# Similarity Search Token-Based Tree Distances

#### 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

Similarity Search

WS 2025/26

Token-based Tree Distances

#### Tokens for Trees – Intuition

- *q*-Grams for strings:
  - split string into substrings (q-grams) of length q
  - strings with many common substrings are similar
- Tokens for trees:
  - split tree into small subunits (tokens) of the same shape
  - tokens may be individual nodes, subtrees, or subgraphs
  - trees with many common tokens are similar
- Example: the so-called pq-gram tokens are besom-shaped subtrees with p + q nodes



#### Outline

- Token-based Tree Distances
- 2 Binary Branches
- g pq-Grams
- Conclusion

Augsten (Univ. Salzburg)

Augsten (Univ. Salzburg)

Token Profile and Label Tuples

• Token profile P(T): set of all tokens of tree T

Token-based Tree Distances

- a token may be a subtree or a subgraph of the tree
- the token profile P(T) of a tree T is the set of all its tokens
- A linear encoding of a token traverses all its nodes in preorder:

$$v_1$$

$$\vdots$$

$$v_p = (v_1, \dots, v_p, v_{p+1}, \dots, v_{p+q})$$

$$v_{p+1} \dots v_{p+q}$$

• Label tuple  $\lambda(t)$ : tuple of the nodes labels  $\lambda(v_i)$  of token  $t = (v_1, v_2, \dots, v_k)$  in preorder:

$$\lambda(t) = (\lambda(\mathsf{v}_1), \lambda(\mathsf{v}_2), \dots, \lambda(\mathsf{v}_k))$$

Similarity Search

WS 2025/26

Token-based Tree Distances

#### Token Index

#### Definition (Token Index)

Let P(T) be a token profile of tree T. The token index,  $\mathcal{I}$ , of tree T is the bag of all label tuples of T,

$$\mathcal{I}(\mathsf{T}) = \biguplus_{g \in \mathsf{P}_\mathsf{T}} \lambda(g)$$

- Note:
  - tokens consist of nodes and are unique within a tree
  - but: different tokens may yield identical label tuples
  - thus the token index may contain duplicates

Augsten (Univ. Salzburg

WS 2025/26

Similarity Search

#### Token-based Tree Distances

# Storing the Token Index Efficiently

- Problem: How to store node labels efficiently?
  - Long labels: large storage overhead
  - Varying label length: in a relational database, the inefficient VARCHAR type must be used instead of the efficient CHAR type
- Solution: Hashing
  - compute fingerprint hash for labels
  - store concatenation of the hashed labels
- Fingerprint hash function (e.g., Karp-Rabin [KR87]):
  - maps a string s to a hash value h(s)
  - h(s) is of fixed length
  - h(s) is unique with high probability (for two different strings  $s_1 \neq s_2$ ,  $h(s_1) \neq h(s_2)$  with high probability)

#### Token-based Tree Distances

#### Token-Based Distance

#### Definition (Token-Based Distance)

The token-based distance between two trees, T and T', with token indexes  $\mathcal{I}(\mathsf{T})$  and  $\mathcal{I}(\mathsf{T}')$ , respectively, is defined as

$$\delta(\mathsf{T},\mathsf{T}') = |\mathcal{I}(\mathsf{T}) \uplus \mathcal{I}(\mathsf{T}')| - 2|\mathcal{I}(\mathsf{T}) \cap \mathcal{I}(\mathsf{T}')|$$

- Metric normalization to [0..1]:  $\delta_g'(\mathsf{T},\mathsf{T}') = \frac{\delta_g(\mathsf{T},\mathsf{T}')}{|\mathcal{I}(\mathsf{T}) \uplus \mathcal{I}(\mathsf{T}')| |\mathcal{I}(\mathsf{T}) \uplus \mathcal{I}(\mathsf{T}')|}$
- Pseudo-metric properties hold for normalization [ABG10]:
  - ✓ self-identity:  $x = y \neq \Rightarrow \delta_{\sigma}(x, y) = 0$
  - $\checkmark$  symmetry:  $\delta_g(x,y) = \delta_g(y,x)$
  - ✓ triangle inequality:  $\delta_{\sigma}(x,z) < \delta_{\sigma}(x,y) + \delta_{\sigma}(y,z)$

Token-based Tree Distances

• Different trees may have identical indexes.

Augsten (Univ. Salzburg)

WS 2025/26

### Overview: Token Index

• Token profile: (so-called pq-grams in the example, p = 2, q = 3)



Hashing: map tokens to integers:

Note: labels may be strings of arbitrary length!

• Token index: bag of hashed tokens 33470, 33700, 37000, 36000, 34000, 37000}

Intuition: similar trees have similar token indexes.

Binary Branches

# Binary Tree

- In a binary tree
  - each node has at most two children:
  - left child and right child are distinguished:
     a node can have a right child without having a left child;
- Notation:  $T_B = (N, E_I, E_r)$ 
  - T<sub>B</sub> denotes a binary tree
  - N are the nodes of the binary tree
  - $E_l$  and  $E_r$  are the edges to the left and right children, respectively
- Full binary tree:
  - binary tree
  - each node has exactly zero or two children.

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

9/4

Binary Branches

# Binary Representation of a Tree

- Binary tree transformation:
  - (i) link all neighboring siblings in a tree with edges
  - (ii) delete all parent-child edges except the edge to the first child
- Transformation maintains
  - label information
  - structure information
- Original tree can be reconstructed from the binary tree:
  - a left edge represents a parent-child relationships in the original tree
  - a right edge represents a right-sibling relationship in the original tree

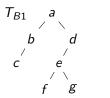
Binary Branches

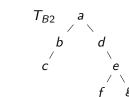
# Example: Binary Tree

• Two different binary trees:  $T_B = (N, E_l, E_r)$ 

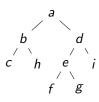
$$T_{B1} = (\{a, b, c, d, e, f, g\}, \{(a, b), (b, c), (d, e), (e, f)\}, \{(a, d), (e, g)\})$$

$$T_{B2} = (\{a, b, c, d, e, f, g\}, \{(a, b), (b, c), (e, f)\}, \{(a, d), (d, e), (e, g)\})$$





• A full binary tree:



Augsten (Univ. Salzburg)

Augsten (Univ. Salzburg)

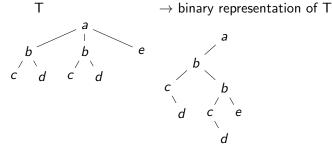
Jiiiiidiity Jea

WS 2025/26

10 /

Example: Binary Tree Transformation

• Represent tree T as a binary tree:



Augsten (Univ. Salzburg) Similarity Search WS 2025/26 11 / 43

Similarity Search

WS 2025/26

### Normalized Binary Tree Representation

- We extend the binary tree with null nodes  $\epsilon$  as follows:
  - a null node for each missing left child of a non-null node
  - a null node for each missing right child of a non-null node
- Note: Leaf nodes get two null-children.
- The resulting normalized binary representation
  - is a full binary tree
  - all non-null nodes have two children
  - all leaves are null nodes (and all null nodes are leaves)

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

13/4

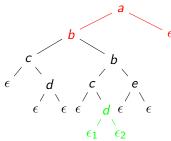
### Binary Branch

- A binary branch BiB(v) is
  - a subtree of the normalized binary tree B(T)
  - consisting of a non-null node v and its two children

Binary Branches

• Example:

$$BiB(a) = (\{a, b, \epsilon\}, \{(a, b)\}, \{(a, \epsilon)\})$$
  
 $BiB(d) = (\{d, \epsilon_1, \epsilon_2\}, \{(d, \epsilon_1)\}, \{(d, \epsilon_2)\})^{-1}$ 



<sup>1</sup>Although the two null nodes have identical labels ( $\epsilon$ ), they are different nodes. We emphasize this by showing their IDs in subscript.

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

15 / 43

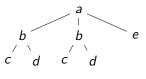
Binary Branches

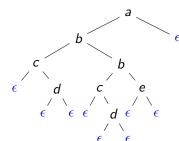
# Example: Normalized Binary Tree

• Transforming T to the normalized binary tree B(T):

Γ







Augsten (Univ. Salzburg)

Similarity Search

Binary Branches

WS 2025/26

14 / 4

# Binary Branches of Trees and Datasets

- Binary branches can be serialized as label tuples:
  - $BiB(v) = (\{v, a, b\}, \{(v, a)\}, \{(v, b)\}) \rightarrow \lambda(v) \circ \lambda(a) \circ \lambda(b)$
- Binary branch profile and index:
  - $P_{bb}(T)$  is the set of all binary branches of T
  - $\mathcal{I}_{bb}(\mathsf{T})$  is the multiset of all binary branch label tuples of  $\mathsf{T}$
- Note:
  - nodes are unique in the tree, thus binary branches are unique
  - labels are *not* unique, thus the label tuples are *not* unique
- Binary branch distance: The binary branch distance between two trees T<sub>1</sub> and T<sub>2</sub> is defined as:

$$\delta_{bb}(\mathsf{T}_1,\mathsf{T}_2) = |\mathcal{I}_{bb}(\mathsf{T}_1) \uplus \mathcal{I}_{bb}(\mathsf{T}_2)| - 2|\mathcal{I}_{bb}(\mathsf{T}_1) \cap \mathcal{I}_{bb}(\mathsf{T}_2)|$$

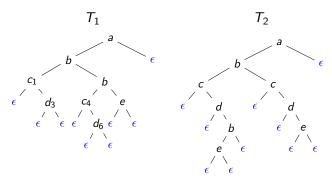
Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

Binary Branches

### Example: Binary Branches and Label Tuples



- $BiB(c_1) \neq BiB(c_4)$ :
  - $BiB(c_1) = (\{c_1, \epsilon_2, d_3\}, \{(c_1, \epsilon_2)\}, \{(c_1, d_3)\})$
  - $BiB(c_4) = (\{c_4, \epsilon_5, d_6\}, \{(c_4, \epsilon_5)\}, \{(c_4, d_6)\})$
- Serialization of both,  $BiB(c_1)$  and  $BiB(c_2)$ , is identical: ' $c \in d$ '

Augsten (Univ. Salzburg)

Sillillarity Searc

WS 2025/2

17 / 43

Binary Branches

Lower Bound Theorem

Theorem (Lower Bound)

Let  $T_1$  and  $T_2$  be two trees. If the tree edit distance between  $T_1$  and  $T_2$  is  $\delta_t(T_1, T_2)$ , then the binary branch distance between them satisfies

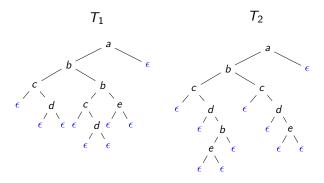
$$\delta_{bb}(\mathsf{T}_1,\mathsf{T}_2) \leq 5 \times \delta_t(\mathsf{T}_1,\mathsf{T}_2).$$

Proof (Sketch — Full Proof in [YKT05]).

- Each node v appears in at most two binary branches.
- Rename: Renaming a node causes at most two binary branches in each tree to mismatch. The sum is 4.
- Similar rational for *insert* and its complementary operation *delete* (at most 5 binary branches mismatch).

Binary Branches

# Example: Binary Branch Distance



$$\mathcal{I}_{bb}(T_1) = \{ab\epsilon, bcb, c\epsilon d, d\epsilon \epsilon, bce, c\epsilon d, d\epsilon \epsilon, e\epsilon \epsilon\}$$

$$\mathcal{I}_{bb}(T_2) = \{ab\epsilon, bcc, c\epsilon d, d\epsilon b, be\epsilon, e\epsilon \epsilon, c\epsilon d, d\epsilon e, e\epsilon \epsilon\}$$

$$\delta_{bb}(T_1, T_2) = 17 - 2 \cdot 4 = 9$$

Binary Branches

Augsten (Univ. Salzburg)

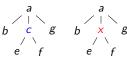
Similarity Searcl

WS 2025/26

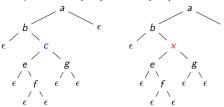
18 / 4

Proof Sketch: Illustration for Rename

• transform  $T_1$  to  $T_2$ : ren(c, x)



• binary trees  $B(T_1)$  and  $B(T_2)$ 



- Two binary branches ( $b\epsilon c$ , ceg) exist only in  $B(T_1)$
- Two binary branches  $(b\epsilon x, xeg)$  exist only in  $B(T_2)$
- $\delta_t(\mathsf{T}_1,\mathsf{T}_2) = 1$  (1 rename)

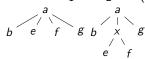
Augsten (Univ. Salzburg)

•  $\delta_{bb}(\mathsf{T}_1,\mathsf{T}_2)=4$  (4 binary branches different)

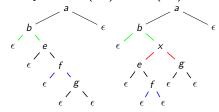
Binary Branches

#### Proof Sketch: Illustration for Insert

• transform  $T_1$  to  $T_2$ : ins(x, a, 2, 2)



• binary trees  $B(T_1)$  and  $B(T_2)$ 



- Two binary branches (bee,  $f \in g$ ) exist only in  $B(T_1)$
- Tree binary branches ( $b\epsilon x$ ,  $f\epsilon\epsilon$ , xeg) exist only in  $B(T_2)$
- $\delta_t(\mathsf{T}_1,\mathsf{T}_2)=1$  (1 insertion)
- $\delta_{bb}(\mathsf{T}_1,\mathsf{T}_2)=5$  (5 binary branches different)

Augsten (Univ. Salzburg)

WS 2025/26

Binary Branches

# Complexity: Binary Branch Distance

- Generating the binary branches: O(n) time and space  $(n = \max\{|T_1|, |T_2|\})$ 
  - the binary branches are formed in a single traversal of the tree
  - for each node of a tree a single binary branch is formed
- Computing the distance:  $O(n \log n)$  time and O(n) space
  - sort binary branch indexes to compute intersection:  $O(n \log n)$
  - alternative: average case O(n) runtime complexity
    - 1. build hash map for index  $\mathcal{I}_{bb}(T_2)$
    - 2. probe label tuples of  $\mathcal{I}_{bb}(T_1)$  to compute size of intersection

Binary Branches

#### **Proof Sketch**

- In general it can be shown that
  - Rename changes at most 4 binary branches
  - Insert changes at most 5 binary branches
  - Delete changes at most 5 binary branches
- Each edit operation changes at most 5 binary branches, thus

$$\delta_{bb}(\mathsf{T}_1,\mathsf{T}_2) \leq 5 \times \delta_t(\mathsf{T}_1,\mathsf{T}_2).$$

Augsten (Univ. Salzburg)

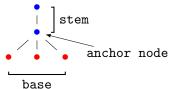
Similarity Search

pq-Grams

WS 2025/26

pq-Grams

• The shape of a pq-gram (p=2, q=3):



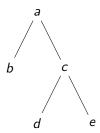
- p nodes (anchor node and p-1 ancestors) form the stem
- q nodes (q consecutive children of the anchor node) form the base

Augsten (Univ. Salzburg) WS 2025/26 23 / 43 WS 2025/26 Similarity Search Augsten (Univ. Salzburg) Similarity Search

pq-Grams

### pq-Extended Tree

• Problem: How can we split the following tree T into 2, 3-grams?



- Solution: Extend tree T with dummy nodes (•):
  - p-1 ancestors to the root node
  - $\bullet$  q-1 children before the first and after the last child of each non-leaf
  - q children for each leaf
- The result is the pq-extendd tree  $T^{pq}$ .

Augsten (Univ. Salzburg)

Sillillarity Search

WS 2025/26

25 / 4

pg-Grams

# Definition: pq-Gram [ABG05]

#### Definition (pq-Gram)

Let T be a tree,  $T^{p,q}$  the respective extended tree, p > 0, q > 0. A subtree of  $T^{p,q}$  is a pq-gram g of T iff

- (a) g has q leaf nodes and p non-leaf nodes,
- (b) all leaf nodes of g are children of a single node  $a \in N(g)$  with fanout g, called the anchor node,
- (c) the leaf nodes of g are consecutive siblings in  $T^{p,q}$ .
- Stem: anchor node and its ancestors in the pq-gram.
- Base: children of the anchor node in the pq-gram.

#### Definition (pq-Gram Profile)

The pq-gram profile,  $P_T$ , of a tree T is the set of all its pq-grams.

pq-Gran

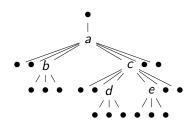
# Example: Extended Tree

• An example tree T and its extended tree  $T^{pq}$  (p=2, q=3):

Τ

2, 3-extended tree  $T^{2,3}$ 





Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

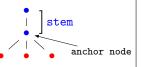
26 / 4

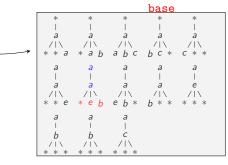
# Example: Systematically Split Tree

• pq-Gram: small subtree with stem and base Example: p = 2, q = 3

• Systematically split tree into pq-grams

• pq-Gram profile: set of all pq-grams of a tree.





Augsten (Univ. Salzburg

Similarity Search

WS 2025/26

pg-Grams

### Label Tuples

 Linear encoding of a pq-gram g with anchor node v<sub>p</sub>: (traverse pq-gram in preorder)

$$v_1$$
 $v_p$ 
 $v_{p+1} \dots v_{p+q}$ 
 $v_{p+q}$ 
 $v_{p+q}$ 
 $v_{p+q}$ 

• Label tuple: tuple of the pq-gram's node labels

$$\lambda(g) = (\lambda(\mathsf{v}_1), \dots, \lambda(\mathsf{v}_{p+q}))$$

for the pq-gram  $g = (v_1, \dots, v_{p+q})$ .

Augsten (Univ. Salzburg)

Jillilarity Searci

WS 2025/26

29 / 43

#### Size of the pq-Gram Index

### Theorem (Size of the pq-Gram Index)

Let T be a tree of size n = l + i with l leaves and i non-leaves. The size of the pq-gram index of T is linear in the tree size:

$$|\mathcal{I}^{pq}(\mathsf{T})| = 2I + qi - 1 = O(n)$$

#### Proof.

- 1. We count all pq-grams whose leftmost leaf is a dummy node: Each leaf is the anchor node of exactly one pq-gram whose leftmost leaf is a dummy node, giving l pq-grams. Each non-leaf is the anchor of q-1 pq-grams whose leftmost leaf is a dummy, giving i(q-1) pq-grams.
- 2. We count all pq-grams whose leftmost leaf is not a dummy node: Each node of the tree except the root is the leftmost leaf of exactly one pq-gram, giving l+i-1 pq-grams.

Overall number of pq-grams: l + i(q - 1) + (l + i - 1) = 2l + qi - 1.

pq-Grams

# pq-Gram Index

#### Definition (pq-Gram Index)

Let T be a tree with profile  $P_T$ , p>0, q>0. The pq-gram index,  $\mathcal{I}$ , of tree T is the **bag of all label tuples** of T,

$$\mathcal{I}(\mathsf{T}) = \biguplus_{g \in \mathsf{P}_\mathsf{T}} \lambda(g)$$

- Note:
  - pq-grams are unique within a tree
  - but: different pq-grams may yield identical label tuples
  - thus the pq-gram index may contain duplicates

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

30 / 4

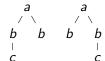
# The pq-Gram Distance

#### Definition (pq-Gram Distance)

The pq-gram distance between two trees, T and T', is defined as

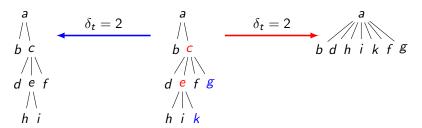
$$\delta_{g}(\mathsf{T},\mathsf{T}') = |\mathcal{I}(\mathsf{T}) \uplus \mathcal{I}(\mathsf{T}')| - 2|\mathcal{I}(\mathsf{T}) \cap \mathcal{I}(\mathsf{T}')|$$

- $\bullet \ \, \text{Metric normalization to [0..1]:} \ \, \delta_g'(\mathsf{T},\mathsf{T}') = \tfrac{\delta_g(\mathsf{T},\mathsf{T}')}{|\mathcal{I}(\mathsf{T}) \uplus \mathcal{I}(\mathsf{T}')| |\mathcal{I}(\mathsf{T}) \uplus \mathcal{I}(\mathsf{T}')|}$
- Pseudo-metric properties hold for normalization [ABG10]:
  - ✓ self-identity:  $x = y \not\leftarrow \Rightarrow \delta_g(x, y) = 0$
  - ✓ symmetry:  $\delta_g(x,y) = \delta_g(y,x)$
  - ✓ triangle inequality:  $\delta_g(x,z) \leq \delta_g(x,y) + \delta_g(y,z)$



• Different trees may have identical indexes:

# Motivation: Unit Cost Model Not Always Intuitive



- Unit cost edit distance:
  - no difference between leaves and non-leaves
  - may lead to non-intuitive results
- Conclusion: Non-leaves should have more weight than leaves.

Augsten (Univ. Salzburg)

Similarity Search

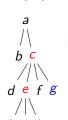
WS 2025/26

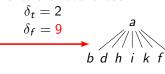
33 / 4

# Example: Fanout-Weighted Tree Edit Distance

- Fanout-Weighted Tree Edit Distance:
  - ullet leaf changes have small cost (c=1) in the example
  - non-leaf changes cost proportional to the node fanout

$$\delta_t = 2$$
 $\delta_f = 2$ 
 $\delta_f = 2$ 
 $\delta_f = 1$ 
 $\delta_f = 2$ 
 $\delta_f = 2$ 
 $\delta_f = 2$ 
 $\delta_f = 1$ 
 $\delta_f = 2$ 





pg-Grams

### Fanout Weighted Tree Edit Distance

#### Definition (Fanout Weighted Tree Edit Distance)

Let T and T' be two trees,  $w \in N(T)$  a node with fanout f,  $w' \in N(T')$  a node with fanout f', c > 0 a constant. The fanout weighted tree edit distance,  $\delta_f = (T, T')$ , between T and T' is defined as the tree edit distance with the following costs for the edit operations:

- Delete:  $\alpha(\mathsf{w} \to \epsilon) = f + c$
- Insert:  $\alpha(\epsilon \to \mathsf{w}') = f' + c$
- Rename:  $\alpha(\mathbf{w} \to \mathbf{w}') = (f + f')/2 + c$
- Cost of changing a non-leaf node: proportional to its fanout.
- Cost of changing a leaf node: constant c.

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

34 /

# pg-Gram Distance Lower Bound

#### Theorem

Let p=1 and  $c \ge \max(2q-1,2)$  be the cost of changing a leaf node. The pq-gram distance provides a lower bound for the fanout weighted tree edit distance, i.e., for any two trees,  $\mathsf{T}$  and  $\mathsf{T}'$ ,

$$\frac{\delta_{\mathbf{g}}(\mathsf{T},\mathsf{T}')}{2} \leq \delta_{\mathbf{f}}(\mathsf{T},\mathsf{T}').$$

#### Proof.

See [ABG10] (ACM Transactions on Database Systems).

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

35 / 43

Augsten (Univ. Salzburg)

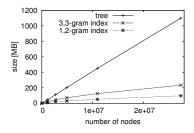
Similarity Search

WS 2025/26

pg-Grams

# Size of the pq-Gram Index

- pq-Gram index size: linear in the tree size
- Experiment:
  - compute pq-gram index for trees with different number of nodes
  - compare tree and index size



[Trees created with xmlgen.]

Why is the *pq*-gram index smaller than the tree?

- hash values are smaller than labels
- duplicate *pq*-grams of a tree are stored only once

Augsten (Univ. Salzburg)

Similarity Search

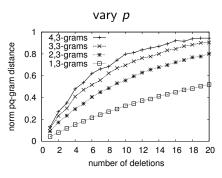
WS 2025/26

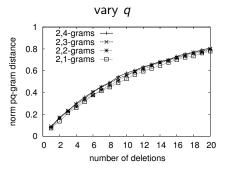
37 / 43

pg-Grams

# Sensitivity to Structure Change — Non-Leaf

- Cost for non-leaf change  $\rightarrow$  controlled by p
- Experiment:
  - delete non-leaf nodes
  - measure normalized pq-gram distance





(Artificial tree with 144 nodes, 102 leaves, fanout 2-6 and depth 6. Average over 100 runs.)

Augsten (Univ. Salzburg)

Similarity Search

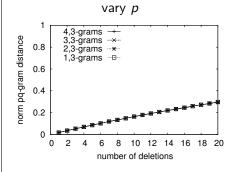
WS 2025/26

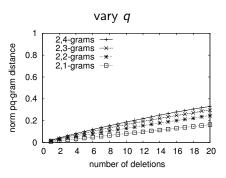
39 / 43

pg-Grams

# Sensitivity to Structure Change — Leaf

- Cost of leaf change  $\rightarrow$  depends only on q
- Experiment:
  - delete leaf nodes
  - measure normalized pq-gram distance





(Artificial tree with 144 nodes, 102 leaves, fanout 2-6 and depth 6. Average over 100 runs.)

Augsten (Univ. Salzburg)

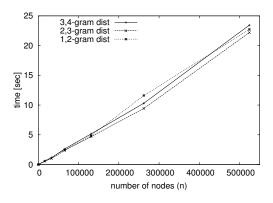
Similarity

WS 2025/26

38 / 4

# Influence of p and q on Scalability

- Scalability (almost) independent of p and q.
- Experiment: For pair of trees
  - compute pq-gram distance for varying p and q
  - vary tree size: up 10<sup>6</sup> nodes
  - measure wall clock time



Augsten (Univ. Salzburg) Simila

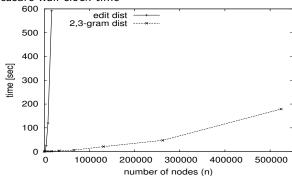
/ Search

WS 2025/26

pq-Gram

### Scalability to Large Trees

- pq-gram distance  $\rightarrow$  scalable to large trees
- compare with edit distance
- Experiment: For pair of trees
  - compute tree edit distance and pq-gram distance
  - vary tree size: up  $5 \times 10^5$  nodes
  - measure wall clock time



Augsten (Univ. Salzburg)

Onimarity Octares

Conclusion

WS 2025/26

41 / 43

# Summary

- Binary Branch Distance
  - lower bound of the unit cost tree edit distance
  - trees are split into binary branches (small subgraphs)
  - similar trees have many common binary branches
  - complexity  $O(n \log n)$  time and (n) space
- pq-Gram Distance
  - lower bound for the fanout weighted tree edit distance
  - trees are split into pq-grams (small subtrees)
  - ullet similar trees have many common pq-grams
  - complexity  $O(n \log n)$  time and O(n) space

pg-Grams

# pq-Grams vs. other Edit Distance Approximations

Effectiveness: pq-grams outperform all other approximations

Experiment: two sets of address trees (299 and 302 trees)

- compute distances between all tree pairs
- find matches (symmetric nearest neighbor)

Distance	Correct	Recall	Precision	f-Measure	Runtime
fanout edit dist	259	86.6%	98.5%	0.922	19 min
unit edit dist	247	82.6%	96.5%	0.890	14 min
node intersection	197	65.9%	93.8%	0.774	4.3s
p,q-grams	236	78.9%	98.7%	0.877	8.1s
tree-embedding	206	68.9%	96.3%	0.803	7.1s
binary branch	193	64.5%	93.2%	0.763	7.4s
bottom-up	148	49.6%	92.5%	0.645	67.0s

Augsten (Univ. Salzburg)

Similarity Search

WS 2025/26

42 / 43

43 / 43

- Nikolaus Augsten, Michael Böhlen, and Johann Gamper.

  Approximate matching of hierarchical data using pq-grams.

  In Proceedings of the International Conference on Very Large

  Databases (VLDB), pages 301–312, Trondheim, Norway, September 2005. ACM Press.
- 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.
- Richard M. Karp and Michael O. Rabin.
  Efficient randomized pattern-matching algorithms.

  IBM Journal of Research and Development, 31(2):249–260, March 1987.
  - Similarity evaluation on tree-structured data.

    In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 754–765, Baltimore, Maryland, USA, June 2005. ACM Press.

Rui Yang, Panos Kalnis, and Anthony K. H. Tung.

Augsten (Univ. Salzburg) Similarity Search WS 2025/26 43 / 43 Augsten (Univ. Salzburg) Similarity Search WS 2025/26