# Similarity Search The q-Gram Distance

#### Nikolaus Augsten

nikolaus.augsten@plus.ac.at Department of Computer Science University of Salzburg



WS 2025/26

Version November 6, 2025

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

1 / 44

/ 44

#### Filters for the Edit Distance

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

2 / 44

Filters for

Filters for the Edit Distance Motivation

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- 3 Conclusion

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 3 / 44

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

4 / 44

Filters for the Edit Distance Motivation

## Application Scenario

- Scenario:
  - A company offers a number of services on the Web.
  - You can subscribe for each service independently.
  - Each service has its own database (no unique key across databases).
- Example: customer tables of two different services:

Α			В			
ID	name		ID	name		
1023	Frodo Baggins		948483	John R. R. Tolkien		
21	J. R. R. Tolkien		153494	C. S. Lewis		
239	C.S. Lewis		494392	Fordo Baggins		
863	Bilbo Baggins		799294	Biblo Baggins		

Task: Created unified customer view!

Augsten (Univ. Salzburg)

WS 2025/26

Filters for the Edit Distance Motivation

## Effectiveness and Efficiency of the Approximate Join

• Effectiveness: Join result for k = 3:

1		
name	ID	name
Frodo Baggins	494392	Fordo Baggins
J. R. R. Tolkien	948483	John R. R. Tolkien
C.S. Lewis	153494	C. S. Lewis
Bilbo Baggins	799294	Biblo Baggins
	Frodo Baggins J. R. R. Tolkien C.S. Lewis	Frodo Baggins 494392 J. R. R. Tolkien 948483 C.S. Lewis 153494

- $\Rightarrow$  very good (100% correct)
- Efficiency: How does the DB evaluate the query?
  - (1) compute  $A \times B$
  - (2) evaluate UDF on each tuple  $t \in A \times B$
- Prohibitive runtime!

Filters for the Edit Distance Motivation

## The Join Approach

• Solution: Join customer tables on name attribute (Q1):

```
SELECT * FROM A,B
WHERE A.name = B.name
```

- Exact Join: Does not work!
- Similarity Join: Allow k errors...
  - (1) Register UDF (User Defined Function) for the edit distance:

returns the union cost edit distance between the strings x and y.

(2) Rewrite query Q1 as similarity join (Q2):

```
SELECT * FROM A,B
WHERE ed(A.name, B.name) <= k
```

Augsten (Univ. Salzburg)

WS 2025/26

Filters for the Edit Distance Motivation

## Using a Filter for Search Space Reduction

- Search space:  $A \times B$  ( $\Rightarrow |A| \cdot |B|$  edit distance computations)
- Filtering (Pruning): Remove tuples that can not match, without actually computing the distance.

Augsten (Univ. Salzburg) WS 2025/26 WS 2025/26 Similarity Search Augsten (Univ. Salzburg) Similarity Search

Filters for the Edit Distance Lower Bound Filters

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- Conclusion

Augsten (Univ. Salzburg)

Similarity Search

Filters for the Edit Distance Lower Bound Filters

WS 2025/26

9 / 44

#### Lower Bound Filters

• Lower bound (lb) for distance dist(x, y):

$$dist(x, y) \ge lb_{dist}(x, y)$$

• Query Q3 with Lower Bound Filter:

- 1b(A.name, B.name) is a cheap function
- database will optimize query: compute ed(A.name, B.name) only if  $lb(A.name, B.name) \leq k$
- No false negatives!

Filters for the Edit Distance Lower Bound Filters

## Filter Properties

• Error Types:

#### Correct Result

		positive	negative
Filter	positive	true positive	false positive
Test	negative	false negative	true negative

- Example: "Are x and y within edit distance k?"
  - Correct result: compute edit distance and test ed(x, y) < k
  - Filter test: give answer without computing edit distance
  - False negatives: x and y are pruned although ed(x, y) < k.
  - False positives: x and y are not pruned although  $ed(x, y) \leq k$ .
- Good filters have
  - no false negatives (i.e., miss no correct results)
  - few false positive (i.e., avoid unnecessary distance computations)

Augsten (Univ. Salzburg)

Similarity Search

Filters for the Edit Distance Length Filter

WS 2025/26

10 / 44

12 / 44

Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- The g-Gram Distance
- Conclusion

11 / 44 Augsten (Univ. Salzburg) Similarity Search WS 2025/26 Filters for the Edit Distance Length Filter

## Length Filtering

#### Theorem (Length Filtering [GIJ+01])

If two strings x and y are within edit distance k, their lengths cannot differ by more than k:

$$\operatorname{ed}(x,y) \ge \operatorname{abs}(|x| - |y|)$$

- Proof: At least abs(|x|-|y|) inserts are needed to bring x and y to the same length.
- Query Q4 with Length Filtering:

```
SELECT * FROM A,B
WHERE ABS(LENGTH(A.name)-LENGTH(B.name)) <= k AND
      ed(A.name, B.name) <= k
```

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Filters for the Edit Distance q-Grams: Count Filter

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- Conclusion

Filters for the Edit Distance Length Filter

## Example: Length Filtering

• Execute query without/with length filter (k = 3):

Α			В			
ID	name		ID	name		
1023	Frodo Baggins <sub>13</sub>		948483	John R. R. Tolkien <sub>18</sub>		
21	J. R. R. Tolkien <sub>16</sub>		153494	C. S. Lewis <sub>11</sub>		
239	C.S. Lewis <sub>10</sub>		494392	Fordo Baggins <sub>13</sub>		
863	Bilbo Baggins <sub>13</sub>		799294	Biblo Baggins <sub>13</sub>		

- Without length filter: 16 edit distance computations
- With length filter (k = 3): 12 edit distance computations
  - $\bullet$  J. R. R. Tolkien  $\leftrightarrow$  C. S. Lewis is pruned
  - all pairs (..., John R. R. Tolkien) except (J. R. R. Tolkien, John R. R. Tolkien) are pruned

Augsten (Univ. Salzburg)

Similarity Search

Filters for the Edit Distance q-Grams: Count Filter

WS 2025/26

What is a q-Gram?

- Intuition:
  - slide window of length q over string  $x \in \Sigma^*$
  - characters covered by window form a q-gram
  - where window extends string: fill with dummy character #  $\notin \Sigma$
- Example: x = Frodo, q = 3

```
extended: ##Frodo##
q-grams: ##F
         #Fr
          Fro
            rod
             odo
```

• q-Gram Profile  $G_x$ : bag of all q-grams of x

d o# 0##

• Profile size:  $|G_x| = |x| + q - 1$ 

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

15 / 44

#### Single Edit Operations and Changing *q*-Grams

- Intuition: Strings within small edit distance share many q-grams.
- How many q-grams (q = 3) change/remain?

X		$ G_x $	y	$ G_y $	$ G_x \cap G_y $
pete	er	7	meter	7	4
pete	er	7	peters	8	5
pete	er	7	peer	6	4

 $\bullet$  ed $(x, y) = 1 \Rightarrow |G_x \cap G_y| = \max(|G_x|, |G_y|) - a$ 

Augsten (Univ. Salzburg)

Filters for the Edit Distance q-Grams: Count Filter

WS 2025/26

# Theorem (Count Filtering [GIJ+01])

Consider two strings x and y with the q-gram profiles  $G_x$  and  $G_y$ , respectively. If x and y are within edit distance k, then the cardinality of the a-gram profile intersection is at least

$$|G_x \cap G_y| \ge \max(|G_x|, |G_y|) - kq$$

- Proof (by induction):
  - true for k = 1:  $|G_x \cap G_y| \ge \max(|G_x|, |G_y|) q$
  - $k \rightarrow k + 1$ : each additional edit operation changes at most q *q*-grams.

## Multiple Edit Operations and Changing q-Grams

- $\bullet$  ed $(x, y) = 1 \Rightarrow |G_x \cap G_y| = \max(|G_x|, |G_y|) q$
- What if ed(x, y) = k > 1?

X		$ G_x $	y	$ G_y $	$ G_x \cap G_y $
pete	er	7	meters	8	2
pete	er	7	petal	7	3

• Multiple edit operations may affect the same q-gram:

$$peter \rightarrow G_x = \{ \#p, \#pe, pet, ete, ter, er\#, r\#\# \}$$

$$petal \rightarrow G_y = \{ \#p, \#pe, pet, eta, tal, al\#, l\#\# \}$$

• Each edit operation affects at most q q-grams.

Augsten (Univ. Salzburg)

Filters for the Edit Distance q-Grams: Count Filter

WS 2025/26

Implementation of q-Grams

- Given: tables A and B with schema (id, name)
  - *id* is the key attribute
  - name is string-valued
- Compute auxiliary tables QA and QB with schema (id, qgram):
  - each tuple stores one *q*-gram
  - string x of attribute name is represented by its |x| + q 1 q-grams
  - QA.id is the key value (A.id) of a tuple with A.name = x
  - QA.qgram is one of the q-grams of x
- Example:

QA		
am		
F		
r		
J		

Filters for the Edit Distance q-Grams: Count Filter

#### Count Filtering Query

Query Q5 with Count Filtering:

```
SELECT
        A.id, B.id, A.name, B.name
FROM
         A, QA, B, QB
WHF.R.F.
         A.id = QA.id AND
         B.id = QB.id AND
         QA.qgram = QB.qgram AND
         ABS(LENGTH(A.name)-LENGTH(B.name)) <= k
GROUP BY A.id, B.id, A.name, B.name
        COUNT(*) >= LENGTH(A.name)-1-(k-1)*q AND
HAVING
         COUNT(*) >= LENGTH(B.name)-1-(k-1)*q AND
         ed(A.name,B.name) <= k
```

Augsten (Univ. Salzburg)

Filters for the Edit Distance q-Grams: Count Filter

WS 2025/26

#### Fixing Count Filtering Query

- Fix query to avoid false negatives [GIJ+03]:
  - Join pairs (x, y) with  $kq \ge \max(|G_x|, |G_y|)$  using only length filter.
  - Union results with results of previous query Q5.
- Query Q6 without false negatives (extends previous query Q5):

UNION SELECT A.id, B.id, A.name, B.name FROM A, B WHERE LENGTH(A.name)+ $q-1 \le k*q$  AND  $LENGTH(B.name)+q-1 \le k*q AND$ ABS(LENGTH(A.name) - LENGTH(B.name)) <= k AND ed(A.name,B.name) <= k

• Note: We omit this part in subsequent versions of the query since it remains unchanged.

Filters for the Edit Distance q-Grams: Count Filter

## Problem with Count Filtering Query

- Previous guery Q5 works fine for  $kq < \max(|G_x|, |G_y|)$ .
- However: If  $kq \ge \max(|G_x|, |G_v|)$ , no q-grams may match even if  $ed(x, y) \le k$ .
- Example (q = 3, k = 2):

WHERE-clause prunes x and y, although  $ed(x, y) \le k$ 

- False negatives:
  - short strings with respect to edit distance (e.g., |x| = 3, k = 3)
  - even if within given edit distance, matches tend to be meaningless (e.g., abc and xyz are within edit distance k = 3)

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

24 / 44

Filters for the Edit Distance q-Grams: Position Filtering

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- Conclusion

Augsten (Univ. Salzburg) WS 2025/26 23 / 44 WS 2025/26 Similarity Search Augsten (Univ. Salzburg) Similarity Search

Filters for the Edit Distance q-Grams: Position Filtering

#### Positional q-Grams

- Enrich *q*-grams with position information:
  - extended string: prefix and suffix string x with q-1 characters #
  - slide window of length q over extended string x'
  - characters covered by window after shifting it *i* times form the *q*-gram at position i + 1
- Example: x = Frodo

```
extended string:
                        ##Frodo##
positional q-grams: (1, \# F)
                      (2, \# Fr)
                        (3.Fro)
                          (4,rod)
                           (5, odo)
                             (6,d o \#)
                              (7.0 # #)
```

Augsten (Univ. Salzburg)

Similarity Search

Filters for the Edit Distance q-Grams: Position Filtering

WS 2025/26

## Corresponding *q*-Grams

- Corresponding *q*-gram:
  - Given: positional q-grams (i, g) of x
  - transform x to y applying edit operations
  - (i,g) "becomes" (i,g) in y
  - We define: (i,g) corresponds to (i,g)
- Example:
  - x' = #abaZabaabaaba##, y' = #abaabaabaabaaba##
  - edit distance is 1 (delete Z from x)
  - (7, aba) in x corresponds to (6, aba) in y
  - ... but not to (9, aba)

Filters for the Edit Distance q-Grams: Position Filtering

## Computing Positional q-Grams in SQL

- Given: table N
  - N has a single attribute i
  - N is filled with numbers from 1 to max (max is the maximum string length plus q-1)
- Positional q-grams for table A in SQL (Q7):

```
CREATE TABLE QA AS
       SELECT A.id, N.i AS pos,
           SUBSTRING (CONCAT (
               SUBSTRING('\#..\#', 1, q - 1),
               LOWER(A.name),
               SUBSTRING('#..#', 1, q - 1)),
           N.i, q) AS ggram
       FROM A, N
       WHERE N.i <= LENGTH(A.name) + q - 1
```

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Filters for the Edit Distance q-Grams: Position Filtering

# Position Filtering

#### Theorem (Position Filtering [GIJ+01])

If two strings x and y are within edit distance k, then a positional q-gram in one cannot correspond to a positional q-gram in the other that differs from it by more then k positions.

- Proof:
  - each increment (decrement) of a position requires an insert (delete);
  - a shift by k positions requires k inserts/deletes.

Augsten (Univ. Salzburg) WS 2025/26 WS 2025/26 Similarity Search Augsten (Univ. Salzburg) Similarity Search

Filters for the Edit Distance q-Grams: Position Filtering

## Position Filtering

• Query Q8 with Count and Position Filtering:

A.id, B.id, A.name, B.name

FROM A, QA, B, QB

WHERE A.id = QA.id AND

B.id = QB.id AND

QA.qgram = QB.qgram AND

ABS(LENGTH(A.name)-LENGTH(B.name)) <= k AND

ABS(QA.pos-QB.pos)<=k

GROUP BY A.id, B.id, A.name, B.name

COUNT(\*) >= LENGTH(A.name)-1-(k-1)\*q ANDHAVING

COUNT(\*) >= LENGTH(B.name)-1-(k-1)\*q AND

ed(A.name, B.name) <= k

Augsten (Univ. Salzburg)

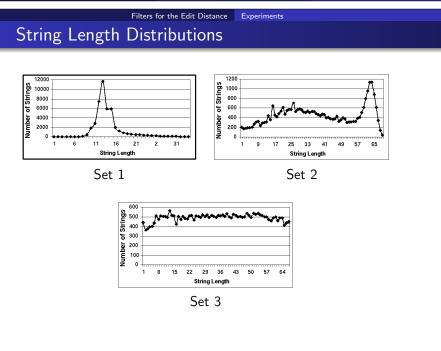
Similarity Search Filters for the Edit Distance Experiments WS 2025/26

## Experimental Data

- All experimental results taken from [GIJ<sup>+</sup>01]
- Three string data sets:
  - set1: 40K tuples, average length: 14 chars
  - set2: 30K tuples, average length: 38 chars
  - set3: 30K tuples, average length: 33 chars

Filters for the Edit Distance Experiments Outline Filters for the Edit Distance Motivation Lower Bound Filters Length Filter • q-Grams: Count Filter • q-Grams: Position Filtering Experiments Conclusion

Similarity Search



Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

31 / 44

Augsten (Univ. Salzburg)

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

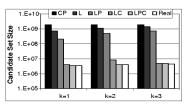
WS 2025/26

32 / 44

Filters for the Edit Distance Experiments

#### Candidate Set Size

- Question: How many edit distances do we have to compute?
- Show candidate set size for different filters (small is good).
- q = 2
- Caption:
  - CP: cross product
  - L: length filtering, P: position filtering, C: count filtering
  - Real: number of real matches



Set 1

Augsten (Univ. Salzburg)

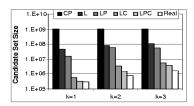
WS 2025/26

#### Candidate Set Size

- Question: How many edit distances do we have to compute?
- Show candidate set size for different filters (small is good).

Filters for the Edit Distance Experiments

- q = 2
- Caption:
  - CP: cross product
  - L: length filtering, P: position filtering, C: count filtering
  - Real: number of real matches

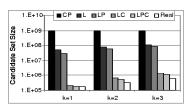


Set 3

Filters for the Edit Distance Experiments

#### Candidate Set Size

- Question: How many edit distances do we have to compute?
- Show candidate set size for different filters (small is good).
- q = 2
- Caption:
  - CP: cross product
  - L: length filtering, P: position filtering, C: count filtering
  - Real: number of real matches



Set 2

Augsten (Univ. Salzburg

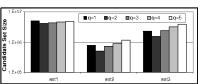
Filters for the Edit Distance Experiments

WS 2025/26

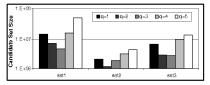
33 / 44

## Various q-Gram Lengths

- Question: How does the choice of *q* influence the filter effectiveness?
- Show candidate set size for different q values (small is good).



Edit Distance Threshold k = 2



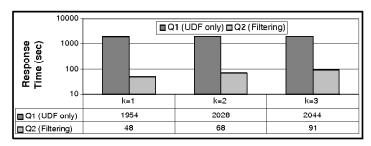
Edit Distance Threshold k = 3

Augsten (Univ. Salzburg) WS 2025/26 33 / 44 Augsten (Univ. Salzburg) WS 2025/26 34 / 44 Similarity Search Similarity Search

Filters for the Edit Distance Experiments

## Response Time

- Approximate self-join on small sample of 1000 tuples (set 1) (full dataset > 3 days without filters!)
- Measure response time (small is good).
- Caption:
  - k: edit distance threshold
  - Q1: edit distance without filters
  - Q2: edit distance with filters



Augsten (Univ. Salzburg)

WS 2025/26

The g-Gram Distance

The g-Gram Distance

WS 2025/26

The q-Gram Distance

## The q-Gram Distance

#### Definition (q-Gram Distance [Ukk92])

Let  $G_x$  and  $G_y$  be the q-gram profiles of the strings x and y, respectively. The *q-gram distance* between two strings is the number of *q*-grams in  $G_x$ and  $G_v$  that have no match in the other profile,

$$\mathsf{dist}_q(x,y) = |G_x \uplus G_y| - 2|G_x \cap G_y|.$$

• Example: q = 2, x = abab, y = abcab

$$G_x = \{\text{\#a}, \text{ab}, \text{ba}, \text{ab}, \text{b\#}\}$$

$$G_{v} = \{ \text{#a, ab, bc, ca, ab, b#} \}$$

$$G_x \uplus G_y = \{\text{\#a,ab,ba,ab,b\#,\#a,ab,bc,ca,ab,b\#}\}$$

$$G_x \cap G_y = \{\text{#a, ab, ab, b#}\}$$

$$dist_{a}(x, y) = |G_{x} \uplus G_{y}| - 2|G_{x} \cap G_{y}| = 11 - 2 \cdot 4 = 3$$

#### Outline

- Filters for the Edit Distance
  - Motivation
  - Lower Bound Filters
  - Length Filter
  - q-Grams: Count Filter
  - q-Grams: Position Filtering
  - Experiments
- 2 The q-Gram Distance
- Conclusion

Augsten (Univ. Salzburg)

Pseudo Metric q-Gram Distance

• The q-gram distance is a pseudo metric: For all  $x, y, z \in \Sigma^*$ 

- $\operatorname{dist}_q(x,y) + \operatorname{dist}_q(y,z) \ge \operatorname{dist}_q(x,z)$  (triangle inequality)
- $\operatorname{dist}_{a}(x, y) = \operatorname{dist}_{a}(y, x)$  (symmetric)
- $\operatorname{dist}_{a}(x, y) = 0 \Leftarrow x = y$
- Note: Identity condition relaxed: dist<sub>a</sub> $(x, y) = 0 \Rightarrow x = y$ i.e., the q-gram distance between two different strings can be 0
- Example:

$$\begin{aligned} & \mathsf{dist}_q(\mathsf{axybxycxyd}, \mathsf{axycxybxyd}) = 0 \\ & G_x = G_y = \{ \texttt{\##a}, \texttt{\#ax}, \mathsf{axy}, \mathsf{xyb}, \mathsf{ybx}, \mathsf{bxy}, \mathsf{xyc}, \mathsf{ycx}, \mathsf{cxy}, \mathsf{xyd}, \mathsf{yd\#}, \mathsf{d\#\#} \} \end{aligned}$$

Augsten (Univ. Salzburg) WS 2025/26 WS 2025/26 38 / 44 Similarity Search Augsten (Univ. Salzburg) Similarity Search

The g-Gram Distance

## Distance Normalization (1/3)

• What is a good threshold?

```
ed(International Business Machines Corporation,
    International Bussiness Machine Corporation) = 2
ed(IBM, BMW) = 2
ed(Int. Business Machines Corp.,
    International Business Machines Corporation) = 17
```

- Problem: Absolute numbers not always meaningful...
- Solution: Compute error relative to string length!

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

39 / 4

The q-Gram Distance

#### Distance Normalization (3/3)

Normalized edit distance:

```
\label{eq:corporation} \begin{array}{ll} \text{norm-ed(International Business Machines Corporation}, \\ & \text{International Bussiness Machine Corporation}) = 0.047 \\ \text{norm-ed(IBM, BMW)} = 0.66 \\ \text{norm-ed(Int. Business Machines Corp.}, \\ & \text{International Business Machines Corporation}) = 0.4 \\ \end{array}
```

• Normalized q-gram distance (q = 3):

```
\begin{array}{c} \mathsf{norm\text{-}dist}_q(\mathsf{International}\ \mathsf{Business}\ \mathsf{Machines}\ \mathsf{Corporation}, \\ \mathsf{International}\ \mathsf{Bussiness}\ \mathsf{Machine}\ \mathsf{Corporation}) = 0.089 \\ \mathsf{norm\text{-}dist}_q(\mathsf{IBM}, \mathsf{BMW}) = 1.0 \\ \mathsf{norm\text{-}dist}_q(\mathsf{Int}.\ \mathsf{Business}\ \mathsf{Machines}\ \mathsf{Corp.}, \\ \mathsf{International}\ \mathsf{Business}\ \mathsf{Machines}\ \mathsf{Corporation}) = 0.36 \\ \end{array}
```

The q-Gram Distance

## Distance Normalization (2/3)

- Normalize distance such that  $\delta(x, y) \in [0..1]$
- Edit Distance:
  - non-normalized:  $0 \le \operatorname{ed}(x, y) \le \max(|x|, |y|)$
  - normalized edit distance: 0 < norm-ed(x, y) < 1

$$\mathsf{norm\text{-}ed}(x,y) = \frac{\mathsf{ed}(x,y)}{\mathsf{max}(|x|,|y|)}$$

- q-Gram Distance:
  - non-normalized:  $0 \le \operatorname{dist}_q(x, y) \le |G_x \uplus G_y| |G_x \cap G_y|$
  - normalized q-gram distance:  $0 \le \text{norm-dist}_q(x, y) \le 1$

$$\mathsf{norm\text{-}dist}_q(x,y) = \frac{\mathsf{dist}_q(x,y)}{|\mathit{G}_x \uplus \mathit{G}_y| - |\mathit{G}_x \cap \mathit{G}_y|} = 1 - \frac{|\mathit{G}_x \cap \mathit{G}_y|}{|\mathit{G}_x \uplus \mathit{G}_y| - |\mathit{G}_x \cap \mathit{G}_y|}$$

<sup>1</sup>Jaccard normalization. Dividing by  $|G_x \uplus G_y|$  (Dice normalization) also normalizes to [0..1], but the metric properties (triangle inequality) get lost [ABG10].

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

40 / 44

The q-Gram Distance

## Edit Distance vs. q-Gram Distance

• Edit distance can not handle block-moves well:

```
x = \text{Nikolaus Augsten} y = \text{Augsten Nikolaus} norm-ed(x, y) = 1.0 norm-dist_q(x, y) = 0.39 (q = 3)
```

• q-Gram distance may be too strict:

```
x = +39-06-46-74-22 y = (39\ 06\ 467422)

norm-ed(x,y) = 0.4

norm-dist_q(x,y) = 1.0 (q = 3)
```

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 41/44 Augsten (Univ. Salzburg) Similarity Search WS 2025/26 42/44



Nikolaus Augsten, Michael Böhlen, and Johann Gamper.
The *pq*-gram distance between ordered labeled trees.

ACM Transactions on Database Systems (TODS), 35(1):1–36, 2010.

Luis Gravano, Panagiotis G. Ipeirotis, H. V. Jagadish, Nick Koudas, S. Muthukrishnan, and Divesh Srivastava.

Approximate string joins in a database (almost) for free.

In Proceedings of the International Conference on Very Large Databases (VLDB), pages 491–500, Roma, Italy, September 2001. Morgan Kaufmann Publishers Inc.

Luis Gravano, Panagiotis G. Ipeirotis, H. V. Jagadish, Nick Koudas, S. Muthukrishnan, and Divesh Srivastava.

Approximate string joins in a database (almost) for free — Erratum. Technical Report CUCS-011-03, Department of Computer Science, Columbia University, 2003.

Esko Ukkonen.

Approximate string-matching with q-grams and maximal matches.

Approximate join with edit distance inefficient.
Edit distance filters speed up join:

Length filter: based on the string length
Count filter: based on q-Grams
Position filter: based on positional q-Grams

Theoretical Computer Science, 92(1):191–211, January 1992.

Similarity Search

Augsten (Univ. Salzburg)

WS 2025/26

44 / 44

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 44 / 44

WS 2025/26