# String Similarity Join

Dong Deng

# Real World Data is Dirty

- Misspellings of the query "*britney spears*" on Google

| | | |
|---|---|---|
| 488941 britney spears | 29 britent spears | 9 brinttany spears |
| 40134 brittany spears | 29 brittnany spears | 9 britanay spears |
| 36315 brittney spears | 29 britttany spears | 9 britinany spears |
| 24342 britany spears | 29 btiney spears | 9 britn spears |
| 7331 britny spears | 26 birttney spears | 9 britnew spears |
| 6633 briteny spears | 26 breitney spears | 9 britneyn spears |
| 2696 britteny spears | 26 brinity spears | 9 britrney spears |
| 1807 briney spears | 26 britenay spears | 9 brtiny spears |
| 1635 brittny spears | 26 britneyt spears | 9 brtittney spears |
| 1479 brintey spears | 26 brittan spears | 9 brtny spears |
| 1479 britanny spears | 26 brittne spears | 9 brytny spears |
| 1338 britiny spears | 26 btittany spears | 9 rbitney spears |
| 1211 britnet spears | 24 beitney spears | 8 birtiny spears |
| 1096 britiney spears | 24 birteny spears | 8 bithney spears |
| 991 britaney spears | 24 brightney spears | 8 brattany spears |
| 991 britnay spears | 24 brintiny spears | 8 breitny spears |
| 811 brithney spears | 24 britanty spears | 8 breteny spears |
| 811 brtiney spears | 24 britenny spears | 8 brighty spears |
| 664 birtney spears | 24 britini spears | 8 brintay spears |
| 664 brintney spears | 24 britnwy spears | 8 brinttey spears |
| 664 briteney spears | 24 brittni spears | 8 briotney spears |

# Set Similarity Functions

- Overlap Size

$$overlap(x, y) = |x \cap y|$$

- Jaccard Similarity

$$J(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

| |
|---|
| $x$ = {A,B,C,D,E} <br> $y$ = {B,C,D,E,F} |
| overlap(x, y) = 4 <br><br> J(x, y) = 4/6 = 0.67 |

# String Similarity Function

- Edit Distance ED(*r, s*): the minimum number of edit operations (insertion/deletion/substitution) needed to transform *r* to *s*.

- For example:   ED(*hilton,* huston) = 2

h*i*lton

*substitute i with u*

hu*l*ton

*substitute l with s*

hu**s**ton

- Edit Similarity: $EDS(r, s) = 1 - \frac{ED(r,s)}{\max(|r|,|s|)}$

# Calculating Edit Distance

*Calculating ED(r,s)*

*Let $r_n$ and $s_m$ be the last characters in r and s.*

a) Match $r_n$ and $s_m$

ED(r,s)= 0 + ED(r[1,n-1],s[1,m-1])

b) Substitute $r_n$ with $s_m$

ED(r,s)= 1 + ED(r[1,n-1],s[1,m-1])

c) Delete $r_n$

ED(r,s)= 1 + ED(r[1,n-1],s[1,m])

d) Insert $s_m$

ED(r,s)= 1 + ED(r[1,n],s[1,m-1])

# Calculating Edit Distance

| | $\$$ | $s_1$ | ... | $s_j$ | ... | $s_m$ |
|---|---|---|---|---|---|---|
| $\$$ | | | | | | |
| $r_1$ | | | | | | |
| $\vdots$ | | | | | | |
| $r_i$ | | | | | | |
| $\vdots$ | | | | | | |
| $r_n$ | | | | | | |

ED($\$$, s[1,j]) = j

*match/substitute*

*insert*

*delete*

ED(r[1,i],s[1,j])

ED(r, s)

# Calculating Edit Distance

|   | $ | b | r | i | t | n | e | y |
|---|---|---|---|---|---|---|---|---|
| $ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| b | 1 |   |   |   |   |   |   |   |
| t | 2 |   |   |   |   |   |   |   |
| i | 3 |   |   |   |   |   |   |   |
| n | 4 |   |   |   |   |   |   |   |
| e | 5 |   |   |   |   |   |   |   |
| y | 6 |   |   |   |   |   |   |   |

# Calculating Edit Distance

|   | $ | b | r | i | t | n | e | y |
|---|---|---|---|---|---|---|---|---|
| $ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| b | 1 | 0 |   |   |   |   |   |   |
| t | 2 |   |   |   |   |   |   |   |
| i | 3 |   |   |   |   |   |   |   |
| n | 4 |   |   |   |   |   |   |   |
| e | 5 |   |   |   |   |   |   |   |
| y | 6 |   |   |   |   |   |   |   |

# Calculating Edit Distance

|     | $   | b   | r   | i   | t   | n   | e   | y   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $   | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
| b   | 1   | 0 → 1 |   |     |     |     |     |     |
| t   | 2   |     |     |     |     |     |     |     |
| i   | 3   |     |     |     |     |     |     |     |
| n   | 4   |     |     |     |     |     |     |     |
| e   | 5   |     |     |     |     |     |     |     |
| y   | 6   |     |     |     |     |     |     |     |

# Calculating Edit Distance

|     | $   | b   | r   | i   | t   | n   | e   | y   |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| $   | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
| b   | 1   | 0 → 1 → 2 |     |     |     |     |     |     |
| t   | 2   |     |     |     |     |     |     |     |
| i   | 3   |     |     |     |     |     |     |     |
| n   | 4   |     |     |     |     |     |     |     |
| e   | 5   |     |     |     |     |     |     |     |
| y   | 6   |     |     |     |     |     |     |     |

# Calculating Edit Distance

|     | $   | b   | r   | i   | t   | n   | e   | y   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $   | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
| b   | 1   | 0   | 1   | 2   | 3   | 4   | 5   | 6   |
| t   | 2   |     |     |     |     |     |     |     |
| i   | 3   |     |     |     |     |     |     |     |
| n   | 4   |     |     |     |     |     |     |     |
| e   | 5   |     |     |     |     |     |     |     |
| y   | 6   |     |     |     |     |     |     |     |

# Calculating Edit Distance

|     | $   | b   | r   | i   | t   | n   | e   | y   |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **$** | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
| **b** | 1   | 0   | 1   | 2   | 3   | 4   | 5   | 6   |
| **t** | 2   | 1   | 1   | 2   | 3   | 4   | 5   | 6   |
| **i** | 3   |     |     |     |     |     |     |     |
| **n** | 4   |     |     |     |     |     |     |     |
| **e** | 5   |     |     |     |     |     |     |     |
| **y** | 6   |     |     |     |     |     |     |     |

# Calculating Edit Distance

|     | $ | b | r | i | t | n | e | y |
|-----|---|---|---|---|---|---|---|---|
| $   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| b   | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| t   | 2 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| i   | 3 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| n   | 4 | 3 | 3 | 2 | 2 | 2 | 3 | 4 |
| e   | 5 | 4 | 4 | 3 | 3 | 3 | 2 | 3 |
| y   | 6 | 5 | 5 | 4 | 4 | 4 | 3 | 2 |

# Calculating Edit Distance

|     | $   | b   | r   | i   | t   | n   | e   | y   |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **$** | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| **b** | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| **t** | 2 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| **i** | 3 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| **n** | 4 | 3 | 3 | 2 | 2 | 2 | 3 | 4 |
| **e** | 5 | 4 | 4 | 3 | 3 | 3 | 2 | 3 |
| **y** | 6 | 5 | 5 | 4 | 4 | 4 | 3 | **2** |

# Calculating Edit Distance

|   | $ | b | r | i | t | n | e | y |
|---|---|---|---|---|---|---|---|---|
| **$** | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| **b** | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| **t** | 2 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| **i** | 3 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| **n** | 4 | 3 | 3 | 2 | 2 | 2 | 3 | 4 |
| **e** | 5 | 4 | 4 | 3 | 3 | 3 | 2 | 3 |
| **y** | 6 | 5 | 5 | 4 | 4 | 4 | 3 | 2 |

*substitute*

*delete*

From s='*btiney*' to r='*britney*'

Match s[1]='b' → *b*

Substitute s[2]='r' with 't'→*bt*

Match s[3] = 'i' → *bti*

Delete s[4]='t'→ *bti*

Match s[5]='n' → *btin*

Match s[6]='e' → *btine*

Match s[7]='y' → *btiney* = r

# What's the difference?

| Functions (normalization) | Edit Distance Edit Similarity | Overlap Size Jaccard Similarity |
|---|---|---|
| Input | Sequences | Sets |
| Example Representations | DNA, String, Time Series | Image, Document, Vector, Friend List |

# String Similarity Join

- Input:

  - A collection of strings $S$

  - A threshold $\tau$

- Output:

  - All string pairs $(s, r) \in S \times S \text{ such that } ED(s, r) \leq \tau$

# String Similarity Join

- Give threshold $\tau = 3$

| ID | Strings |
|----|---------|
| $s_1$ | vankatesh |
| $s_2$ | avataresha |
| $s_3$ | kaushic chaduri |
| $s_4$ | kaushik chakrab |
| $s_5$ | kaushuk chadhui |
| $s_6$ | caushik chakrabar |

$ED(s_1, s_2) = 5$   $ED(s_1, s_3) = 13$   $ED(s_1, s_4) = 12$   $ED(s_1, s_5) = 12$
$ED(s_1, s_6) = 14$   $ED(s_2, s_3) = 12$   $ED(s_2, s_4) = 12$   $ED(s_2, s_5) = 12$
$ED(s_2, s_6) = 14$   $ED(s_3, s_4) = 5$   $ED(s_3, s_5) = 4$   $ED(s_3, s_6) = 8$
$ED(s_4, s_5) = 4$   **$ED(s_4, s_6) = 3$**   $ED(s_5, s_6) = 8$

# Data Cleaning & Integration

## Relation with Duplicates

| ID | name | ZIP | Income |
|----|------|------|--------|
| P1 | Green | 51519 | 30k |
| P2 | Green | 51518 | 32k |
| P3 | Peter | 30528 | 40k |
| P4 | Peter | 30528 | 40k |
| P5 | Gree | 51519 | 55k |
| P6 | Chuck | 51519 | 30k |

# Challenges

$$O(n^2) \; pairs \; of \; strings.$$
$$1 \; million \; strings \; result \; in \; 1 \; trillion \; pairs \; !$$

$$O(|r||s|) \; time \; to \; calculate \; ED(r, s)$$

# Filter-and-Refine Framework

- Basic idea

  - **Filter** a large number of dissimilar string pairs

  - **Verify** the remaining potentially similar pairs


- Good Filter Condition

  - Efficient to check

  - Effective for pruning dissimilar pairs

# Length Filter

- ED($r, s$): The minimum number of edit operations (insertion/deletion/substitution) needed to transform $r$ to $s$.

- *Property:*  $\text{ED}(\boldsymbol{r}, \boldsymbol{s}) \geq \big||r|-|s|\big|$

  – it needs at least $\big||r|-|s|\big|$ deletions to just make $r$ and $s$ have the same length

# Applying the Length Filter

- Give threshold $\tau = 3$

- Pruning Condition:

$$\left| |s_i| - |s_j| \right| > 3$$

| ID | Strings | Length |
|----|---------|--------|
| $s_1$ | vankatesh | 9 |
| $s_2$ | avataresha | 10 |
| $s_3$ | kaushic chaduri | 15 |
| $s_4$ | kaushik chakrab | 15 |
| $s_5$ | kaushuk chadhui | 15 |
| $s_6$ | caushik chakrabar | 17 |

ED($s_1$ ,$s_2$)=5    ~~ED($s_1$ ,$s_3$)=13~~  ~~ED($s_1$ ,$s_4$)=12~~  ~~ED($s_1$ ,$s_5$)=12~~
~~ED($s_1$ ,$s_6$)=14~~  ~~ED($s_2$ ,$s_3$)=12~~  ~~ED($s_2$ ,$s_4$)=12~~  ~~ED($s_2$ ,$s_5$)=12~~
~~ED($s_2$ ,$s_6$)=14~~  ED($s_3$ ,$s_4$)=5    ED($s_3$ ,$s_5$)=4    ED($s_3$ ,$s_6$)=8
ED($s_4$ ,$s_5$)=4    **ED($s_4$ ,$s_6$)=3**    ED($s_5$ ,$s_6$)=8

# Partitioning Filter

- Give threshold $\tau = 1$

*hi*lton

**1**

huston

*hi* dose not appear in huston and needs at least 1 edit

# Partitioning Filter

- Give threshold $\tau = 1$

hi*lton*

**1**

huston

*minimum # of edit operations is 2. Prune!*

# Partitioning Filter

- Threshold $\tau$

*split r to $\tau$ +1 disjoint segments*



- String r

- String s

Is there any substring of $s$ matching a segment of $r$ ?

Yes        No

<r, s> are a candidate pair      <r, s> are dissimilar, prune

# How to Partition?

- Give threshold $\tau = 1$

*hil**ton***

**Match**

hus**ton**          *Candidate!*

# Partition Scheme

- Even Partition Scheme

    – Given $\tau = 3$, *"avataresha"* $\rightarrow$ {*"av"*, *"at"*, *"are"*, *"sha"*}

- Other Schemes

    – Select good partition strategies.

    – Adaptive partition scheme [Deng et al. 2012a].

# Challenge

*there are $|s|^2$ substrings in s*

how to reduce the number of substrings
to compare with for each segment?

# Length-based Method

- For each segment, only compare to the substrings with the same length as the segment.

$$r = \text{``vankatesh''} \quad \tau = 3 \quad s = \text{``avataresha''}$$

$r^1$ va $\longleftarrow$ av va at ta ar re es sh ha

$r^2$ nk $\longleftarrow$ av va at ta ar re es sh ha

$r^3$ at $\longleftarrow$ av va at ta ar re es sh ha

$r^4$ esh $\longleftarrow$ ava vat ata tar are res esh sha

# Length-based Method

- For any strings *r*, *s*, and $\tau$, the number of comparisons:

$$(\tau+1)(|s|+1)-|r|$$

- For $r = $ "vankatesh" and $s = $ "avataresha", the number is 35

# Shift-based Method



*Just transform $r_l$ to $s_l$ needs at least $\left||s_l|-|r_l|\right|$ edit operations*

**Pruning Condition:** $\left||s_l|-|r_l|\right| > \tau \longrightarrow |p - p_i| > \tau$

- For each segment $r^{\,i}$ with the start position $p_i$, only compares to the substrings with start positions $p$ in $[p_i - \tau, p_i + \tau]$

# Shift-based Method

- For each segment $r^i$ with the start position $p_i$, only compares to the substrings with start positions $p$ in $[p_i - \tau, p_i + \tau]$

$$r = \text{``vankatesh''} \quad \tau = 3 \quad s = \text{``avataresha''}$$

| | | |
|---|---|---|
| $r^1$ | va | ⟵ av va at ta |
| $r^2$ | nk | ⟵ av va at ta ar re |
| $r^3$ | at | ⟵ va at ta ar re es sh |
| $r^4$ | esh | ⟵ tar are res esh sha |

# Shift-based Method

- For any strings $r$, $s$, and $\tau$, the number of comparisons:

$$(\tau + 1)(2\tau + 1)$$

- For $r = $ "vankatesh" and $s = $ "avataresha", the number is 22.

# Position-aware Method

$r = $ "vank[a]tesh" $\implies \{$va, nk, (at), esh$\}$

$r_l \qquad r_r$

$s = $ "av[at]aresha"

$s_l \qquad s_r$

*Transform $r_l$ to $s_l$, match "at", and transform $r_r$ to $s_r$*

$$\Big| |s_l| - |r_l| \Big| + \Big| |s_r| - |r_r| \Big| = 2 + 3 > \tau = 3$$

# Position-aware Method



*Transform $r_l$ to $s_l$, match $r_m$ and $s_m$ then transform $r_r$ to $s_r$*

**Pruning Condition:** $\big||s_l|-|r_l|\big| + \big||s_r|-|r_r|\big| > \tau$

- For each segment $r^i$ with the start position $p_i$, only compare to the substrings with start position in $[P_i - \big\lfloor\frac{\tau-\Delta}{2}\big\rfloor, P_i + \big\lfloor\frac{\tau+\Delta}{2}\big\rfloor]$ where $\Delta = |s| - |r|$

# Position-aware Method

- For each segment $r^i$ with the start position $p_i$, only compare to the substrings with start position in $[P_i - \lfloor \frac{\tau-\Delta}{2} \rfloor, P_i + \lfloor \frac{\tau+\Delta}{2} \rfloor]$ where $\Delta = |s| - |r|$

$\tau = 3$  $r = $ "vankatesh"  $s = $ "avataresha"

$r^1$  va  ⟵  av  va  at

$r^2$  nk  ⟵  va  at  ta  ar

$r^3$  at  ⟵  ta  ar  re  es

$r^4$  esh ⟵  res  esh  sha

# Position-aware Method

- For any strings $r$, $s$, and $\tau$, the number of comparisons:

$$(\tau + 1)^2$$

- For $r =$ "vankatesh" and $s =$ "avataresha", the number is 14.
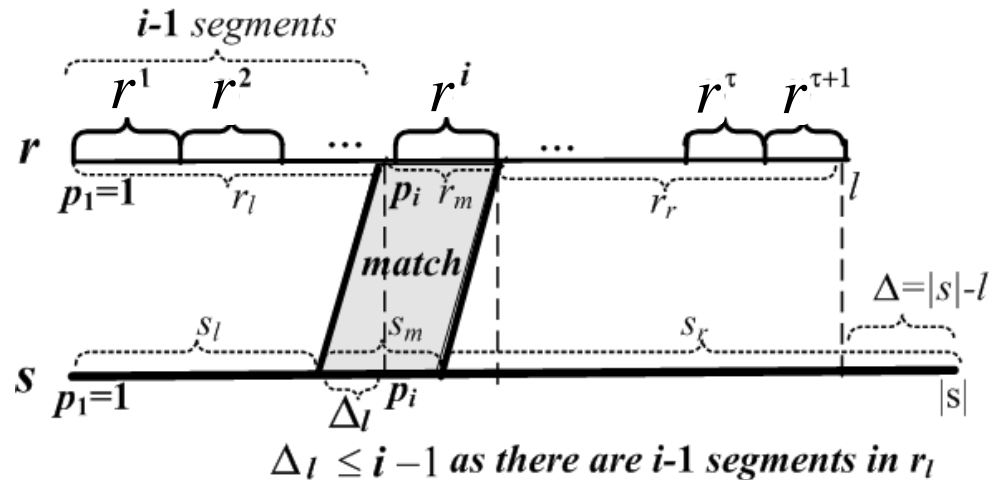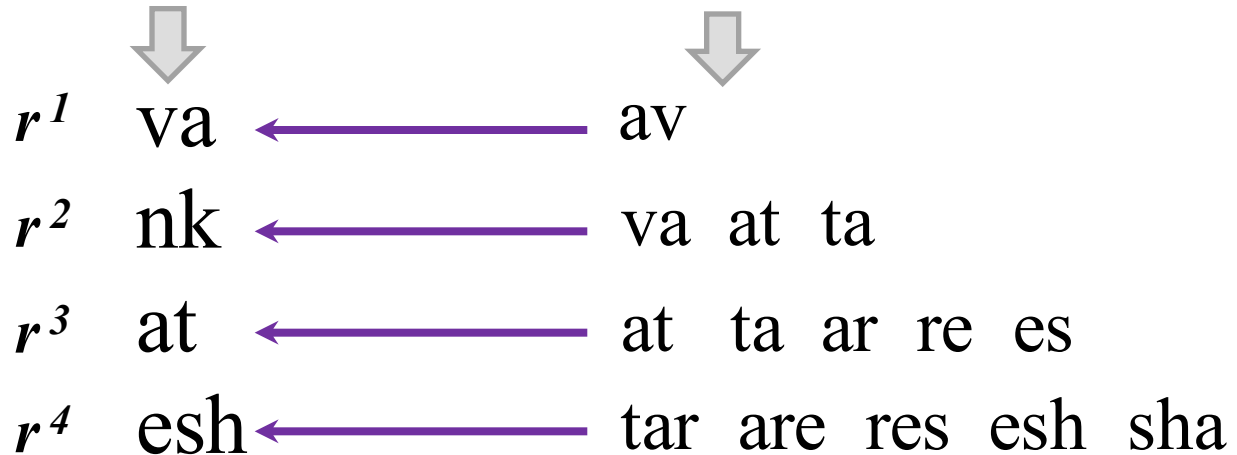
# Multi-match-aware Method

## -- Left-side Perspective

$r = \text{``vankatesh''} \implies \{va, nk, at, esh\}$

$r_l = \text{`` ''}$

$\tau = 3$

3 unvisited segments

$s = \text{``avataresha''}$

$s_l = \text{``a''}$

*If none of the 3 unvisited segments having matching substring in s, then r and s cannot be similar*

$1 \quad 1 \quad 1$

$+$

$||s_l| - |r_l|| = 1$

$= 4 > \tau$

*Thus we can safely skip the current matching segment and look for the next*

# Multi-match-aware Method

## *-- Left-side Perspective*



$\Delta_l \leq i-1$ **as there are i-1 segments in $r_l$**

**Pruning Condition:** $\left| |s_l| - |r_l| \right| +$ (# of unvisited segments) $> \tau$

- For each segment $r^i$ with the start position $p_i$, only compare to the substrings with start position in $[P_i - (i-1), P_i + (i-1)]$

# Multi-match-aware Method

## -- Left-side Perspective

- For each segment $r^i$ with the start position $p_i$, only compare to the substrings with start position in $[P_i - (i - 1), P_i + (i - 1)]$

$\tau = 3$ $r = $ "vankatesh" $s = $ "avataresha"

| | | |
|---|---|---|
| $r^1$ | va ⟵ | av |
| $r^2$ | nk ⟵ | va at ta |
| $r^3$ | at ⟵ | at ta ar re es |
| $r^4$ | esh ⟵ | tar are res esh sha |

# Multi-match-aware Method

### -- Left-side Perspective

- For any strings $r$, $s$, and $\tau$, the number of comparisons:

$$\tau^2 + 2\tau$$

- For $r = $ "vankatesh" and $s = $ "avataresha", the number is 14.

# Multi-match-aware Method

## *-- Right-side Perspective*



$$\Delta_r \le \tau + 1 - i \text{ as there are } \tau + 1 - i \text{ segments in } r_r$$

**Pruning Condition:** (# unvisited segments)$+\Big|\,|s_r|\text{-}|r_r|\,\Big| > \tau$

- For each segment $r^i$ with the start position $p_i$, only compare to the substrings with start position in $[P_i + \Delta - (\tau + 1 - i), P_i + \Delta + (\tau + 1 - i)]$

# Multi-match-aware Method

- Interestingly, we can apply the multi-match-aware method from left- and right-side perspectives simultaneously.

- For each segment $r^i$ with start position $p_i$, only compare to the substrings with start position in

$$[max(P_i - (i - 1), P_i + \Delta - (\tau + 1 - i)),$$
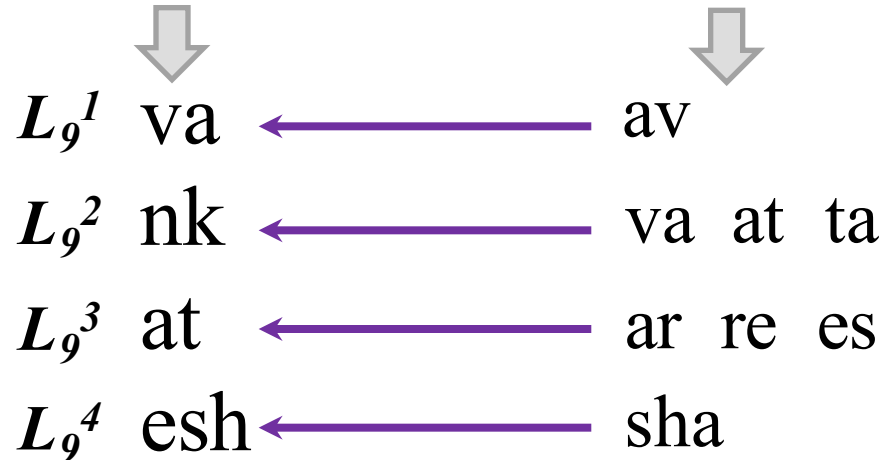$$min(P_i + (i - 1), P_i + \Delta + (\tau + 1 - i))]$$

# Multi-match-aware Method

- For each segment $r^i$ with start position $p_i$, only compare to the substrings with start position in $[max(P_i - (i-1), P_i + \Delta - (\tau + 1 - i)), min(P_i + (i-1), P_i + \Delta + (\tau + 1 - i))]$

$\tau = 3$   $r =$ "vankatesh"   $s =$ "avataresha"

| | | |
|---|---|---|
| $L_9^1$ va | ⟵ | av |
| $L_9^2$ nk | ⟵ | va  at  ta |
| $L_9^3$ at | ⟵ | ar  re  es |
| $L_9^4$ esh | ⟵ | sha |

# Multi-match-aware Method

- For any strings $r$, $s$, and $\tau$, the number of comparisons:

$$\left\lceil \left| \frac{\tau^2 - \Delta^2}{2} \right| \right\rceil + \tau + 1$$

- For $r = $ "vankatesh" and $s = $ "avataresha", the number is 8.

# Theoretical Results

- The number of comparisons by the multi-match-aware method is **minimum** while guarantees completeness

- For any $s,\ r$ and $\tau$,

$$W_{multi\text{-}match}(s,\ r,\ \tau)\ \subseteq\ W_{position}(s,\ r,\ \tau)\ \subseteq\ W_{shift}(s,\ r,\ \tau)\ \subseteq\ W_{length}(s,\ r,\ \tau)$$

# *Applying Partitioning Filter*

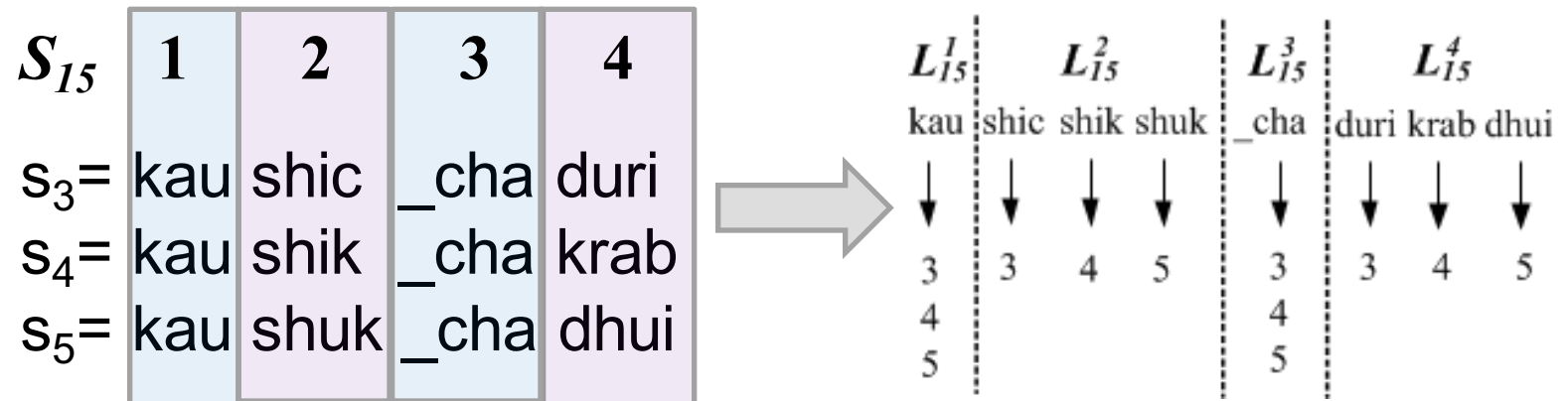# Applying Partitioning Filter

- 1. Group all the strings by length: **S**

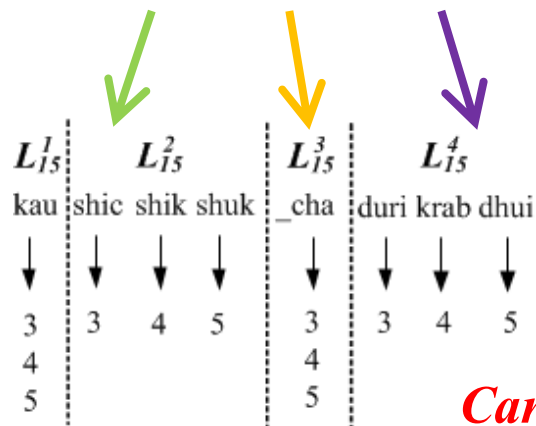| ID | Strings | Length | |
|----|---------|--------|----|
| $s_1$ | vankatesh | 9 | $S_9$ |
| $s_2$ | avataresha | 10 | $S_{10}$ |
| $s_3$ | kaushic chaduri | 15 | |
| $s_4$ | kaushik chakrab | 15 | $S_{15}$ |
| $s_5$ | kaushuk chadhui | 15 | |
| $s_6$ | caushik chakrabar | 17 | $S_{17}$ |

# Applying Partitioning Filter

- 2. For each group $S_l$, partition its strings into $\tau + 1$ segments and build $\tau + 1$ inverted indexes $L_l^i$

# Applying Partitioning Filter

- 3. For each string **s** and index $L_l^i$, select substrings from **s** based on the partitioning filter to get candidates:

$s_6$= cau shik _cha krabar

$$[max(P_i - (i - 1), P_i + \Delta - (\tau + 1 - i)),$$
$$min(P_i + (i - 1), P_i + \Delta + (\tau + 1 - i))]$$



*the segments in $L_l^i$ have the same start positions $p_i$ and same string lengths l*

**Candidates:** $\langle s_3, s_6 \rangle \langle s_4, s_6 \rangle \langle s_5, s_6 \rangle$

# Applying Partitioning Filter

- 4. Verify the candidates

*Candidates*:  $\langle s_3, s_6 \rangle \langle s_4, s_6 \rangle \langle s_5, s_6 \rangle$

$ED(s_3, s_6) > 3$ ✗

$ED(s_4, s_6) = 3$ ✓

$ED(s_5, s_6) > 3$ ✗

# *Improving Verification*

# Calculating Edit Distance

|   | $ | b | r | i | t | n | e | y |
|---|---|---|---|---|---|---|---|---|
| **$** | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| **b** | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| **t** | 2 | 1 | 1 | 2 | 3 | 4 | 5 | 6 |
| **i** | 3 | 2 | 2 | 1 | 2 | 3 | 4 | 5 |
| **n** | 4 | 3 | 3 | 2 | 2 | 2 | 3 | 4 |
| **e** | 5 | 4 | 4 | 3 | 3 | 3 | 2 | 3 |
| **y** | 6 | 5 | 5 | 4 | 4 | 4 | 3 | **2** |

# Verification

|     | $   | b   | r   | i   | t   | n   | e   | y   |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| $   | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   |
| b   | 1   | 0   | 1   | 2   | 3   | 4   | 5   | 6   |
| t   | 2   | 1   | 1   | 2   | 3   | 4   | 5   | 6   |
| i   | 3   | 2   | 2   | 1   | 2   | 3   | 4   | 5   |
| n   | 4   | 3   | 3   | 2   | 2   | 2   | 3   | 4   |
| e   | 5   | 4   | 4   | 3   | 3   | 3   | 2   | 3   |
| y   | 6   | 5   | 5   | 4   | 4   | 4   | 3   | 2   |

$$ED(r[1,i],s[1,j]) \geq |i - j|$$

# Verification

|   | $ | b | r | i | t | n | e | y |
|---|---|---|---|---|---|---|---|---|
| $ | 0 | 1 | 2 |   |   |   |   |   |
| b | 1 | 0 | 1 | 2 |   |   |   |   |
| t | 2 | 1 | 1 | 2 | 3 |   |   |   |
| i |   | 2 | 2 | 1 | 2 | 3 |   |   |
| n |   |   | 3 | 2 | 2 | 2 | 3 |   |
| e |   |   |   | 3 | 3 | 3 | 2 | 3 |
| y |   |   |   |   | 4 | 4 | 3 | 2 |

# **Verification**

*only need to calculate a band of width 2τ+1*

|     | **$** | **b** | **r** | **i** | **t** | **n** | **e** | **y** |
|-----|-------|-------|-------|-------|-------|-------|-------|-------|
| **$** | 0 | 1 | 2 |   |   |   |   |   |
| **b** | 1 | 0 | 1 | 2 |   |   |   |   |
| **t** | 2 | 1 | 1 | 2 | 3 |   |   |   |
| **i** |   | 2 | 2 | 1 | 2 | 3 |   |   |
| **n** |   |   | 3 | 2 | 2 | 2 | 3 |   |
| **e** |   |   |   | 3 | 3 | 3 | 2 | 3 |
| **y** |   |   |   |   | 4 | 4 | 3 | 2 |

# Early Termination



*all derived cells must have values larger than $\tau$.*

# Length-aware Verification



*total length difference is 4 > τ. Thus no need to calculate M[2][5].*

# Length-aware Verification



*only need to calculate a band of width $\tau$*

# Extension-base Verification

# Extension-base Verification

- We can verify a candidate pair using tighter thresholds:
  - For the left parts we can set $\tau_l = i - 1$
  - For the right parts we can set $\tau_r = \tau + 1 - i$

*band widths are $\tau_l < \tau$ and $\tau_r < \tau$*

# *Alternative: Prefix Filter*

# *q*-gram

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

$$youtbecom$$

2-grams $\begin{cases} yo \\ ou \\ ut \\ tb \\ be \\ ec \\ co \\ om \end{cases}$

# *q*-gram

- 1 edit operation destroies at most $q$ grams.

$$youtdecom$$

$$yo$$
$$ou$$
$$ut$$
$$td$$
$$de$$
$$ec$$
$$co$$
$$om$$

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

# Prefix Filter

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



*Prefix(r)*

*Prefix size: qτ+1*

The sorted *q*-grams

*Prefix(s)*

*Prefix Filter:*
*If prefix(r) ∩ prefix(s) = ϕ, ED(r,s) > τ*

# Disjoint q-gram
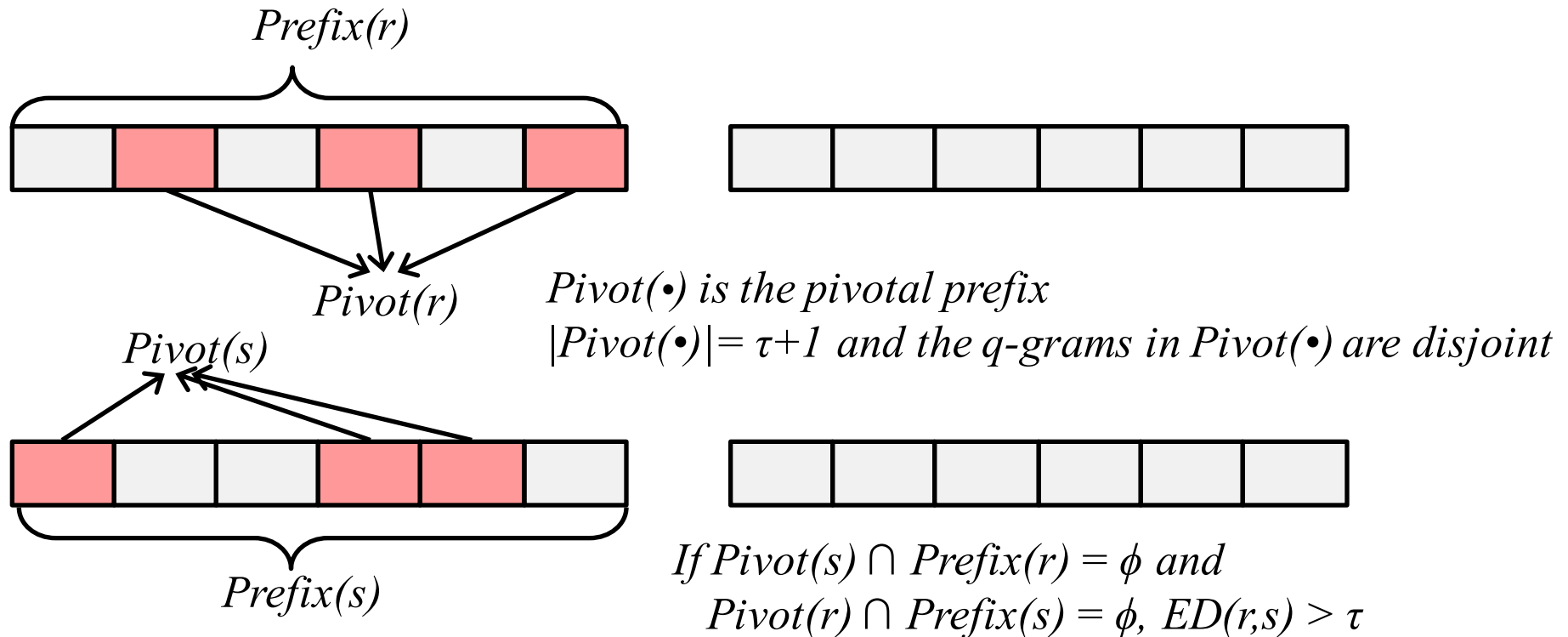
- One edit operation destroies at most *1* disjoint gram.

$$yout\ decom$$

$$yo$$

$$ut$$

$$de$$

$$om$$

- $\tau$ edit operations destroy *at most* $\tau$ disjoint q-grams.

- if r and s have more than $\tau$ mismatch disjoint q-grams, ED(r, s)> $\tau$

# Pivotal Prefix Filter

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

*Prefix(r)*



*Pivot(r)*

*Pivot(s)*

*Prefix(s)*

*Pivot(•) is the pivotal prefix*
*$|Pivot(•)| = \tau+1$ and the q-grams in Pivot(•) are disjoint*

*If Pivot(s) $\cap$ Prefix(r) = $\phi$ and*
*Pivot(r) $\cap$ Prefix(s) = $\phi$, ED(r,s) > $\tau$*

# Takeaways

(1)  The partitioning filter for Edit Distance

(2)  The multi-match-aware method

(3)  The extension-based verification

(4)  The prefix filters for Edit Distance

# References

- Pass-Join: A Partition based Method for Similarity Joins. G. Li, D. Deng, J. Wang, J. Feng. VLDB 2012.

- A Pivotal Prefix Based Filtering Algorithm for String Similarity Search. D. Deng, G. Li, J. Feng. SIGMOD 2014.

- Ed-Join: An Efficient Algorithm for Similarity Joins with Edit Distance Constraints. C. Xiao, W. Wang, X. Lin. VLDB 2008