

Faster Weighted and Unweighted Tree Edit Distance and APSP Equivalence

Jakob Nogler¹ Adam Polak² Barna Saha³ Virginia Vassilevska Williams⁴ Yinzhan Xu³ Christopher Ye³

¹ETH Zurich ²Bocconi University ³UC San Diego ⁴MIT

Problem Definition

(String) Edit Distance

Input: Strings S, S' .

Output: Cheapest transformation of S into S' .

Tree Edit Distance (TED)

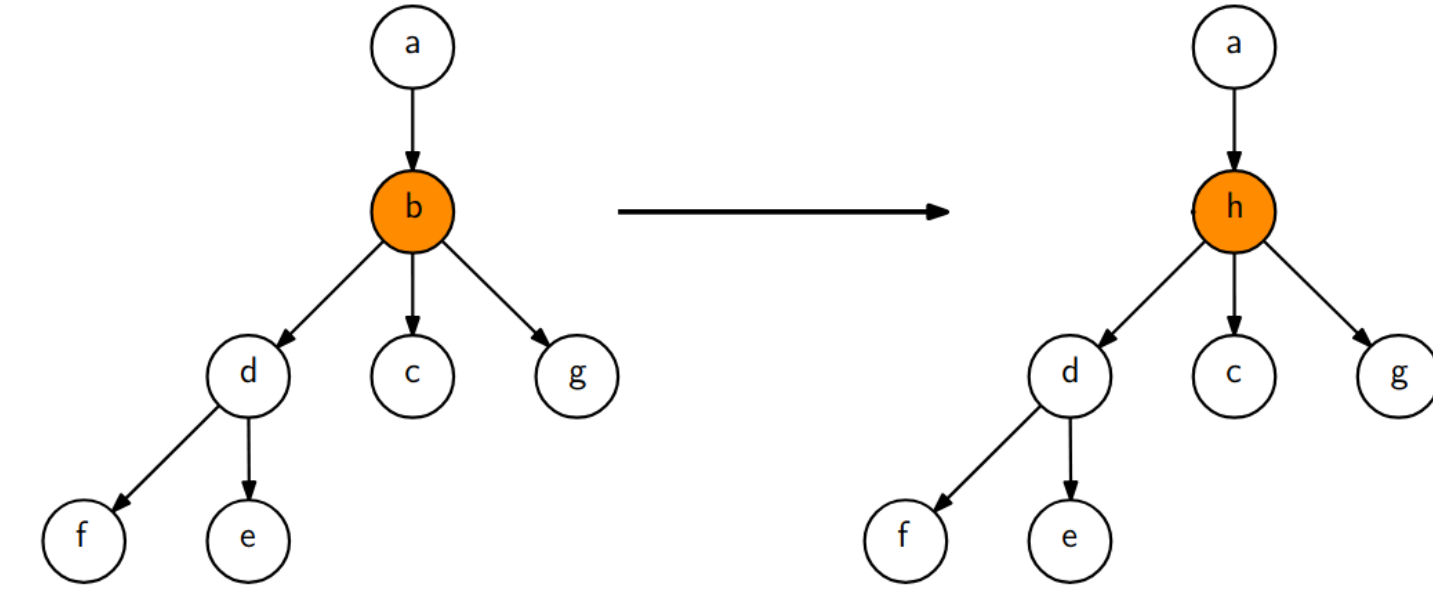
Input: Labeled, left-to-right ordered trees T, T' .

Output: Cheapest transformation of T into T' .

Allowed Operations:

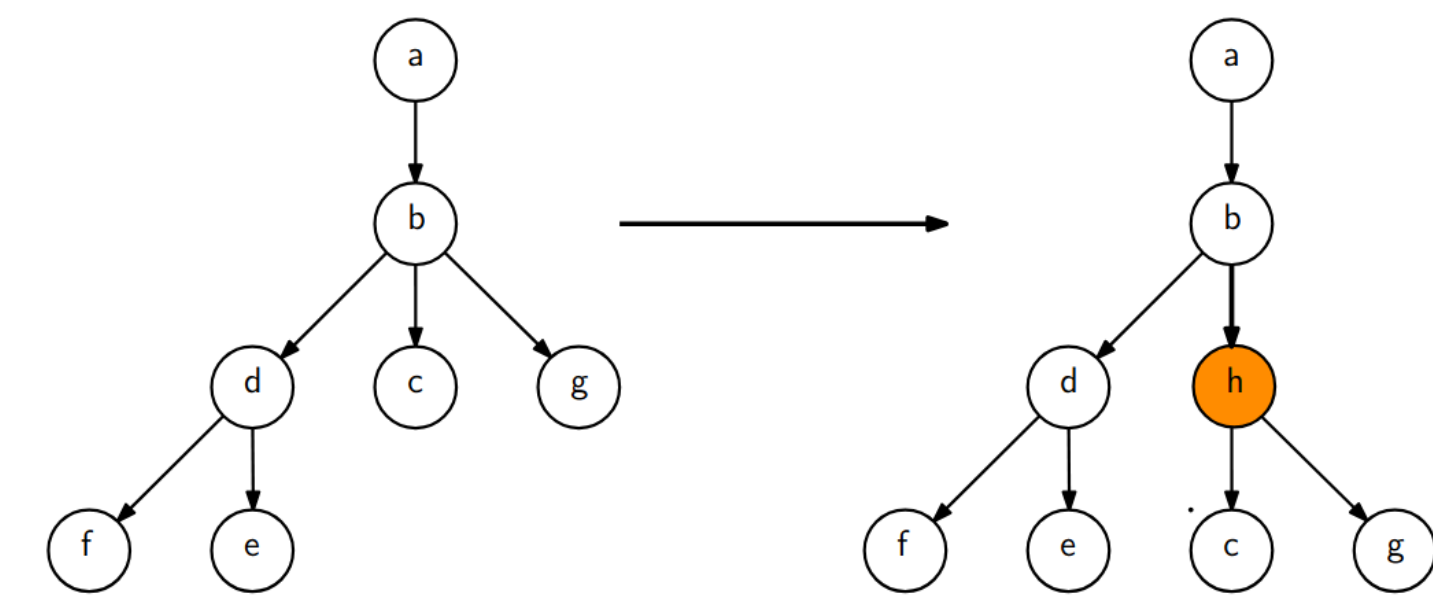
① Substitution

abc**x**ef \longrightarrow abc**y**ef



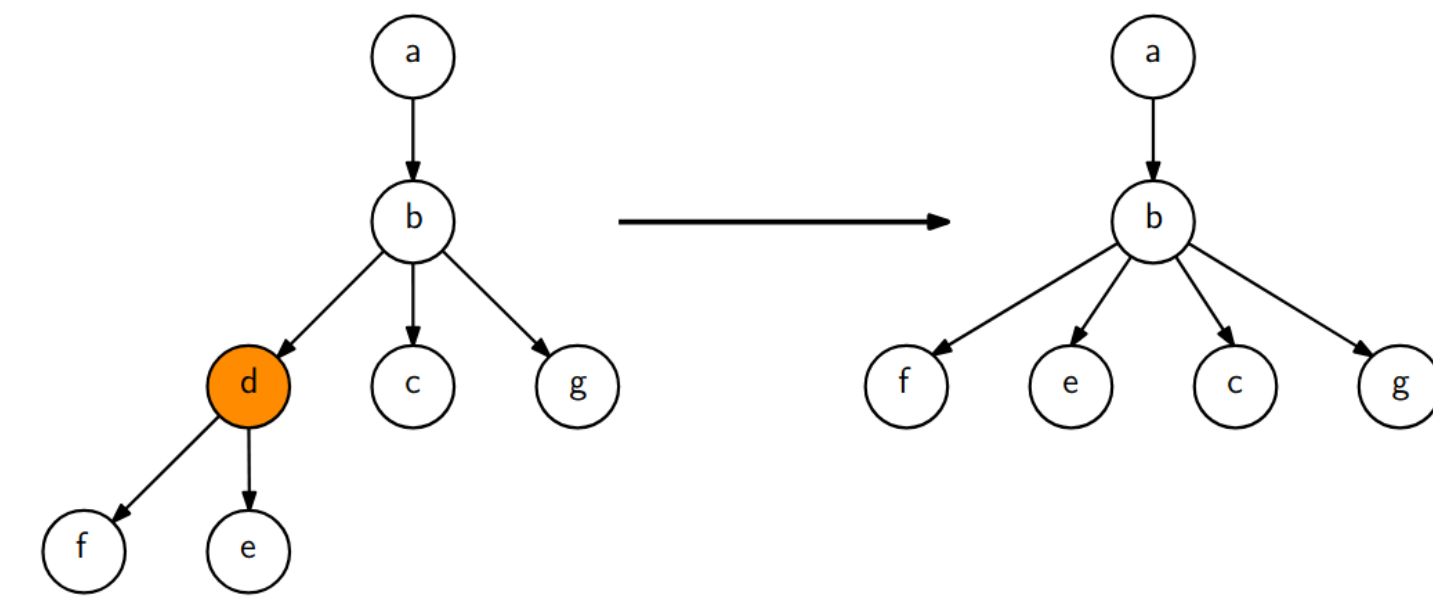
② Insertion

abce**f** \longrightarrow abc**x**ef



③ Deletion

abc**d**ef \longrightarrow abdef



Right hand side images taken from [SS19].

Unweighted: All operations have cost 1.

Weighted: Cost of operations is specified by a cost function δ , i.e. $\delta(c, c')$, $\delta(c, \varepsilon)$ and $\delta(\varepsilon, c)$ specify substitution, deletion and insertion cost of character c .

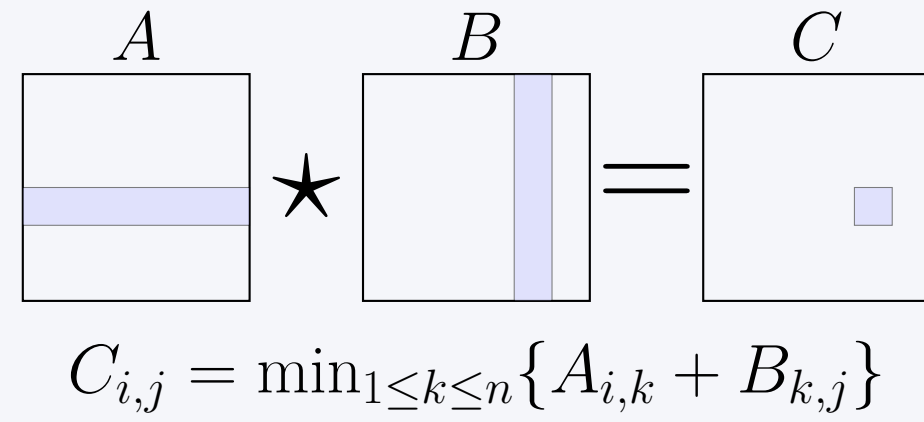
Previous Works and Our Results

Year	Work	Setting	Complexity
1979	Tai	weighted	$\mathcal{O}(n^6)$
1989	Shasha, Zhang	weighted	$\mathcal{O}(n^4)$
1998	Klein	weighted	$\mathcal{O}(n^3 \log n)$
2007	Demaine, Mozes, Rossman, Weimann	weighted	$\mathcal{O}(n^3)$
2020	Bringmann, Gawrychowski, Mozes, Weinmann	weighted	no $\mathcal{O}(n^{3-\epsilon})$ algo under APSP
2024	This work	weighted	$n^3/2^{\Omega(\sqrt{\log n})}$
2022	Mao	unweighted	$\mathcal{O}(n^{2.9546})$
2023	Dürr	unweighted	$\mathcal{O}(n^{2.9148})$
2024	This work	unweighted	$\mathcal{O}(n^{2.687})$

More details on our results:

- We show how to reduce TED to APSP. Combined with the APSP lower bound, this gives equivalence between TED and APSP.
- Via the state-of-the-art APSP Algorithm [Williams'18] we get faster weighted TED.
- Previous truly subcubic algorithm by Mao and Dürr, reduce (unweighted) TED to monotone min-plus product, but lose some polynomial factors on the way. We show how to get a tight reduction. By doing so, via the state-of-the-art algorithm for this product [CDXZ'22], we obtain a $\mathcal{O}(n^{(\omega+3)/2})$ -time algorithm.

Monotone Min-plus Product



$$C_{i,j} = \min_{1 \leq k \leq n} \{A_{i,k} + B_{k,j}\}$$

- B is row monotone:

$$\forall i, j \quad B_{i,j} \leq B_{i,j+1}.$$

- A, B are bounded:

$$\forall i, j \quad A_{i,j}, B_{i,j} = \mathcal{O}(n).$$

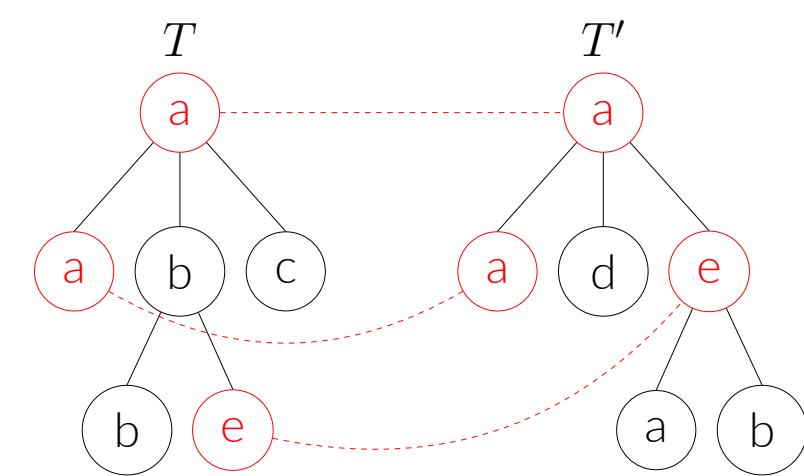
Problem Re-Definition

Rather than finding we finding a least-cost transformation, we aim to a maximize a matching.

We define the cost $\eta(c, c')$ of matching two characters:

$$\eta(c, c') := \delta(c, \varepsilon) + \delta(\varepsilon, c') - \delta(c, c').$$

m i t t a g s s e m i n a r
c h r i s t m a s



String Similarity: maximum total weight matching, under the condition that matching do not cross.

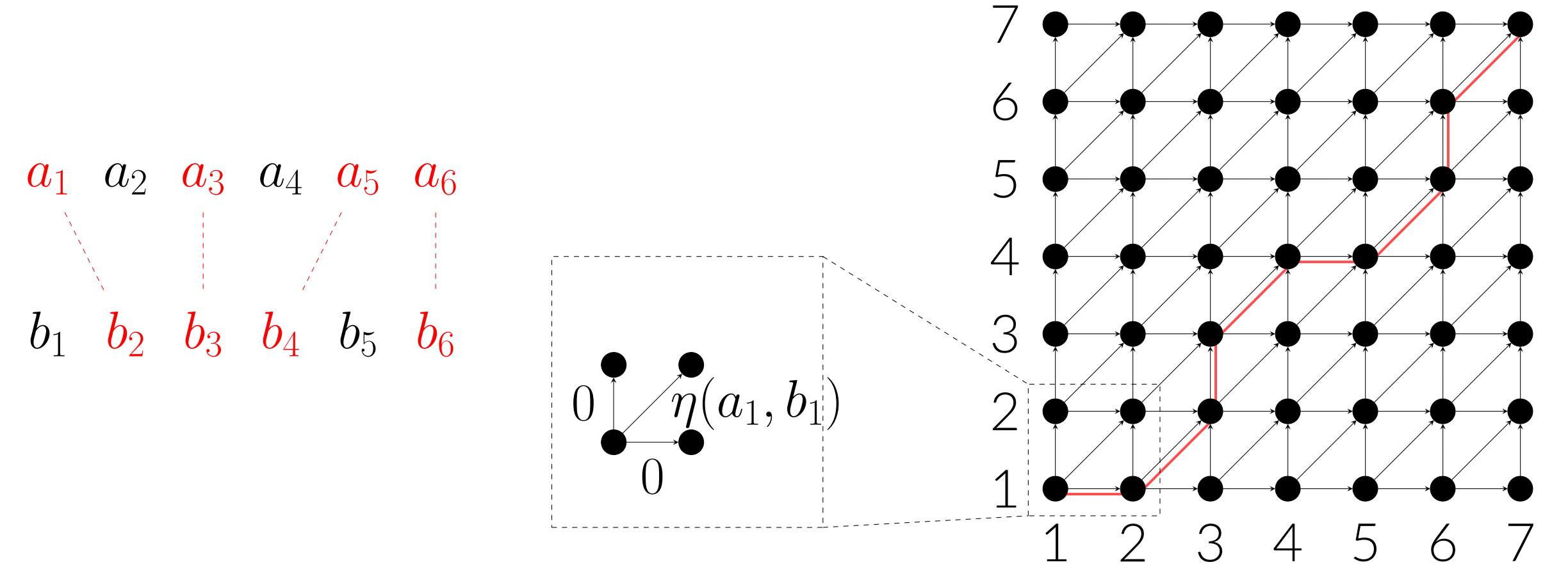
Tree Similarity: maximum total weight matching, under the condition that matching respect same ancestry and pre-order orderings of the matched nodes.

The following relation holds for string/tree distance and similarity:

$$\text{sim}(A, B) = \sum_{a \in A} \delta(a, \varepsilon) + \sum_{b \in B} \delta(\varepsilon, b) - \text{ed}(A, B).$$

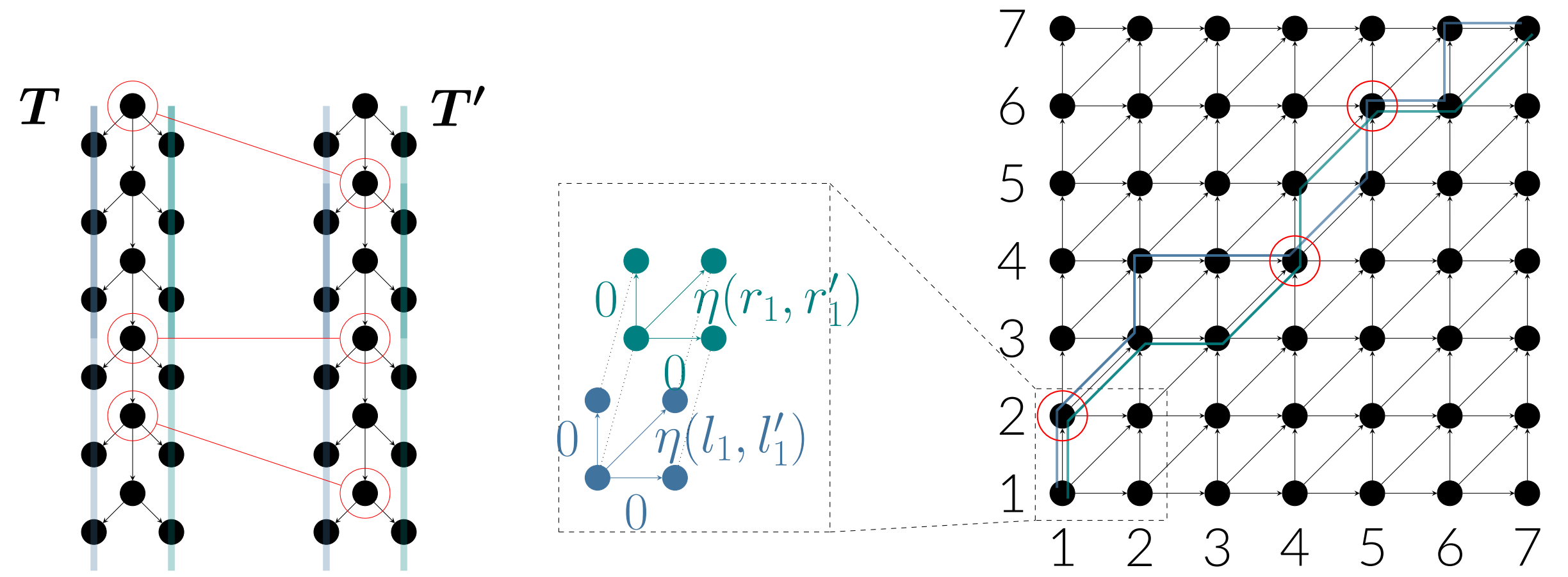
Similarity Computation on a Grid: Strings vs Caterpillar Trees

① Strings



N.B. Border-to-border distances in such graphs can be computed in $\mathcal{O}(n^2)$ time.

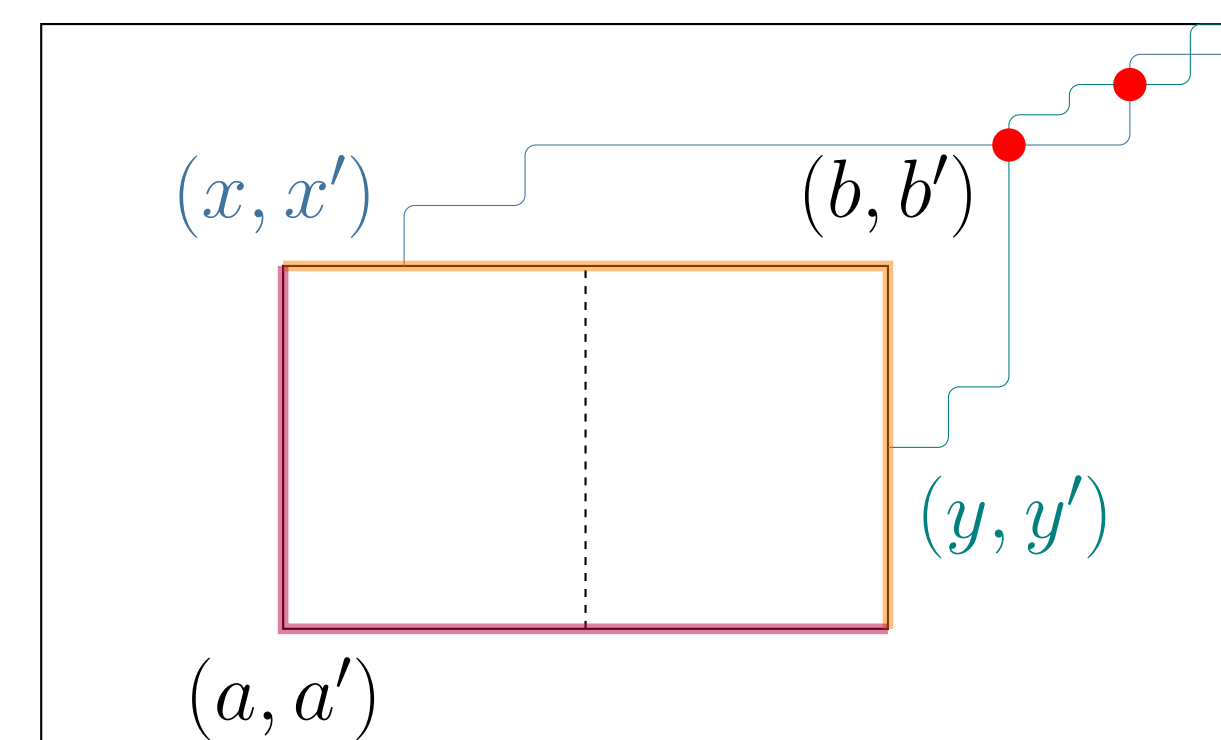
② Caterpillar Trees



N.B. Here we assume the left, right and spine nodes $l_1, l_2, \dots, r_1, r_2, \dots$, and c_1, c_2, \dots of \mathbf{T} only match with the left, right and spine nodes $l'_1, l'_2, \dots, r'_1, r'_2, \dots$, and c'_1, c'_2, \dots of \mathbf{T}' .

A Divide and Conquer Scheme for Caterpillar Tree Similarity

$(n+1, n+1)$



$(1, 1)$

Input:

- The lower-left corner (a, a') and upper-right corner (b, b') of a rectangle.
- $\text{sim}((x, x'), (y, y')) \quad \forall (x, x'), (y, y') \in ([a..b] \times \{b'\}) \cup (\{b\} \times [a'..b'])$.

Output:

- $\text{sim}((x, x'), (y, y')) \quad \forall (x, x'), (y, y') \in ([a..b] \times \{a'\}) \cup (\{a\} \times [a'..b'])$.

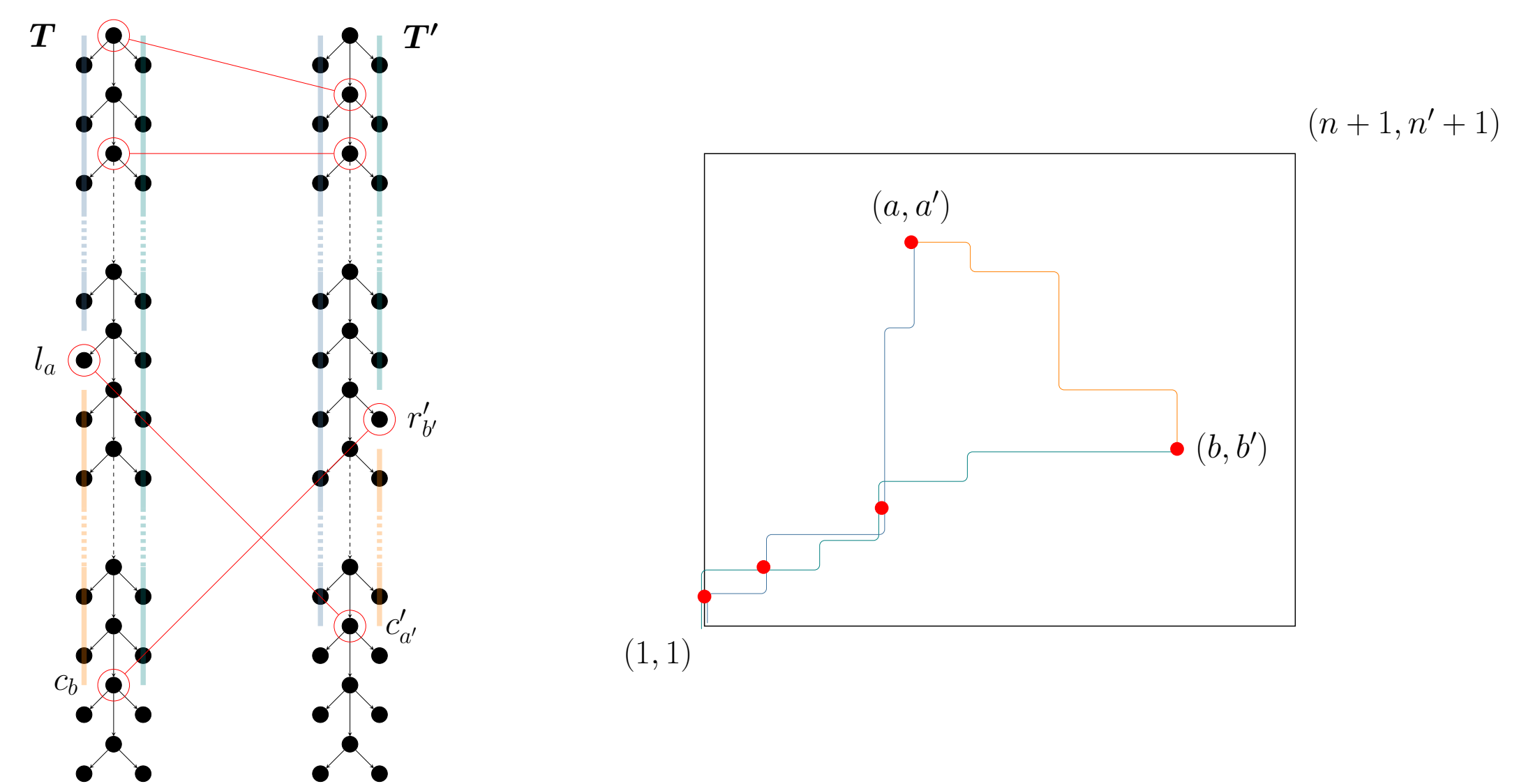
$\text{sim}((x, x'), (y, y'))$ equals to the maximum achievable sum of:

- the weight of a path from (x, x') to $(n+1, n+1)$ in the alignment graph of $\text{sim}(L, L')$;
- the weight of a path from (y, y') to $(n+1, n+1)$ in the alignment graph of $\text{sim}(R, R')$; and
- values $\eta(c_i, c'_i)$ for (i, i') where the two paths intersect (each c_i and c'_i appears at most once).

We obtain the recurrence $T(n) = 4T(n/2) + \mathcal{O}(T_{\text{APSP}}(n))$. Thus, $T(n) = \mathcal{O}(T_{\text{APSP}}(n))$.

How to Extend to General Caterpillars

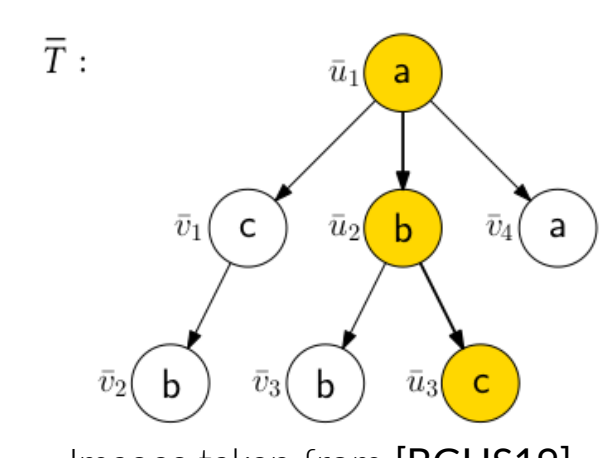
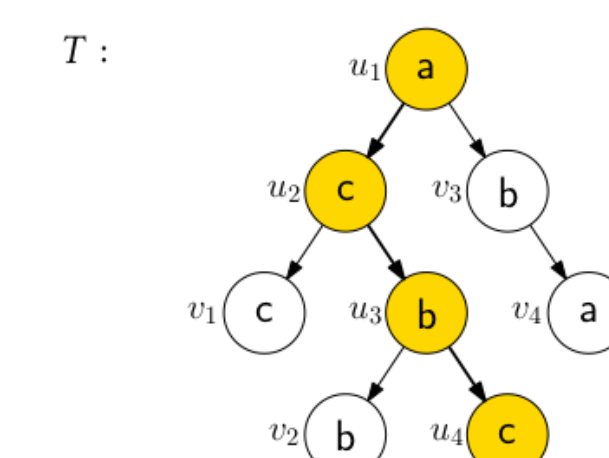
Dropping the assumption on the mapping adds one more path in a string alignment graph.



How to Extend to General Trees

We generalize TED on caterpillar trees to **Spine Edit Distance**.

Input: trees T, T' , root-to-leaf paths $S \subseteq T, S' \subseteq T'$, and $\text{sim}(\text{sub}(v), \text{sub}(v')) \quad \forall (v, v') \in (T \times T') \setminus (S \times S')$.



Images taken from [BGHS19]

Output: $\text{sim}(\text{sub}(v), \text{sub}(v')) \quad \forall (v, v') \in S \times S'$.

Why?

- There are still spine, left, and right nodes.
- The underlying paths are not in string alignment graphs anymore but in *forest alignment graphs*.
- TED on arbitrary trees is fine-grained equivalent to Spine Edit Distance.