





#### **Optimal Algorithms for Bounded Weighted Edit Distance**

Alejandro Cassis Saarland University MPI-INF, SIC Tomasz Kociumaka MPI-INF. SIC

Philip Wellnitz

MPI-INF, SIC

1



- Bounded Weighted Edit Distance is harder than Bounded Unweighted Edit Distance
- (Weighted) Edit Distance is Shortest Paths on Planar Graphs
- APSP-based Lower Bounds via Grid Construction and Dynamic Intermediate Problem





#### How similar are two strings X and Y?







#### How similar are two strings X and Y?









#### How similar are two strings X and Y?





















ed(0PIN1CIV, OPINION) = 5













#### ed(0PIN1CIV, OPINION) = 5

ed(0PIN1CIV, PICNIC) = 5



# Background



 $w(0,0) := 1 \quad w(1,1) := 1 \quad w(C,0) := 1 \quad w(*,*) := 2 \quad w(*,\varepsilon) := 1 \quad w(\varepsilon,*) := 10$ 









# Background



 $w(0,0) := 1 \quad w(1,1) := 1 \quad w(C,0) := 1 \quad w(*,*) := 2 \quad w(*,\varepsilon) := 1 \quad w(\varepsilon,*) := 10$ 



0 P I N I 0 N \* | | | \* \* × × 0 P I N 1 C I V

ed<sup>w</sup>(0PIN1CIV, PICNIC) ≤ 14

ed<sup>w</sup>(0PIN1CIV, 0PINION) = 6





# Background



 $w(0,0) := 1 \quad w(1,1) := 1 \quad w(C,0) := 1 \quad w(*,*) := 2 \quad w(*,\varepsilon) := 1 \quad w(\varepsilon,*) := 10$ 



0 P I N I 0 N \* | | | \* \* \* × 0 P I N 1 C I V

ed<sup>w</sup>(0PIN1CIV, PICNIC) = 8

ed<sup>w</sup>(0PIN1CIV, 0PINION) = 6





# Computing (Weighted) Edit Distance

# How fast can we compute the (weighted) edit distance of two strings? $\rightsquigarrow O(n^2)$ (you know this!)





# Computing (Weighted) Edit Distance

# How fast can we compute the (weighted) edit distance of two strings? $\rightsquigarrow O(n^2)$ (you know this!)

#### What if the (weighted) edit distance is small (at most k)?







# How fast can we compute the (weighted) edit distance of two strings? $\rightsquigarrow O(n^2)$ (you know this!)

# What if the (weighted) edit distance is small (at most *k*)?

→→ Bounded (Weighted) Edit Distance





# Algorithms for (Bounded) Edit Distance



Existing algorithms for Edit Distance ed(X, Y), where  $|X|, |Y| \le n$ 



Cassis, Kociumaka, **Wellnitz** Optimal Algorithms for Bounded Weighted Edit Distance

7-1

# Algorithms for (Bounded) Edit Distance



Existing algorithms for Edit Distance ed(X, Y), where  $|X|, |Y| \le n$ 



Cassis, Kociumaka, **Wellnitz** 7-2

# Algorithms for (Bounded) Edit Distance



Existing algorithms for Edit Distance ed(X, Y), where  $|X|, |Y| \le n$ 



Cassis, Kociumaka, **Wellnitz** 7-3

# Algorithms for (Bounded) Weighted Edit Distance

#### What about Weighted Edit Distance?





## Algorithms for (Bounded) Weighted Edit Distance



Existing algorithms for Weighted Edit Distance  $ed^{w}(X, Y)$ , where  $|X|, |Y| \le n$ 





## Algorithms for (Bounded) Weighted Edit Distance



Existing algorithms for Weighted Edit Distance  $ed^{w}(X, Y)$ , where  $|X|, |Y| \le n$ 



## Algorithms for (Bounded) Weighted Edit Distance



Existing algorithms for Weighted Edit Distance  $ed^{w}(X, Y)$ , where  $|X|, |Y| \le n$ 





### hms for (Poundad) Waightad Edit Dis

### Algorithms for (Bounded) Weighted Edit Distance



Existing algorithms for Weighted Edit Distance  $ed^{w}(X, Y)$ , where  $|X|, |Y| \le n$ 





### Main Theorem 1 (Upper Bound)

Strings X, Y each of length at most n Oracle access to (normalized) weight function w Can compute  $k = ed^{w}(X, Y)$  in time  $O(n + \sqrt{nk^3} \log^3 n)$ 







### Main Theorem 1 (Upper Bound)

Strings X, Y each of length at most n Oracle access to (normalized) weight function w Can compute  $k = ed^{w}(X, Y)$  in time  $O(n + \sqrt{nk^3} \log^3 n)$ 

## Main Theorem 2 (Lower Bound)

Assuming the APSP Hypothesis and for  $\sqrt{n} \le k \le n$ , Main Theorem 1 is tight (up to  $n^{o(1)}$ -factors)











• Justified Assumption: w is normalized,  $w(x, y) \ge 1$  for all  $x \ne y$ .





◆ Justified Assumption: w is normalized,  $w(x, y) \ge 1$  for all  $x \ne y$ .  $\rightarrow$  Have (via  $O(n + k^2)$  algo): Alignment  $A : X \rightarrow Y$  of unweighted cost  $\le k$ 





Justified Assumption: w is normalized, w(x, y) ≥ 1 for all x ≠ y.
 → Have (via O(n + k<sup>2</sup>) algo): Alignment A : X → Y of unweighted cost ≤ k
 → A aligns edit-free most of X with Y [DGHKS23]





 $\rightsquigarrow$  A aligns edit-free most of X with Y [DGHKS23]

Theorem (Universal Kernel)

Can trim X and Y to length- $O(k^4)$  strings X', Y' with  $\operatorname{ed}_{\leq k}^{W}(X, Y) = \operatorname{ed}_{\leq k}^{W}(X', Y')$ .



Cassis, Kociumaka, **Wellnitz**Optimal Algorithms for Bounded Weighted Edit Distance
10-5

[DGHKS23]







Tool 1: Alignment Graphs and Multiple-Source Shortest Path

• X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each





- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and Y









- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and Y  $\rightsquigarrow$  ed<sup>w</sup>(X, Y) is distance (0, 0)  $\rightsquigarrow$  (|X|, |Y|)









- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and  $Y \rightarrow ed^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>2</sup>: Trim AG to O(k) diagonals  $\rightarrow ed_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$









- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and  $Y \rightarrow ed^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>2</sup>: Trim AG to O(k) diagonals  $\rightarrow ed_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$ 
  - $\rightsquigarrow$  AG has  $O(k^5)$  vertices  $\rightsquigarrow$  Dijkstra yields  $\tilde{O}(k^5)$  algo









- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and  $Y \rightarrow ed^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>2</sup>: Trim AG to O(k) diagonals  $\rightarrow ed_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>3</sup>: Split AG according to structure of X and Y
   Compute all b-to-b dist [Kleino5] + stitch together results ((min, +)-product [SMAWK87])









- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and  $Y \rightarrow ed^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>2</sup>: Trim AG to O(k) diagonals  $\rightsquigarrow \operatorname{ed}_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightsquigarrow (|X|, |Y|)$
- Idea<sup>3</sup>: Split AG according to structure of X and Y
   Compute all b-to-b dist [Kleino5] + stitch together results ((min, +)-product [SMAWK87])
   → k<sup>2</sup> · Õ(k<sup>3</sup>) for periodic pieces + k<sup>2</sup> · Õ(k<sup>2</sup>) for stitching







# Upper Bounds

#### ••••••••••••

Tool 1: Alignment Graphs and Multiple-Source Shortest Path

- X and Y consist in  $O(k^2)$  pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use alignment graph AG of X and  $Y \rightarrow ed^{W}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>2</sup>: Trim AG to O(k) diagonals  $\rightarrow ed_{\leq k}^{w}(X, Y)$  is distance  $(0, 0) \rightarrow (|X|, |Y|)$
- Idea<sup>3</sup>: Split AG according to structure of X and Y
   Compute all b-to-b dist [Kleino5] + stitch together results ((min, +)-product [SMAWK87]) for periodic pieces: fast exponentiation;

 $\rightsquigarrow k^2\cdot \tilde{O}(k^2)$  for periodic pieces +  $k^2\cdot \tilde{O}(k^2)$  for stitching









- X and Y consist in  $O(k^2)$  periodic pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use AG, trimmed to O(k) diags  $\rightsquigarrow \operatorname{ed}_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightsquigarrow (|X|, |Y|)$
- Idea<sup>3</sup>: Compute all b-to-b dist for periodic pieces [Kleino5] + fast exponentiation; stitch together results using min-plus product [SMAWK87]
- Idea<sup>4</sup>: Use Divide-and-Conquer to reduce number of periodic pieces to O(k)





- X and Y consist in  $O(k^2)$  periodic pieces of length  $O(k^2)$  and with period O(k) each
- Idea: Use AG, trimmed to O(k) diags  $\rightsquigarrow \operatorname{ed}_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightsquigarrow (|X|, |Y|)$
- Idea<sup>3</sup>: Compute all b-to-b dist for periodic pieces [Kleino5] + fast exponentiation; stitch together results using min-plus product [SMAWK87]
- Idea<sup>4</sup>: Use Divide-and-Conquer to reduce number of periodic pieces to O(k) $\rightarrow k \cdot \tilde{O}(k^2)$  for periodic pieces (and padding) +  $k \cdot \tilde{O}(k^2)$  for stitching





# Tool 3: Compressibility instead of Periodicity

- X and Y consist in O(k) pieces of length  $O(k^2)$  each
- Idea: Use AG, trimmed to O(k) diags  $\rightsquigarrow \operatorname{ed}_{\leq k}^{W}(X, Y)$  is distance  $(0, 0) \rightsquigarrow (|X|, |Y|)$
- Idea<sup>3</sup>: Compute all b-to-b dist for periodic pieces [Kleino5] + fast exponentiation; stitch together results using min-plus product [SMAWK87]
- ◆ Idea<sup>4</sup>: Use Divide-and-Conquer to reduce number of periodic pieces to O(k)
- Idea<sup>5</sup>: Use tailor-made compressibility measure instead of periodicity
   + (w)ED algorithms for compressed strings

 $\rightsquigarrow \tilde{O}(n + \sqrt{k^3 n})$  time in total





### Main Theorem 1 (Upper Bound) 🗸

Strings X, Y each of length at most n Oracle access to (normalized) weight function w Can compute  $k = ed^{w}(X, Y)$  in time  $O(n + \sqrt{nk^3} \log^3 n)$ 

## Main Theorem 2 (Lower Bound)

Assuming the APSP Hypothesis and for  $\sqrt{n} \le k \le n$ , Main Theorem 1 is tight (up to  $n^{o(1)}$ -factors)



14







































- ◆ Step 1: Compute matrix-vector product A ⊕ b as weighted ED of two gadget strings
   → encode matrix/vector as substitution costs
   → use telescoping sums
- Step 2: Compute A ⊕ b' by replacing single character of X







#### $\bullet \bullet \circ \circ \circ$

Tool 1: Matrix-Vector Min-Plus Multiplication via wED

### **Min-Plus Product Minimum**

Given 3 matrices A, B, C, check if min diag entry of (min, +)-product  $A \oplus B \oplus C$  is negative

- ◆ Step 1: Compute matrix-vector product A ⊕ b as weighted ED of two gadget strings
   → encode matrix/vector as substitution costs
   → use telescoping sums
- Step 2: Compute A ⊕ b' by replacing single character of X
   → Allows to compute A ⊕ B; encode C similarly to B
   → Use selection gadget for ed<sup>w</sup>(X, Y) ≈ min(A ⊕ B ⊕ C)
  - $\leadsto$  Lower bound for dynamic problem







Tool 2: Minimum Gadget via Intermediate Strings

◆ Have: Lower bound for dynamic problem
 → turn into LB for static problem





#### **\*\*\*\*\*\*\*\*\*\*\*\*\*\***

Tool 2: Minimum Gadget via Intermediate Strings

- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)



...

X<sub>1</sub>

 $X_{2}$ 

X

Χ,

γ

γ

γ

Y



- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings  $(Y, \hat{X}_{(i-1)i})$







#### ••••••

- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost







- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost







#### ••••••

- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost
- Step 5: Forbid inserting or deleting Y (simplification)
   (w(ε, Y) := ∞, w(Y, ε) := ∞)
   → Have to align exactly one X; to Y







#### •••••••••••••••

- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost
- Step 5: Forbid inserting or deleting Y (simplification)
   (w(ε, Y) := ∞, w(Y, ε) := ∞)
   → Have to align exactly one X; to Y







- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost
- Step 5: Forbid inserting or deleting Y (simplification)
   (w(ε, Y) := ∞, w(Y, ε) := ∞)
   → Have to align exactly one X; to Y







#### **◆◆◆◆◆◆◆◆◆◆◆◆◆◆**◆◇

- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost
- Step 5: Forbid inserting or deleting Y (simplification)
   (w(ε, Y) := ∞, w(Y, ε) := ∞)
   → Have to align exactly one X; to Y







- ◆ Have: Lower bound for dynamic problem
   → turn into LB for static problem
- Step 3: Take snapshots (X<sub>i</sub>, Y)
- Step 4: Create intermediate strings (Y, X̂<sub>(i-1)i</sub>)
   → may align X̂<sub>i</sub> with either X<sub>i</sub> or X<sub>i+1</sub> for the same cost
- Step 5: Forbid inserting or deleting Y (simplification)
   (w(ε, Y) := ∞, w(Y, ε) := ∞)
  - $\rightsquigarrow$  Have to align exactly one  $X_i$  to Y
  - $\rightsquigarrow$  Minimum gadget:  $ed^{w}(\hat{X}, \hat{Y}) \approx \min ed^{w}(X_{i}, Y)$
  - $\rightsquigarrow$  Static LB







### Main Theorem 1 (Upper Bound) ✓

Strings X, Y each of length at most n Oracle access to (normalized) weight function w Can compute  $k = ed^w(X, Y)$  in time  $O(n + \sqrt{nk^3} \log^3 n)$ 

## Main Theorem 2 (Lower Bound) 🗸

Assuming the APSP Hypothesis and for  $\sqrt{n} \le k \le n$ , Main Theorem 1 is tight (up to  $n^{o(1)}$ -factors)



18

# Key Messages and Open Problems

Key Messages

- Bounded Weighted Edit Distance is harder than Bounded Unweighted Edit Distance
- (Weighted) Edit Distance is Shortest Paths on Planar Graphs
- APSP-based Lower Bounds via Grid Construction and Dynamic Intermediate Problem



# Key Messages and Open Problems

Key Messages

- Bounded Weighted Edit Distance is harder than Bounded Unweighted Edit Distance
- (Weighted) Edit Distance is Shortest Paths on Planar Graphs
- APSP-based Lower Bounds via Grid Construction and Dynamic Intermediate Problem

**Open Questions** 

- What is the true complexity of  $n^{1/3} \le k \le n^{1/2}$ ? Must be between  $n + k^{2.5}$  and  $\sqrt{nk^3}$ .
- Is the problem easier for small (constant-sized) alphabets?
- Is there any easy class of weight functions?







