
Set Similarity Join

Dong Deng

Similarity Query Processing: Re-cap

- String/Sequence, Edit Distance (first lecture) - *Partition*
- Set, Overlap (second lecture) - *Subset Enumeration*
- (Weighted) Set, Jaccard/Cosine Similarity and more (this lecture) - *Prefix Filtering*
- Real-value Vectors, Euclidian Distance and more (next lecture) - *Approximation*

Programming Assignment I

Near Duplicate Data



Home > Sports > Top Stories

Federer bests Sampras in exhibition at Madison Square Garden

03/11/2008 | 11:28 AM

[Email this](#) | [Email the Editor](#) | [Print](#) | [Digg this](#) | [Add to del.icio.us](#)

NEW YORK – Past and present stood across the net from each other during a third-set tiebreaker at Madison Square Garden.

On one end, a winded Pete Sampras tried to summon enough energy to give the New York fans another memorable win to talk about it on the subway ride home. On the other side, Roger Federer wore a sly grin like he knew age was about to catch up to the former world No. 1 – the man who owns the record of 14 Grand Slams he wants.

Youth is served, indeed.

Current No. 1 Federer emerged with a 6-3, 6-7 (4), 7-6 (6) victory Monday night in an exhibition that featured a little bit of everything – some laughter, some stellar shots, uneven play and compelling tennis.

There was even a Tiger Woods sighting.

"It was a great night for tennis," Sampras said.

There were moments when, if you squinted a bit, you would have sworn that was the Sampras of old, rather than an old Sampras. There were moments when, if you listened to the whip of

On one end, a winded Pete Sampras tried to summon enough energy to give the New York fans another memorable win to talk about it on the subway ride home. On the other side, Roger Federer wore a sly grin like he knew age was about to catch up to the former world No. 1 – the man who owns the record of 14 Grand Slams he wants.

03/11/2008 | 11:28 AM



Roger Federer beats Pete Sampras in sold-out exhibition at Madison Square Garden

By JAY COHEN, AP Sports Writer
Mar 11, 4:23 am EDT

[Buzz Up](#) [Print](#)

NEW YORK (AP)—Past and present stood across the net from each other during a third-set tiebreaker at Madison Square Garden.

On one end, a winded Pete Sampras tried to summon enough energy to give the New York fans another memorable win to talk about it on the subway ride home. On the other side, Roger Federer wore a sly grin like he knew age was about to catch up to the former world No. 1 – the man who owns the record of 14 Grand Slams he wants.

Youth is served, indeed.
Current No. 1 Federer emerged with a 6-3, 6-7 (4), 7-6 (6) victory Monday night in an exhibition that featured a little bit of everything – some laughter, some stellar shots, uneven play and compelling tennis.

There was even a Tiger Woods sighting.



"It was a great night for tennis," Sampras said.

There were moments when, if you squinted a bit, you would have sworn that was the Sampras of old, rather than an old Sampras. There were moments when, if you listened to the whip of the racket through the air, you would have been absolutely sure Federer was giving it his all.

And then there were moments when, as you watched Sampras throw his racket to the ground in mock disgust or saw Federer raise an index finger to celebrate four

By JAY COHEN, AP Sports Writer
Mar 11, 4:23 am EDT

Fuzzy Matching: A use case :-)

- Find speeches similar to Melania Trump's RNC speech



Fuzz Matching: A use case :-)



**Melania
Trump**

RNC speech | July 18, 2016



From a young age, my parents impressed on me the values that you work hard for what you want in life, that your word is your bond and you do what you say and keep your promise, that you treat people with respect. They taught and showed me values and morals in their daily lives. That is a lesson that I continue to pass along to our son.

And we need to pass those lessons on to the many generations to follow. Because we want our children in this nation to know that the only limit to your achievements is the strength of your dreams and your willingness to work for them."



**Michelle
Obama**

DNC speech | August 25, 2008



And Barack and I were raised with so many of the same values: that you work hard for what you want in life; that your word is your bond and you do what you say you're going to do; that you treat people with dignity and respect, even if you don't know them, and even if you don't agree with them.

And Barack and I set out to build lives guided by these values, and to pass them on to the next generation. Because we want our children -- and all children in this nation -- to know that the only limit to the height of your achievements is the reach of your dreams and your willingness to work for them."

Set Similarity Functions

- Set Similarity Function

- $Jaccard(X, Y) = \frac{|X \cap Y|}{|X \cup Y|}$

- $Cosine(X, Y) = \frac{|X \cap Y|}{\sqrt{|X||Y|}}$

- $Dice(X, Y) = \frac{2|X \cap Y|}{|X| + |Y|}$

$X = \{A, B, C, D, E\}$

$Y = \{B, C, D, E, F\}$

$4/6 = 0.67$

$4/5 = 0.8$

$8/10 = 0.8$

Problem Definition

- Input:
 - A collection of sets R
 - A set similarity function Sim
 - A similarity threshold δ
- Output
 - All pairs $(X, Y) \in R \times R$ such that $Sim(X, Y) \geq \delta$

An Example

- Input:

- R

<i>id</i>	<i>The records</i>	
1	\mathcal{X}_1	$\{x_1, x_2, x_5, x_6, x_7, x_{10}, x_{11}, x_{13}, x_{14}\}$
2	\mathcal{X}_2	$\{x_2, x_4, x_5, x_6, x_9, x_{11}, x_{13}, x_{14}, x_{15}\}$
3	\mathcal{X}_3	$\{x_1, x_3, x_6, x_7, x_9, x_{10}, x_{11}, x_{13}, x_{14}\}$
4	\mathcal{X}_4	$\{x_3, x_4, x_5, x_7, x_8, x_{10}, x_{12}, x_{13}, x_{14}\}$
5	\mathcal{X}_5	$\{x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_{10}, x_{11}, x_{13}, x_{14}\}$

- $Sim = \text{Jaccard Similarity}$

- $\delta = 0.73$

- Output

- $\text{Jaccard}(\mathcal{X}_1, \mathcal{X}_5) = 0.82 \geq \delta$

Applications

- For Web search engines:
 - Perform focused crawling
 - Increase the quality and diversity
 - Identify spams
- For Web mining:
 - Perform document clustering
 - Find replicate Web collections
 - Detect plagiarism

Q. What are the advantages of RAID5 over RAID4?

A. 1. Several write requests could be processed in parallel, since the bottleneck of a **unique** check disk has been eliminated. 2. Read requests have a higher level of parallelism. Since the data is distributed over all disks, read requests involve all disks, whereas in systems with a **dedicated check disk** the check disk never participates in read.

<NUMBER> drew lucky star winning numbers

<NUMBER> which consequently won in the 2ND

Q. What are the advantages of RAID5 over RAID4?

A. 1. Several write requests could be processed in parallel, since the bottleneck of a **single** check disk has been eliminated. 2. Read requests have a higher level of parallelism **on RAID5**. Since the data is distributed over all disks, read requests involve all disks, whereas in systems with a **check disk** the check disk never participates in read.

Set Similarity Functions

- The above similarity functions can be equivalently converted to overlap

- Overlap:

$$O(X, Y) = |X \cap Y|$$

- Jaccard:

$$J(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} \geq \theta \quad \leftrightarrow \quad O(X, Y) \geq \frac{\theta}{1 + \theta} (|X| + |Y|) = t$$

- We can only consider the **overlap size**, as others can be transformed to the overlap similarity.

Example

- Overlap threshold $t = 3$
- Result (RID pairs): $\{(1,3), (3,4), (3,5)\}$

RID	Name
1	Database System Concepts
2	Database Concepts Techniques
3	Database System Programming Concepts Oracle Linux
4	Database Programming Concepts Illustrated
5	System Programming Concepts Techniques Oracle Linux

Naïve Solutions

- Enumerate all string pairs
- Verify whether the strings in the pair are similar

Naïve Algorithm

need to compare all $O(n^2)$ pairs !!

RID	Name
1	Database System Concepts
2	Database Concepts Techniques
3	Database System Programming Concepts Oracle Linux
4	Database Programming Concepts Illustrated
5	System Programming Concepts Techniques Oracle Linux

$t = 3$

Result:
 (3, 1)
 (4, 3)
 (5, 4)

overlap of (i, j)	i=1	2	3	4	5
j=1		2	3	2	2
2			2	2	2
3				3	5
4					2

Framework

- for each $S_j \in S$
 - Candidates = ϕ
 - Candidates = **getCandidates**(S_j)
 - **Verify**(S_j , Candidates)

Verify($x, \{y_1, y_2, \dots\}$)

for each y_i
 if $\text{overlap}(x, y_i) \geq t$
 output($\langle x, y_i \rangle$)
end

- Algorithms differ in **getCandidates**()
 - naïve alg: return $\{S_i \mid i < j\}$

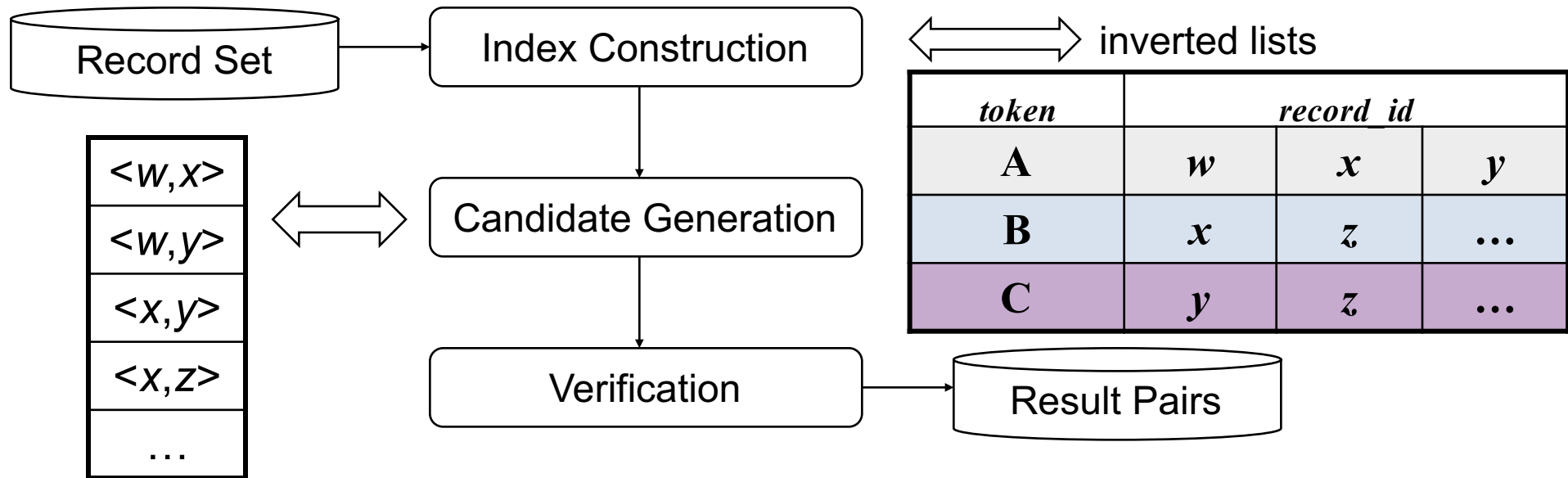
Not-so-naïve Algorithm

- Observation: (*, 6) should not be compared
 - If x and y have no common token, they won't be in the result
- Idea: Use *inverted index* to consider promising candidate pairs only

RID	Name
1	Database System Concepts
2	Database Concepts Techniques
3	Database System Programming Concepts Oracle Linux
4	Database Programming Concepts Illustrated
5	System Programming Concepts Techniques Oracle Linux
6	Harry Potter and the Sorcerer's Stone

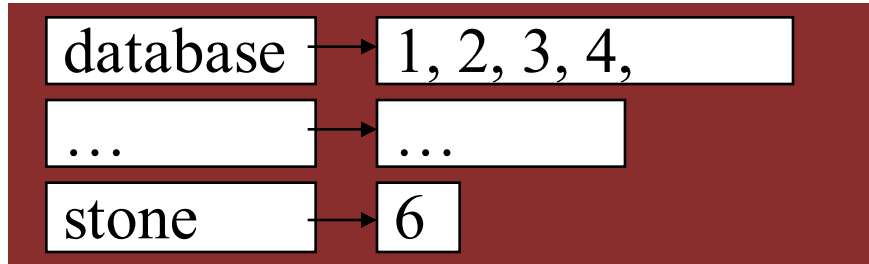
Index-Based Algorithms

- Index-based Algorithm



Inverted Index

- Conceptually, an inverted index has an inverted list for each token to be indexed from the document collection
 - An inverted list is just an sorted array of document identifiers (in our case, RIDs) such that the token appears in the corresponding document



RID	Name
1	Database System Concepts
...
6	Harry Potter and the Sorcerer's Stone

Probe Count

RID	Name
1	Database System Concepts
2	Database Concepts Techniques
3	Database System Programming Concepts Oracle Linux
4	Database Programming Concepts Illustrated
5	System Programming Concepts Techniques Oracle Linux

token	Inverted list
Database	{1, 2, 3, 4}
System	{1, 3, 5}
Concepts	{1, 2, 3, 4, 5}
Techniques	{2, 5}
Programming	{3, 4, 5}
Oracle	{3, 5}
Linux	{3, 5}
Illustrated	{4}

Pro

RID	Name
1	Database System Concepts
2	Database Concepts Techniques

3	term	RID-list
4	Database	{1, 2, 3, 4}
5	System	{1, 3, 5}
	Concepts	{1, 2, 3, 4, 5}
	Techniques	{2, 5}
	Programming	{3, 4, 5}
	Oracle	{3, 5}
	Linux	{3, 5}
	Illustrated	{4}

→ {1, 2, 3, 4, 5}

1 needs to be compared
with {2, 3, 4, 5}

1. Consider the terms in the current record and retrieve their inverted lists through the index
2. Merge these RID-lists
3. Repeat step 1 & 2 for every record

Note that RIDs in the inverted list are sorted

Pro

RID	Name	
1	Database System Concepts	
2	Database Concepts Techniques	
3	term	RID-list
4	Database	{1, 2, 3, 4}
5	System	{1, 3, 5}
	Concepts	{1, 2, 3, 4, 5}
	Techniques	{2, 5}
	Programming	{3, 4, 5}
	Oracle	{3, 5}
	Linux	{3, 5}
	Illustrated	{4}

} {1/3, 2/2, 3/3, 4/2, 5/2}

→

RID / count. e.g., overlap(1,2)=2

1. Consider each term, t , in the current record and retrieve their inverted lists through the index
2. Count Filtering: Merge and filter these RID-lists
 - prune those whose count $\leq t$
3. Repeat steps 1 & 2 for every record

Note that RIDs in the inverted list are sorted

Framework

- for each $S_j \in S$
 - Candidates = ϕ
 - Candidates = getCandidates(S_j)
 - Verify(S_j , Candidates)
- Algorithms differ in getCandidates()
 - naïve alg: return $\{S_i \mid i < j\}$
 - index-based alg: return $\{S_i \mid i < j \wedge S_i \cap S_j \neq \phi\}$
// use inverted indexes

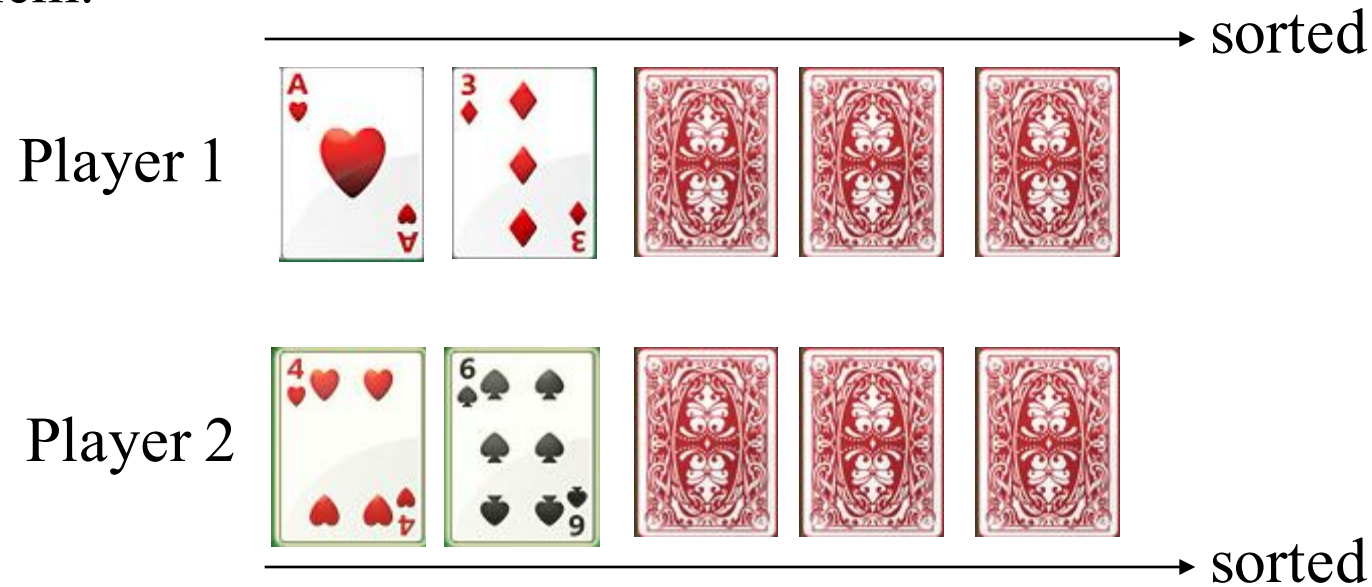
Problems

- Still too many comparisons
 - tokens that appears in many documents not only results in large inverted lists, but also slow down the computation

RID	Name
1	Database System Concepts
2	Database Concepts Techniques
3	Database System Programming Concepts Oracle Linux
4	Database Programming Concepts Illustrated
5	System Programming Concepts Techniques Oracle Linux
6	Database of Respiratory Diseases

Prefix Filtering

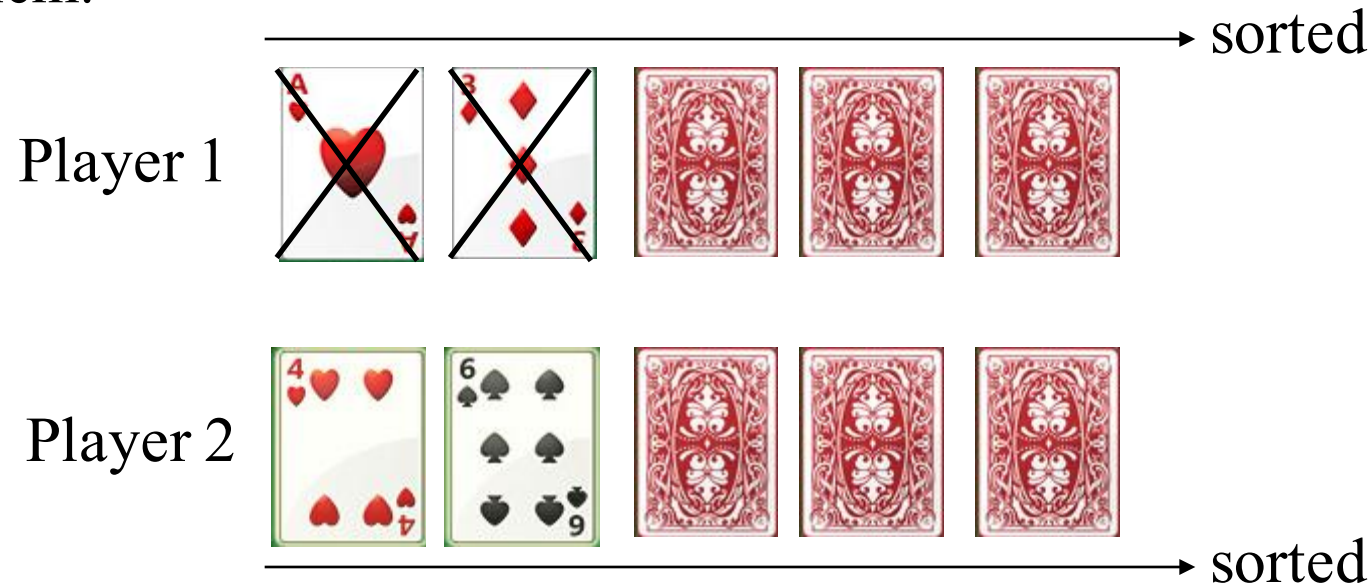
- Establish an upper bound of the overlap between two sets by looking at only *part of* them.



*what's the maximum possible number of cards held by both players
(denomination not considered) ?*

Prefix Filtering

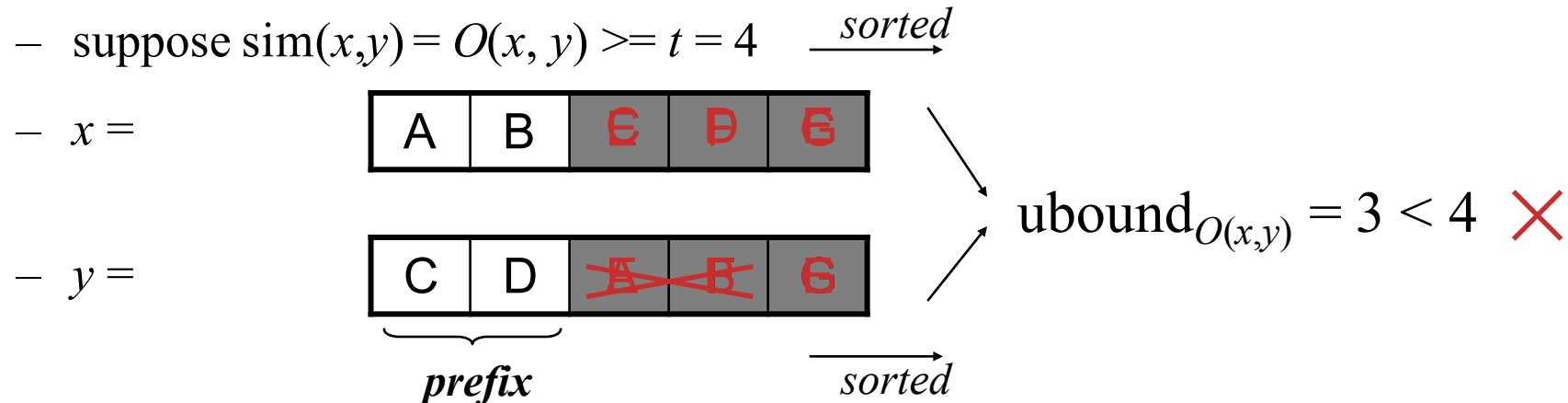
- Establish an upper bound of the overlap between two sets by looking at only *part of* them.



*what's the maximum possible number of cards held by both players
(denomination not considered) ?*

Prefix Filter

- Sort the tokens by a **global ordering**
 - increasing order of frequency
- Index the first few tokens (prefix) for each record
- Example:



- Must share at least one token in prefix to be a **candidate pair**

Prefix Filter

- For x (or y)
 - $O(x,y) \geq t \rightarrow \text{prefix length} = |x| - (t - 1)$
 - $J(x,y) \geq \delta \rightarrow O(x,y) \geq \delta |x| \rightarrow \text{prefix length} = \lfloor (1 - \delta) |x| + 1 \rfloor$
- Example: suppose $\text{sim}(x,y) = J(x,y) \geq \delta = 0.8$
 - $w = \{\underline{C}, D, E, F\}$
 - $x = \{\underline{B}, \underline{C}, D, E, F\}$
 - $y = \{\underline{A}, \underline{B}, C, D, F\}$
 - $z = \{\underline{G}, \underline{A}, B, E, F\}$

	Candidate Pairs	Results
	$\langle w,x \rangle, \langle x,y \rangle, \langle y,z \rangle$	$\langle w,x \rangle$

Prefix Filtering for Overlap

- Constraint: $\text{Overlap}(S_1, S_2) \geq t$
- Pre-requisite
 - All sets sorted in a global order
 - Usually the increasing **frequency** order \leftarrow facilitates the join
- Preprocessing
 - Set $S \rightarrow \text{prefix}(S)$, s.t., $|\text{prefix}(S)| = |S| - (t - 1)$
- Filter
 - If $\text{prefix}(S_1) \cap \text{prefix}(S_2) = \emptyset$, then $\text{Overlap}(S_1, S_2) < t$
 - i.e., (S_1, S_2) can be filtered
- Refine: verify the survived pairs

Framework

- for each $S_j \in S$
 - Candidates = ϕ
 - Candidates = **getCandidates**(S_j)
 - Verify(S_j , Candidates)
- Algorithms differ in **getCandidates()**
 - naïve alg: return $\{S_i \mid i < j\}$
 - index-based alg: return $\{S_i \mid i < j \wedge S_i \cap S_j \neq \phi\}$ // use inverted index
 - **prefix-filtering-base alg: return $\{S_i \mid i < j \wedge \text{prefix}(S_i) \cap \text{prefix}(S_j) \neq \phi\}$ using inverted index on prefix set**

Pro

RID	Name
1	Database System Concepts
2	Database Concepts Techniques
3	Database System Programming Concepts Oracle Linux
4	Database Programming Concepts Illustrated
5	System Programming Concepts Techniques Oracle Linux

token	Inverted list	df	idf	Order
Database	{1, 2, 3, 4}	4	1/4	7
System	{1, 3, 5}	3	1/3	6
Concepts	{1, 2, 3, 4, 5}	5	1/5	8
Techniques	{2, 5}	2	1/2	4
Programming	{3, 4, 5}	3	1/3	5
Oracle	{3, 5}	2	1/2	3
Linux	{3, 5}	2	1/2	2
Illustrated	{4}	1	1/1	1

Pro

RID	Name	
1	<i>System</i>	Database Concepts
2	<i>Techniques</i>	Database Concepts
3	<i>Linux Oracle Programming System</i>	Database Concepts
4	<i>Illustrated Programming</i>	Database Concepts
5	<i>Linux Oracle Techniques Programming</i>	System Concepts

Order: Illustrated, Linux, Oracle, Techniques, Programming, System, Database, Concepts

$$t=3 \rightarrow \text{prefix-len} = |S| - 2$$

j=1

token	Inverted list
System	{1}

j=3

token	Inverted list
System	{1,3}
Techniques	{2}
Linux	{3}
Oracle	{3}
Programming	{3}

j=2

token	Inverted list
System	{1}
Techniques	{2}

Pro

RID	Name	
1	<i>System</i>	Database Concepts
2	<i>Techniques</i>	Database Concepts
3	<i>Linux Oracle Programming System</i>	Database Concepts
4	<i>Illustrated Programming</i>	Database Concepts
5	<i>Linux Oracle Techniques Programming</i>	System Concepts

$t=3 \rightarrow \text{prefix-len} = |S| - 2$

$j=3$

token	Inverted list
System	{1,3}
Techniques	{2}
Linux	{3}
Oracle	{3}
Programming	{3}

$j=4$

token	Inverted list
System	{1,3}
Techniques	{2}
Linux	{3}
Oracle	{3}
Programming	{3, 4}
Illustrated	{4}

$j=5$

token	Inverted list
System	{1,3}
Techniques	{2, 5}
Linux	{3, 5}
Oracle	{3, 5}
Programming	{3, 4}
Illustrated	{4}

Framework

- for each $S_j \in S$
 - Candidates = ϕ
 - Candidates = getCandidates(S_j)
 - Verify(S_j , Candidates)
- getCandidates()
 - prefix-filtering-base algorithm:
 - Return $\{S_i \mid i < j \wedge \text{prefix}(S_i) \cap \text{prefix}(S_j) \neq \phi\}$
 - Using inverted lists

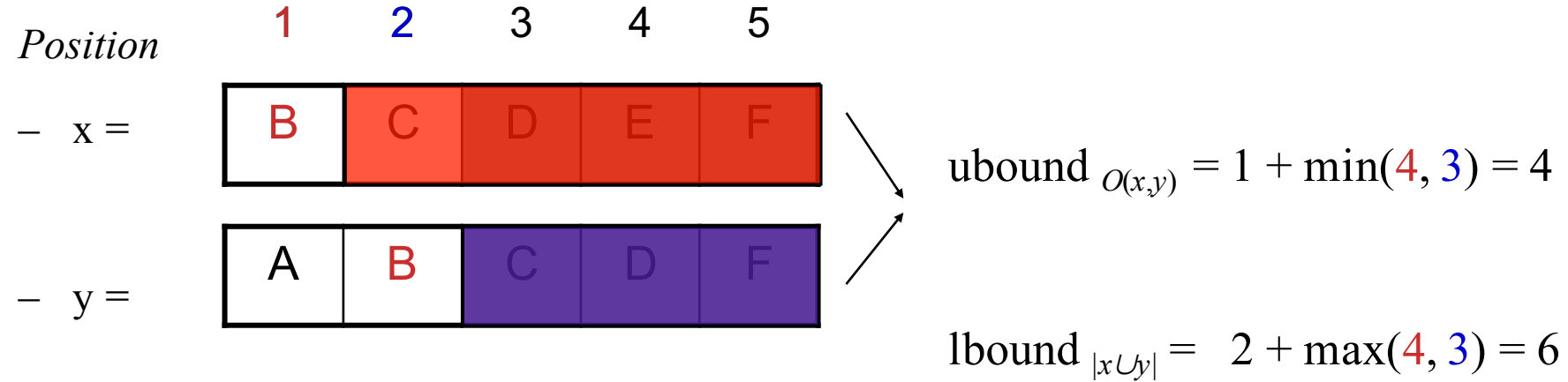
Supporting two datasets

- Join R, S
- Generate prefix sets for R and S
- Build inverted lists for prefix sets of R
- for each $S_j \in S$
 - Candidates = ϕ
 - Candidates = `getCandidates(S_i)`
 - `Verify(S_j , Candidates)`
- Algorithms differ in `getCandidates()`
 - **prefix-filtering-base algorithm:**
 - Return $\{S_i \mid \text{prefix}(S_i) \cap \text{prefix}(R_j) \neq \phi\}$
 - Using inverted lists

Pruning Techniques for Jaccard Similarity

Positional Filter within Prefix

- Index both tokens and their positions

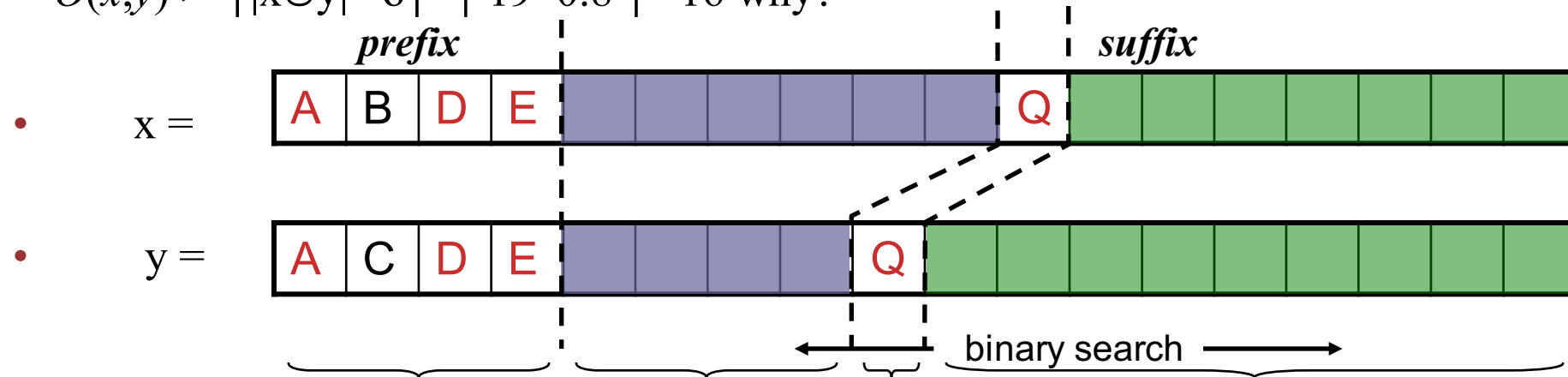


$$\text{ubound}_{J(x,y)} = \frac{|x \cap y|}{|x \cup y|} = \frac{4}{6} < \delta = 0.8 \quad \times$$

- $\text{ubound}_{O(x,y)} = p_x \cap p_y + \min(|x| - |p_x|, |y| - |p_y|)$ $p_x = B$
- $\text{lbound}_{x \cup y} = p_x \cup p_y + \max(|x| - |p_x|, |y| - |p_y|)$ $p_y = A \ B$

Positional Filter within Suffix

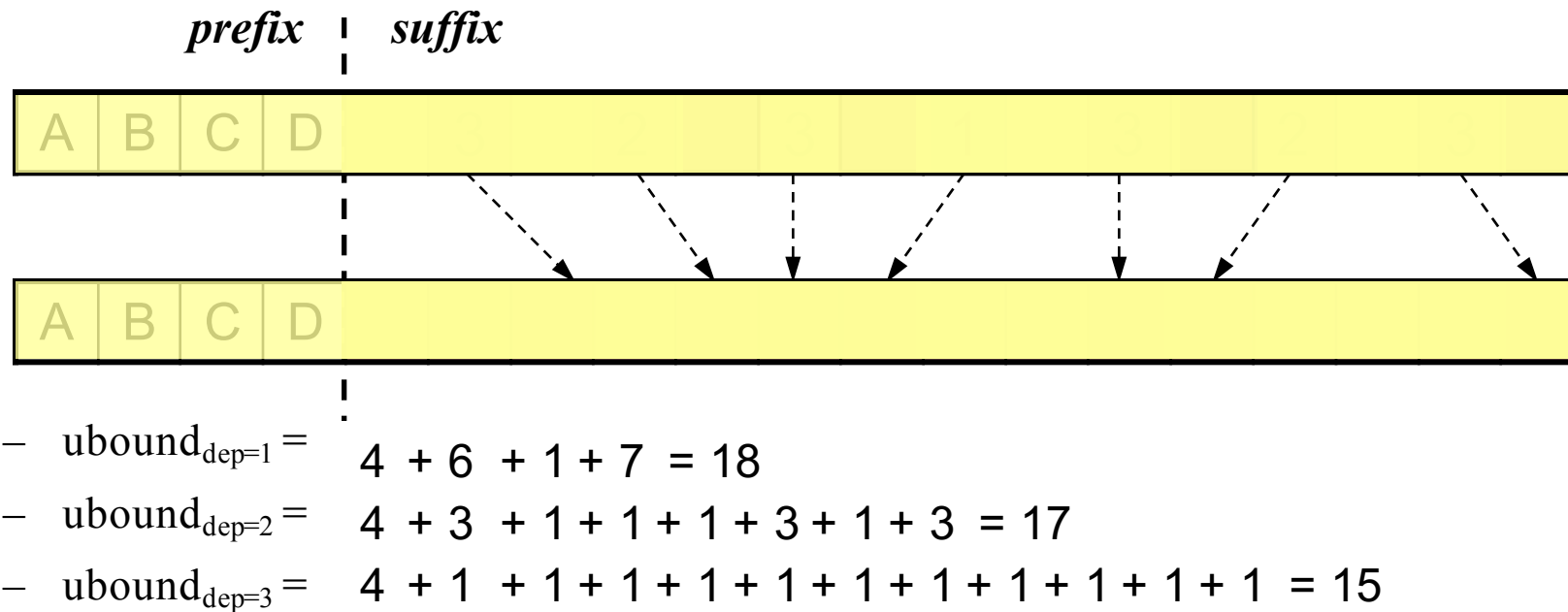
- probe tokens in suffix, and compare their positions
- suppose $\text{sim}(x,y) = J(x,y) \geq \delta = 0.8$
- $|x| = |y| = 18$, $|\text{prefix}| = \lfloor (1 - \delta) * x + 1 \rfloor = \lfloor (1 - 0.8) * 18 + 1 \rfloor = 4$ why?
- $O(x,y) \geq \lceil |x \cup y| * \delta \rceil = \lceil 19 * 0.8 \rceil = 16$ why?



- $\text{ubound}_{O(x,y)} = 3 \quad \times \quad +4 \quad +1 \quad +7$
 $= 15 < 16$

Positional Filter within Suffix

- Divide and Conquer



- probe suffix recursively, until either candidate pair is pruned, or reach max-depth

Effect of Filters

- $\text{sim}(x,y) = J(x,y) \geq \delta = 0.8$
 - $u = \{\underline{C}, D, E, F\}$
 - $v = \{\underline{B}, \underline{C}, D, E, F\}$
 - $w = \{\underline{A}, \underline{B}, C, D, F\}$
 - $x = \{\underline{G}, \underline{A}, B, E, F\}$
 - $y = \{\underline{A}, \underline{B}, D, E, F\}$
 - $z = \{\underline{G}, \underline{A}, C, D, E, F\}$
- after prefix filter:
 - $\langle u,v \rangle, \langle v,w \rangle, \langle v,y \rangle, \langle w,x \rangle, \langle w,y \rangle, \langle w,z \rangle, \langle x,y \rangle, \langle x,z \rangle, \langle y,z \rangle$
- after prefix position filter:
 - $\langle u,v \rangle, \langle w,y \rangle, \langle w,z \rangle, \langle x,z \rangle, \langle y,z \rangle$
- after suffix position + (max-depth = 1):
 - $\langle u,v \rangle, \langle x,z \rangle, \langle y,z \rangle$
- real result:
 - $\langle u,v \rangle$

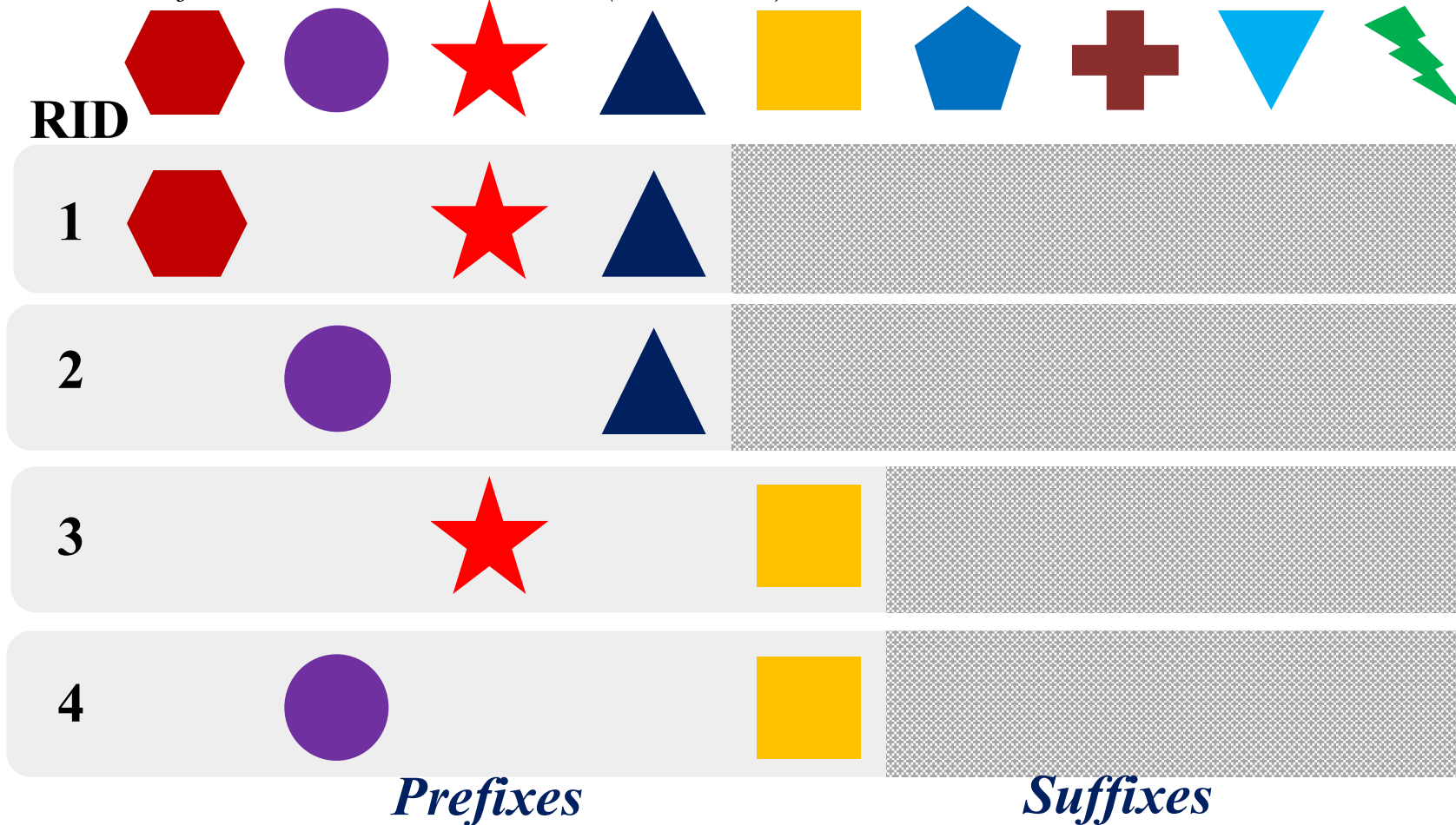
Similarity Join for Jaccard

- Generate tokens
- Compute frequency of token
- Filter: Generate Candidate Pairs
 - Using inverted lists
 - Using prefix sets
 - Using prefix sets with prefix-position pruning
 - Using suffix pruning
- Verify Candidate Pairs

A Partition-based Method for Set Similarity Join

Prefix Filter Framework

The list of all elements in order (universe)



Prefix Filter: $\text{Sim}(X, Y) \geq \delta$ only if $\text{Prefix}(X) \cap \text{Prefix}(Y) \neq \emptyset$

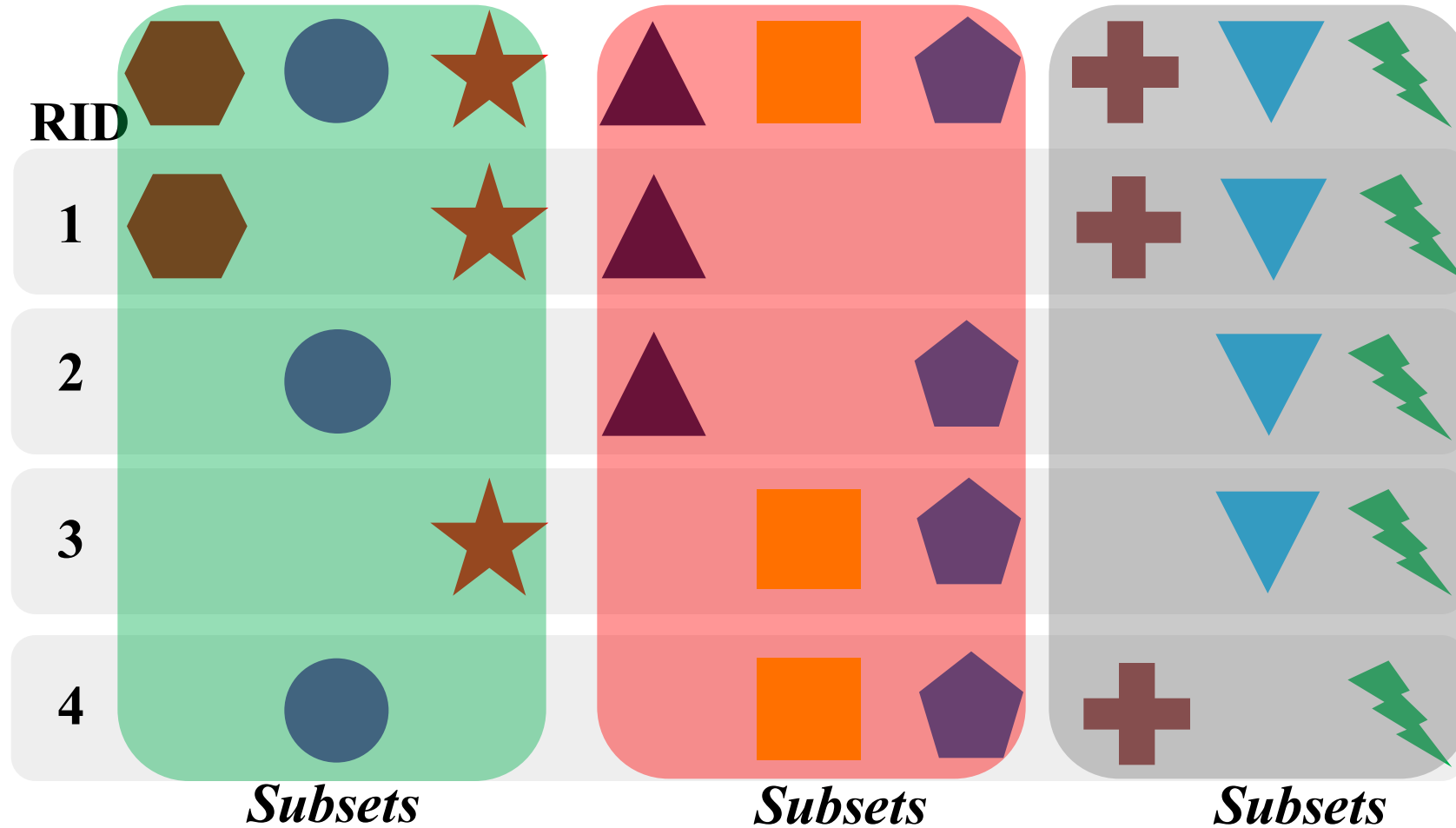
Prefix Filter Framework

the pruning power is limited!

*two dissimilar sets are a candidate
if they share 1 element in their prefixes*

Partition-based Framework

The list of all elements in order (universe)



$\text{Sim}(X, Y) \geq \delta$ only if $\text{Subsets}(X) \cap \text{Subsets}(Y) \neq \emptyset$

Partition-based Framework

What is the minimum number of partitions that can guarantee completeness?

The Number of Partitions

Intuition:

1: Deduce an overlap lower bound based on the similarity function and the threshold

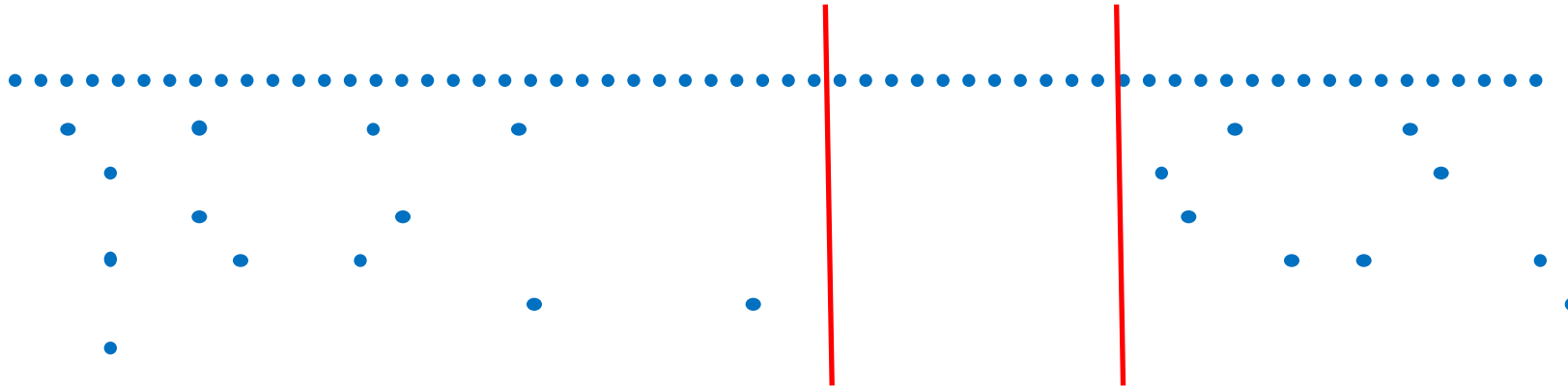
$$\text{Sim}(X, Y) \geq \delta \rightarrow |X \cap Y| \geq m$$

2: Partition them into $m + 1$ subsets

Then two similar sets must share at least 1 subset

Element Skew Problem

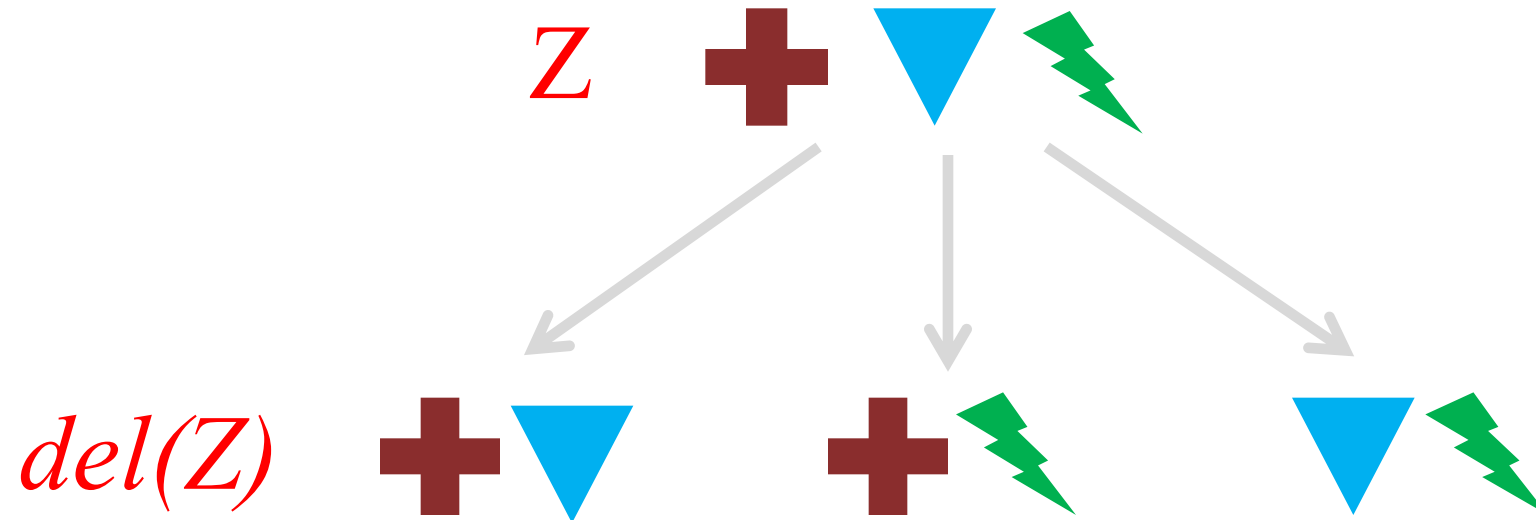
- Some subsets have limited number of elements
 - The ‘empty’ subsets yield quadratic candidates



- Solution: Add some flexibility
 - Skip the subsets with less elements
 - Select more signatures from subsets with more elements to guarantee completeness

Signatures: 1-deletion neighborhoods

- Given a non-empty set Z , its 1-deletion neighborhoods are its subsets with size of $|Z| - 1$, denote as $del(Z)$



Using 1-deletion neighborhood

$$U \neq V, V \notin \text{del}(U), U \notin \text{del}(V) \rightarrow |U \Delta V| \geq 2$$



$$|U \Delta V| = 2$$

Skip x subsets & select 1-deletion neighborhoods from another x subsets

guaranteed completeness !!

Optimal Allocation Strategy

$$v_i = \begin{cases} 0: & \text{skip the } i^{\text{th}} \text{ subset} \\ 1: & \text{only use the } i^{\text{th}} \text{ subset as signature} \\ 2: & \text{use both the } i^{\text{th}} \text{ subset and 1-deletions} \end{cases}$$



Constraint: $\sum_{i=1}^{m+1} v_i = m + 1$

Object: minimize $\sum_{i=1}^{m+1} c_{v_i}^i$

$c_0^i = 0$ c_1^i : the # of sets sharing the i^{th} subset

c_2^i : the # of sets sharing (subset or 1-deletion) signatures

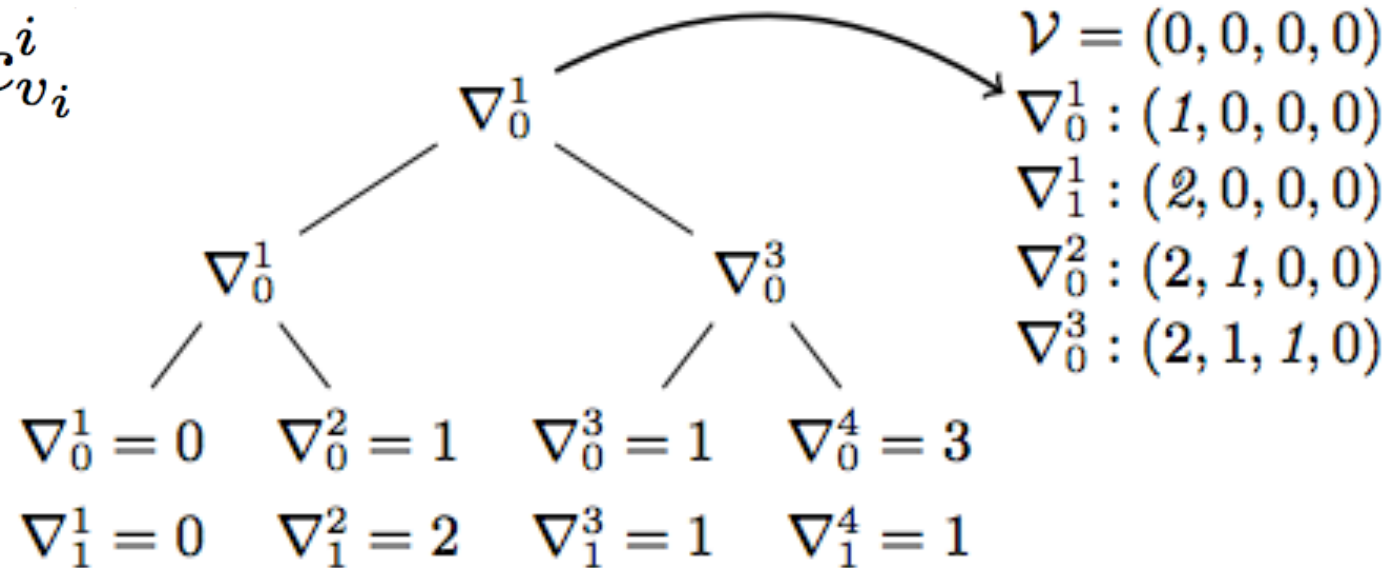
Time Complexity

- Dynamic Programming
 - Optimal: # of candidates
 - $O(s^2)$ time complexity as $m = O(s)$ where s is the set size
- Each set is partitioned $s - \delta s + 1 = O(s)$ times
- Allocation time complexity is $O(s^3)$ for each set
- Next we reduce it to $O(s \log s)$

Greedy Method for Allocation Selection

- Heap-based Method

$$\nabla_{v_i}^i = c_{v_i+1}^i - c_{v_i}^i$$



- 2-approximation algorithm
- Time complexity: $O(s^2) \rightarrow O(s \log s)$

Adaptive Grouping



Adaptive Grouping

- The k -th group includes all the sets with size within $[\frac{l_{min}}{\alpha^{k-1}}, \frac{l_{min}}{\alpha^k})$
where $\alpha \in [\frac{1}{2}, 1]$
- The size range becomes more and more ‘broader’
- The partition times is bounded by $\log_{\alpha} \delta + 1 = O(1)$ for any set
- The time complexity is $O(s \log s \log_{\alpha} \delta) = O(s \log s)$

Prefix Filtering for Edit Distance

q-gram Based Filtering

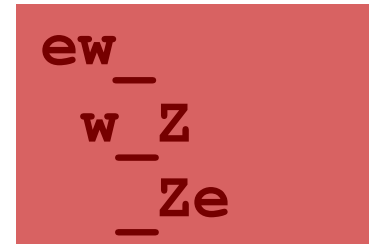
- Naïve algorithm
 - compute edit distance: $O(n^2)$ time complexity, where n is string length
 - do this for $N^2/2$ pairs, where N is the number of strings
- q-gram based filtering
 - q-gram is a substring of length q
 - filter-and-refine
 - length filtering
 - $|\text{len}(s) - \text{len}(t)| \leq \tau$

Matching q-grams

- count filtering
 - at least $LB(s, t)$ common q -grams
 - $LB(s, t) = \max(|s|, |t|) - q + 1 - q * \tau$
- Bottleneck: generating candidate pairs which share at least $LB(s, t)$ matching q -grams

New_Zealand

New



Zea

eal

ala

lan

and

Prefix Length

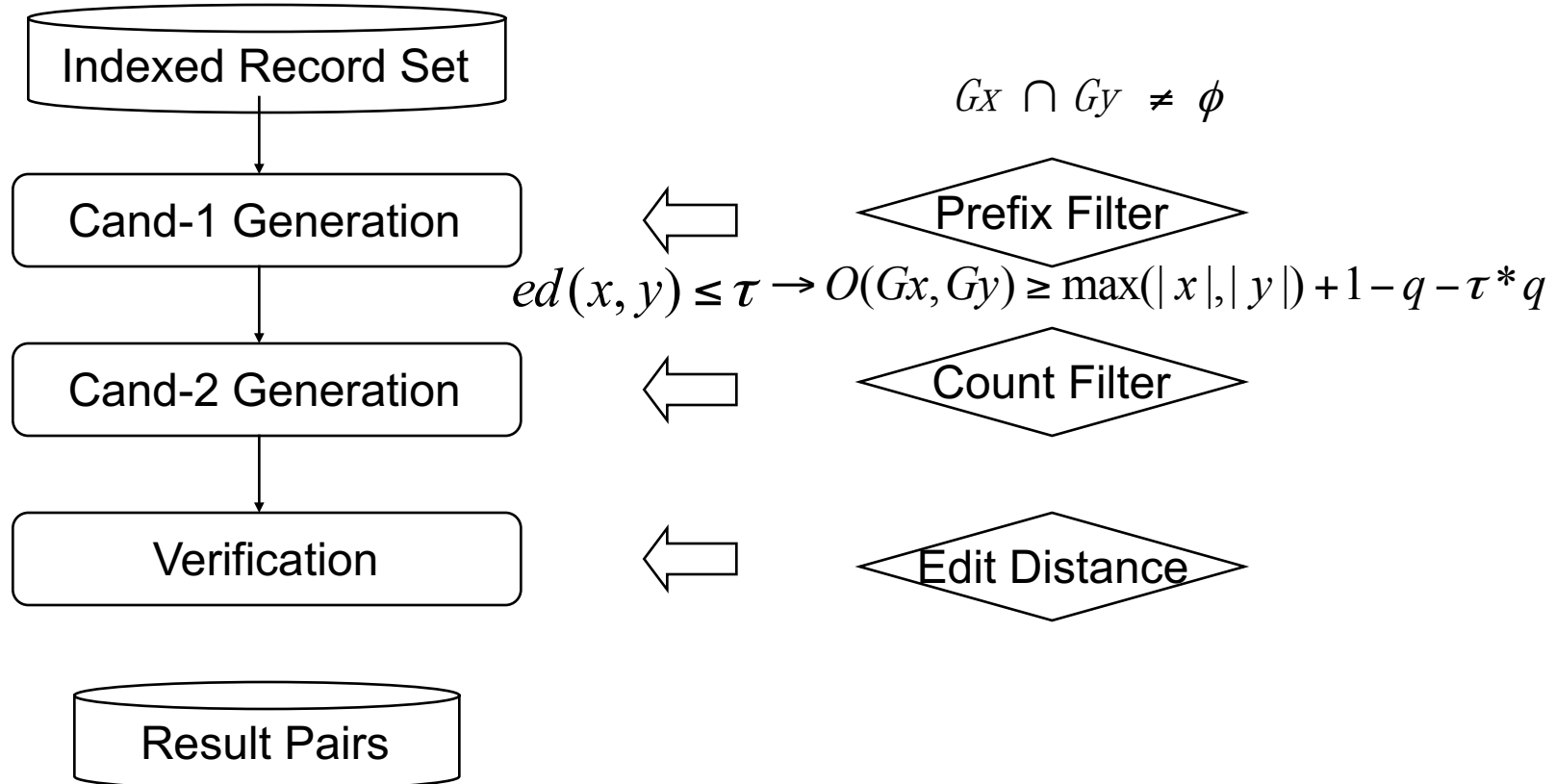
- Edit distance: $ed(x, y) \leq \tau$

$$\rightarrow O(Gx, Gy) \geq \max(|x|, |y|) + 1 - q - \tau * q$$

Gx and Gy are the gram sets of x and y respectively

$$\text{Prefix length} = |x| + 1 - q - (0 - 1) = \tau * q + 1$$

Similarity-Join Algorithm for ED



Prefix Filter

- Prefix Filter
 - sort q-grams by global ordering, e.g., *frequency*

$$\begin{array}{c}
 \underbrace{q^* \tau + 1} \quad \underbrace{|s| + 1 - q - q^* \tau - 1} \\
 - \quad Q_s = \begin{array}{|c|c|c|c|c|} \hline q_a & q_b & & & \\ \hline \end{array}
 \end{array}$$

$$\begin{array}{c}
 - \quad Q_t = \begin{array}{|c|c|c|c|c|} \hline q_x & q_y & & & \\ \hline \end{array} \\
 |t| + 1 - q - q^* \tau - 1
 \end{array}$$

Example – All-Pairs-Ed

- $\tau=1, q=2$ \longrightarrow $\text{prefix_len} = q * \tau + 1 = 3$
 - $a = \text{'Austria'}$ $Qa = \{\text{ri, Au, us, ...}\}$
 - $b = \text{'Australia'}$ $Qb = \{\text{ra, li, Au, ...}\}$
 - $c = \text{'Australiana'}$ $Qc = \{\text{na, ra, li, ...}\}$
 - $d = \text{'New_Zealand'}$ $Qd = \{\text{_Z, Ze, Ne, ...}\}$
 - $e = \text{'New_Sealand'}$ $Qe = \{\text{_S, Se, Ne, ...}\}$
- after prefix filter: $\langle a, b \rangle \langle b, c \rangle \langle d, e \rangle$
- after count filter: $\langle b, c \rangle \langle d, e \rangle$
- after edit distance verification: $\langle d, e \rangle$

Edit-Distance Join (Ed-Join)

- Idea
 - mismatching q -grams provide useful information

Location-Based Filtering

- Idea: reduce prefix length

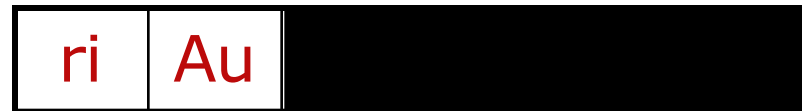
- Example, $\tau=1$, $q=2$

– $s = \text{'Austria'}$

– $t = \text{'Australia'}$

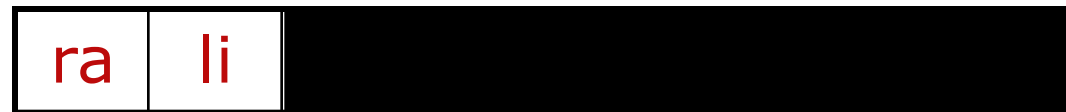
location 5 1

– $Q_s =$



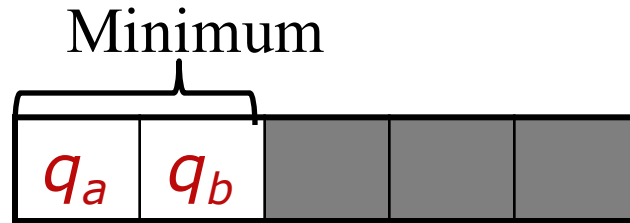
– $Q_t =$

location 5 7



Find Minimum Prefix Length

- Given a prefix set of q -grams, find the minimum prefix length that needs at least $\tau + 1$ edit operations to destroy them



- $Q_s =$

- Example

– $s = \text{'Austria'}$

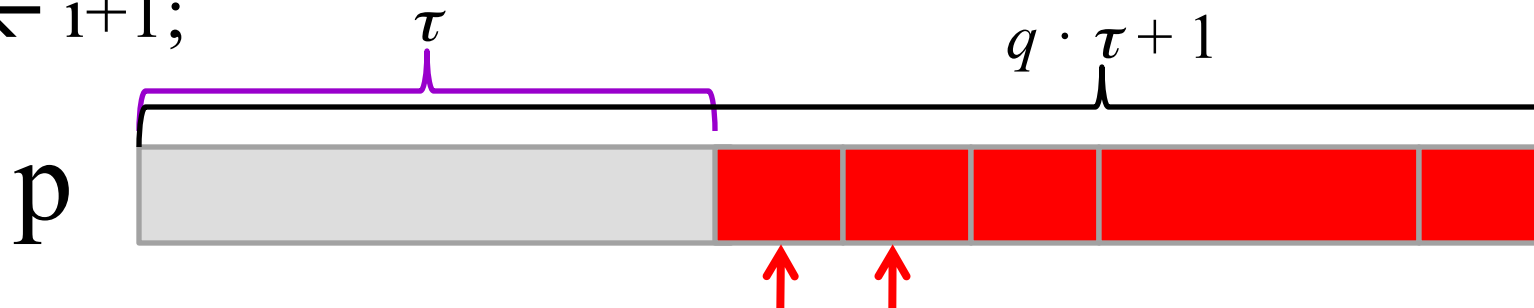
– $t = \text{'Australia'}$

– $p_s = \text{ri Au us (5 1 2)}$ → ri Au

– $P_t = \text{ra li Au (5 7 1)}$ → ra li

Algorithms

- Find the minimum prefix length of p
- $i \leftarrow \tau + 1$; $\text{end} \leftarrow q \cdot \tau + 1$;
- for $i \leq \text{end}$ do
- $\text{err} \leftarrow \text{MinEditErrors}(p[1 \dots i])$;
- if $\text{err} > \tau$ then
- break;
- $i \leftarrow i + 1$;

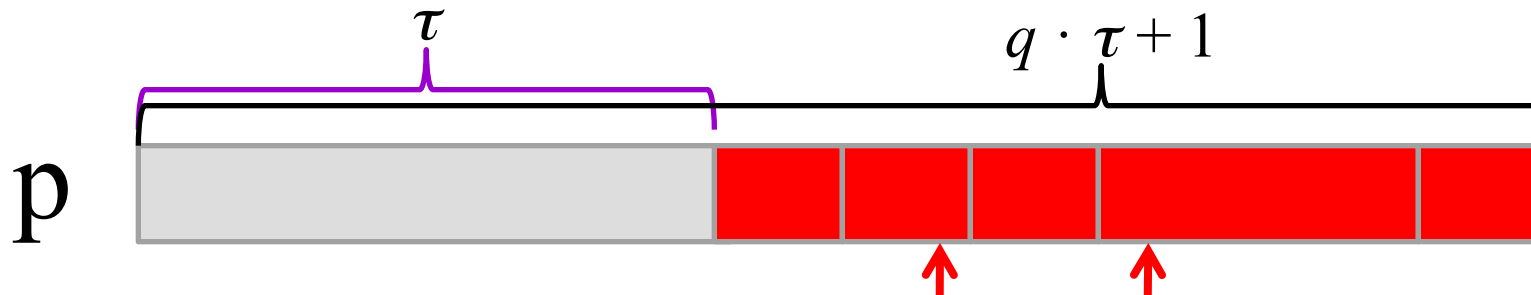


Algorithms

- Find Minimum Edit Errors of a set of q-grams
- Sort q-grams in Q in increasing order of locations;
- $\text{cnt} \leftarrow 0; \text{loc} \leftarrow 0;$
- for $i = 1$ to $|Q|$ do
- if $Q[i].\text{loc} > \text{loc}$ then
- $\text{cnt} \leftarrow \text{cnt} + 1;$
- $\text{loc} \leftarrow Q[i].\text{loc} + q - 1;$
- return cnt
- $p_s = \text{ri Au us } (5 \ 1 \ 2) \rightarrow 1 \ 2 \ 5 \rightarrow 2$
- $P_t = \text{ra li Au } (5 \ 7 \ 1) \rightarrow 1 \ 5 \ 7 \rightarrow 3$

Improved Algorithms

- Find the minimum prefix length
- $\text{left} \leftarrow \tau + 1; \text{right} \leftarrow q \cdot \tau + 1;$
- while $\text{left} < \text{right}$ do
- $\text{mid} \leftarrow (\text{left} + \text{right})/2;$
- $\text{err} \leftarrow \text{MinEditErrors}(x[1 \dots \text{mid}]);$
- if $\text{err} \leq \tau$ then $\text{left} \leftarrow \text{mid} + 1;$
- else $\text{right} \leftarrow \text{mid};$
- return left



Pivotal Prefix Filter

q-gram

- *q*-gram is the substring of length *q*

2-grams { *yo*
ou
ut
tb
be
ec
co
om

q-gram

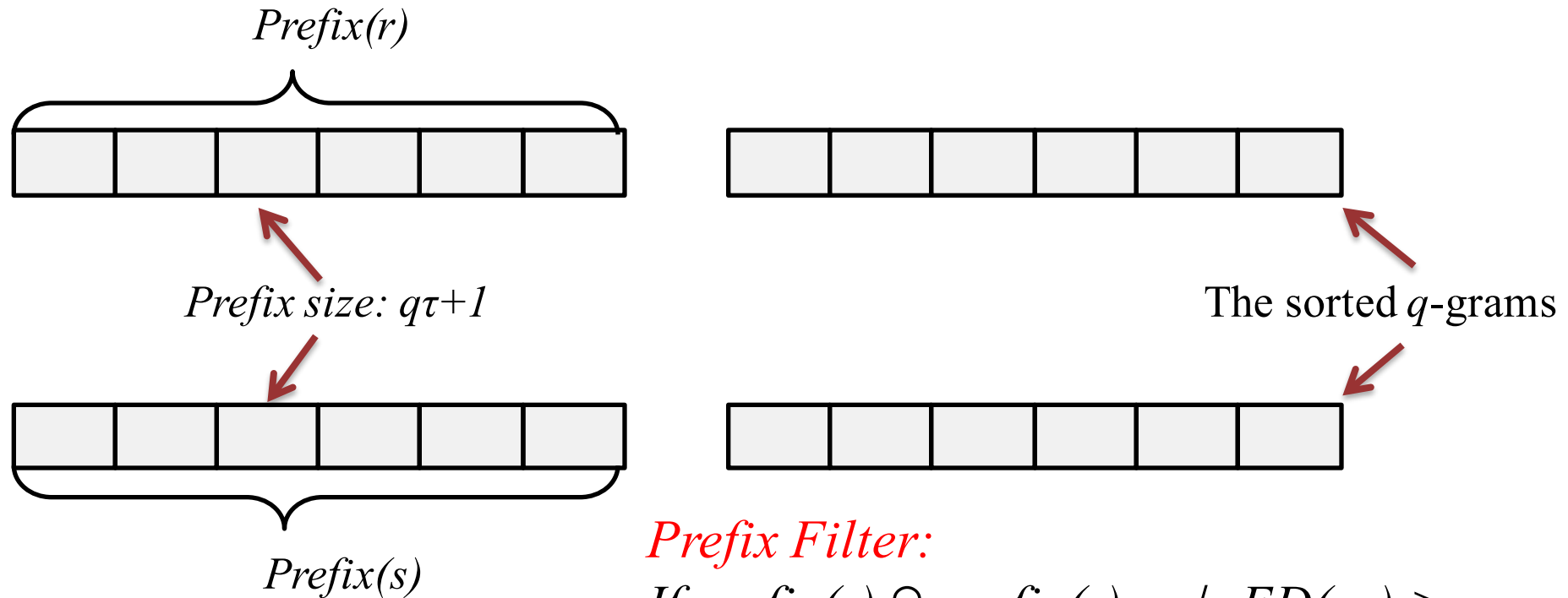
- 1 edit operation destroys at most *q* grams.

youtdecom
yo
ou
ut
td
de
ec
co
om

- τ edit operations destroy *at most* *q* τ q-grams.
- if *r* and *s* have more than *q* τ mismatch q-grams, $ED(r, s) > \tau$

Prefix Filter

Sort all q -grams by global ordering, such as *alphabetic order*



Prefix Filter:

If $prefix(r) \cap prefix(s) = \phi$, $ED(r,s) > \tau$

Disjoint q-gram

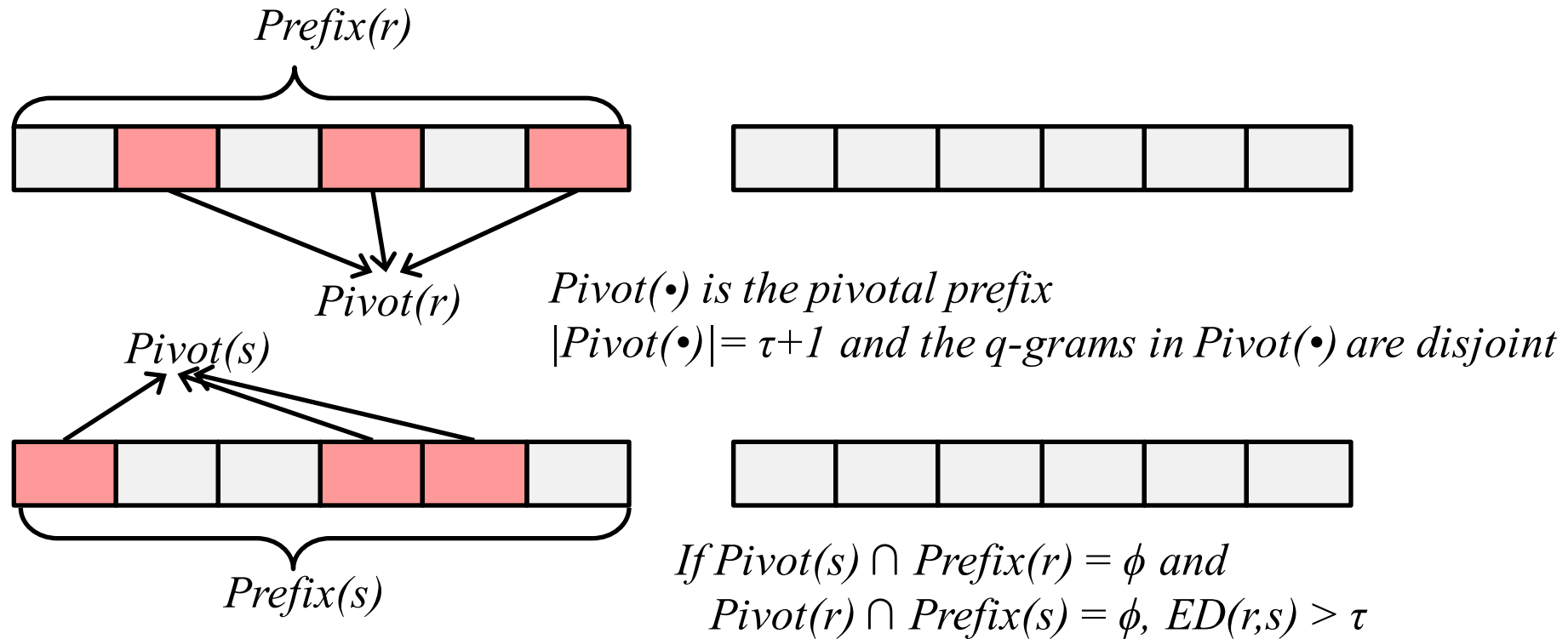
- One edit operation destroys at most *1 disjoint* gram.

yout decom
yo
ut
de
om

- τ edit operations destroy *at most* *τ disjoint* q-grams.
- if r and s have more than τ mismatch *disjoint* q-grams, $ED(r, s) > \tau$

Pivotal Prefix Filter

Sort all q-grams by global ordering, such as *alphabetical order*



Summary

- Prefix Filtering for Set Similarity Join
- Position Filter
- Suffix Filter
- Partition-based Framework
- 1-Deletion Neighborhood
- Optimal Allocation Algorithm
- 2-approximation Greedy Algorithm
- Adaptive Grouping Mechanism
- Prefix Filtering for Edit Distance
- Pivotal Prefix Filter