# A Scalable Index for Top-k Subtree Similarity Queries

Presented by: Aliza Subedi, Bishal Bashyal



#### Problem Statement

#### <root>

<item id="1">Hello</item> <item id="1">Hello</item> <item id="2">World</item> <item id="3">World</item> <item id="4">Print</item> </root>

#### **Document T**



#### <item id="1">Hello</item>









#### Tree Edit Distance

#### Minimum number of node edit operations that transform one tree to another.





#### Tree Edit Distance

#### Minimum number of node edit operations that transform one tree to another.



Space Complexity: O(n^2)

Time Complexity: O(n^3)



#### **Cost Estimate and Bounds**



#### • Label Lower Bound (IIb): how many labels are not similar in both trees ?

#### • Size Lower Bound (slb): Absolute size difference between T and Q.



#### **Cost Estimate and Bounds**



- Query labels =  $\{A, C, D, E, F\}$
- Tree labels =  $\{A, C, D, E, B\}$
- LLB = 1 (since F is missing)

#### • Label Lower Bound (IIb): how many labels are not similar in both trees ?

#### • Size Lower Bound (slb): Absolute size difference between T and Q.



#### **Cost Estimate and Bounds**



- Query labels =  $\{A, C, D, E, F\}$
- Tree labels =  $\{A, C, D, E, B\}$
- LLB = 1 (since F is missing)

- Number of Query nodes = 5
- Number of Tree nodes = 5
- SLB = 0 ( size of both trees are same)

#### • Label Lower Bound (IIb): how many labels are not similar in both trees ?

#### • Size Lower Bound (slb): Absolute size difference between T and Q.





	Tasm-Postorder	
Index	Index-Free method	
Querying	Slow	
Memory Usage	Low	
Index Updates	_	





# **Bounding Condition**

# LowerBound(Q, $T_i$ ) $\geq \delta(R'[k])$

Avoids unnecessary Tree Edit Distance calculations by eliminating subtrees early if they cannot match the query within the allowed edit distance.



#### Inverted Index



(a) Example of Query Tree (Q) and Document Tree (T)

(b) Inverted List index of T





# MergeAll ALgorithm

• Start from the S0 stripe whith SLB = 0 (possible top ranking).

 Compute IIb(Q,Ti) aand if IIb(Q,Ti) != slb(Q,Ti), cache that Ti to a bucket b[j] where j is the stripe number.

• If IIb = slb, compute TED and update the ranking list.

 When evaluating a new stripe, evaluate bucket b[j] first.

• As processing continues, check if |R| = k and  $j \ge \delta(R[k])$  for terminating the filtering process.





- •
- :

# **Cone Filtering**

• Mergeall could be inefficient by processing large partitions within a stripe.

 Cone Accesses short, sorted lists first, often terminating before reaching large partitions.
Processes only a subset of partitions unlike mergeall.

Uses min = |Q| - B, max = |Q| + B - nml (Ti) values in each round.

• Terminates if |R| = k and  $B \ge \delta(R[k])$ .

• Quadratic Space Complexity.



#### Slim Inverted list



(a) Slim inverted list index.



### SlimCone



## Algorithm

• Reconstruct subtrees on the fly by unveiling.

- Round Based with B=0, and increases in each round.
- Generate possible candidates by calculating llb(Q,Ti).
- If(IIb(Q,Ti) = B), evauate TED, else cache.
- Check if  $|\mathbf{R}| = k \wedge \mathbf{B} \ge \delta$  ( $\mathbf{R}[\mathbf{k}]$  for terminating filtering algorithm.
- Works with a Linear space (slim) index.
- Implements sliminvertedlist as a Binary Search tree to support index updates.



# **Supporting Updates**



(a) Slim inverted list index.

### Dynamic Node Index

- Each node entry contains:Ibl: The node's label (e.g., a,b,x a, b, x a,b,x).
- size: The size of the subtree rooted at the node.
- par: The parent node's identifier.
- cl: The identifier of the first child.
- sib: The identifier of the next (right) sibling.
- No need to reconstruct the entire index.

![](_page_15_Picture_10.jpeg)

# Comparision

	MERGEALL	CONE
Handles large datasets efficiently	×	
Space-efficient index (linear size)	×	×
Faster Index Build	×	×
Supports dynamic updates	×	×

![](_page_16_Figure_2.jpeg)

![](_page_16_Picture_3.jpeg)

#### Datasets

Name	Size T	Size [No	
	[MB]	T	
XMark1	112	$3.6 \cdot 10^{6}$	
XMark2	223	$7.2 \cdot 10^{6}$	
XMark4	447	$14.4 \cdot 10^{6}$	
XMark8	895	$28.9 \cdot 10^{6}$	
XMark16	1,790	$57.8 \cdot 10^6$	
TreeBank	83	$3.8\cdot 10^6$	
DBLP	2,161	$126.5 \cdot 10^{6}$	
SwissProt	6,137	$479.3 \cdot 10^{6}$	

des] avg.  $|T_i|$ 6.2 6.2 6.2 6.2 6.2 8.4 3.4 5.1

![](_page_17_Picture_3.jpeg)

![](_page_18_Figure_0.jpeg)

![](_page_18_Figure_1.jpeg)

 $\star$ 

memory [MB]

 $10^{4}$ 

10<sup>3</sup>

#### (a) Build time, XMark.

![](_page_18_Figure_3.jpeg)

#### (b) Index size, XMark.

![](_page_18_Picture_5.jpeg)

![](_page_19_Figure_0.jpeg)

![](_page_19_Picture_1.jpeg)

# Slimcone

![](_page_20_Picture_1.jpeg)

### Linear Space indexing

#### Scalable

### Updatable

![](_page_20_Picture_5.jpeg)

# Conclusion

- Proposed a scalable and novel solution for top-k subtree similarity query retrieval.
- Presented baseline algorithm, Mergeall and cone and builds slim on top of it.
- Introduced slim index reducing memory complexity from quadratic to linear.
- Extended Slimcone to Slim-Dyn for dynamic document updates.
- Validated performance through extensive evaluation on synthetic and real-world datasets.

![](_page_21_Picture_11.jpeg)

# Thank you!

![](_page_22_Picture_1.jpeg)

![](_page_22_Picture_2.jpeg)