

# **An Efficient Similarity Join Algorithm with Cosine Similarity Predicates**

---

Dongjoo Lee, Jaehui Park, Sang-goo Lee

Seoul National University, Korea

Junho Shim

Sookmyung Women's University, Korea

# Agenda

---

- ☐ Introduction
- ☐ Previous Filtering Techniques
- ☐ **MMJoin**
- ☐ Experimental evaluation
- ☐ Extension to Dice and Tanimoto Similarity Joins
- ☐ Conclusion & Future Work

# Similarity Join

---

- Similarity Join
  - Finds all similar pairs of objects
  - $\text{simjoin}(\mathcal{D}, \text{sim}(), t) = \{(x, y, \text{sim}(x, y)) \mid x, y \in \mathcal{D} \wedge \text{sim}(x, y) \geq t\}$ 
    - $\text{sim}()$  : similarity function
    - $t$  : similarity threshold
  
- Similarity functions
  - Jaccard coefficient, Dice coefficient
  - Cosine similarity, Tanimoto coefficient
  - Edit distance (similarity)

# Application Area

- Mining
  - Clustering
  - Bio informatics
- Near Duplicate Detection (NDD)
- Data integration & cleansing



Won Kim Data Mining: Promises, Reality, and Future.  
Won Kim Data Mining: Promises, Reality, and Future - Part 2.  
Won Kim Data Mining: Promises, Reality, and Future - Part 3.

Bon Jovi, Livin' on a player  
Bon Jovi, Living on a player

# Two Issues

---

- How to measure similarity between objects
  - Various similarity functions
  
- **Computation time**
  - $O(n^2)$ ,  $n$  is the total number of objects
  - Efficient algorithms

# Motivation

---

- There are many application domains in which finding similar object pairs is fundamental operation.
- Considering importance of features of each object is also important.
- Weight vector is a good representation for various objects to consider importance of features, but **there is no efficient similarity join algorithm over weight vectors.**
  - ppjoin+ for sets (binary vectors)
  - EdJoin for sequences
  - All-pairs for weight vectors but not efficient

# Contribution

---

- ❑ We propose new similarity upper bounds that is computable with little overhead.
- ❑ We propose an efficient algorithm for cosine similarity joins by exploiting the proposed similarity upper bounds.
- ❑ We extended the algorithm to Dice and Tanimoto similarity joins.
- ❑ We show that our algorithm outperforms a state-of-the-art algorithm using empirical evaluation with large scale datasets.

# Problem Statement

---

- **Self**-similarity join.
  
- Similarity between objects is quantified by **cosine similarity**.
  
- Each object is represented as a vector.
  - Feature (token) universe
    - $\mathcal{U} = \{t_1, t_2, \dots, t_i, \dots, t_m\}$
  - Weight Vector
    - $\vec{x} = \langle x[1], x[2], \dots, x[i], \dots, x[m] \rangle$
    - $\vec{y} = \langle y[1], y[2], \dots, y[i], \dots, y[m] \rangle$
    - $x[i]$  and  $y[i]$  are weights for the  $i$ th token of  $\vec{x}$  and  $\vec{y}$



# Cosine Similarity Join

---

## □ Cosine Similarity between Vectors

$$C(x, y) = \frac{\sum_{i=1}^m x[i] \cdot y[i]}{\sqrt{\sum_{i=1}^m x[i]^2} \sqrt{\sum_{i=1}^m y[i]^2}} = \frac{\text{dot}(x, y)}{\|x\| \|y\|}$$

## □ Vector normalization

- Original vector:  $x_0 = \langle x_0[1], x_0[2], \dots, x_0[i], \dots, x_0[m] \rangle$
- Normalized vector:  $x, x[i] = \frac{x_0[i]}{\|x_0\|}$

## □ For normalized vectors, $C(x, y) = \text{dot}(x, y)$

$$\text{dot}(x, y) = \sum_{i=1}^m x[i] \cdot y[i]$$

# Calculating Dot-Product

We want to find pairs whose cosine similarity is above 0.9 from a dataset.

$\text{simjoin}(\mathcal{D}, C(), 0.9)$

Input

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
o1	0.46	0.31			0.31			0.31		0.62	0.31				0.15
o2	0.46	0.31		0.31	0.31			0.31		0.62					0.15
o3		0.33		0.17		0.33	0.33		0.33	0.50	0.50		0.17		
o4		0.55	0.18	0.18	0.37	0.18		0.37	0.18		0.37		0.18		0.37
o5	0.15		0.30		0.30			0.61				0.46		0.46	

Output

	o1	o2	o3	o4	o5
o1		0.91	0.57	0.57	0.35
o2			0.47	0.51	0.35
o3				0.55	0.00
o4					0.39
o5					

# Calculating Dot-Product

---

- Naïve pair-wise computation.
  - NestedLoopJoin
  
- Using inverted index.
  - InvertedIndexJoin
  - Merits?
    - We can avoid unnecessary computation using filtering techniques.

# NestedLoopJoin

---

---

## Algorithm 1 NestedLoopJoin( $\mathcal{D}, t$ )

---

**Input:** a set of vectors  $\mathcal{D}$ , similarity threshold  $t$

**Output:**  $\{(x, y) | x, y \in \mathcal{D} \text{ and } \text{sim}(x, y) \geq t\}$

```
1:  $O \leftarrow \phi$ 
2: for each  $x \in \mathcal{D}$  do
3:   for each  $y \in \mathcal{D}$  such that  $x \succ y$  do
4:     if  $\text{sim}(x, y) \geq t$  then
5:        $O \leftarrow O \cup \{(x, y)\}$ 
6:     end if
7:   end for
8: end for
9: return  $O$ 
```

---

# InvertedIndexJoin

---

**Algorithm 2** InvertedIndexJoin( $\mathcal{D}$ ,  $t$ )

---

**Input:** a set of vectors  $\mathcal{D}$ , similarity threshold  $t$

**Output:**  $\{(x, y) | x, y \in \mathcal{D} \text{ and } \text{sim}(x, y) \geq t\}$

```
1:  $O \leftarrow \phi$ 
2:  $I_1, \dots, I_m \leftarrow \phi$ 
3: for each  $x \in \mathcal{D}$  do
4:    $C \leftarrow$  empty map from id to weight
5:   for  $i = 1$  to  $m$  such that  $x[i] > 0$  do
6:     for each  $(y, y[i]) \in I_i$  do
7:        $C[y] \leftarrow C[y] + x[i] \cdot y[i]$ 
8:     end for
9:      $I_i \leftarrow I_i \cup \{(x, x[i])\}$ 
10:  end for
11:  for each  $y \in C$  do
12:    if  $C[y] \geq t$  then
13:       $O \leftarrow O \cup \{(x, y)\}$ 
14:    end if
15:  end for
16: end for
17: return  $O$ 
```

---

$$\text{dot}(x, y) = \sum_{i=1}^m x[i] \cdot y[i]$$

# InvertedIndexJoin

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
o1	0.46	0.31			0.31			0.31		0.62	0.31				0.15
o2	0.46	0.31		0.31	0.31			0.31		0.62					0.15
o3		0.33		0.17		0.33	0.33		0.33	0.50	0.50		0.17		
o4		0.55	0.18	0.18	0.37	0.18		0.37	0.18		0.37		0.18		0.37
o5	0.15		0.30		0.30			0.61				0.46		0.46	

	o1	o2	o3	o4	o5
o1		0.91	0.57		
o2			0.47		
o3					
o4					
o5					

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
o1 0.46	o1 0.31		o2 0.31	o1 0.31	o3 0.33	o3 0.33	o1 0.31	o3 0.33	o1 0.62	o1 0.31		o3 0.17		o1 0.15
o2 0.46	o2 0.31		o3 0.17	o2 0.31			o2 0.31		o2 0.62	o3 0.50				o2 0.15
	o3 0.33								o3 0.50					

# Input

# Output

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
o1	0.46	0.31			0.31	?									
o2	0.46	0.31		0.31	0.31										
o3		0.33		0.17											
o4		0.55	0.18	0.18	0.37										
o5	0.15		0.30		0.30			0.61				0.46		0.46	

	o1	o2	o3	o4	o5
o1		0.91	0.57	0.57	0.16
o2			0.47	0.51	0.16
o3				0.55	0.00
o4					0.17
o5					

1 p Incremental computation m

$$\text{dot}(o_1, o_5) = 0.16 + 0.61 \cdot o_1[H] + 0.46 \cdot o_1[L] + 0.46 \cdot o_1[N] < t = 0.9$$

$$\text{dot}(o_2, o_5) = 0.16 + 0.61 \cdot o_2[H] + 0.46 \cdot o_2[L] + 0.46 \cdot o_2[N] < t = 0.9$$

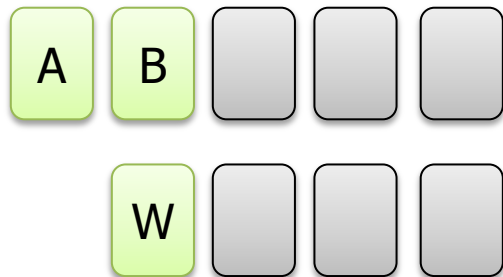
$$\text{dot}(o_3, o_5) = 0.00 + 0.61 \cdot o_3[H] + 0.46 \cdot o_3[L] + 0.46 \cdot o_3[N] < t = 0.9$$

$$\text{dot}(o_4, o_5) = 0.17 + 0.61 \cdot o_4[H] + 0.46 \cdot o_4[L] + 0.46 \cdot o_4[N] < t = 0.9$$

Can we know these? If we can, how?

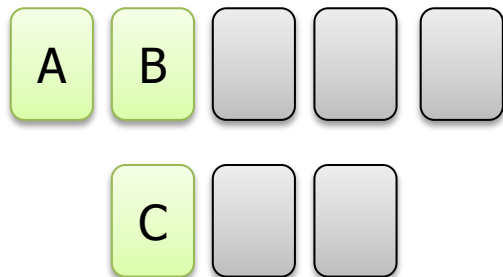
# Card Example

- From Wei Wang. Efficient Exact Similarity Join Algorithms. Seminar at University of Technology, Sydney, Oct 2009.

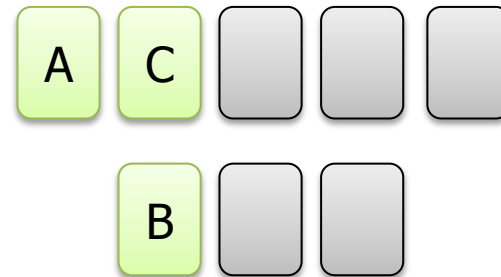


Possible maximum overlap is 3

- If we want to know whether the overlap between two cards lists is more than 4, we do not need to see more cards.
- Other examples



Possible maximum overlap is 3



Possible maximum overlap is 2



# Preliminary: Prefix and Suffix

---

□  $x = \langle x[1], \dots, x[p], x[p+1], \dots, x[m] \rangle$

**Prefix** at  $p : x'_p = \langle x[1], \dots, x[p], 0, \dots, 0 \rangle$

**Suffix** at  $p : x''_p = \langle 0, \dots, 0, x[p+1], \dots, x[m] \rangle$

□ Properties

$$x = x'_p + x''_p$$

$$|x| = |x'_p| + |x''_p|$$

$$\|x\|^2 = \|x'_p\|^2 + \|x''_p\|^2$$

$$\text{dot}(x'_p, y) = \text{dot}(x, y'_p) = \text{dot}(x'_p, y'_p)$$

$$\text{dot}(x'_p, y'_q) = \text{dot}(x'_{\min(p,q)}, y) = \text{dot}(x, y'_{\min(p,q)})$$

# Filtering Principle

- We can avoid unnecessary similarity computation through calculating upper-bounds given **known information** with little overhead.

$$\begin{aligned}\text{dot}(x, y) &= \sum_{i=1}^m x[i] \cdot y[i] \\ &= \sum_{i=1}^p x[i] \cdot y[i] + \sum_{i=p+1}^m x[i] \cdot y[i] \\ &= \text{dot}(x'_p, y'_p) + \text{dot}(x''_p, y''_p) \\ &\leq \text{dot}(x'_p, y'_p) + \text{ubdot}(x''_p, y''_p)\end{aligned}$$

$$\begin{aligned}\text{dot}(x'_p, y'_p) + \text{ubdot}(x''_p, y''_p) < t &\Rightarrow \text{dot}(x, y) < t \\ \text{ubdot}(x, y) < t &\Rightarrow \text{dot}(x, y) < t\end{aligned}$$

Filtering Statements

# Upper-Bounds with Maximum Weight (1)

---

- Let  $\max w_i(\mathcal{D})$  be the maximum weight  $y[i]$  for all  $y$  in  $\mathcal{D}$
- $\mathcal{M}$  is a vector whose  $i$ th weight is  $\max w_i(\mathcal{D})$ , then

$$\begin{aligned}\text{dot}(x_p'', y) &= \sum_{i=p+1}^m x[i] \cdot y[i] \\ &\leq \sum_{i=p+1}^m x[i] \cdot \max w_i(\mathcal{D}) \\ &= \text{dot}(x'', \mathcal{M})\end{aligned}$$

Bayardo et al., WWW 2007

## Upper-Bounds with Maximum Weight (2)

---

- Let  $\text{maxw}(x)$  be the maximum weight  $x[i]$  for all  $i$  over  $1\dots m$ .
- $\mathcal{M}(x)$  is a vector whose  $i$ th weight is  $\min(\text{maxw}(x), \text{maxw}_i(\mathcal{D}))$
- If  $\text{maxw}(x) \geq \text{maxw}(y)$ , then

$$\begin{aligned}\text{dot}(x_p'', y) &= \sum_{i=p+1}^m x[i] \cdot y[i] \\ &\leq \sum_{i=p+1}^m x[i] \cdot \min(\text{maxw}(x), \text{maxw}_i(\mathcal{D})) \\ &= \text{dot}(x'', \mathcal{M}(x))\end{aligned}$$

Bayardo et al., WWW 2007

# Upper-Bounds with Vector Size

---

$$\begin{aligned}\text{dot}(x, y) &= \sum_{i=1}^p x[i] \cdot y[i] + \sum_{i=p+1}^m x[i] \cdot y[i] \\ &\leq \text{dot}(x'_p, y'_p) + \min(|x''_p|, |y''_p|) \cdot \maxw(x''_p) \cdot \maxw(y''_p) \\ &\leq \min(|x|, |y|) \cdot \maxw(x) \cdot \maxw(y)\end{aligned}$$

Bayardo et al., WWW 2007

# Upper-Bounds with Vector Length

---

$$\begin{aligned}\text{dot}(x_p'', y_p'') &= \sum_{i=p+1}^m x[i] \cdot y[i] \\ &\leq \sum_{i=p+1}^m \frac{x[i]^2 + y[i]^2}{2} = \frac{1}{2} \sum_{i=p+1}^m x[i]^2 + \frac{1}{2} \sum_{i=p+1}^m y[i]^2 \\ &= \frac{1}{2} \|x_p''\|^2 + \frac{1}{2} \|y_p''\|^2 \\ &\leq \frac{1}{2} \|x_p''\|^2 + \frac{1}{2}\end{aligned}$$

# Filtering Statements

---

## □ Prefix filtering

$$\text{dot}(x'_p, y) = 0 \wedge \text{dot}(x''_p, \mathcal{M}) < t \Rightarrow \text{dot}(x, y) < t$$

$$\text{dot}(x'_p, y) = 0 \wedge \text{maxw}(x) \geq \text{maxw}(y) \wedge \text{dot}(x'', \mathcal{M}(x)) < t \Rightarrow \text{dot}(x, y) < t$$

$$\text{dot}(x'_p, y) = 0 \wedge \frac{1}{2} \|x''_p\|^2 + \frac{1}{2} < t \Rightarrow \text{dot}(x, y) < t$$

## □ Size filtering

$$\min(|x|, |y|) \cdot \text{maxw}(x) \cdot \text{maxw}(y) < t \Rightarrow \text{dot}(x, y) < t$$

$$|y| \cdot \text{maxw}(x) \cdot \text{maxw}(y) < t \Rightarrow \text{dot}(x, y) < t$$

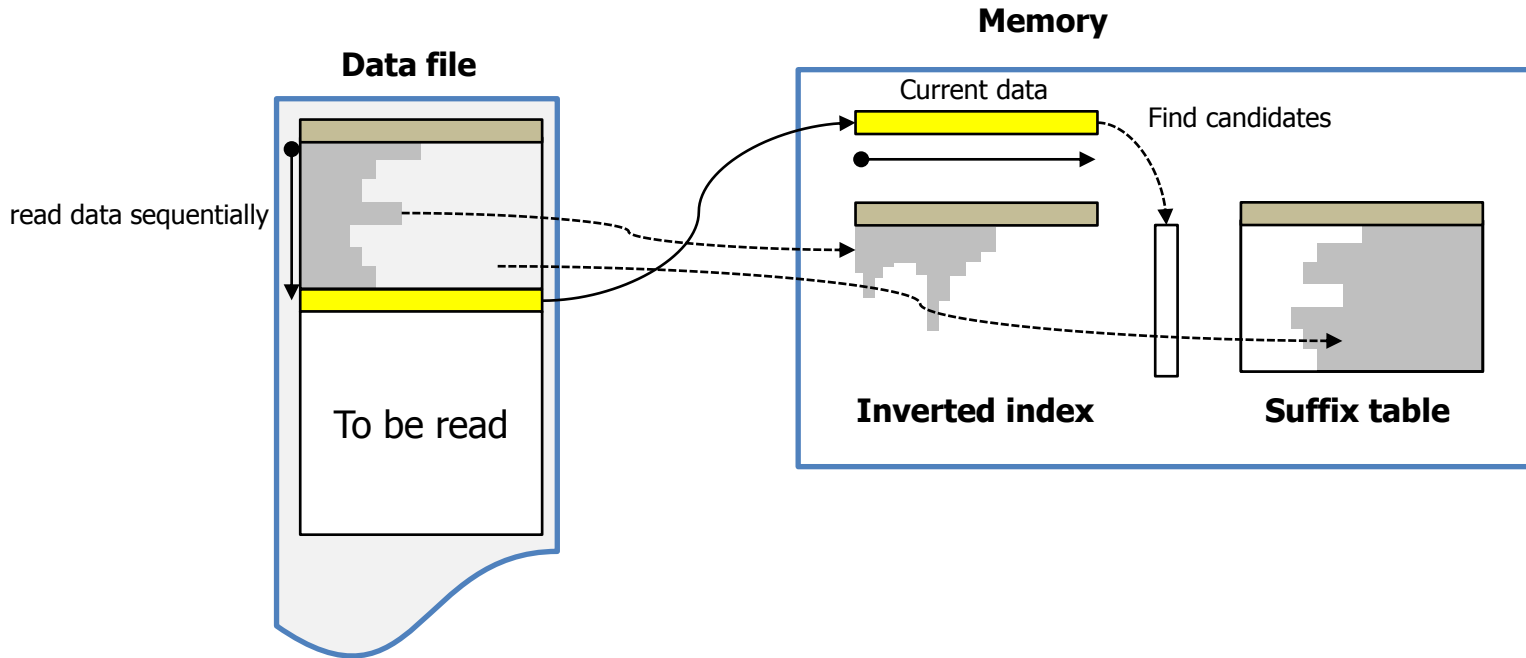
$$|y| \cdot \text{maxw}(y) < \frac{t}{\text{maxw}(x)} \Rightarrow \text{dot}(x, y) < t$$

$$\text{dot}(x, y'_q) + \min(|x|, |y''_q|) \cdot \text{maxw}(x) \cdot \text{maxw}(y''_q) \Rightarrow \text{dot}(x, y) < t$$

## □ Length filtering

$$\text{dot}(x'_p, y'_p) + \frac{1}{2} \|x''_p\|^2 + \frac{1}{2} \|y''_p\|^2 < t \Rightarrow \text{dot}(x, y) < t$$

# General Architecture of Filtering-based Methods



## Preprocessing

- 1) Store  $\max w(x)$  and  $\max w_i(D)$  for all vectors and dimensions
- 2) Sort dimensions as least nonblank entries come first
- 3) Sort vectors in descending order of  $\max w(x)$



# MMJoin

---

---

## Algorithm 2 MMJoin( $\mathcal{D}, t$ )

---

**Input:**  $\mathcal{D} = \{o_1, o_2, \dots, o_n\}$ , similarity threshold  $t$

**Output:**  $\{(x, y, \text{dot}(x, y)) | x, y \in \mathcal{D} \wedge \text{dot}(x, y) \geq t\}$

- 1: Reorder the dimension  $1 \dots m$  such that dimension with the least non-zero entries in  $\mathcal{D}$  appear first
  - 2: Denote the max. of  $x[i]$  over all  $x \in \mathcal{D}$  as  $\text{maxw}_i(\mathcal{D})$
  - 3: Denote the max. of  $x[i]$  for  $1 \dots m$  as  $\text{maxw}(x)$
  - 4:  $O \leftarrow \emptyset$
  - 5:  $I_1, I_2, \dots, I_m \leftarrow \emptyset$
  - 6: **for each**  $x \in \mathcal{D}$  in decreasing order of  $\text{maxw}(x)$  **do**
  - 7:    $C \leftarrow \text{Find-Candidates}(x, I_1, I_2, \dots, I_m, t)$
  - 8:    $O \leftarrow O \cup \text{Verify}(x, C, t)$
  - 9:    $\text{Index}(x, I_1, I_2, \dots, I_m, t)$
  - 10: **return**  $O$
-

# Find-Candidates

---

Algorithm 3 Find-Candidates( $x, I_1, I_2, \dots, I_m, t$ )

---

Input: a vector  $x$ , inverted lists  $I_1, I_2, \dots, I_m$ , similarity threshold  $t$

Output:  $\{(y, \text{dot}(x, y'_q)) | (y, y[i], l_y) \in I_i \wedge \text{dot}(x'_p, y'_q) > 0\}$

```
1:  $C \leftarrow$  empty map from id to weight
2:  $b_1 \leftarrow \sum_{i=1}^m x[i] \cdot \text{maxw}_i(\mathcal{D})$ 
3:  $b_2 \leftarrow 1$ 
4:  $l_x \leftarrow \frac{1}{2}$ 
5:  $\text{minsize} \leftarrow \frac{t}{\text{maxw}(x)}$ 
6: for  $i \leftarrow 1$  to  $m$  such that  $x[i] > 0$  do
7:   Remove  $(y, y[i], l_y)$  from  $I_i$  s.t.  $|y| \cdot \text{maxw}(y) < \text{minsize}$ 
8:    $l_x \leftarrow l_x - \frac{1}{2}x[i]^2$ 
9:   for each  $(y, y[i], l_y) \in I_i$  do
10:    if  $\min(b_1, b_2) \geq t$  or  $C[y] > 0$  then
11:       $C[y] \leftarrow C[y] + x[i] \cdot y[i]$ 
12:      if  $C[y] + l_x + l_y < t$  then
13:         $C[y] \leftarrow 0$ 
14:    $b_1 \leftarrow b_1 - x[i] \cdot \text{maxw}_i(\mathcal{D})$ 
15:    $b_2 \leftarrow l_x + \frac{1}{2}$ 
16:    $x.l[i] \leftarrow l_x$ 
17: return  $C$ 
```

---

size filtering

prefix filtering

length filtering

# Index

---

---

**Algorithm 4**  $\text{Index}(x, I_1, I_2, \dots, I_m, t)$ 

---

**Input:** a vector  $x$ , inverted lists  $I_1, I_2, \dots, I_m$ , similarity threshold  $t$

```
1:  $b_1 \leftarrow \sum_{i=1}^m x[i] \cdot \min(\text{maxw}_i(\mathcal{D}), \text{maxw}(x))$ 
2:  $b_2 \leftarrow 1$ 
3:  $l \leftarrow \frac{1}{2}$ 
4: for  $i \leftarrow 1$  to  $m$  such that  $x[i] > 0$  do
5:    $l \leftarrow l - \frac{1}{2}x[i]^2$ 
6:    $I_i \leftarrow I_i \cup \{(x, x[i], l)\}$ 
7:   remove  $x[i]$  and  $x.l[i]$ 
8:    $b_1 \leftarrow b_1 - x[i] \cdot \min(\text{maxw}(x), \text{maxw}_i(\mathcal{D}))$ 
9:    $b_2 \leftarrow b_2 - \frac{1}{2}x[i]^2$ 
10:  if  $\min(b_1, b_2) < t$  then
11:    break;
```

---

# Verify

---

## Algorithm 5 $\text{Verify}(x, C, t)$

---

**Input:** a vector  $x$ , a map from id to weight  $C$ , threshold  $t$

**Output:**  $\{(x, y, \text{dot}(x, y)) \mid C[y] > 0 \wedge \text{dot}(x, y) \geq t\}$

```
1:  $O \leftarrow \emptyset$ 
2: for each  $y \in C$  do
3:   if  $C[y] + \min(|x|, |y''_q|) \cdot \text{maxw}(x) \cdot \text{maxw}(y''_q) \geq t$  then
4:      $\text{unmatched} \leftarrow 0$ 
5:      $d \leftarrow C[y]$ 
6:     for  $i \leftarrow q + 1$  to  $m$  such that  $y[i] > 0$  do
7:       if  $x[i] > 0$  then
8:          $d \leftarrow d + x[i] \cdot y[i]$ 
9:         if  $\text{unmatched} > 1$  then
10:          if  $d + x.l[i] + y.l[i] < t$  then
11:             $d \leftarrow 0$ 
12:            break;
13:           $\text{unmatched} \leftarrow 0$ 
14:       else
15:          $\text{unmatched} \leftarrow \text{unmatched} + 1$ 
16:       if  $d \geq t$  then
17:          $O \leftarrow O \cup \{(x, y, d)\}$ 
18: return  $O$ 
```

size filtering

length filtering

---

# MMJoin

maxw		G	L	N	C	F	I	M	A	D	J	K	O	B	E	H
0.62	o1								0.46		0.62	0.31 0.16	0.15 0.14	0.31 0.10	0.31 0.05	0.31 0.00
0.62	o2								0.46	0.31	0.62		0.15 0.14	0.31 0.10	0.31 0.05	0.31 0.00
0.61	o5		0.46	0.46 0.29	0.30 0.24				0.15 0.23						0.30 0.19	0.61 0.00
0.55	o4				0.18 0.49	0.18 0.47	0.18 0.46	0.18 0.45		0.18 0.44		0.37 0.36	0.37 0.29	0.55 0.14	0.37 0.07	0.37 0.00
0.50	o3	0.33				0.33	0.33	0.17		0.17	0.50	0.50		0.33		
maxw <sub>i</sub>		0.33	0.46	0.46	0.30	0.33	0.33	0.18	0.46	0.31	0.62	0.50	0.37	0.55	0.37	0.61

C[y]	o1	o2	o5	o4	o3
o1		0.91			
o2				0.06	
o5					
o4					
o3					

Candidate map

G	L	N	C	F	I	M	A	D	J	K	O	B	E	H
	o5 0.46 0.39						o1 0.46 0.40	o2 0.31 0.35	o1 0.62 0.20					
							o2 0.46 0.40		o2 0.62 0.16					

Inverted index

$$\text{dot}(x'_D, y'_D) + \frac{1}{2} \|x''_D\|^2 + \frac{1}{2} \|y''_D\|^2 < t \Rightarrow \text{dot}(x, y) < t$$

$$0.06 + 0.44 + 0.35 = 0.85 < 0.9$$

# Experimental Evaluation

---

- Time Performance of Filtering Techniques
  - Average prefix size, candidate size, computation time
  
- Previous filtering techniques (Bayardo et al., WWW 2007)
  - All-Pairs0: Basic prefix filtering
  - All-Pairs1: Exploiting a vector sort order
  - All-Pairs2: Size filtering in the candidate generation phase
  - All-Pairs: Size filtering in the verification phase
  
- Our filtering techniques
  - MMJoin0: Improved prefix selection
  - MMJoin1: Length filtering in the candidate generation phase
  - MMJoin2: Length filtering in the verification phase

# Experimental Setup

---

## □ Environment

- Java 1.6, standard libraries
- 2.83 GHz Intel Core2 Quad, 8 Gbytes of RAM and two 7200 RPM SATA II-IDE hard drives.

## □ Datasets

- DBLP (Bibliographic dataset), TREC filtering dataset (Text REtrieval Conference), LASTFM (Song, title and artist)
- UKBENCH (University of Kentucky, image benchmark dataset)
  - 10200 images (640\*480), 4 picture images for the same scene with different angles or distances
  - Feature vectors are extracted and abstracted as visual words



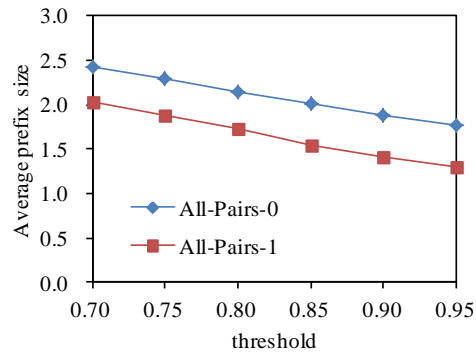
# Experimental Setup

- Tokenization
  - by punctuation.
  - by punctuation and into 4grams.
- Weighting
  - tf-idf weighting scheme
- Dataset Statistics

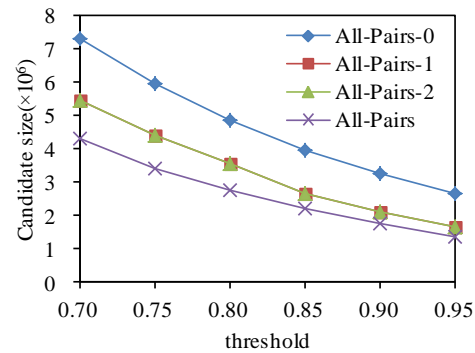
Dataset	$n$	$\text{avg}(\text{size})$	$ U $	$\text{avg}(DF)$
LAST.FM	134,949	4.8	47,295	13.8
LAST.FM 4GRAM		11.2	44,272	34.3
DBLP	1,298,016	8.6	381,450	29.3
DBLP 4GRAM		23.9	135,204	224.5
TREC	348,566	77.1	298,302	90.1
TREC 4GRAM		232.4	121,252	658.8
UKBENCH	10,200	425.7	533,412	6.9



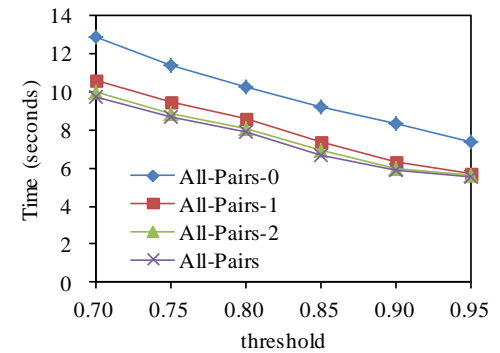
# Effects of Previous Filtering Techniques



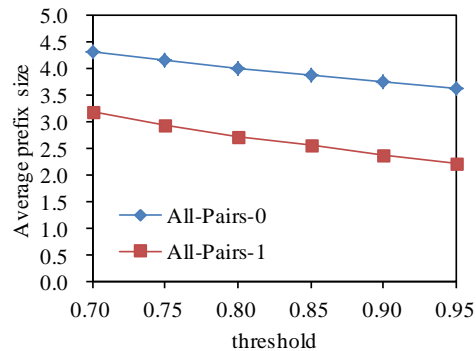
LASTFM, Average prefix size



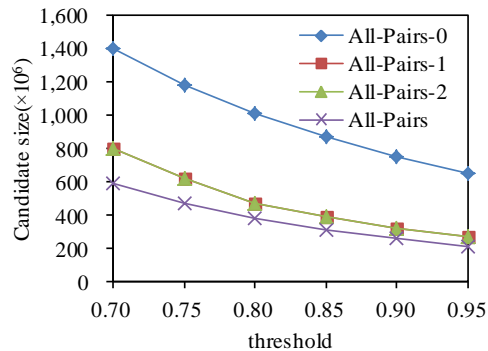
LASTFM, Candidate size



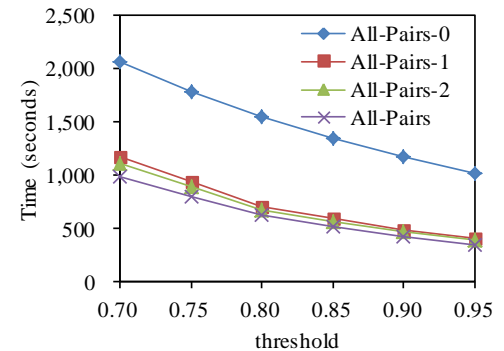
LASTFM, Computation Time



DBLP, Average prefix size

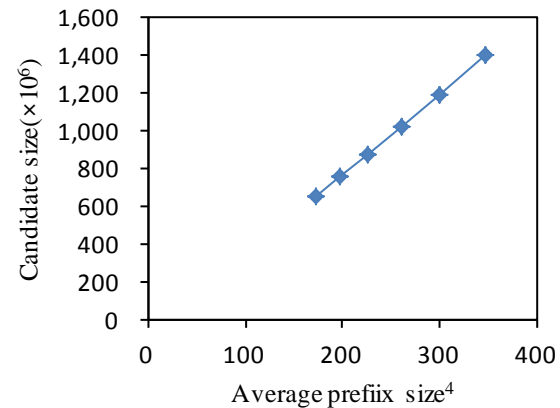
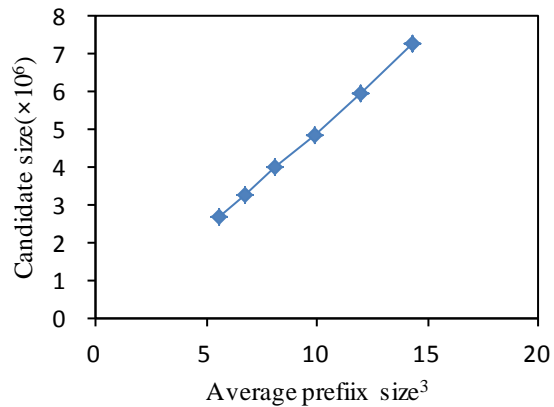


DBLP, Candidate size

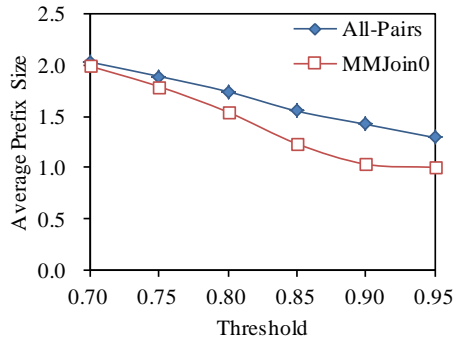


DBLP, Computation Time

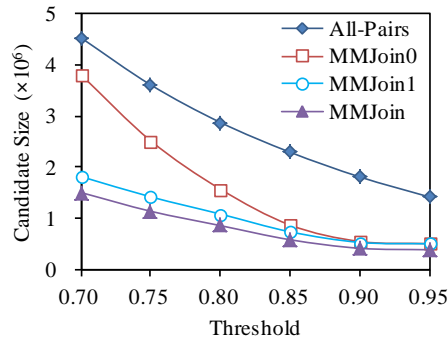
# Relationship between Prefix Size and Candidate Size



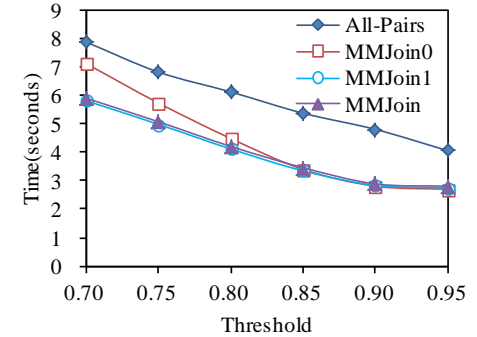
# Effects of Proposed Filtering Techniques (1/3)



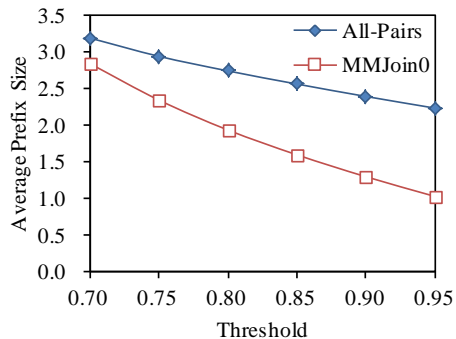
LASTFM, Average prefix size



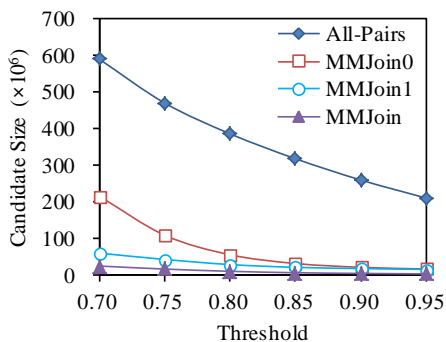
LASTFM, Candidate size



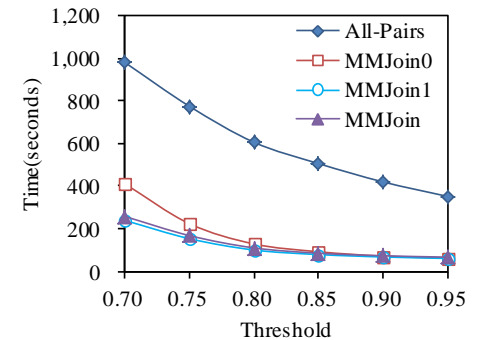
LASTFM, Computation Time



DBLP, Average prefix size

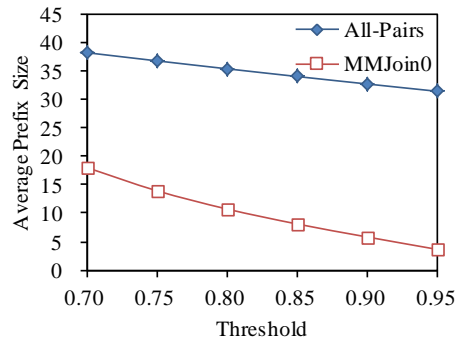


DBLP, Candidate size

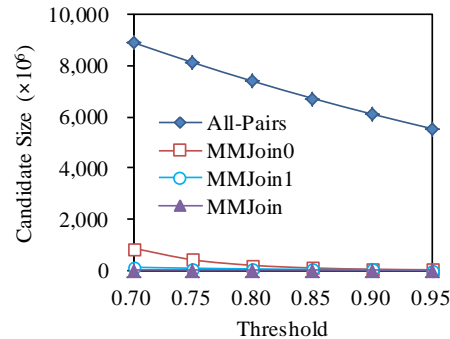


DBLP, Computation Time

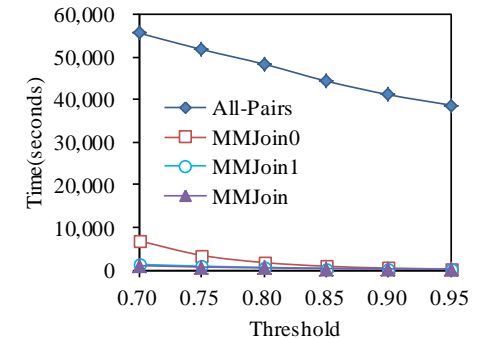
# Effects of Proposed Filtering Techniques (2/3)



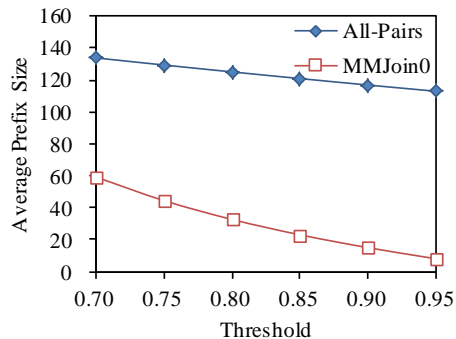
TREC, Average prefix size



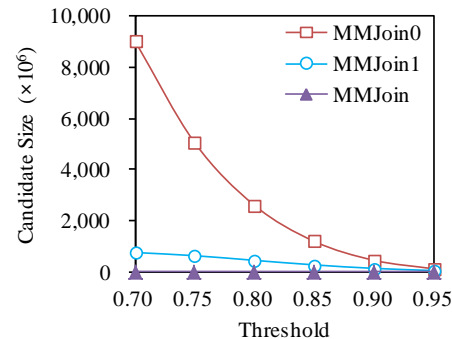
TREC, Candidate size



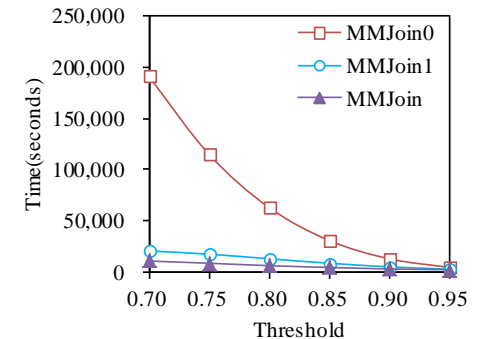
TREC, Computation Time



TREC 4GRAM, Average prefix size

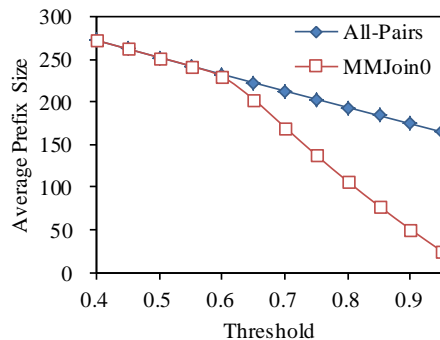


TREC 4GRAM, Candidate size

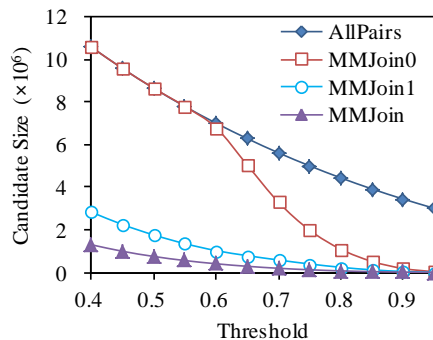


TREC 4GRAM, Computation Time

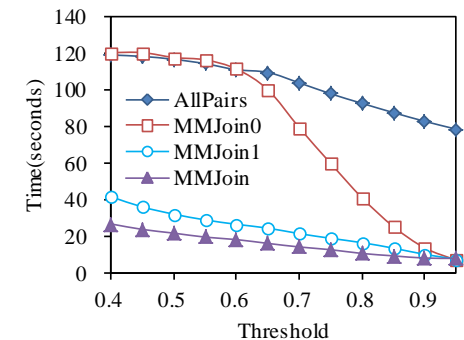
# Effects of Proposed Filtering Techniques (3/3)



UKBENCH, Average prefix size



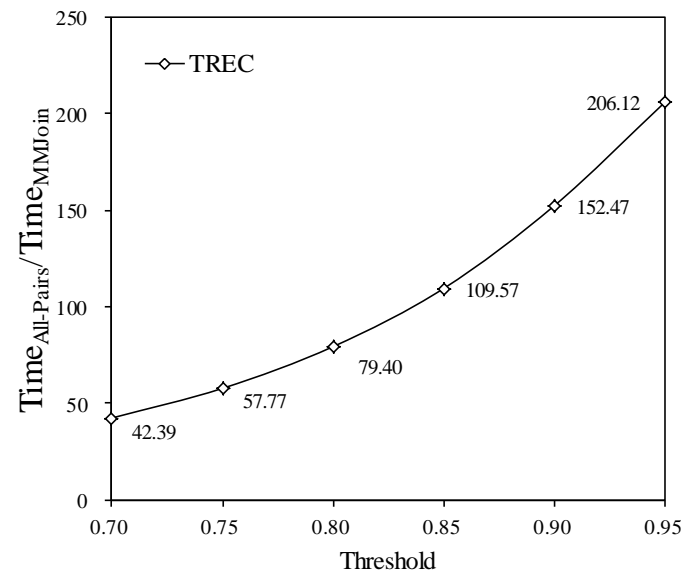
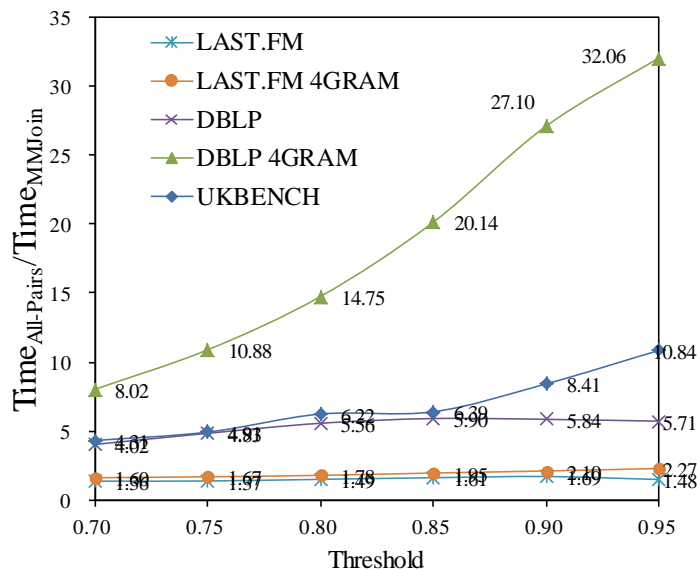
UKBENCH, Candidate size



UKBENCH, Computation Time

# Performance Differences

- MMJoin outperforms All-Pairs in all experimental settings



# Inner Product Family

---

- Inner product

$$\text{dot}(x, y) = \sum_{i=1}^m x[i] \cdot y[i]$$

- Cosine

$$C(x, y) = \frac{\sum_{i=1}^m x[i] \cdot y[i]}{\sqrt{\sum_{i=1}^m x[i]^2} \sqrt{\sum_{i=1}^m y[i]^2}} = \frac{\text{dot}(x, y)}{\|x\| \|y\|}$$

- Tanimoto (extended Jaccard)

$$T(x, y) = \frac{\sum_{i=1}^m x[i] \cdot y[i]}{\sum_{i=1}^m x[i]^2 + \sum_{i=1}^m y[i]^2 - \sum_{i=1}^m x[i] \cdot y[i]} = \frac{\text{dot}(x, y)}{\|x\|^2 + \|y\|^2 - \text{dot}(x, y)}$$

- Dice

$$D(x, y) = \frac{2 \sum_{i=1}^m x[i] \cdot y[i]}{\sum_{i=1}^m x[i]^2 + \sum_{i=1}^m y[i]^2} = \frac{\text{dot}(x, y)}{\frac{1}{2}(\|x\|^2 + \|y\|^2)}$$

# Extension to Dice and Tanimoto Similarity

---

- Tanimoto similarity (Extended Jaccard similarity)

- $$T(x, y) = \frac{\text{dot}(x, y)}{\|x\|^2 + \|y\|^2 - \text{dot}(x, y)}$$

- Dice similarity

- $$D(x, y) = \frac{\text{dot}(x, y)}{\frac{\|x\|^2 + \|y\|^2}{2}}$$



## Observation 1.

---

- Relationship among similarity functions

$$T(x, y) = \frac{\text{dot}(x, y)}{\|x\|^2 + \|y\|^2 - \text{dot}(x, y)} \leq D(x, y) = \frac{\text{dot}(x, y)}{\frac{\|x\|^2 + \|y\|^2}{2}} \leq C(x, y) = \frac{\text{dot}(x, y)}{\|x\| \|y\|}$$

- Cosine similarity constraints is the weakest constraint.

$$T(x, y) \geq t \Rightarrow D(x, y) \geq t \Rightarrow C(x, y) \geq t$$

$$C(x, y) < t \Rightarrow D(x, y) < t \Rightarrow T(x, y) < t$$

## Observation 2.

---

- $x_0$  and  $y_0$  are unnormalized vector, let  $x$  and  $y$  be normalized versions of two vectors, i.e.,  $x[i] = \frac{x_0[i]}{\|x\|}$  and  $y[i] = \frac{y_0[i]}{\|y\|}$ .
- Property of the Dice and Tanimoto similarity
$$T(x_0, y_0) \leq T(x, y)$$
$$D(x_0, y_0) \leq D(x, y) = C(x, y)$$

## Proof of Observation 2.

---

$$T(x_0, y_0) = \frac{\text{dot}(x_0, y_0)}{\|x_0\|^2 + \|y_0\|^2 - \text{dot}(x_0, y_0)}$$

$$= \frac{\frac{\text{dot}(x_0, y_0)}{\|x_0\| \|y_0\|}}{\frac{\|x_0\|^2 + \|y_0\|^2}{\|x_0\| \|y_0\|} - \frac{\text{dot}(x_0, y_0)}{\|x_0\| \|y_0\|}}$$

$$= \frac{\text{dot}(x, y)}{\frac{\|x_0\|^2 + \|y_0\|^2}{\|x_0\| \|y_0\|} - \text{dot}(x, y)}$$

$$\leq \frac{\text{dot}(x, y)}{2 - \text{dot}(x, y)} = T(x, y)$$

$$\frac{\|x_0\|^2 + \|y_0\|^2}{\|x_0\| \|y_0\|} \geq 2$$

# Extended Algorithm

---

---

## Algorithm 6 MMJoin( $\mathcal{D}$ , $\text{sim}()$ , $t$ )

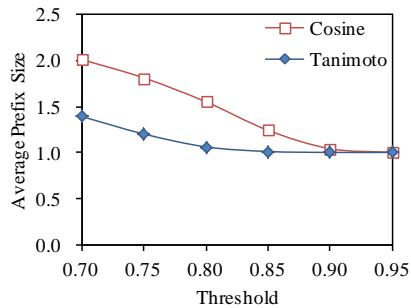
---

**Input:**  $\mathcal{D} = \{o_1, o_2, \dots, o_n\}$ , similarity function  $\text{sim}()$ , similarity threshold  $t$

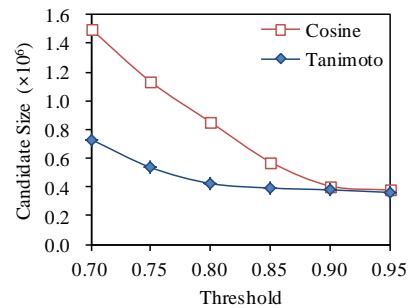
**Output:**  $\{(x_0, y_0, \text{sim}(x_0, y_0)) | x_0, y_0 \in \mathcal{D} \wedge \text{sim}(x_0, y_0) \geq t\}$

- 1: Store the length of each vector and normalize all the vectors
  - 2:  $\alpha \leftarrow \frac{2t}{1+t}$  in case of Tanimoto similarity, otherwise  $\alpha \leftarrow t$
  - 3:  $C \leftarrow \text{MMJoin}(\mathcal{D}, \alpha)$
  - 4:  $O \leftarrow \emptyset$
  - 5: **for each**  $(x, y, \text{dot}(x, y)) \in C$  **do**
  - 6:      $s \leftarrow \text{Calculate}_{\text{sim}}(\|x_0\|, \|y_0\|, \text{dot}(x, y))$
  - 7:     **if**  $s \geq t$  **then**
  - 8:          $O \leftarrow O \cup \{(x_0, y_0, s)\}$
  - 9: **return**  $O$
-

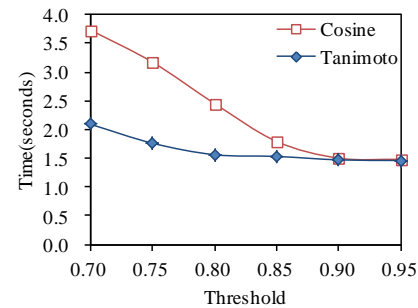
# Experiments



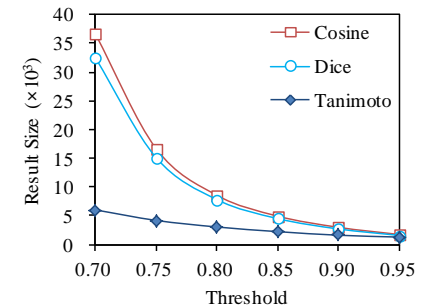
LASTFM, Average prefix size



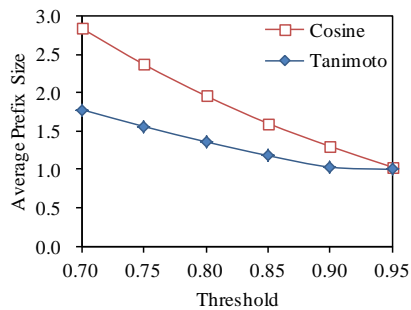
LASTFM, Candidate size



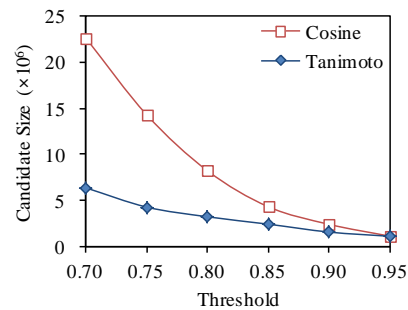
LASTFM, Computation Time



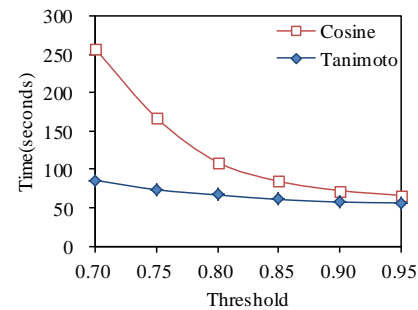
LASTFM, Result Size



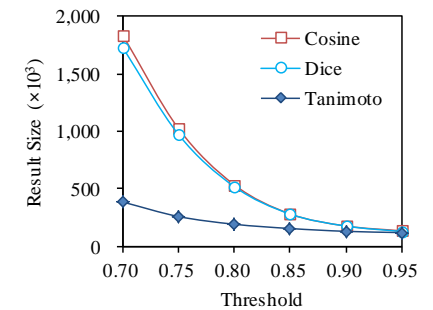
DBLP, Average prefix size



DBLP, Candidate size



DBLP, Computation Time



DBLP, Result Size

# Conclusion

---

- We formalized filtering techniques as **filtering statements** in mathematical way.
- We proposed new cosine similarity upper bounds that is computable with little overhead.
- We proposed an efficient algorithm (MMJoin) for cosine similarity joins by exploiting the proposed similarity upper bounds with little overhead.
- We extended MMJoin to Dice and Tanimoto similarity joins.
- We showed that MMJoin outperforms a state-of-the-art algorithm using empirical evaluation with large scale datasets.

# Future Work

---

- General filtering framework for other similarity functions
  - 8 groups of 38 distance (divergence) and **13 similarity measures** [Cha, 2007]
    - 1) Intersection family
    - 2) Inner Product family
    - 3) Fidelity family (or squared-chord family)
  
- Empirical study on measuring similarity
  - with various feature extraction methods
  - with various weighting schemes
  - with various similarity functions and thresholds