

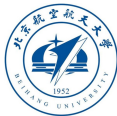
# HyperBlocker: Accelerating Rule-Based Blocking in Entity Resolution Using GPUs

Research track Information Integration and Data Quality III

Xiaoke Zhu

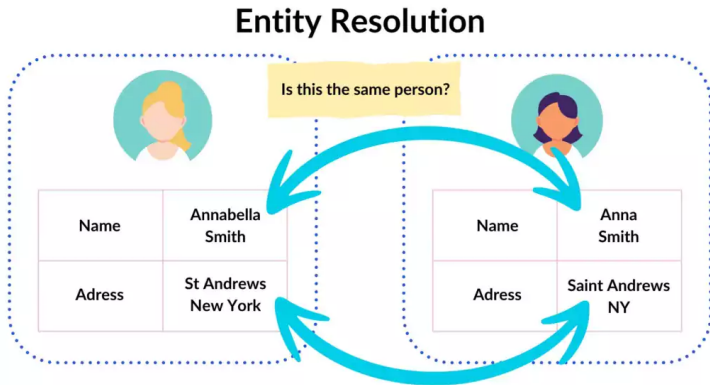
(Beihang University → Shenzhen Institute of Computing Sciences)

Co-authors: Min Xie (Shenzhen Institute of Computing Sciences),  
Ting Deng (Beihang University), Qi Zhang (Meta Platforms)



# What is Entity Resolution?

Problem of identifying and linking/grouping different manifestations of the same real world object.



# What is Entity Resolution?

Problem of identifying and linking/grouping different manifestations of the same real world object.

## What Caused the ER Problem?

- Name/Attribute ambiguity
- Errors due to data entry
- Data integration
- ...

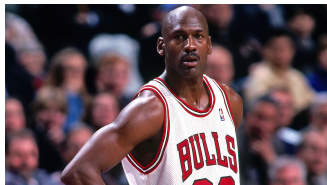


Figure 1: Basketball Player



Figure 2: Professor at UC, Berkeley

## What is Entity Resolution?

Problem of identifying and linking/grouping different manifestations of the same real world object.

### What Caused the ER Problem?

- Name/Attribute ambiguity
- Errors due to data entry
- Data integration
- ...



# What is Entity Resolution?

Problem of identifying and linking/grouping different manifestations of the same real world object.

## What Caused the ER Problem?

- Name/Attribute ambiguity
- Errors due to data entry
- Data integration
- ...



# BigData ER Challenges

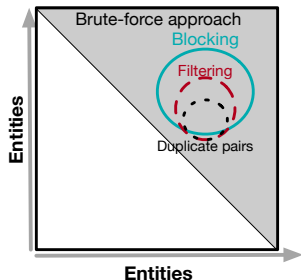
ER suffers from a quadratic time complexity,  $O(n^2)$

1000 business listings each from 1,000 different cities across the world

- 1 trillion comparisons
- 11 days (if each comparison is 1 us)

## Candidate selection step

- Blocking (e.g., hashing/clustering): groups similar entities into blocks.
- Filtering (e.g., similarity join): removes dissimilar entity pairs while potentially keeping some non-matching pairs.



## Rule-based BLOCKING

eid	pno	pname	price	sname	description	color	saddress
$e_1$	$t_1$	Apple Mac Air	\$909	Comp. World	Apple MacBook Air (13-inch, 8GB RAM, 256GB SSD)	Ashen	9 Barton Grove, McCulloughmouth
$e_2$	$t_2$	ThinkPad	-	Smith's Tech	ThinkPad E15, 15.6-inch full HD IPS display, Intel Core i5-1235U processor, (16GB) RAM   512GB PCIe SSD)	Gray	Seg Plaza, Hua qiang North Road
$e_2$	$t_3$	ThinkPad	\$849	Smith's Tech	Lenovo E15 Business ThinkPad, 15.6-inch full HD IPS display, 12 generation Intel Core i5, 16GB RAM, 512GB SSD	Gray	Seg Plaza, Hua qiang North Road
$e_1$	$t_4$	MacBook Air	\$909	Comp. World	Apple 2022 MacBook Air M2 chip 13-inch, 8 GB RAM, 256 GB SSD storage gray	Gray	-
$e_1$	$t_5$	MacBook Air	\$909	Comp. World	-	Gray	Barton Grove, McCulloughmouth

- $\varphi_1 : t.\text{color} = s.\text{color} \wedge t.\text{price} = s.\text{price} \wedge t.\text{sname} = s.\text{sname} \wedge t.\text{pname} \approx_{\text{ED}} s.\text{pname} \rightarrow t.\text{eid} = s.\text{eid}$
- $\varphi_2 : t.\text{sname} = s.\text{sname} \wedge t.\text{description} \approx_{\text{JD}} s.\text{description} \rightarrow t.\text{eid} = s.\text{eid}$
- $\varphi_3 : t.\text{saddress} \approx_{\text{ED}} s.\text{saddress} \wedge t.\text{description} \approx_{\text{JD}} s.\text{description} \rightarrow t.\text{eid} = s.\text{eid}$

## Rule-based BLOCKING

eid	pno	pname	price	sname	description	color	saddress
$e_1$	$t_1$	Apple Mac Air	\$909	Comp. World	Apple MacBook Air (13-inch, 8GB RAM, 256GB SSD)	Ashen	9 Barton Grove, McCulloughmouth
$e_2$	$t_2$	ThinkPad	-	Smith's Tech	ThinkPad E15, 15.6-inch full HD IPS display, Intel Core i5-1235U processor, (16GB) RAM   512GB PCIe SSD	Gray	Seg Plaza, Hua qiang North Road
$e_2$	$t_3$	ThinkPad	\$849	Smith's Tech	Lenovo E15 Business ThinkPad, 15.6-inch full HD IPS display, 12 generation Intel Core i5, 16GB RAM, 512GB SSD	Gray	Seg Plaza, Hua qiang North Road
$e_1$	$t_4$	MacBook Air	\$909	Comp. World	Apple 2022 MacBook Air M2 chip 13-inch, 8 GB RAM, 256 GB SSD storage gray	Gray	-
$e_1$	$t_5$	MacBook Air	\$909	Comp. World	-	Gray	Barton Grove, McCulloughmouth

- $\varphi_1 : t.\text{color} = s.\text{color} \wedge t.\text{price} = s.\text{price} \wedge t.\text{sname} = s.\text{sname} \wedge t.\text{pname} \approx_{\text{ED}} s.\text{pname} \rightarrow t.\text{eid} = s.\text{eid}$ 
  - $t_4$  is a potential match for  $t_5$
- $\varphi_2 : t.\text{sname} = s.\text{sname} \wedge t.\text{description} \approx_{\text{JD}} s.\text{description} \rightarrow t.\text{eid} = s.\text{eid}$ 
  - $t_1$  is a potential match for  $t_4$
- $\varphi_3 : t.\text{saddress} \approx_{\text{ED}} s.\text{saddress} \wedge t.\text{description} \approx_{\text{JD}} s.\text{description} \rightarrow t.\text{eid} = s.\text{eid}$ 
  - $t_2$  is a potential match for  $t_3$



## Entity Resolution on 10+ Million Tuples

**Dataset:** TFACC: 10 million tuples, each described by 16 attributes

### Enviroments:

A cluster of 30 HPC servers, powered with 2.40GHz Intel Xeon Gold CPU  
1 × Nvidia Tesla V100 GPU

### Summary:

Solution	Execution Time	Characteristics
Dedoop	>3 hours	Distributed implementation (30-nodes)
DisDedup	>3 hours	Distributed implementation (30-nodes)
SparkER	>3 hours	Distributed implementation (30-nodes)
Faiss + Ditto	Out of Memory	Similarity Join + Matcher (GPU)
DeepBlocker + Ditto	Out of Memory	Learned Blocker + Matcher (GPU, SOTA)

**Existing methods cannot process 10M+ tuples, even on a 30-node cluster**

## Our Contributions

- The data/rule-aware execution plan for BLOCKING
  - Tree representation of execution plan
  - Cost models for User Defined Function(UDFs)
- The GPU hardware-aware parallelism for BLOCKING
  - Task stealing + parallel sliding window
- Multi-GPUs collaboration

VS Learned Blocker  
(SOTA in Accuracy)



Accuracy  
 $\approx$

VS Learned Blocker  
(SOTA in Accuracy)



Space  
 $\approx 75\%$

VS 30-Node Distributed Blocker  
& GPU-Accelerated Blocker  
(SOTA in Efficiency)



Speed  
80% - 9.1X

## Observation 1

Not all predicates entail equivalent satisfaction or computational costs

eid	pno	pname	price	sname	description	color	saddress
e <sub>1</sub>	t <sub>1</sub>	Apple Mac Air	\$909	Comp. World	Apple MacBook Air (13-inch, 8GB RAM, 256GB SSD)	Ashen	9 Barton Grove, McCulloughmouth
e <sub>2</sub>	t <sub>2</sub>	ThinkPad	-	Smith's Tech	ThinkPad E15, 15.6-inch full HD IPS display, Intel Core i5-1235U processor, (16GB) RAM   512GB PCIe SSD)	Gray	Seg Plaza, Hua qiang North Road
e <sub>2</sub>	t <sub>3</sub>	ThinkPad	\$849	Smith's Tech	Lenovo E15 Business ThinkPad, 15.6-inch full HD IPS display, 12 generation Intel Core i5, 16GB RAM, 512GB SSD	Gray	Seg Plaza, Hua qiang North Road
e <sub>1</sub>	t <sub>4</sub>	MacBook Air	\$909	Comp. World	Apple 2022 MacBook Air M2 chip 13-inch, 8 GB RAM, 256 GB SSD storage gray	Gray	-
e <sub>1</sub>	t <sub>5</sub>	MacBook Air	\$909	Comp. World	-	Gray	Barton Grove, McCulloughmouth

**More discriminative in the “description” attribute**

“description”: Long-text field containing discriminative features for entity resolution

## Observation 1

Not all predicates entail equivalent satisfaction or computational costs

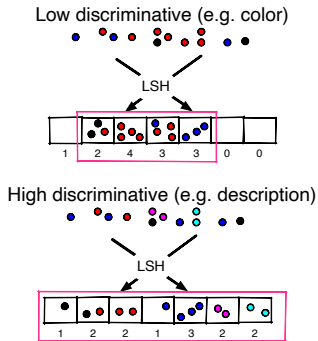
eid	pno	pname	price	sname	description	color	saddress
$e_1$	$t_1$	Apple Mac Air	\$909	Comp. World	Apple MacBook Air (13-inch, 8GB RAM, 256GB SSD)	Ashen	9 Barton Grove, McCulloughmouth
$e_2$	$t_2$	ThinkPad	-	Smith's Tech	ThinkPad E15, 15.6-inch full HD IPS display, Intel Core i5-1235U processor, (16GB) RAM   512GB PCIe SSD)	Gray	Seg Plaza, Hua qiang North Road
$e_2$	$t_3$	ThinkPad	\$849	Smith's Tech	Lenovo E15 Business ThinkPad, 15.6-inch full HD IPS display, 12 generation Intel Core i5, 16GB RAM, 512GB SSD	Gray	Seg Plaza, Hua qiang North Road
$e_1$	$t_4$	MacBook Air	\$909	Comp. World	Apple 2022 MacBook Air M2 chip 13-inch, 8 GB RAM, 256 GB SSD storage gray	Gray	-
$e_1$	$t_5$	MacBook Air	\$909	Comp. World	-	Gray	Barton Grove, McCulloughmouth

**Higher computational efficiency on the “color” attribute.**

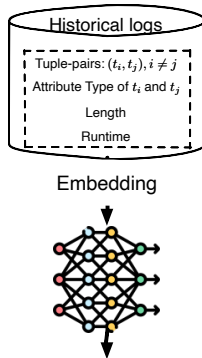
“color”: The computation of short-text similarity necessitates only few operations

## Solution: Query optimization for User Defined Function

Quantifies the likelihood that predicate  $p$  is satisfied on attribute  $a$  of dataset  $D$



Quantifies the efficiency of computing  $p$  on attribute  $a$  of dataset  $D$



Estimated cost of  $p$  on attribute  $a$  of  $D$

A cost model balancing computational cost and selectivity of predicates was proposed for execution plan optimization

## Solution: Query optimization for User Defined Function

Quantifies the likelihood that predicate  $p$  is satisfied when applied to attribute  $a$

- Evaluating selectivity of  $p$  on attribute  $a$  of  $D$

$$sp(p, D) = \text{Norm}\left(\sqrt{\frac{1}{k} \sum_i^k \left(b_i - \frac{|D|}{k}\right)^2}\right)$$

Quantifies the efficiency of computing  $p$  on attribute  $a$

- Cost of predicate  $p$  on attribute  $a$  of  $D$

$$\text{cost}_a(p, D) = \sum_{(t_1, t_2) \in D \times D} T_p(t_1, t_2)$$

- Estimated cost

$$\hat{\text{cost}}_a(p, D) = \text{Norm}\left(\sum_{(t_1, t_2) \in D \times D} \mathcal{N}(p, t_1, t_2)\right)$$

**cost-effectiveness mode**

$$\frac{1 - sp(p, D)}{\hat{\text{cost}}_a(p, D)}$$

## Observation 2 & Solution

### Redundant Computations

$$\varphi_1 : t.\text{color} = s.\text{color} \wedge t.\text{price} = s.\text{price} \wedge t.\text{sname} = s.\text{sname} \wedge t.\text{pname} \approx_{ED} s.\text{pname} \rightarrow t.\text{eid} = s.\text{eid}$$

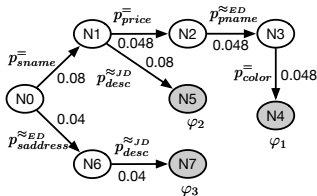
$$\varphi_2 : t.\text{sname} = s.\text{sname} \wedge t.\text{description} \approx_{JD} s.\text{description} \rightarrow t.\text{eid} = s.\text{eid}$$

$$\varphi_3 : t.\text{saddress} \approx_{ED} s.\text{saddress} \wedge t.\text{description} \approx_{JD} s.\text{description} \rightarrow t.\text{eid} = s.\text{eid}$$

### Solution: Execution Tree Representation

$[p_{\text{sname}}^{\approx_{ED}}, p_{\text{price}}^{\approx_{ED}}, p_{\text{saddress}}^{\approx_{ED}}, p_{\text{pname}}^{\approx_{ED}}, p_{\text{desc}}^{\approx_{JD}}, p_{\text{color}}^{\approx_{JD}}]$

(a) order of predicates



(b) Execution Tree

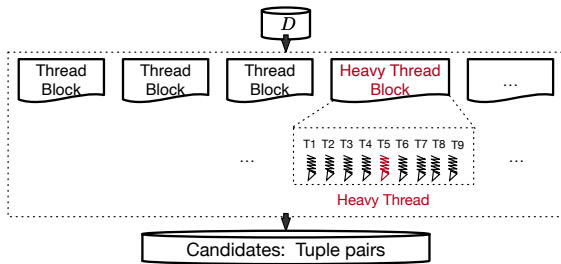
Attribute index      Function pointer

t.sname=s.sname	3	=
t.price=s.price	2	=
t.pname≈ <sub>ED</sub> s.pname	1	≈ <sub>ED</sub>
t.color=s.color	5	=
CP	CP	CP
t.desc≈ <sub>JD</sub> s.desc	4	≈ <sub>JD</sub>
CP	CP	CP
t.saddress≈ <sub>ED</sub> s.saddress	6	≈ <sub>ED</sub>
t.desc≈ <sub>JD</sub> s.desc	4	≈ <sub>JD</sub>
CP	CP	CP

(c) Sequential execution path

## Observation 3

GPU thread workload imbalance occurs during rule deduction on  $D$



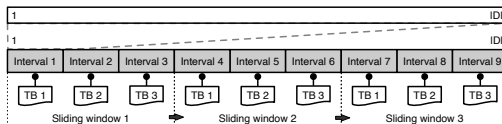
## Unique Challenges for Rule-Based Blocking

- Early termination for a thread (first rule suffice) vs. full traversal for another thread (no early witness).
- Even for the same predicate, the evaluation time on different tuples is different (e.g., long vs. short text)

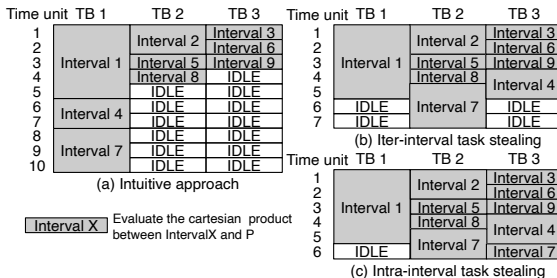


# Solutions

## Parallel Sliding Window



## Task Stealing



# Experimental

## Dataset

Dataset	Domain	#Tuples	Max #Pairs	#GT Pairs	#Attrs	#Rules	#Partitions
Fodors-Zagat	restaurant	866	$1.8 \times 10^4$	112	6	1	1
DBLP-ACM	citation	4591	$6.0 \times 10^6$	2294	4	10	8
DBLP-Scholar	citation	66881	$1.7 \times 10^8$	5348	4	10	8
IMDB	movie	1.5M	$8.1 \times 10^{10}$	0.2M+	6	10	128
Songs	music	0.5M	$2.7 \times 10^{11}$	1.2M	8	10	128
NCV	vote	2M	$1.0 \times 10^{12}$	0.5M+	5	10	512
TPCH	synthetic	4M	$1.6 \times 10^{13}$	#	8	30	512
TFACC	traffic	10M	$1.0 \times 10^{14}$	#	16	50	1024
TFACC <sub>large</sub>	traffic	36M	$1.3 \times 10^{15}$	#	16	50	1024

## Baselines

- Distributed on 30 nodes: SparkER[EDBT'19], DisDedup[VLDB'16], Dedoop[VLDB'12]
- GPU-based: DeepBlocker[VLDB'21], GPUdet[HPIFuture'13]
- Matcher: Ditto[VLDB'20]
- Variants: HyperBlocker<sub>noEPG</sub>: without EPG, HyperBlocker<sub>HO</sub>: disables all hardware optimizations, HyperBlocker<sub>Ditto</sub>: HyperBlocker as blocker and Ditto as the matcher

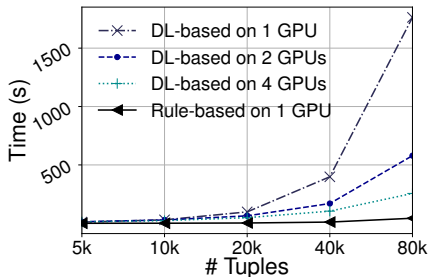
# HyperBlocker is Highly Efficient and Scalable

Achieves speeds  $6.818\times$  faster than baselines, with comparable accuracy

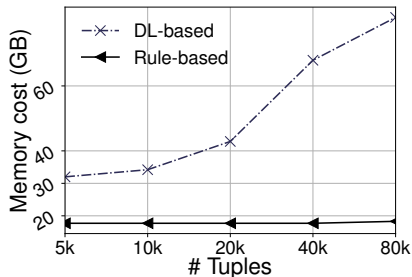
Method	DBLP-ACM		IMDB		Songs		NCV	
	F1-score	Time (s)	F1-score	Time (s)	F1-score	Time (s)	F1-score	Time (s)
SparkER	0.77 (-0.17)	11.0 (13.8 $\times$ )	0.31 (-0.65)	242.9 (6.8 $\times$ )	0.08 (-0.72)	203.4 (15.2 $\times$ )	0.26 (-0.66)	229.3 (49.8 $\times$ )
GPUDet	0.92 (-0.02)	20.1 (25.1 $\times$ )	0.94 (-0.02)	323.8 (9.1 $\times$ )	0.80 (+0)	404.8 (30.2 $\times$ )	0.90 (-0.02)	1252.6 (272.3 $\times$ )
DeepBlocker	0.98 (+0.04)	8.3 (10.4 $\times$ )	/	>3h	/	>3h	/	>3h
HyperBlocker <sub>noEPG</sub>	0.94 (+0)	9.9 (12.4 $\times$ )	/	>3h	0.80 (+0)	1904.1 (142 $\times$ )	0.92 (+0)	2408.6 (523.6 $\times$ )
HyperBlocker <sub>noHO</sub>	0.94 (+0)	9.5 (11.9 $\times$ )	0.96 (+0)	472.6 (13.2 $\times$ )	0.80 (+0)	45.0 (3.4 $\times$ )	0.92 (+0)	35.9 (7.8 $\times$ )
<b>HyperBlocker</b>	<b>0.94</b>	<b>0.8</b>	<b>0.96</b>	<b>35.7</b>	<b>0.80</b>	<b>13.4</b>	<b>0.92</b>	<b>4.6</b>
Dedoop	0.90 (-0.08)	59.4 (9.4 $\times$ )	0.67 (-0.29)	534.0 (15.0 $\times$ )	0.80 (-0.08)	7643.4 (6.5 $\times$ )	/	>3h
DisDedup	0.45 (-0.53)	94.0 (14.9 $\times$ )	0.67 (-0.29)	644.0 (18.0 $\times$ )	0.06 (-0.82)	917.0 (0.8 $\times$ )	/	>3h
Ditto <sub>top2</sub>	0.98 (+0)	9.0 (1.4 $\times$ )	0.79 (-0.17)	6741.2 (188.8 $\times$ )	0.88 (+0)	2308.6 (2.0 $\times$ )	0.97 (+0.03)	381.8 (2.1 $\times$ )
DeepBlocker <sub>Ditto</sub>	0.99 (+0.01)	12.4 (2.0 $\times$ )	/	>3h	/	>3h	/	>3h
<b>HyperBlocker<sub>Ditto</sub></b>	<b>0.98</b>	<b>6.3</b>	<b>*0.96</b>	<b>*35.7</b>	<b>0.88</b>	<b>1179.0</b>	<b>0.94</b>	<b>180.6</b>

HyperBlocker can scale to large dataset: it process 36M tuples within 3h.

## HyperBlocker Reduces Memory Costs



(a) Execution time



(b) Memory cost on a single GPU

Figure 3: DL-based blocking vs. rule-based blocking

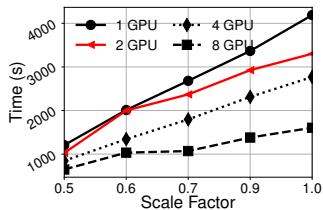
# HyperBlocker Takeaways

- Improves Performance
- Scale up to 36M tuples
- Benefit beyond Blocking in Entity Resolution
- Source code: <https://github.com/hsiaoko/HyperBlocker>

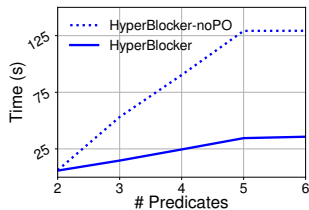
**Thanks**

## **Backup Slides**

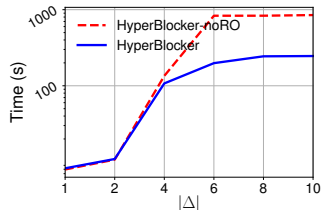
# HyperBlocker is Scalable



(a) TFACC: Varying #GPUs



(b) IMDB: Varying  $|\varphi|$



(c) IMDB: Varying  $|\Delta|$



## EPG VS Existing Query Optimizer

Table 1: Time under different execution plans

Method	Songs	DBLP-ACM	NCV	IMDB
PostgreSQL	40.6 (3.0×)	0.7 (2.3×)	13.5 (2.2×)	13036.0 (449.5×)
TupleX[SIGMOD'21]	14.2 (1.04×)	0.4 (1.3×)	12.4 (2.0×)	36.1 (1.2×)
EPG	13.6	0.3	6.2	29.0

## Before Reachability Querying

Define

Given a graph  $G$ , answer if  $s$  can reach  $t$ , where  $s, t \in V$ .

### High Level Goals

- Index construction time
- Index size
- Query time

### Category

- Label-only: Index size is  $(|V|^2)$ , construction time is  $(|V||E|)$ .
- Label +  $G$ : DFS is needed at run time.
- Label-Free: support dynamic graphs but relative high query time.

## Label + G

### Before Label + G

- $u \rightsquigarrow v$ , if  $Out(v) \subseteq Out(u)$  or  $In(u) \subseteq In(v)$ .
- $u \nrightarrow v$ , if  $Out(v) \not\subseteq Out(u)$  or  $In(u) \not\subseteq In(v)$ .

### Min-wise Independent permutations

- $\min_k\{\pi(\mathcal{A})\} \prec \min_k\{\pi(\mathcal{B})\}$ : if every  $\pi(b) \in \{\min_k\{\pi(\mathcal{B})\} / \min_k\{\pi(\mathcal{A})\}\}$  is greater than the largest number in  $\pi(\mathcal{A})$
- $\mathcal{B} \not\subseteq \mathcal{A}$  if  $\min_k\{\pi(\mathcal{A})\} \succ \min_k\{\pi(\mathcal{B})\}$
- $Out(v) \not\subseteq Out(u)$  if  $\min_k\{\pi(Out(u))\} \succ \min_k\{\pi(Out(v))\}$
- $In(u) \not\subseteq In(v)$  if  $\min_k\{\pi(In(v))\} \succ \min_k\{\pi(In(u))\}$

Reachability querying: An independent permutation labeling approach. In PVLDB 2024.

# Label + G

## Example

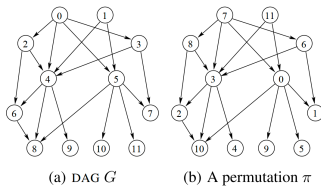


Figure 1: A permutation  $\pi$

Vertex	$\mathcal{L}_{out}$	$\mathcal{L}_{in}$	Vertex	$\mathcal{L}_{out}$	$\mathcal{L}_{in}$
$v_0$	$\{0,1,2,3,4\}$	$\{7\}$	$v_6$	$\{2,10\}$	$\{2,3,6,7,8\}$
$v_1$	$\{0,1,2,3,4\}$	$\{11\}$	$v_7$	$\{1\}$	$\{0,1,6,7,11\}$
$v_2$	$\{2,3,4,8,10\}$	$\{7,8\}$	$v_8$	$\{10\}$	$\{0,2,3,6,7\}$
$v_3$	$\{1,2,3,4,6\}$	$\{6,7\}$	$v_9$	$\{4\}$	$\{3,4,6,7,8\}$
$v_4$	$\{2,3,4,10\}$	$\{3,6,7,8,11\}$	$v_{10}$	$\{9\}$	$\{0,7,9,11\}$
$v_5$	$\{0,1,5,9,10\}$	$\{0,7,11\}$	$v_{11}$	$\{5\}$	$\{0,5,7,11\}$

- $\mathcal{L}_{out} = \min_k \{\pi(\text{Out}(u_i))\}$
- Answer  $u \nrightarrow v$ , if  $\mathcal{L}_{out}(u) \succ \mathcal{L}_{out}(v)$  or  $\mathcal{L}_{in}(v) \succ \mathcal{L}_{out}(u)$
- $\mathcal{Q}(v_2, v_7)$ :  $\mathcal{L}_{out}(v_2) \succ \mathcal{L}_{out}(v_7)$ .
- $\mathcal{Q}(v_3, v_2)$ :  $\mathcal{L}_{out}(v_3) \prec \mathcal{L}_{out}(v_2)$  but  $\mathcal{L}_{in}(v_2) \succ \mathcal{L}_{in}(v_3)$ .
- $\mathcal{Q}(v_1, v_8)$ :  $\mathcal{L}_{out}(v_1) \prec \mathcal{L}_{out}(v_8)$ ,  $\mathcal{L}_{in}(v_1) \succ \mathcal{L}_{in}(v_8)$ . (need DFS)

# Label Free

## Baseline

- Bidirectional BFS.

## Basic Idea

**Input:** The source vertex  $s$ , the destination vertex  $t$ , the teleportation constant  $\alpha$ , the threshold  $\epsilon$

**Output:** Whether  $t$  is reachable from  $s$

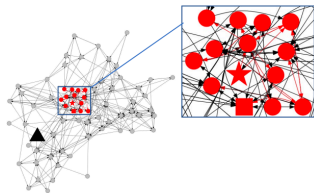
```
1  $\mathbf{r}_s \leftarrow \chi_s$ 
2 while  $\max_{u \in V} \frac{\mathbf{r}_s(u)}{f_{norm}(u)} \geq \epsilon$  do
3   Choose any  $u \in V$  with  $\frac{\mathbf{r}_s(u)}{f_{norm}(u)} \geq \epsilon$ 
4   forall  $u_i \in N_{out}(u)$  do
5     if  $u_i = t$  then
6       return true
7      $\mathbf{r}_s(u_i) \leftarrow \mathbf{r}_s(u_i) + (1 - \alpha) \cdot \frac{\mathbf{r}_s(u)}{f_{dist}(u, u_i)}$ 
8    $\mathbf{r}_s(u) \leftarrow 0$ 
9 return false
```

- PPR with Reachability Querying.
- $r_s$ : The PPR value of vertex  $s$ .

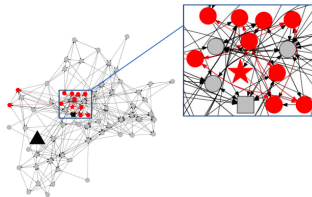
# Label Free

## Optimization

intra-community



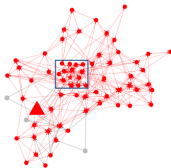
(a) Baseline,  $\epsilon = 3e-2$ , 18 edge accesses



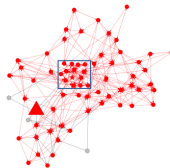
(b) BFS, 18 edge accesses

- ★ source
- intra-community destination
- ▲ inter-community destination

inter-community



(c) Baseline,  $\epsilon = 1e-5$ , 6023 edge accesses



(d) BFS, 244 edge accesses

