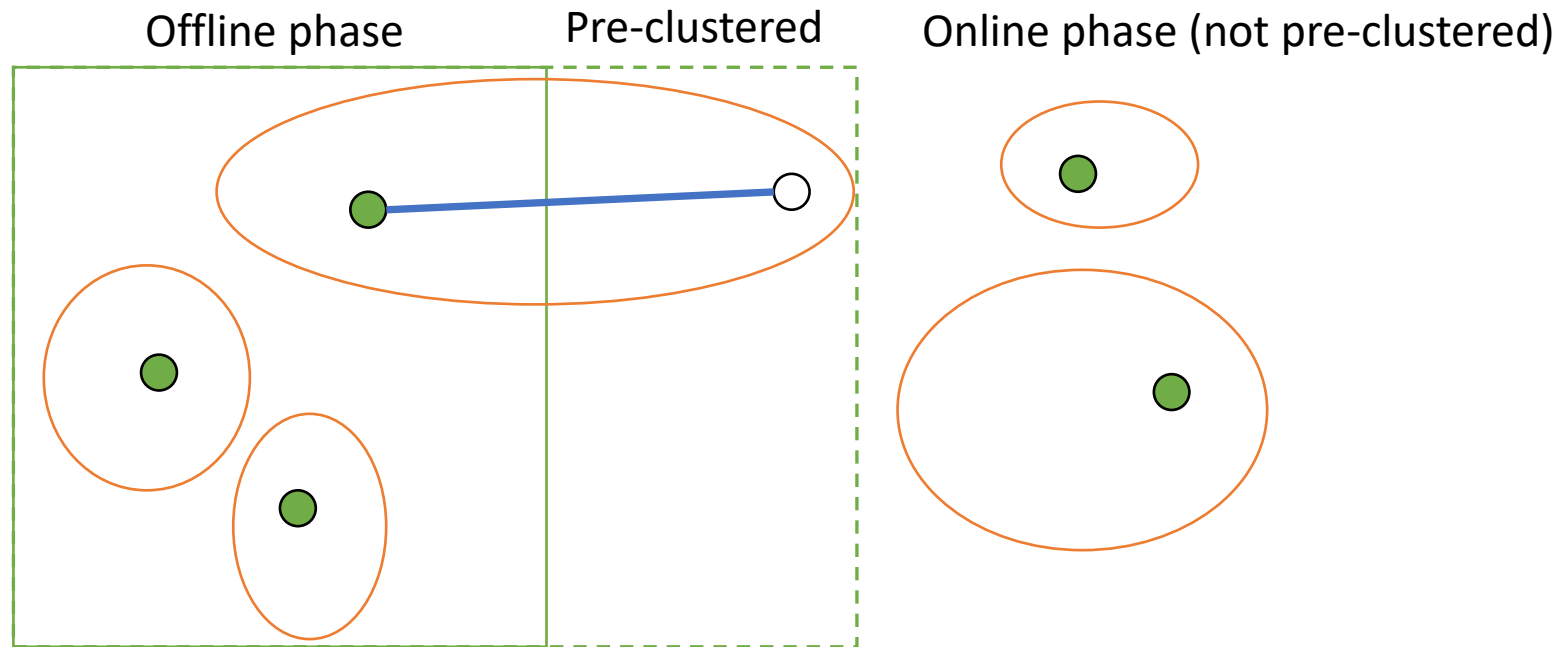
Analysis Overview

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals



Analysis Overview

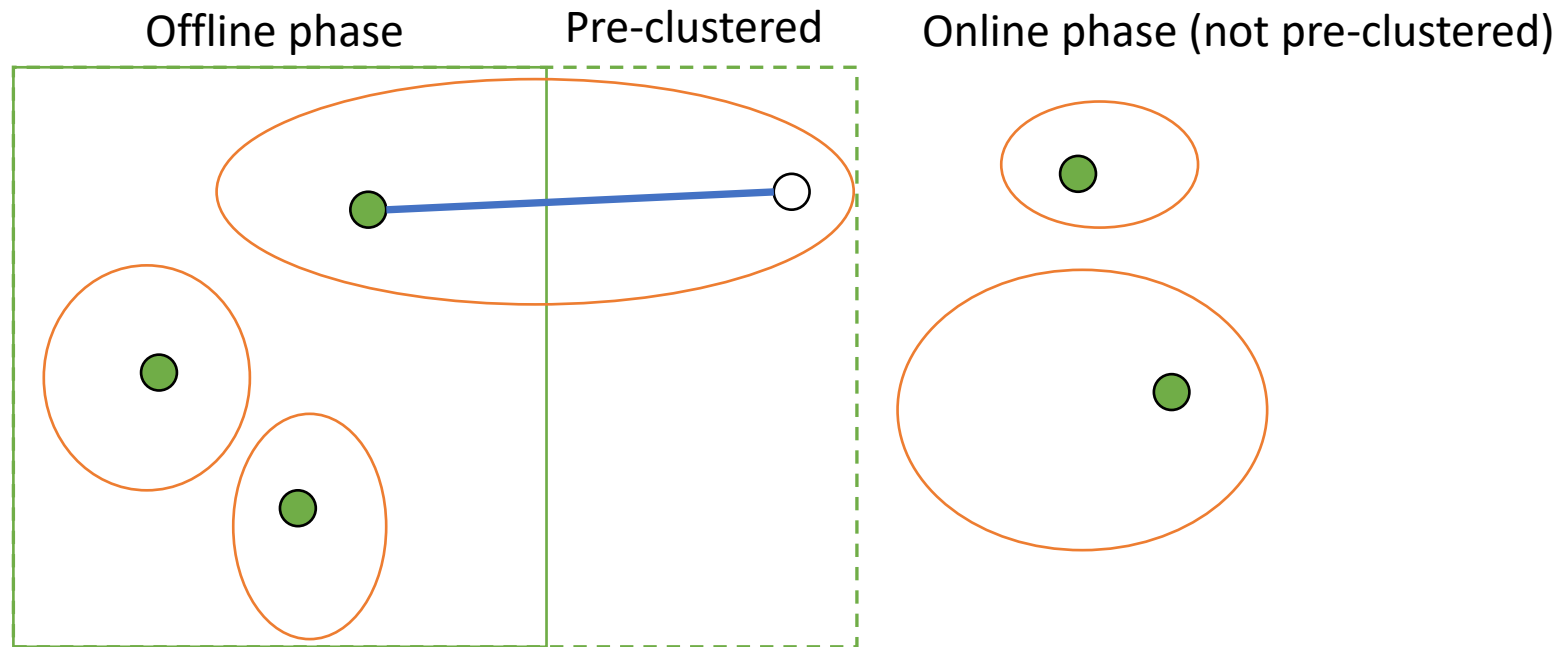
- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals



Analysis Overview

Pivot is $O(1)$ -competitive in random order

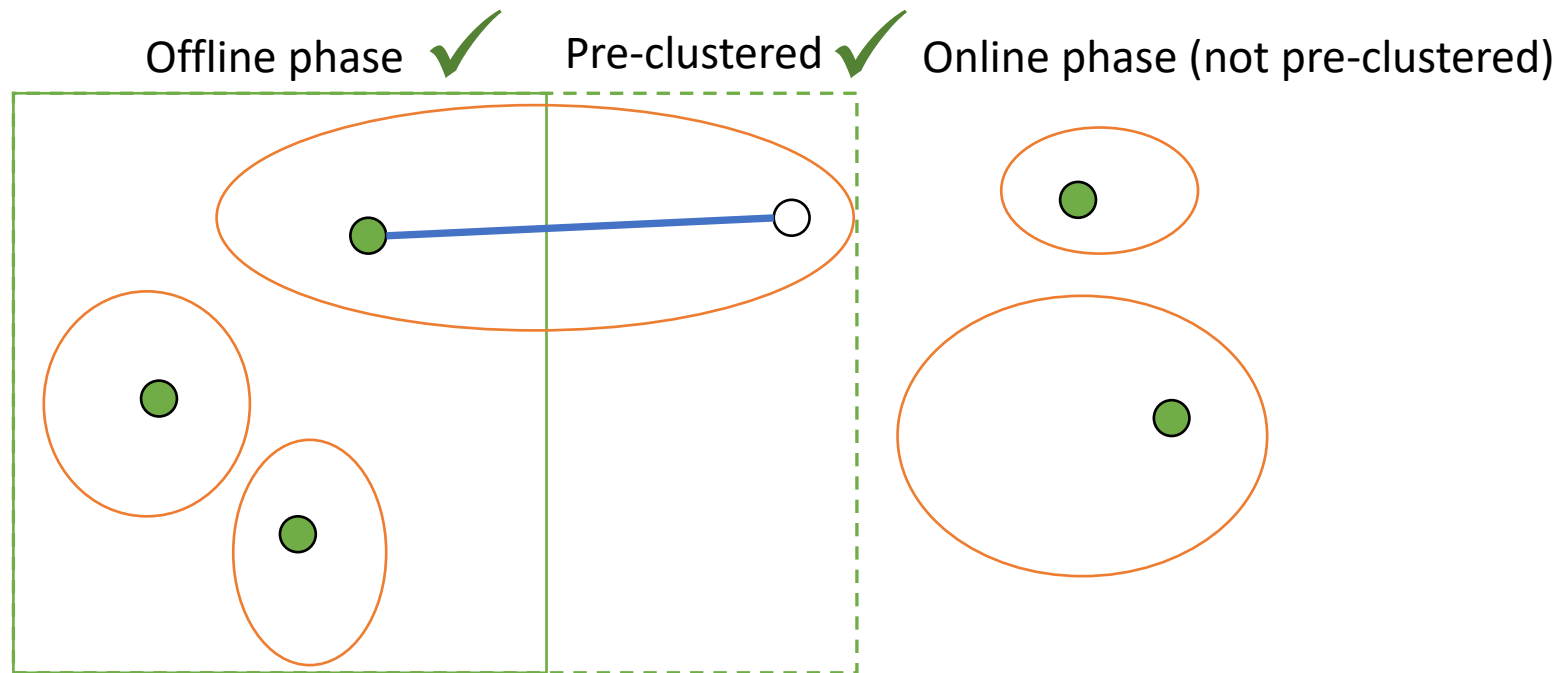
- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals



Pivot is $O(1)$ -competitive in random order

Analysis Overview

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals

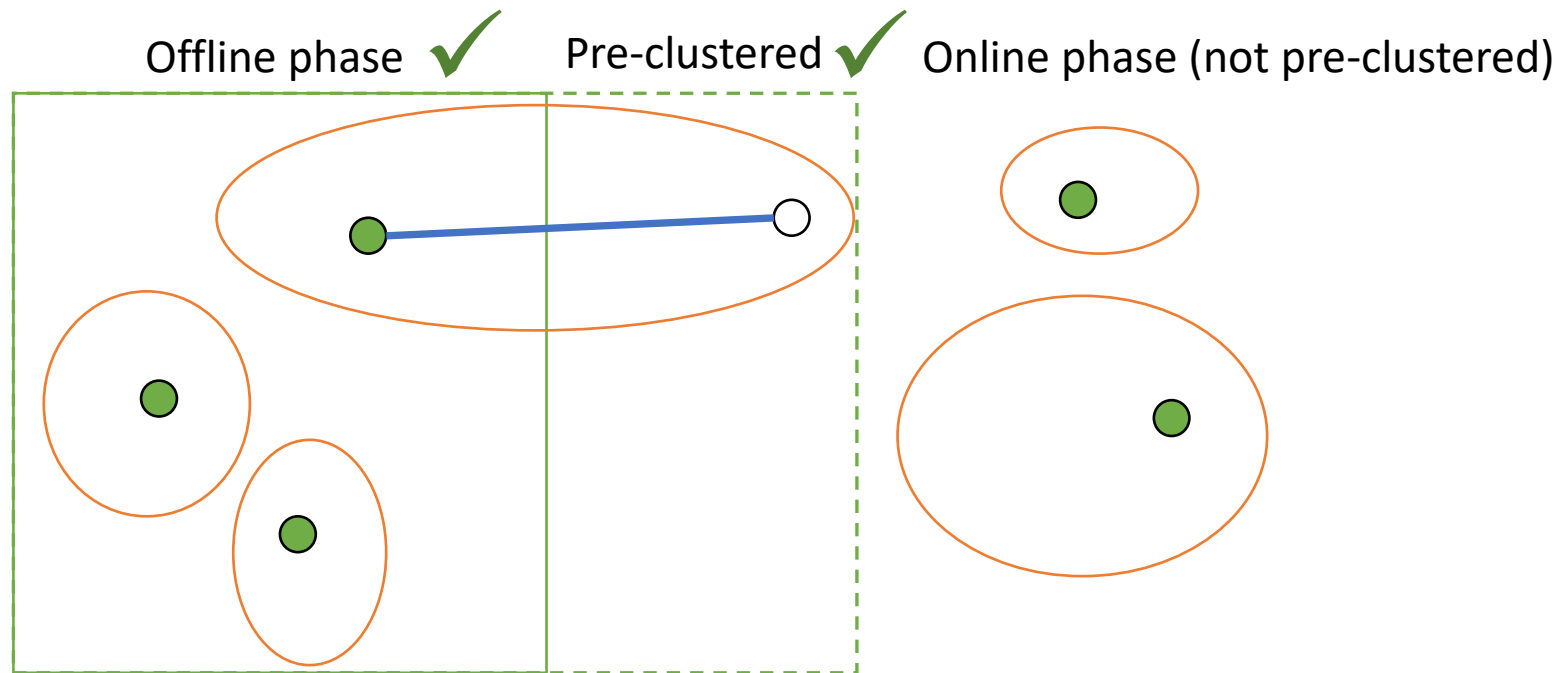


Analysis Overview

Pivot is $O(1)$ -competitive in random order

Not pre-clustered graph is sparse

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals

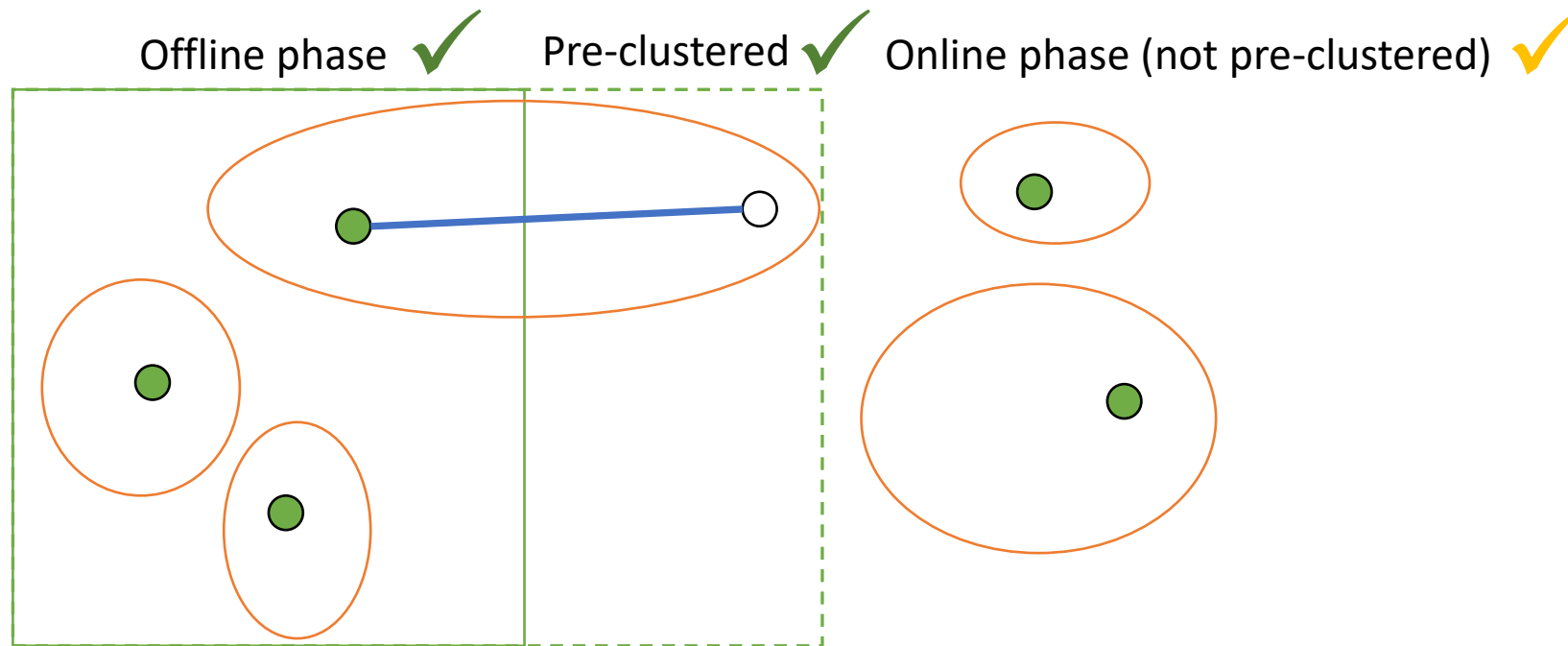


Analysis Overview

Pivot is $O(1)$ -competitive in random order

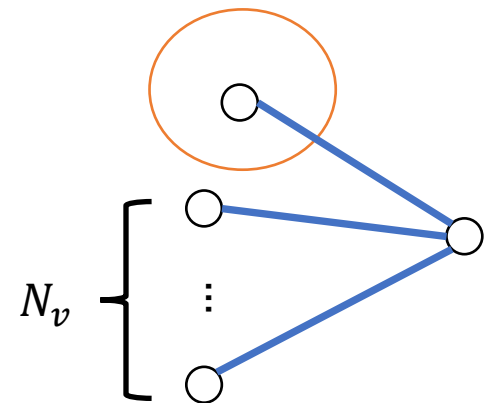
Not pre-clustered graph is sparse

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals



Lemma: For any vertex v , define the random variable N_v such that $N_v = \#(\text{not pre-clustered neighbors of } v)$ if v is not pre-clustered, and $N_v = 0$ otherwise. Then $\mathbb{E} N_v = O(\frac{1}{\epsilon})$.

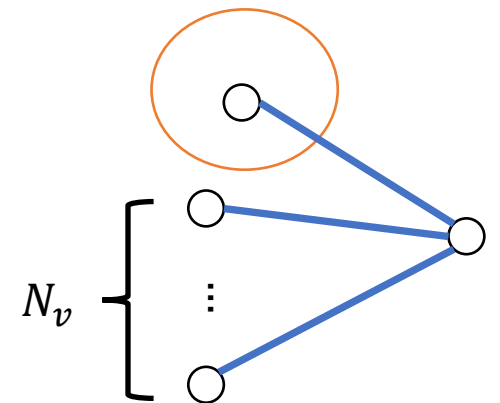
Our Algorithm: Run Pivot in **random order** in offline phase; then continue in **arrival order** in online phase



Lemma: For any vertex v , define the random variable N_v such that $N_v = \#(\text{not pre-clustered neighbors of } v)$ if v is not pre-clustered, and $N_v = 0$ otherwise. Then $\mathbb{E} N_v = O(\frac{1}{\epsilon})$.

- Assume **no corruption**
- $N_v \geq k \Rightarrow$ For each arrival in offline phase, v still has k not pre-clustered neighbors, and none of them arrive next (A_i)

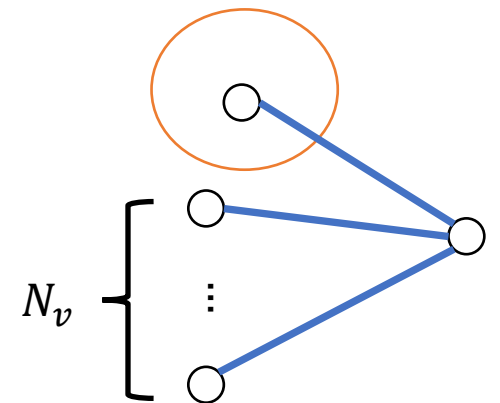
Our Algorithm: Run Pivot in **random order** in offline phase; then continue in **arrival order** in online phase



Lemma: For any vertex v , define the random variable N_v such that $N_v = \#(\text{not pre-clustered neighbors of } v)$ if v is not pre-clustered, and $N_v = 0$ otherwise. Then $\mathbb{E} N_v = O(\frac{1}{\epsilon})$.

- Assume **no corruption**
- $N_v \geq k \Rightarrow$ For each arrival in offline phase, v still has k not pre-clustered neighbors, and none of them arrive next (A_i)
- $\mathbb{P}(A_i \mid A_{i-1}, \dots, A_1) \leq \left(1 - \frac{k}{n}\right) \Rightarrow \mathbb{P}(N_v \geq k) \leq \left(1 - \frac{k}{n}\right)^{\epsilon n} \leq e^{-\epsilon k}$
- Integrating over tail gives $\mathbb{E} N_v = O(\frac{1}{\epsilon})$.

Our Algorithm: Run Pivot in **random order** in offline phase; then continue in **arrival order** in online phase

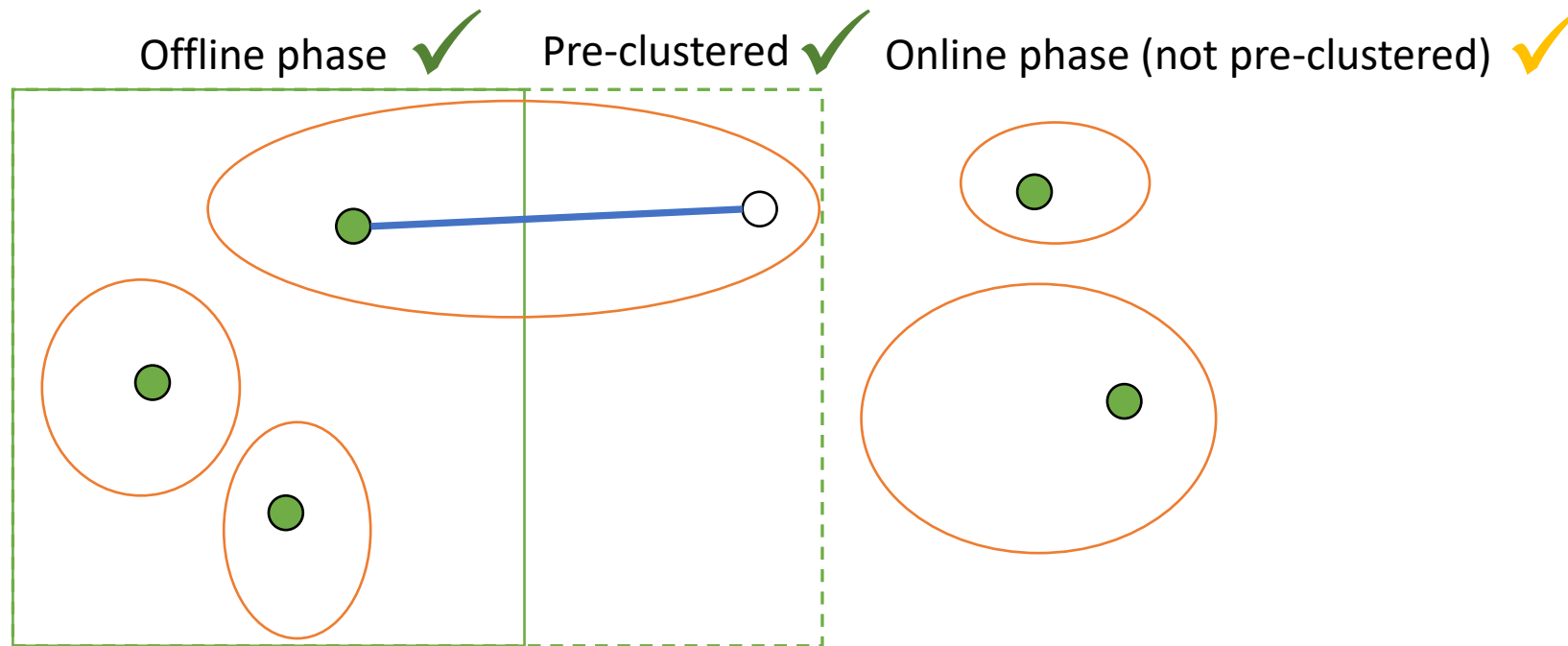


Analysis Overview

Pivot is $O(1)$ -competitive in random order

Not pre-clustered graph is sparse

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals

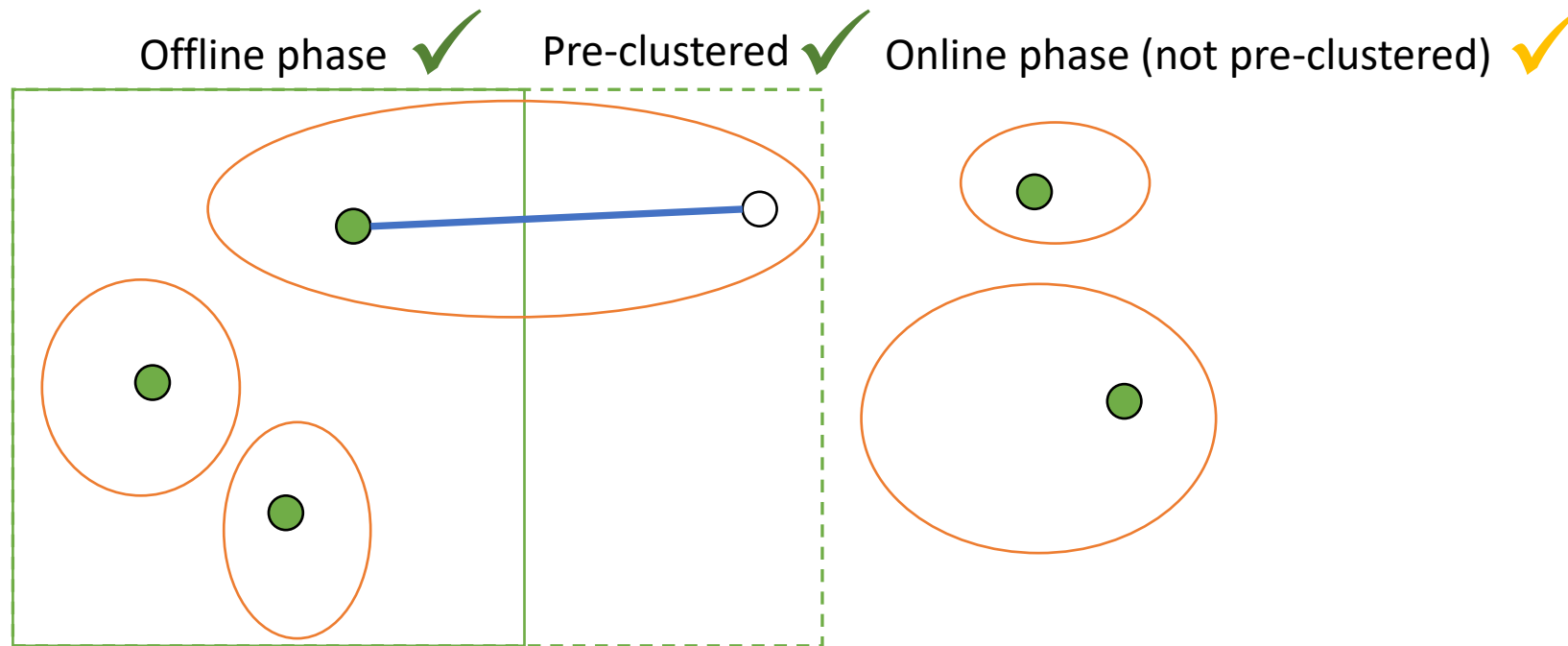


Analysis Overview

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals

Pivot is $O(1)$ -competitive in random order

Not pre-clustered graph has expected degrees all $O(\frac{1}{\epsilon})$

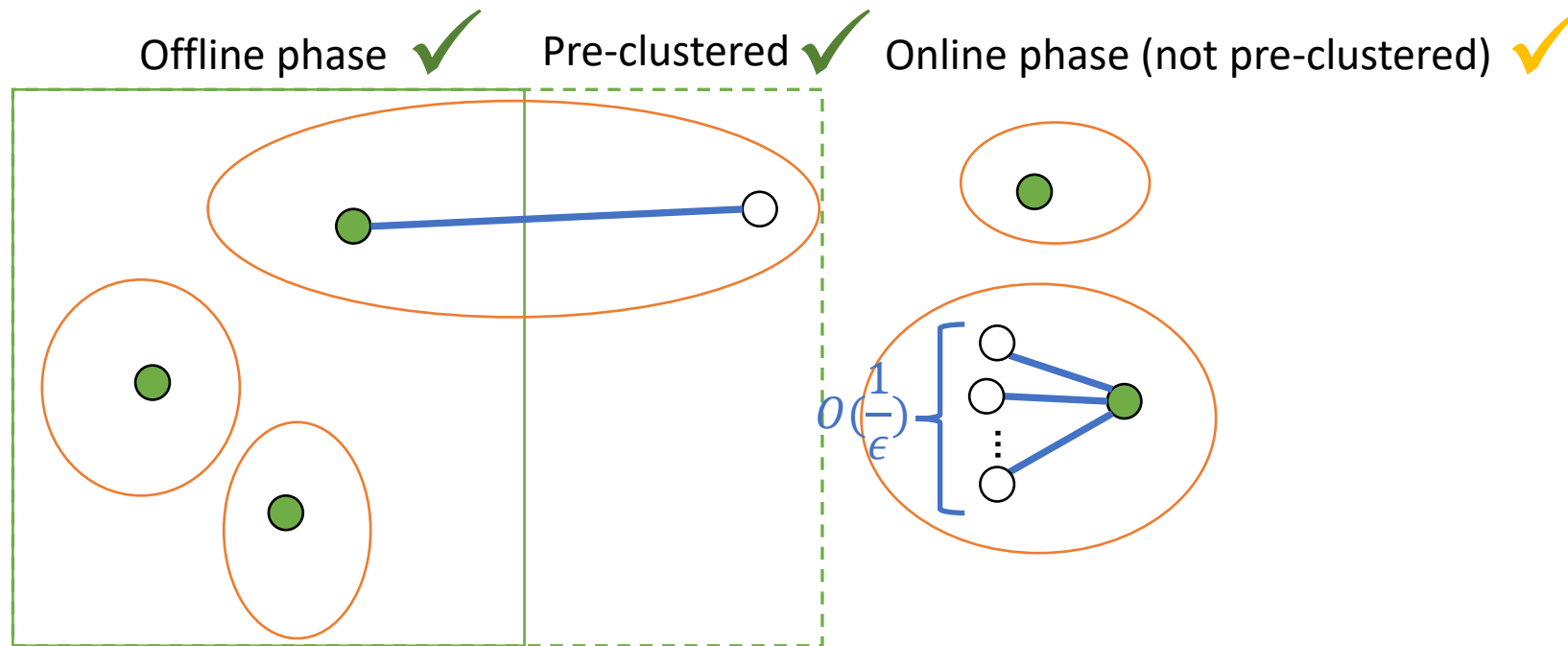


Analysis Overview

- Assume **no corruption**
- Use offline phase to **pre-cluster** online arrivals

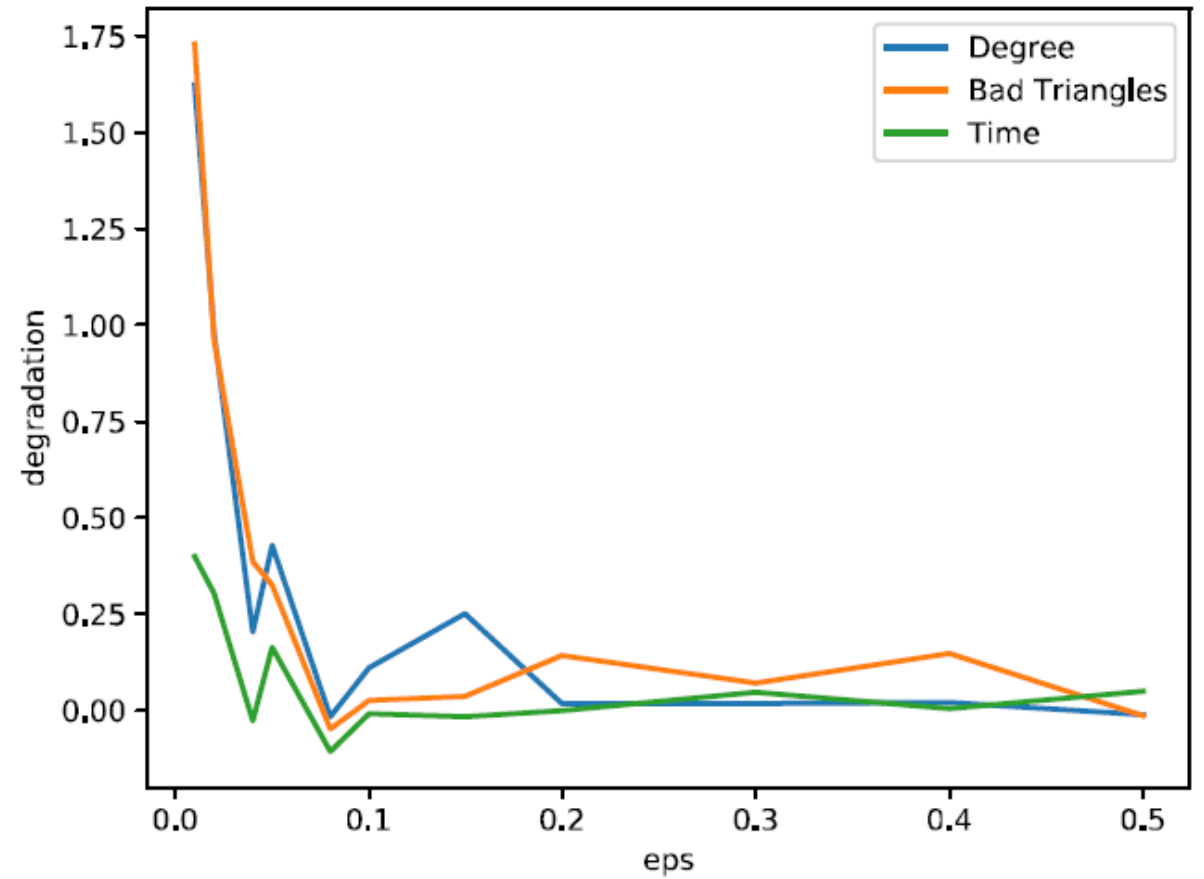
Pivot is $O(1)$ -competitive in random order

Not pre-clustered graph has expected degrees all $O\left(\frac{1}{\epsilon}\right)$



Experimental Results

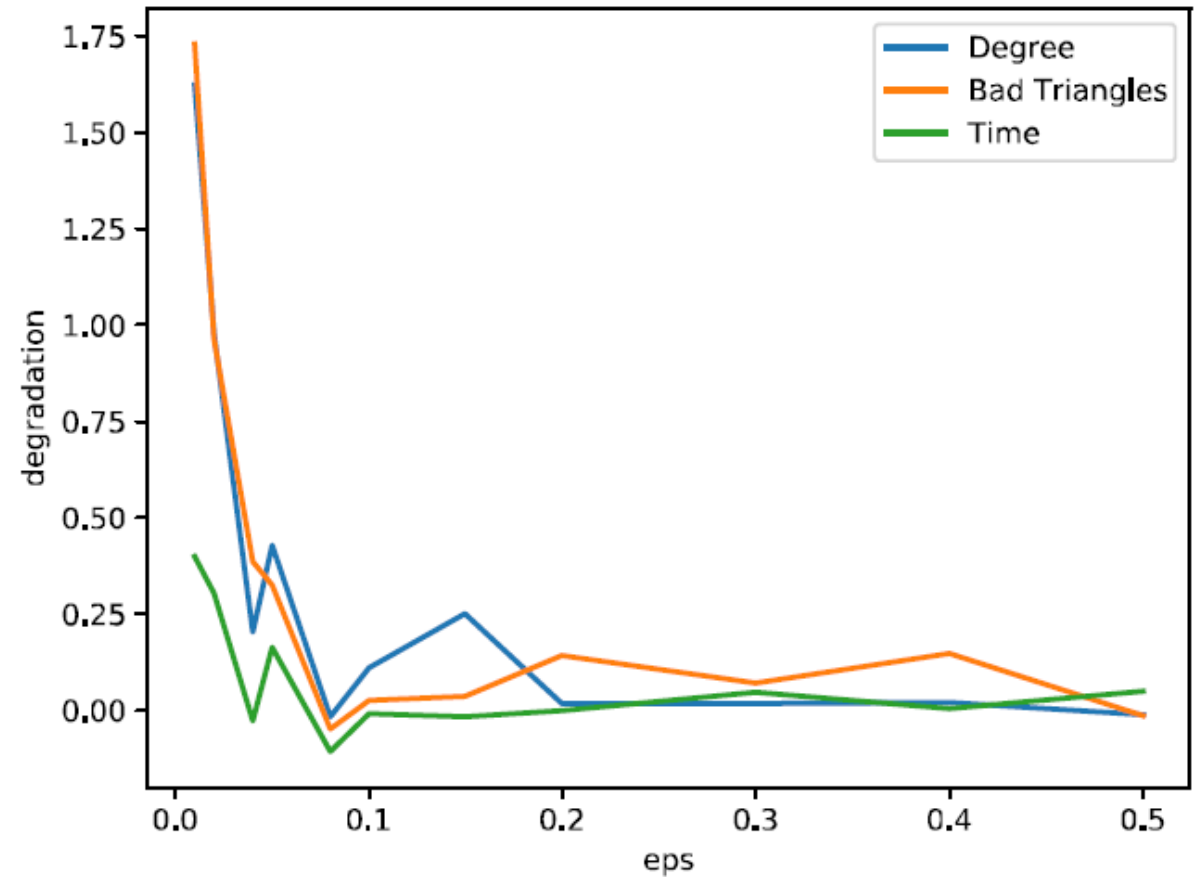
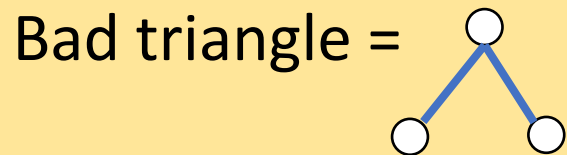
- Algorithm is competitive with offline baseline for moderate ϵ



Experimental Results

- Algorithm is **competitive with offline baseline** for moderate ϵ

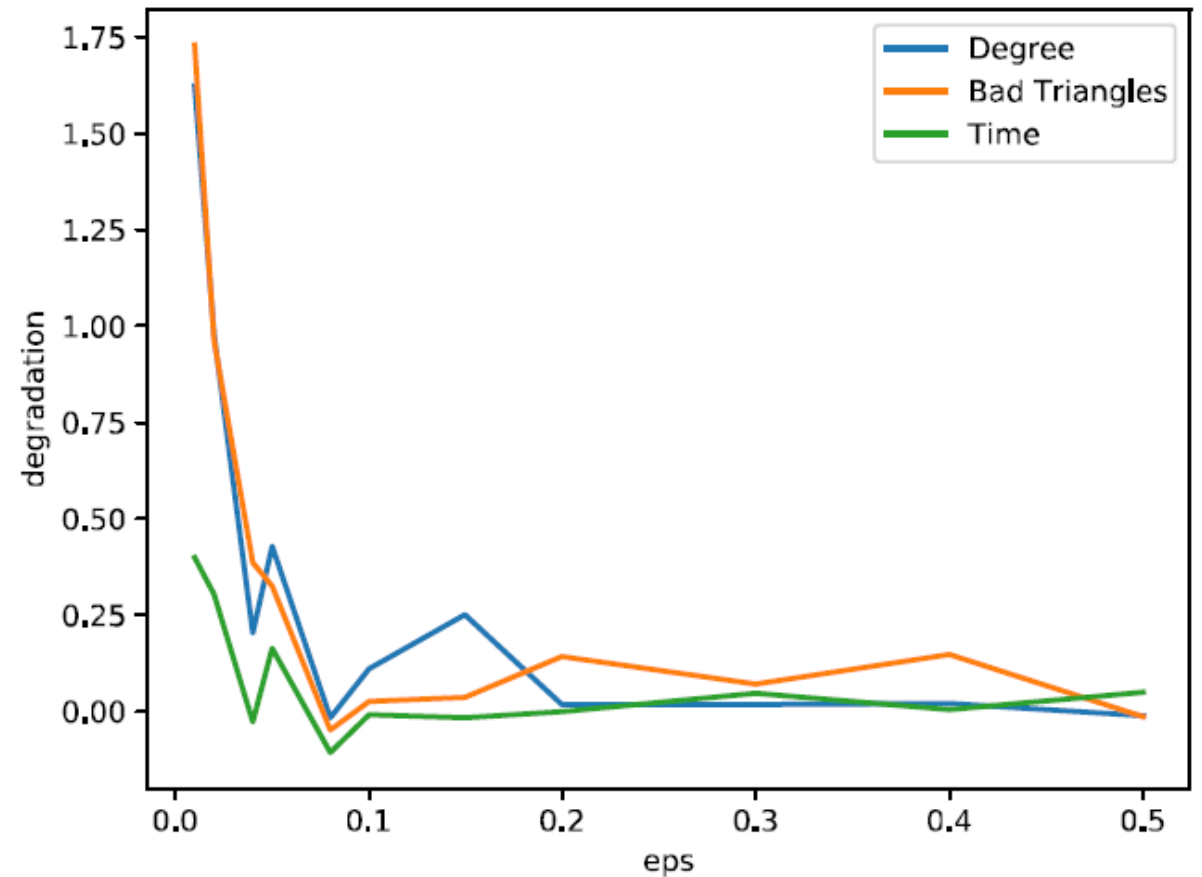
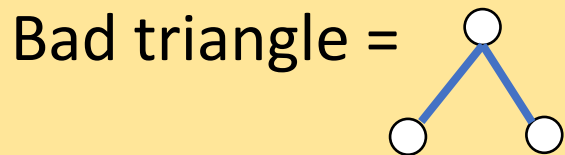
$$\text{Degradation} = \frac{\text{Algorithm} - \text{Baseline}}{\text{Baseline}}$$



Experimental Results

- Algorithm is **competitive with offline baseline** for moderate ϵ
- ... and **robust** to adversarial corruptions
- ... and random sample can be **practically obtained** from past data

$$\text{Degradation} = \frac{\text{Algorithm} - \text{Baseline}}{\text{Baseline}}$$



Summary

- Introduced **semi-online model with adversarial corruptions** ~ **add data-driven decision making to online algorithms**
- Designed novel semi-online algorithm for correlation clustering with tight competitive ratio ~ **best possible way to use historical data**
- $\Omega(\#vertices)$ lower bound online $\Rightarrow O(1)$ -competitive semi-online
- Theory predictive of practice