

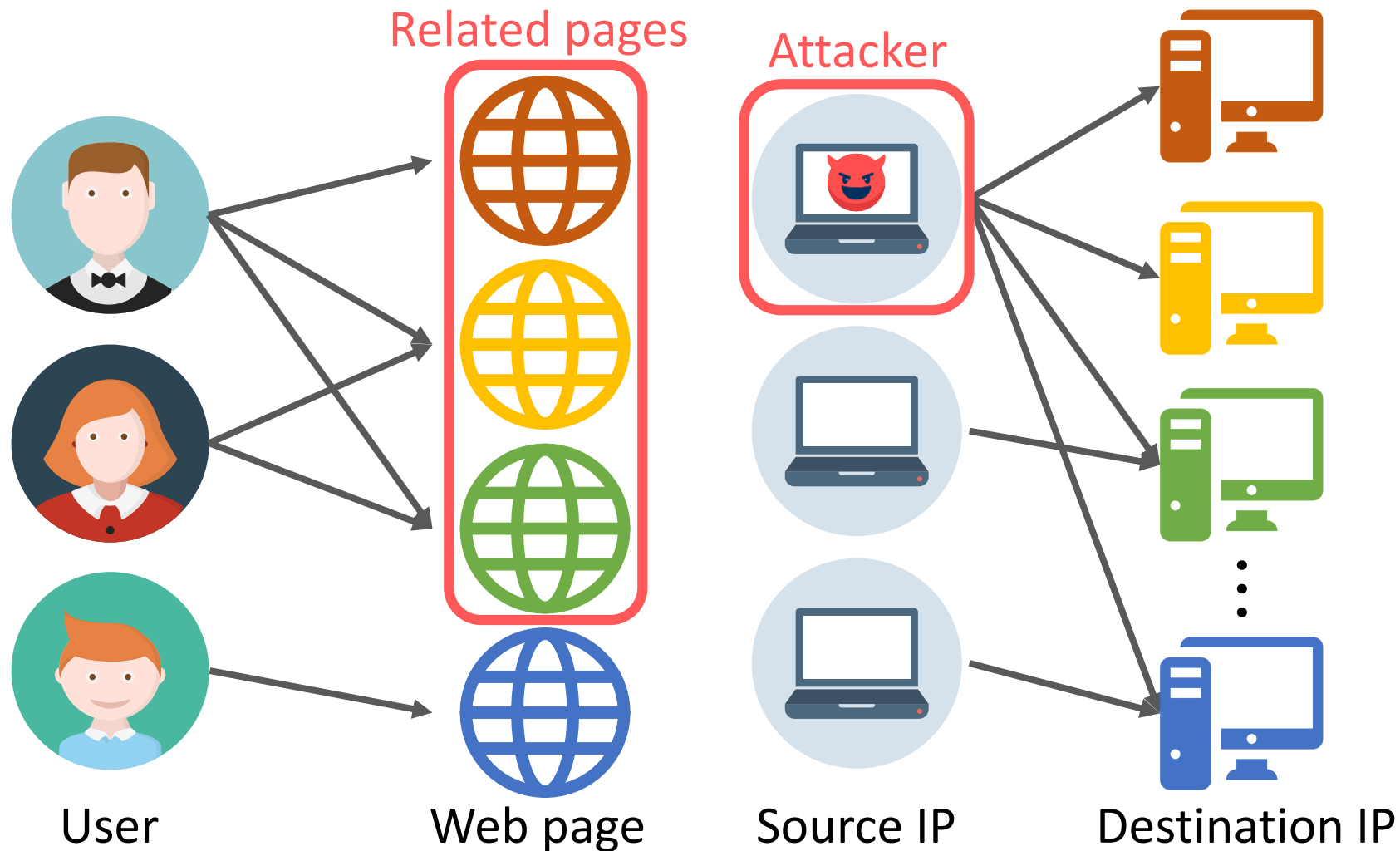
Finding a Concise, Precise, and Exhaustive Set of Near Bi-Cliques in Dynamic Graphs

Hyeonjeong Shin¹, Taehyung Kwon¹, Neil Shah², Kijung Shin¹

¹ KAIST, ² Snap Inc.

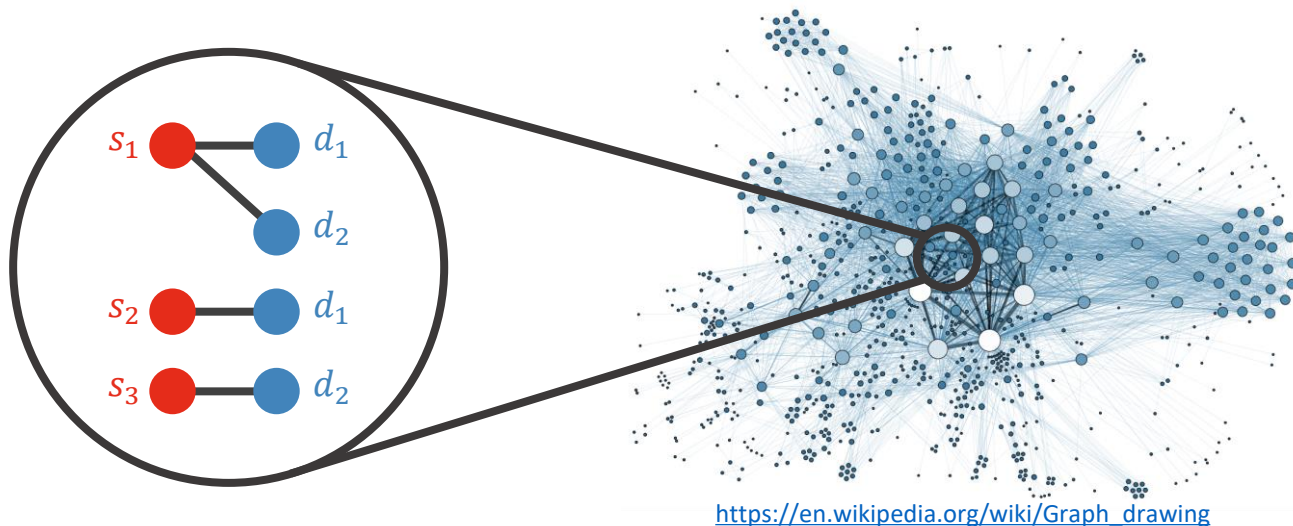


Finding (Near) Bi-Clique Problem



Problem: Too Many (Near) Bi-Cliques

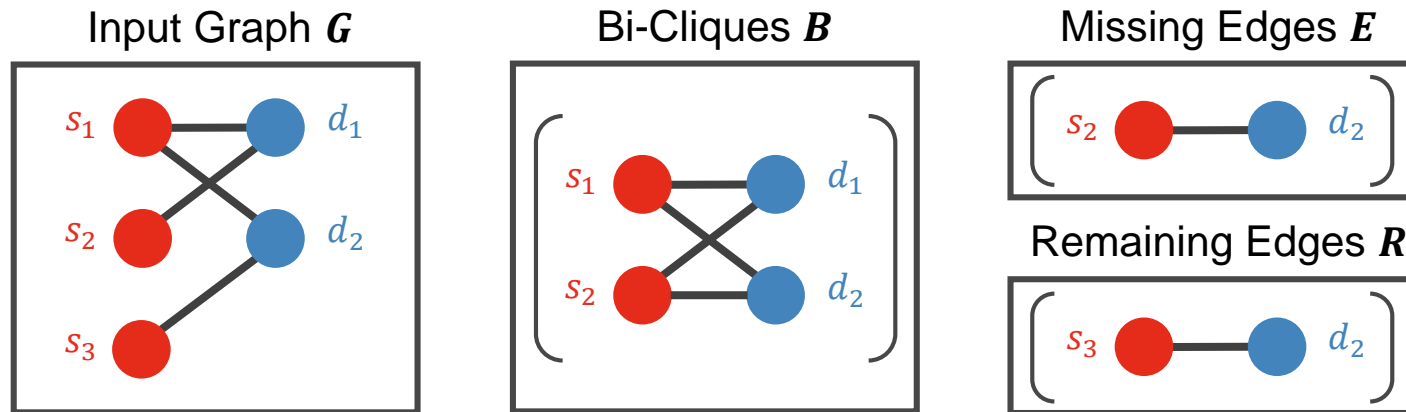
- The number of (near) bi-cliques is exponential in the number of nodes
 - Practically impossible to find all
- Some are too small or highly overlapped
- What is a “good” set of near bi-cliques?



https://en.wikipedia.org/wiki/Graph_drawing

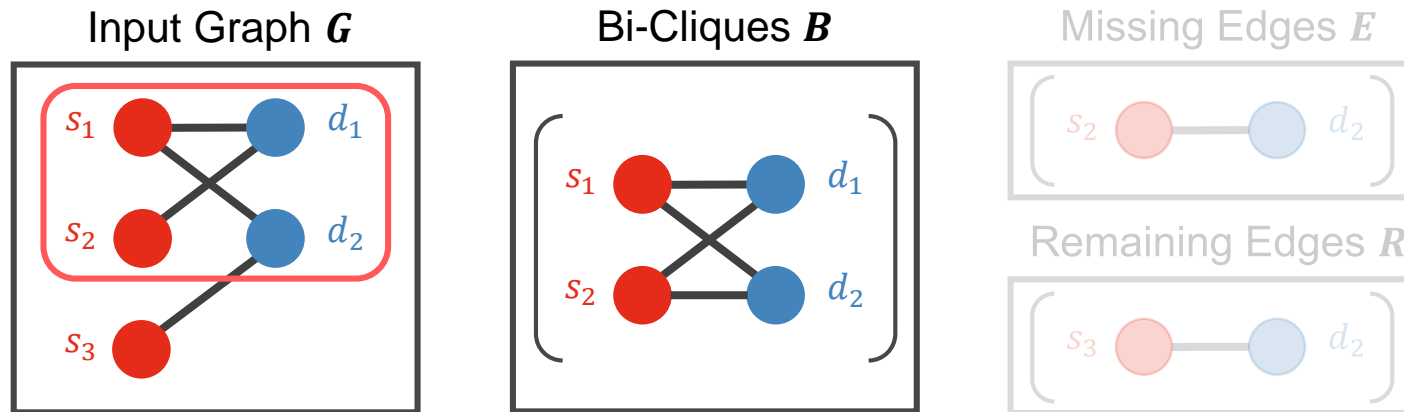
Terminologies

- Exact bi-cliques B
- Missing edges $M := \{e \notin G \mid \exists b \in B \text{ s.t. } e \in b\}$
- Remaining edges $R := \{e \in G \mid \forall b \in B, e \notin b\}$



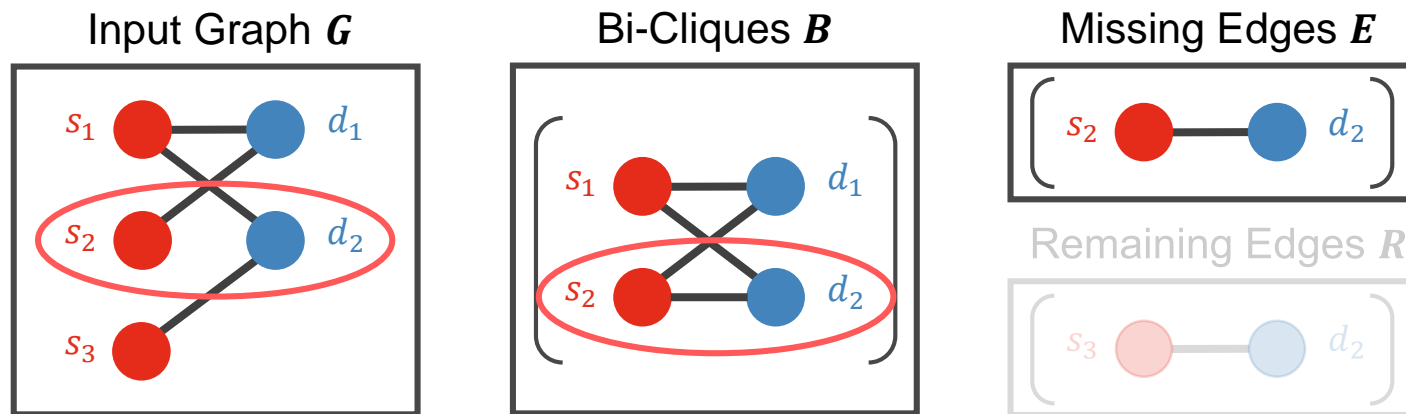
Terminologies

- Exact bi-cliques B
- Missing edges $M := \{e \notin G \mid \exists b \in B \text{ s.t. } e \in b\}$
- Remaining edges $R := \{e \in G \mid \forall b \in B, e \notin b\}$



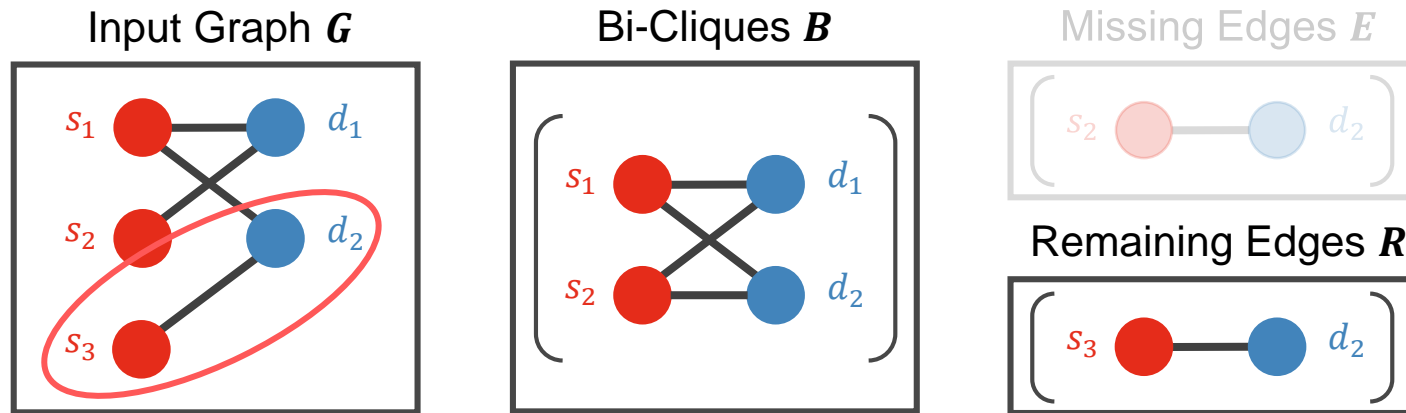
Terminologies

- Exact bi-cliques B
- Missing edges $M := \{e \notin G \mid \exists b \in B \text{ s.t. } e \in b\}$
- Remaining edges $R := \{e \in G \mid \forall b \in B, e \notin b\}$



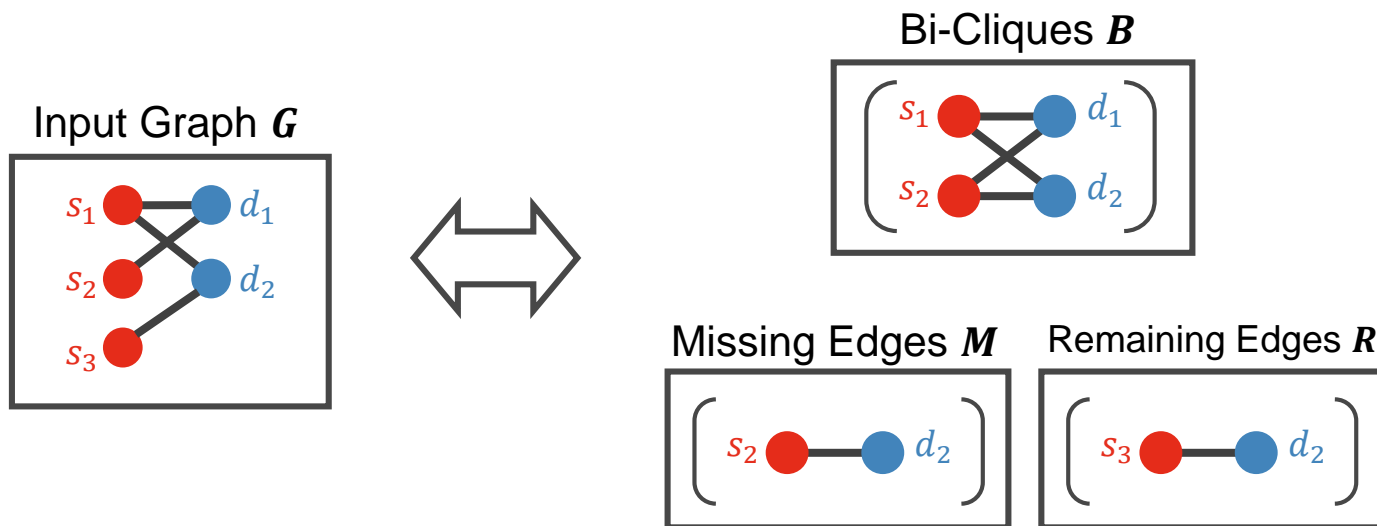
Terminologies

- Exact bi-cliques B
- Missing edges $M := \{e \notin G \mid \exists b \in B \text{ s.t. } e \in b\}$
- Remaining edges $R := \{e \in G \mid \forall b \in B, e \notin b\}$



Problem Definition

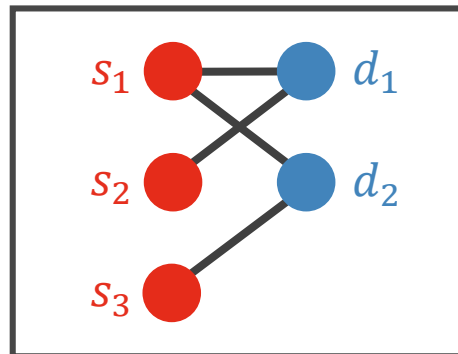
- A dynamic graph can be described using B , M and R
 - By removing $e \in M$ from B and adding $e \in R$ into B



Solution: Find “Good” Near Bi-Cliques

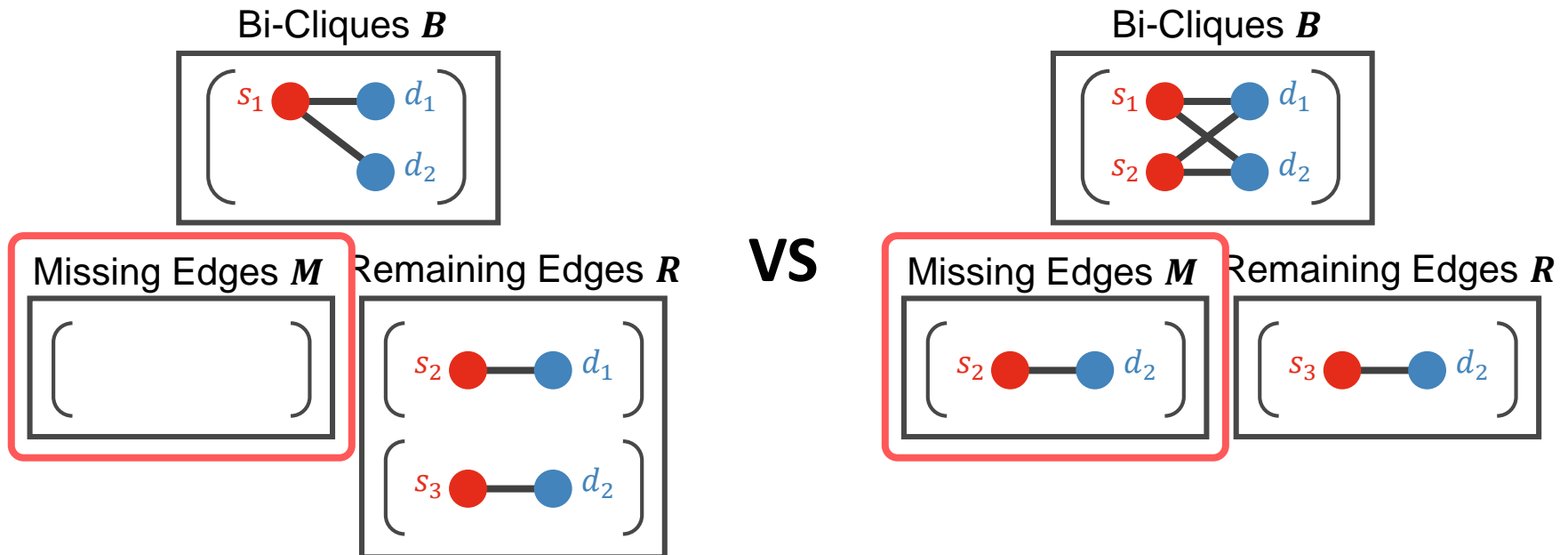
- What are the rules for a “good” set of near bi-cliques?
- Let’s consider the below graph G

Input Graph G



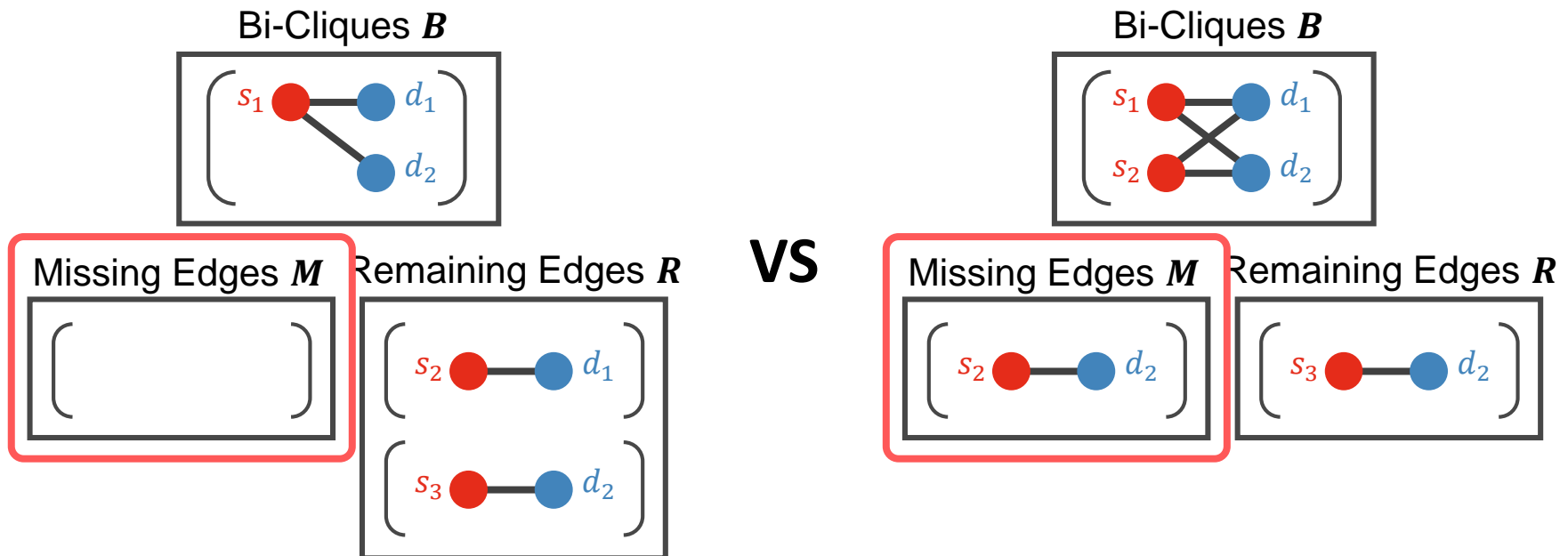
Solution: Find “Good” Near Bi-Cliques

- The left one has less missing edges M
 - i.e., less errors inside B



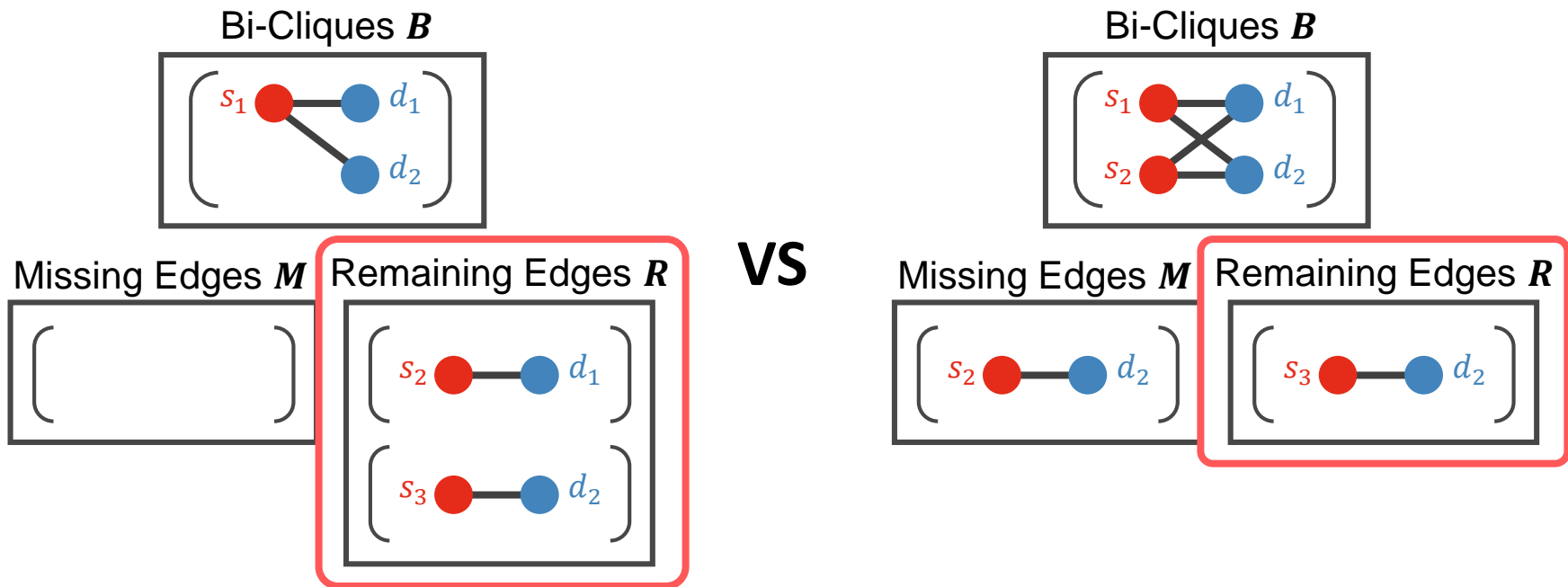
Solution: Find “Good” Near Bi-Cliques

- Rule 1) **Preciseness**: Near bi-cliques should be close to exact bi-cliques while minimizing # of missing edges



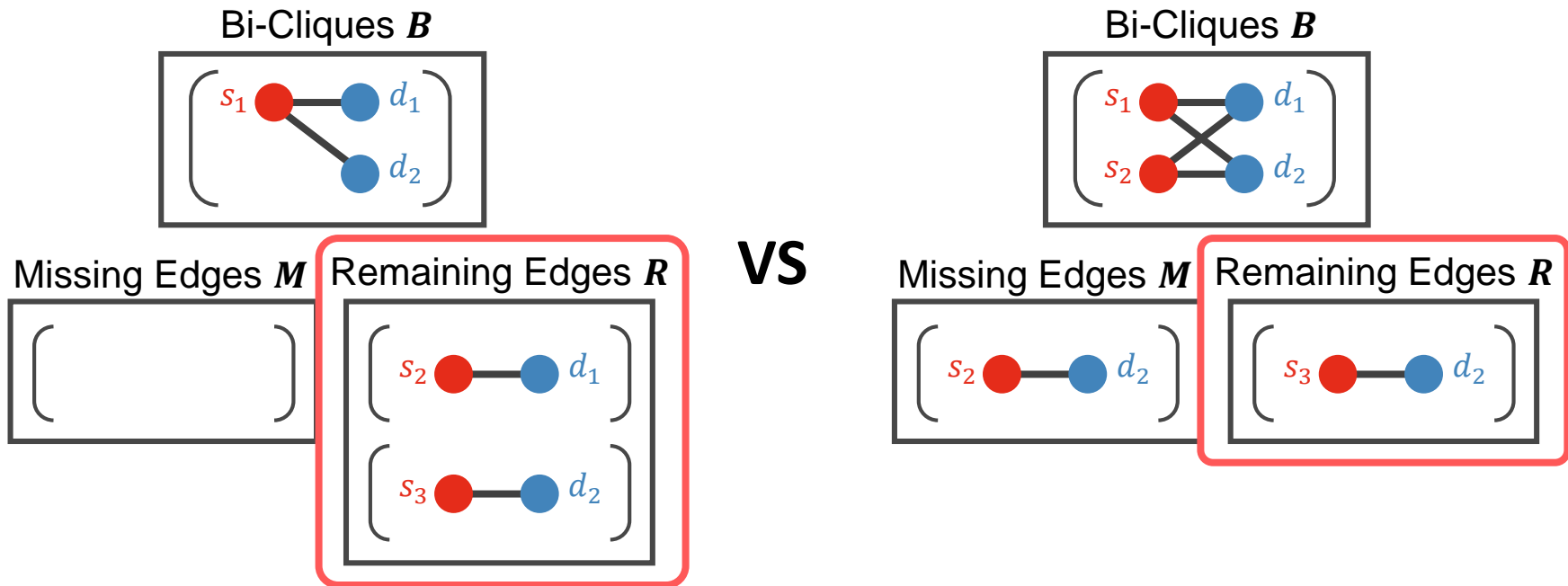
Solution: Find “Good” Near Bi-Cliques

- The right B has less remaining edges, R
 - i.e., covers more edges in G



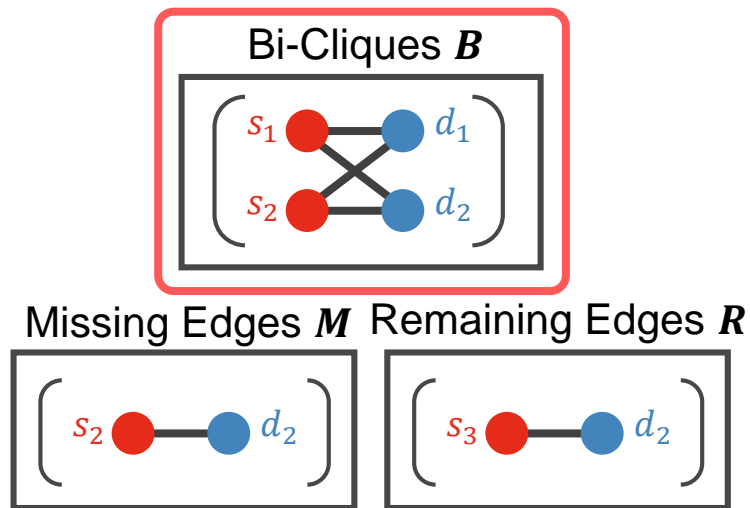
Solution: Find “Good” Near Bi-Cliques

- Rule 2) **Exhaustiveness**: Near bi-cliques should cover a large portion of the given graph G

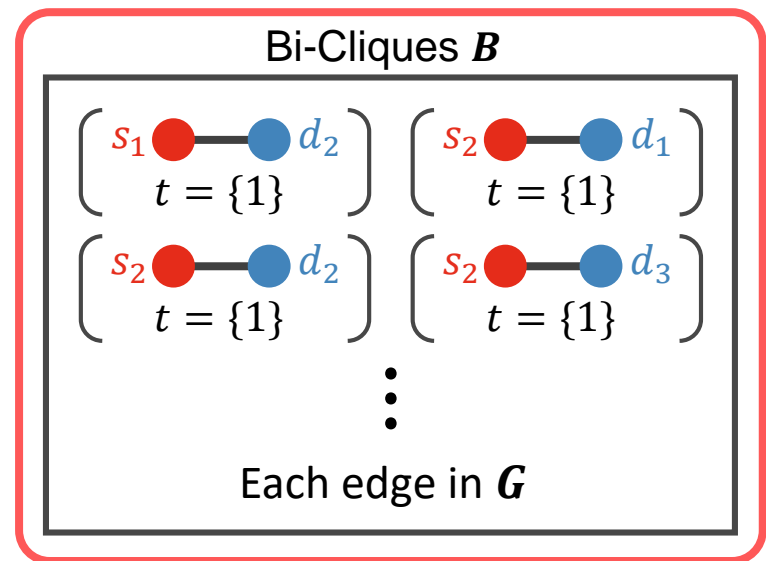


Solution: Find “Good” Near Bi-Cliques

- Rule 3) **Conciseness**: A larger near bi-clique is better than smaller one



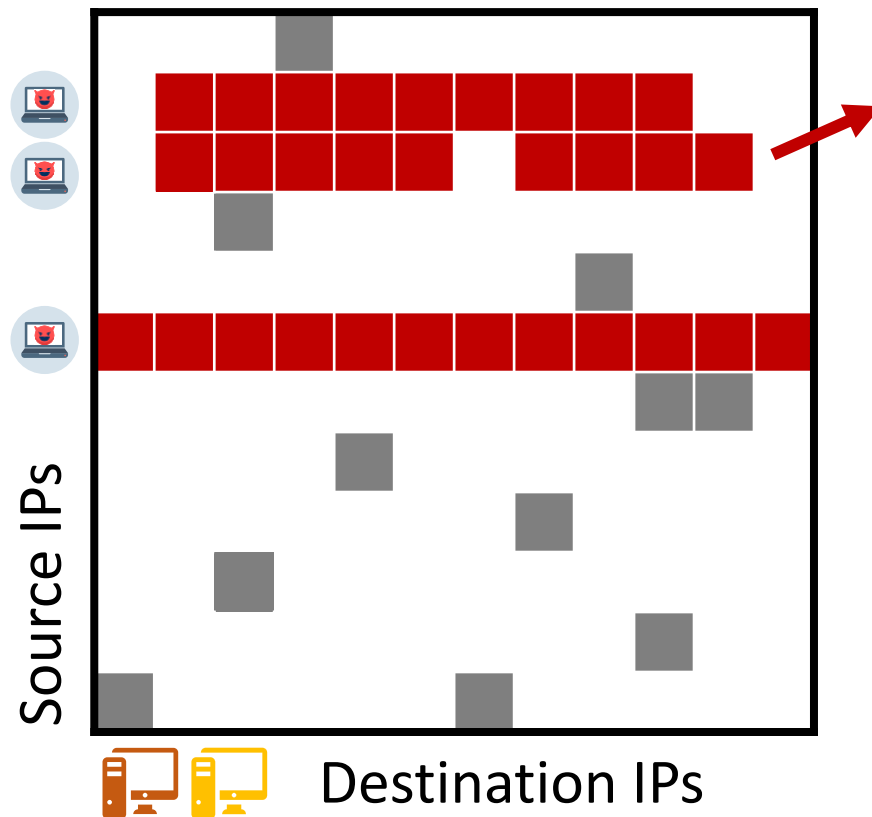
VS



Solution: Find “Good” Near Bi-Cliques

- What is a “good” set of near bi-cliques?

An Adjacency Matrix of a Dynamic Graph



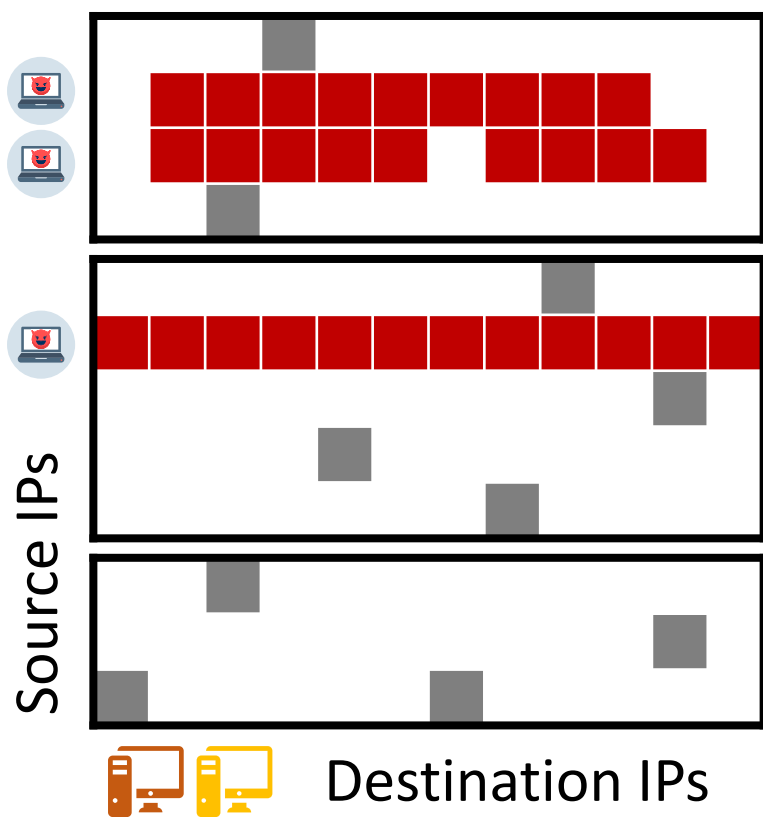
- A set of near bi-cliques
- Covering a large portion
→ Exhaustively
 - Consisting of many nodes
→ Concisely
 - With less error
→ Precisely

Quality Definition

- Each sub-goal minimizes the encoding cost of
 - ✓ Preciseness : missing edges, $\phi_P(\mathbf{M})$
 - ✓ Exhaustiveness : remaining edges, $\phi_E(\mathbf{R})$
 - ✓ Conciseness : exact bi-cliques, $\phi_C(\mathbf{B})$
- Total cost, $\phi(\hat{\mathbf{B}}, \mathbf{G}) := \phi_P(\mathbf{M}) + \phi_E(\mathbf{R}) + \phi_C(\mathbf{B})$
 - $\hat{\mathbf{B}}$: near bi-cliques obtained from \mathbf{B} , \mathbf{R} , and \mathbf{M}
 - Minimizing $\phi(\hat{\mathbf{B}}, \mathbf{G})$ aligns with the Minimum Description Length (MDL) principle to describe \mathbf{G}

Solution: Find “Good” Near Bi-Cliques

- How can we find them?



A detected near bi-clique



Step1. Cut into partitions

Step2. Peel good near bi-cliques from each partition

Problem Definition

- **Given:** a graph G
- **Find:** near bi-cliques \hat{B}
- **To minimize:** the total cost $\phi(\hat{B}, G)$

Search Algorithm

- Our goal is to design an approximation algorithm

- ✓ **High Quality & Speed:** quickly detects high-quality near bi-cliques
- ✓ **Scalability:** runs in near-linear time
- ✓ **Applicability:** can be applied to compression and pattern discovery

- **CutNPeel**, our proposed method, satisfies all the requirements

Finding One Near Bi-Clique

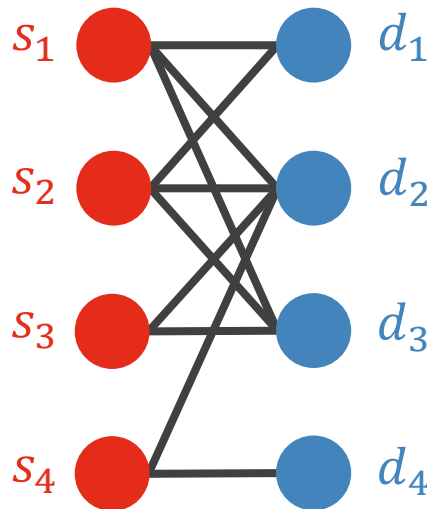
- Find a single near bi-clique \hat{b} with maximum *saving*

$$\text{Saving}(\hat{b}, G) := \phi(\emptyset, G) - \phi(\{\hat{b}\}, G)$$

- Saving approximates $\phi(\hat{B}, G) - \phi(\hat{B} \cup \{\hat{b}\}, G)$
 - the amount of saving in total cost ϕ due to \hat{b}
- Computed in $O(1)$ time with the guarantee in preciseness

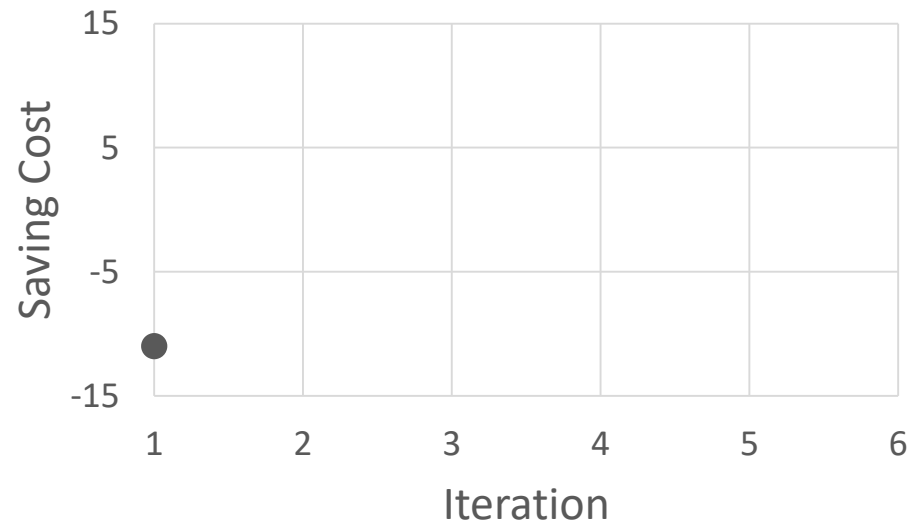
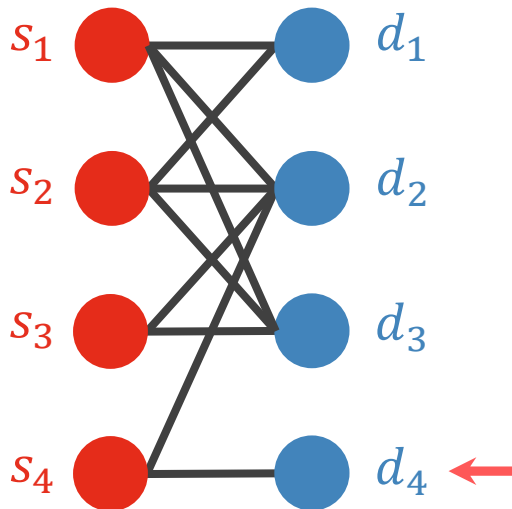
Finding One Near Bi-Clique

- Top-down search
- Removes a node $v \in S \cup D$ with sparsest connectivity
 - S : source nodes
 - D : destination nodes



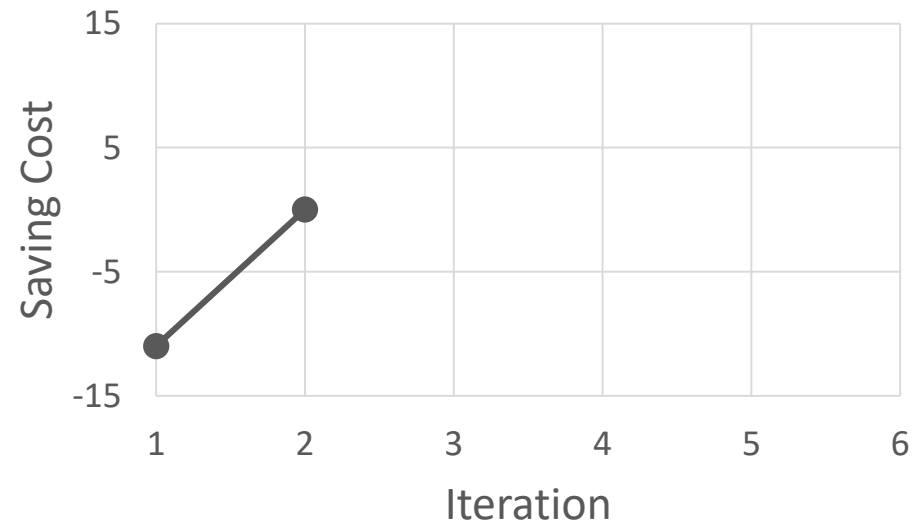
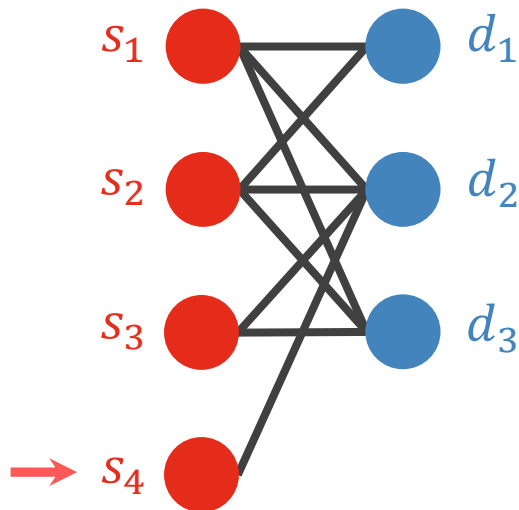
Finding One Near Bi-Clique

- Top-down search
- Removes a node $v \in S \cup D$ with sparsest connectivity



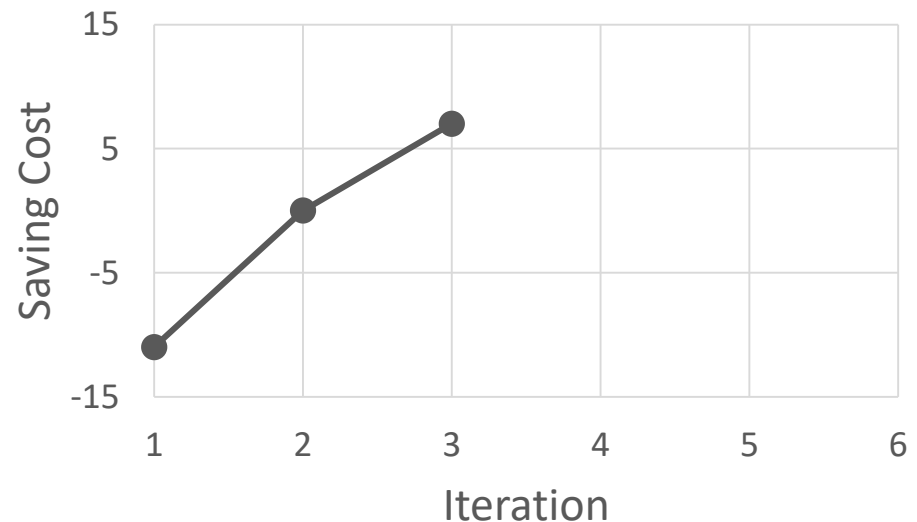
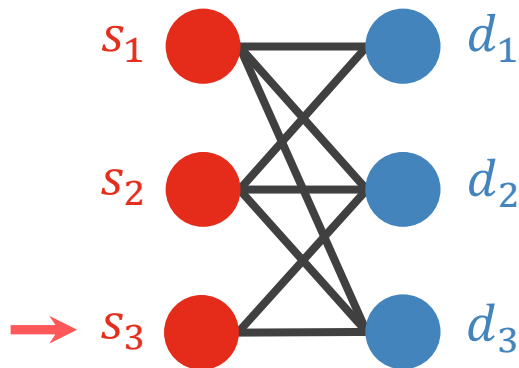
Finding One Near Bi-Clique

- Top-down search
- Removes a node $v \in S \cup D$ with sparsest connectivity



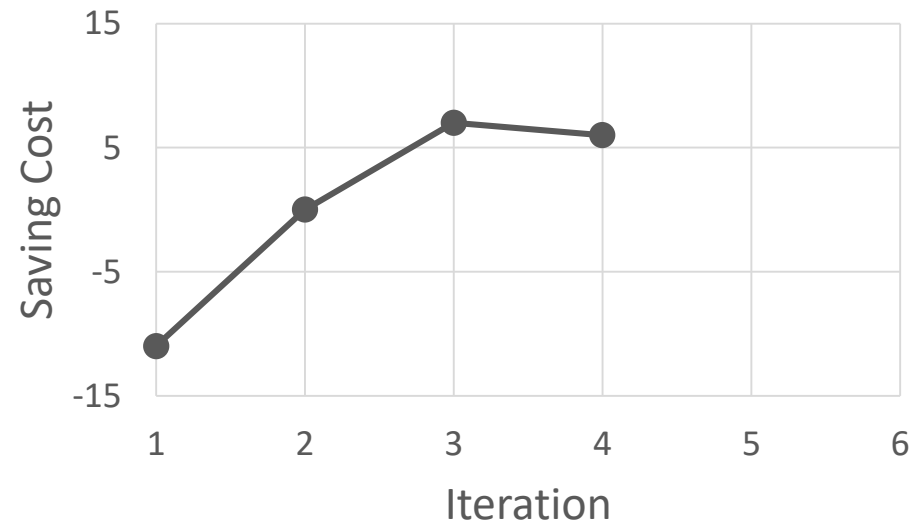
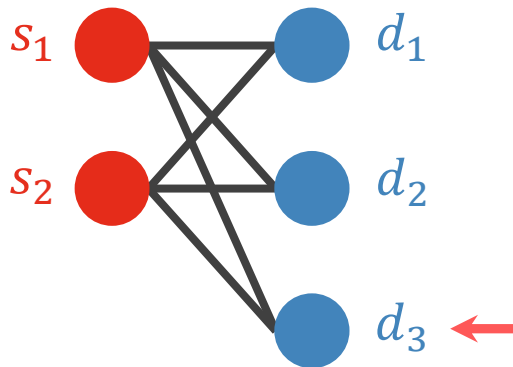
Finding One Near Bi-Clique

- Top-down search
- Removes a node $v \in S \cup D$ with sparsest connectivity



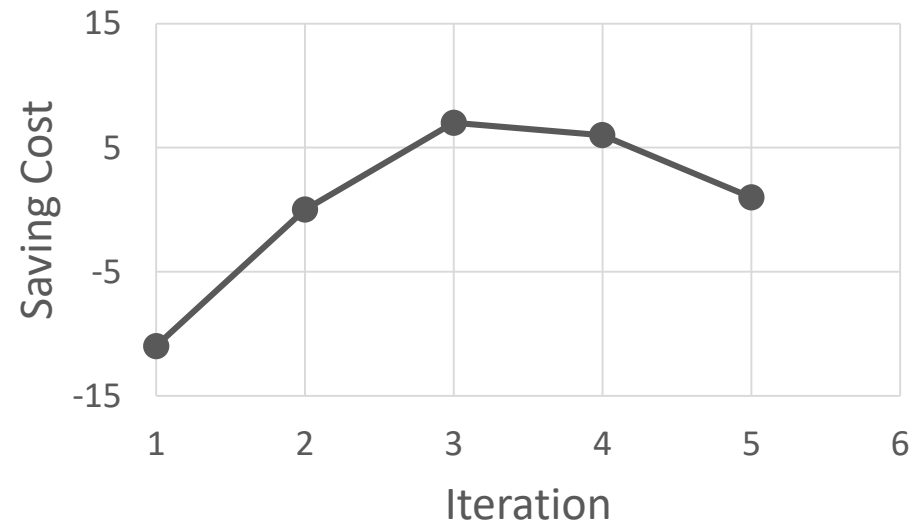
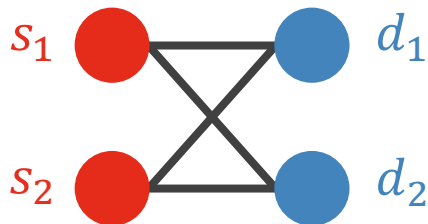
Finding One Near Bi-Clique

- Top-down search
- Removes a node $v \in S \cup D$ with sparsest connectivity



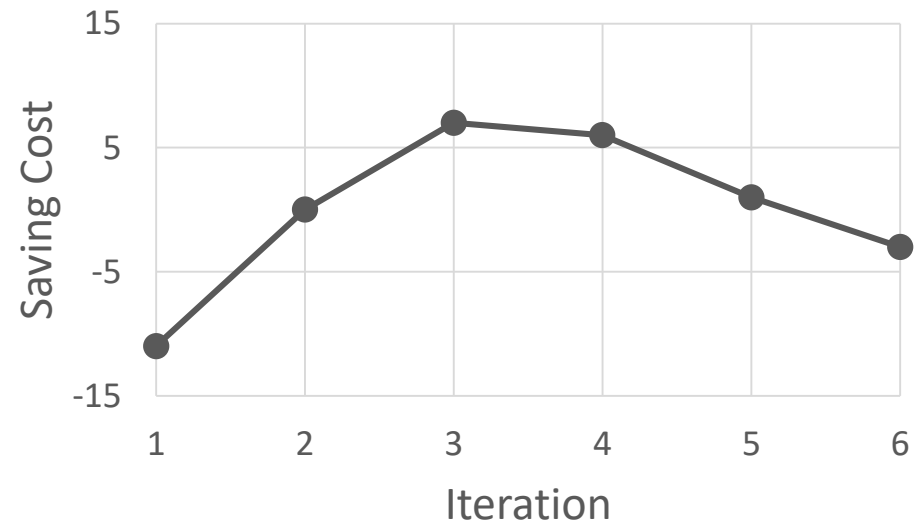
Finding One Near Bi-Clique

- Top-down search
- Removes a node $v \in S \cup D$ with sparsest connectivity



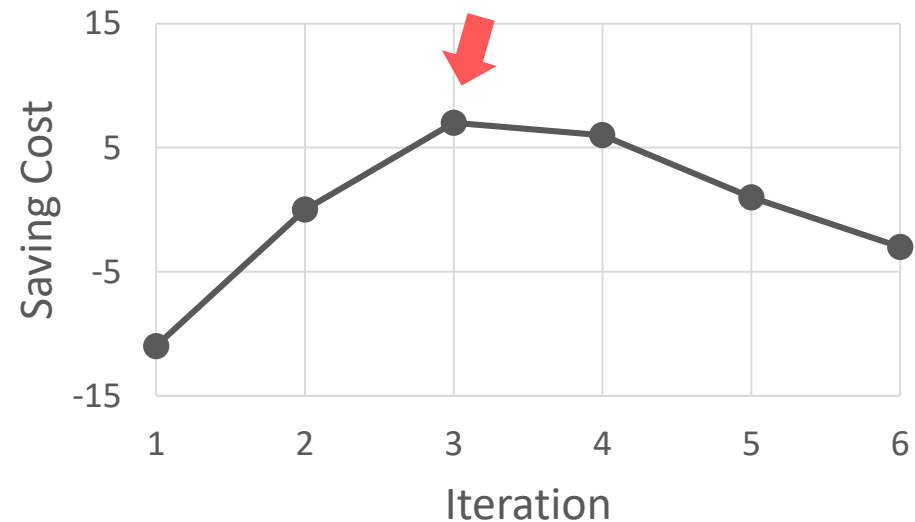
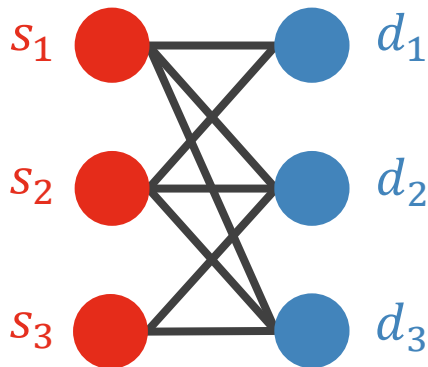
Finding One Near Bi-Clique

- Until all nodes are removed



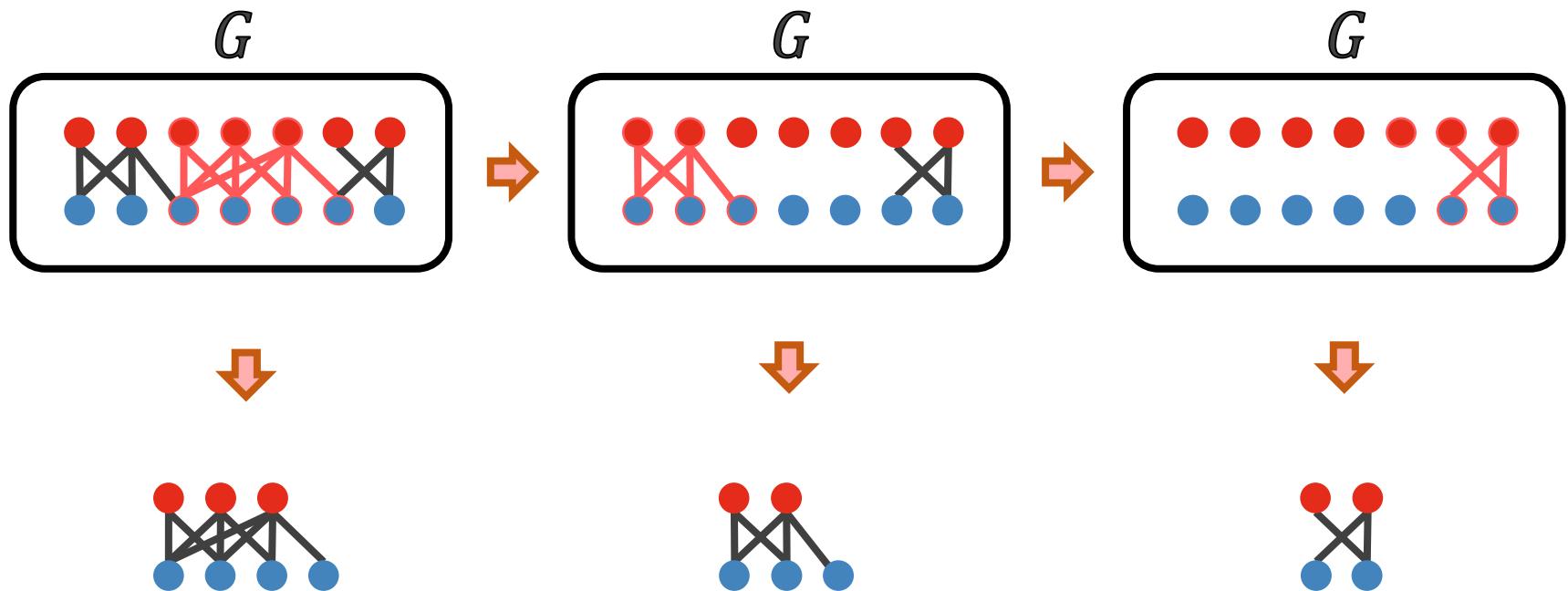
Finding One Near Bi-Clique

- Output: Return the snapshot where the saving cost is maximized



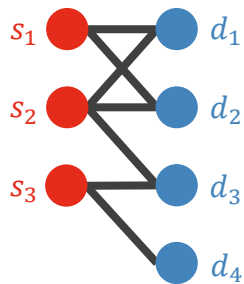
Peel: Preliminary Algorithm

- Remove each detected near bi-clique from the input graph G before finding more bi-cliques

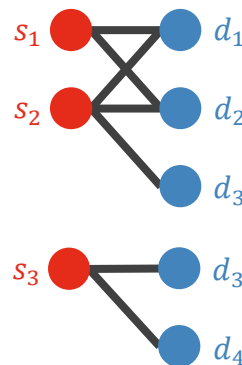


CutNPeel: Search Space Reduction

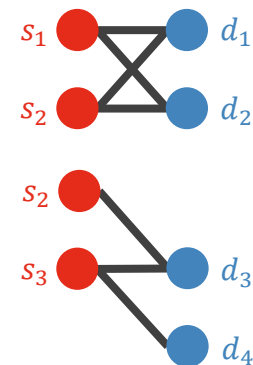
- To reduce the search space, **CutNPeel** uses a partitioning method
- Inspired by min-hashing, **CutNPeel** divides nodes into groups of nodes with similar connectivity
- In each iteration, the leftover graph is re-partitioned



Divide along



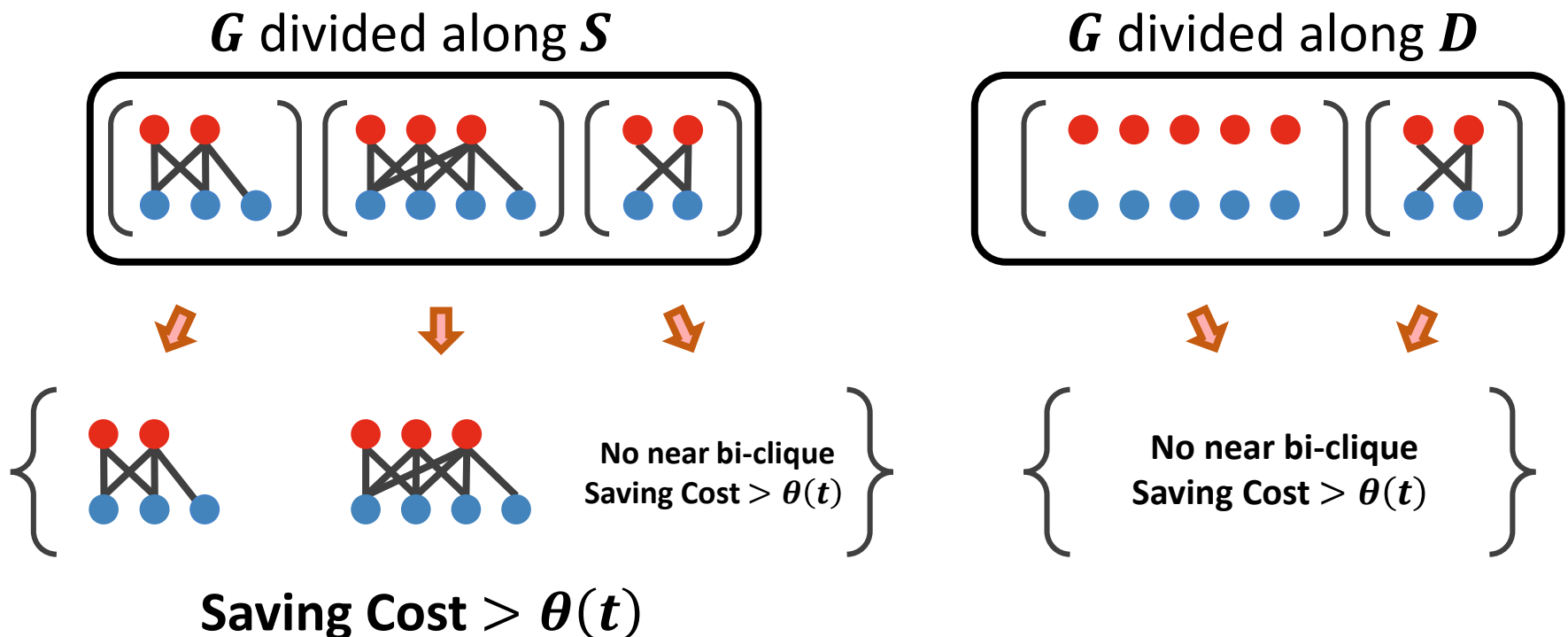
Source S



Destination D

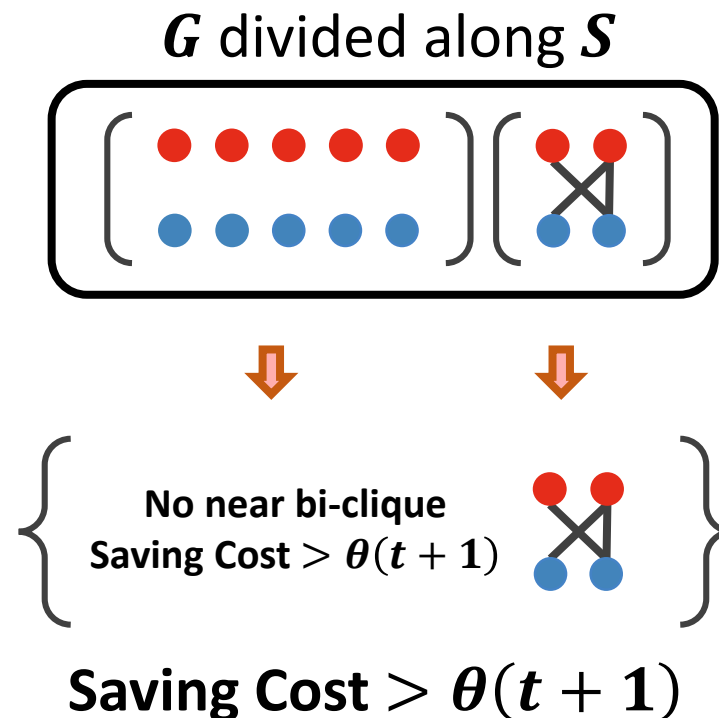
CutNPeel: Quality Guarantee

- In the t -th iteration, **CutNPeel** finds near bi-cliques with saving cost larger than $\theta(t)$ from each partition



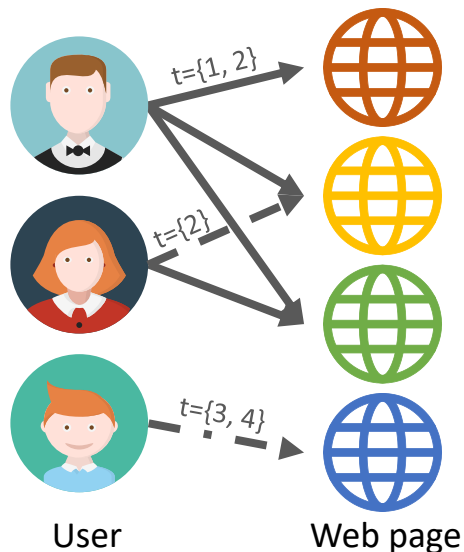
CutNPeel: Quality Guarantee

- In the $(t+1)$ -th iteration, decrease $\theta(t)$ adaptively based on the detected near bi-cliques

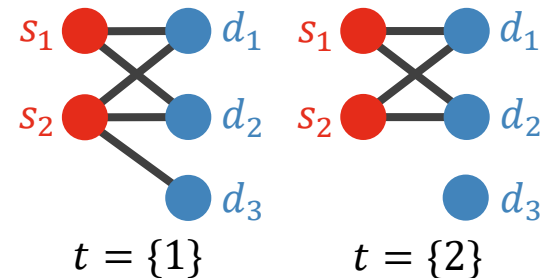


Details in the Paper

- So far, we assumed static graphs for ease of explanation
- In the paper, we considered **dynamic graphs** and **near temporal bi-cliques**



Dynamic Graph



Near Temporal Bi-Cliques

Design of Experiments

- We aim to demonstrate that **CutNPeel**
 - High Quality & Speed:** quickly detects high-quality near bi-cliques
 - Scalability:** runs in near-linear time
 - Applicability:** can be applied to compression and pattern discovery

EXP1. Search Quality & Speed

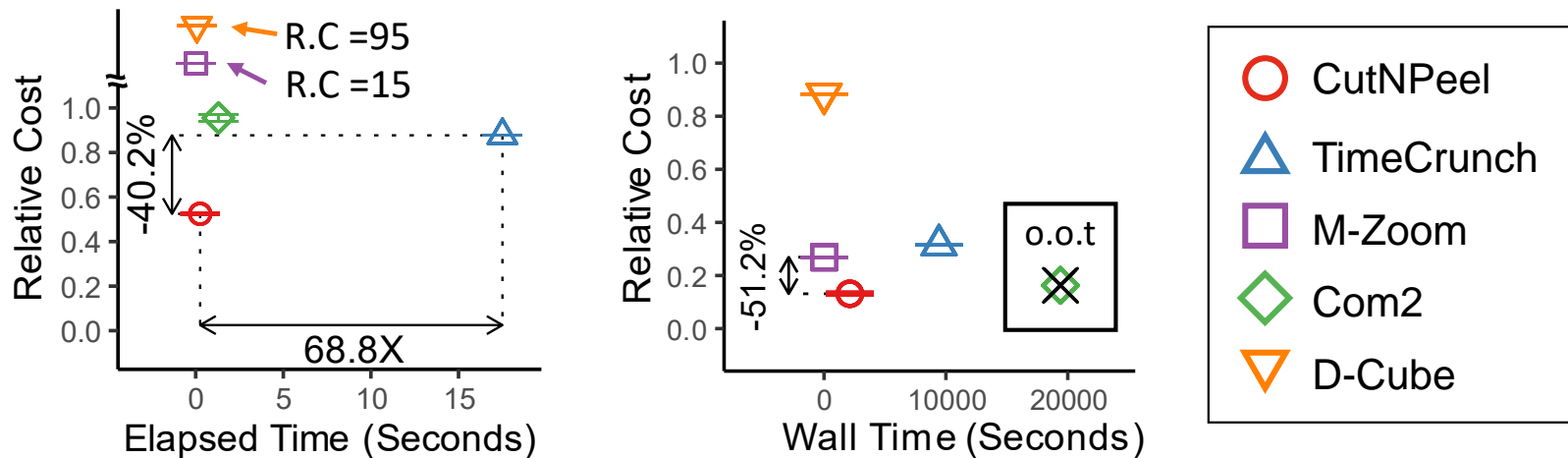
- Q1. Evaluate the **speed** and the **quality** of near bi-cliques
- Quality evaluation metric

$$\textit{Relative Cost} := \frac{\textit{Cost for describing } \mathbf{G} \textit{ using } \mathbf{B}, \mathbf{R} \textit{ and } \mathbf{M}}{\textit{Cost for describing edges in } \mathbf{G} \textit{ individually}}$$

- Methods compared :
 - CutNPeel (Proposed)
 - TimeCrunch
 - M-Zoom
 - Com2
 - D-Cube

EXP1. Search Quality & Speed

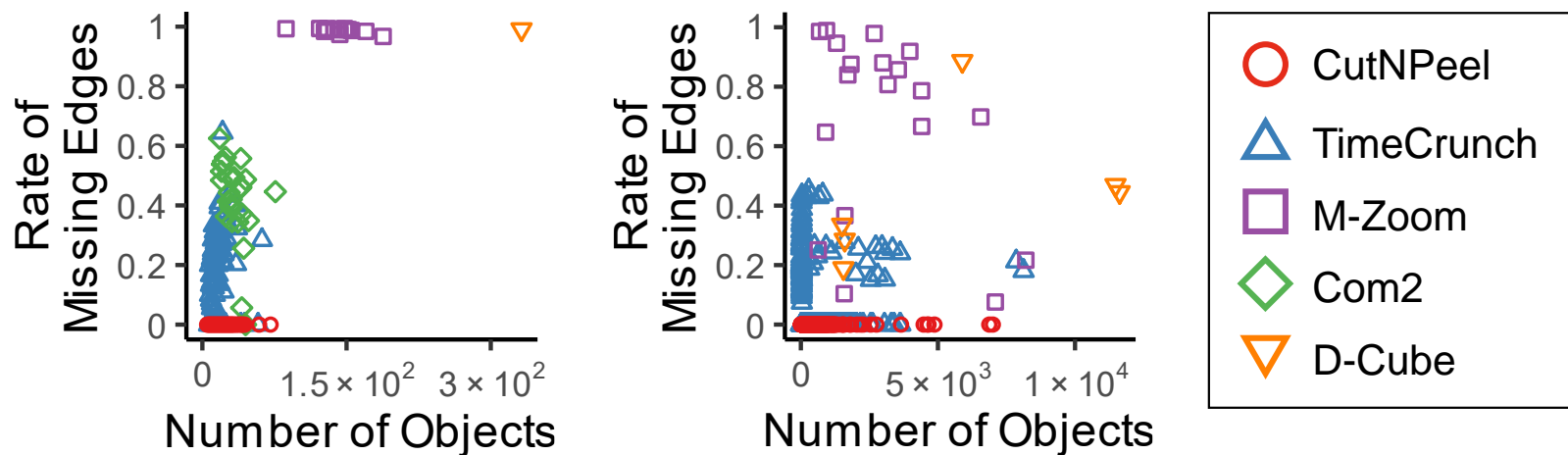
- Six real-world datasets
 - Enron (left) : Sender X Receiver X Time
 - DDoS (right) : Source IP X Destination IP X Time



- **CutNPeel** achieves up to 51.2% better quality up to 68.8X faster compared to the competitors

EXP1. Search Quality & Speed

- Six real-world datasets
 - Enron (left) : Sender X Receiver X Time
 - DDoS (right) : Source IP X Destination IP X Time



- The near bi-cliques detected by **CutNPeel** tend to be precise with smaller ratios of missing edges

EXP1. Search Quality & Speed

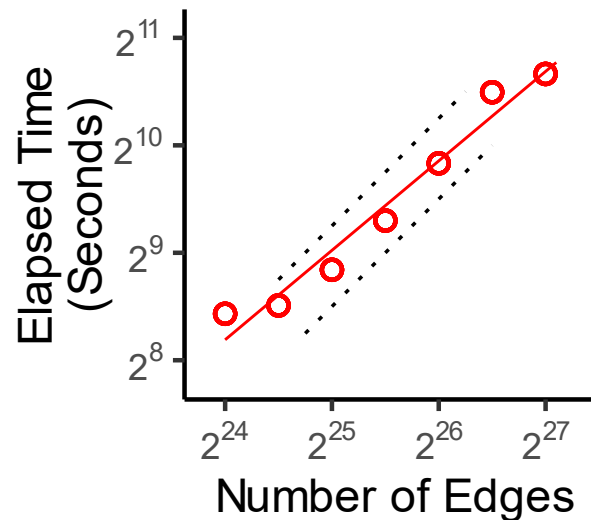
- **CutNPeel** achieves up to 51.2% better quality up to 68.8X faster compared to the competitors
- **CutNPeel** detects subgraphs close to bi-cliques

- High Quality & Speed:** quickly detects high-quality near bi-cliques
- Scalability:** runs in near-linear time
- Applicability:** can be applied to compression and pattern discovery

EXP2. Scalability

- Q2. Measure the **scalability** w.r.t. the size of G
- **CutNPeel** scales almost linearly

○ Actual Running Time ■ ■ ■ Linear Increase



EXP2. Scalability

- High Quality & Speed:** quickly detects high-quality near bi-cliques
- Scalability:** runs in near-linear time
- Applicability:** can be applied to compression and pattern discovery

EXP3. Application: Compression

- Q3. Evaluate the **compression** ability of CutNPeel
- **CutNPeel** achieved the best compression in all considered datasets

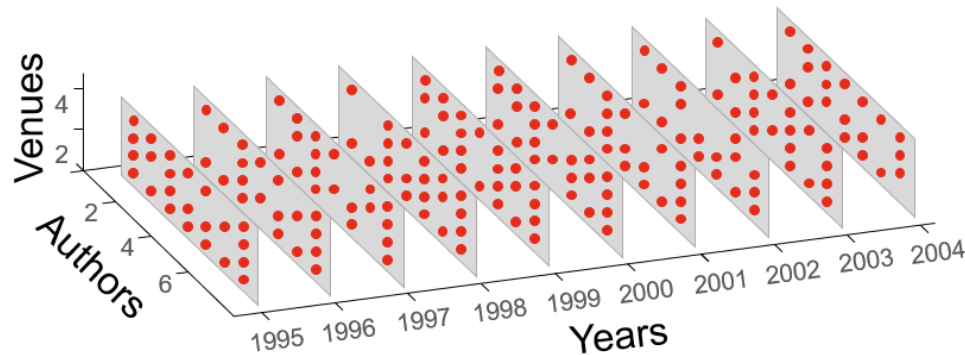
Metric*	Method	Compression Rates in % (the lower, the better)					
		Enron	Darpa	DDoS	DBLP	Yelp	Weeplaces
CutNPeel	CutNPeel	<u>52.5</u>	<u>32.3</u>	<u>13.1</u>	<u>53.7</u>	77.1	<u>55.9</u>
	TimeCrunch	87.7	37.2	31.5	74.2	100.6	o.o.t.
	Com2	95.5	100	o.o.t.	99.9	o.o.t.	o.o.t.
TimeCrunch**	TimeCrunch	86.3	36.7	16.3	79.7	100.0	o.o.t.
Com2	CutNPeel	48.3	27.1	8.0	52.4	<u>78.0</u>	52.4
	TimeCrunch	79.3	34.6	22.8	63.9	100.2	o.o.t.
	Com2	58.3	202.8	o.o.t.	57.5	o.o.t.	o.o.t.

*Different metrics are based on different encoding methods

**The encoding method used in TimeCrunch is not applicable to CutNPeel and Com2

EXP4. Application: Pattern Mining

- **CutNPeel** detects some interesting patterns on the DBLP and DDoS datasets



- ✓ 7 researchers
- ✓ who presented at the 4 same venues
- ✓ from 10 consecutive years from 1995 to 2004

EXP3&4. Application

- ✓ **High Quality & Speed:** quickly detects high-quality near bi-cliques
- ✓ **Scalability:** runs in near-linear time
- ✓ **Applicability:** can be applied to compression and pattern discovery

Conclusion

Considering **the problem of finding a concise, precise and exhaustive set of near bi-cliques** in a dynamic graph,

- We formulated the above problem as an **optimization problem** based on the **MDL principle**
- To solve this problem, we designed **CutNPeel** with the following strengths:

- ✓ **High Quality & Speed**
- ✓ **Scalability**
- ✓ **Applicability**

- Github Link: <https://github.com/hyeonjeong1/cutnpeel>

Finding a Concise, Precise, and Exhaustive Set of Near Bi-Cliques in Dynamic Graphs

Hyeonjeong Shin¹, Taehyung Kwon¹, Neil Shah², Kijung Shin¹

¹ KAIST, ² Snap Inc.

