Approximate matching in grammar-compressed strings

Alexander Tiskin

Department of Computer Science, University of Warwick http://go.warwick.ac.uk/alextiskin

< 4 ► >

ヨト イヨト



- 2 Matrix distance multiplication
- 3 Semi-local string comparison
- 4 Compressed string comparison
- 5 Conclusions and future work

Alexander Tiskin (Warwick)



- Matrix distance multiplication
- 3 Semi-local string comparison
- 4 Compressed string comparison
- 5 Conclusions and future work

Alexander Tiskin (Warwick)

String matching: finding an exact pattern in a string String comparison: finding similar patterns in two strings

イロト イポト イヨト イヨト

String matching: finding an exact pattern in a string String comparison: finding similar patterns in two strings Standard types of string comparison:

- global: whole string vs whole string
- local: substrings vs substrings

Main focus of this work:

• semi-local: whole string vs substrings; prefixes vs suffixes

Closely related to approximate string matching: whole pattern string vs substrings in text string

Standard approach to string comparison: dynamic programming Our approach: the algebra of unit-Monge matrices, a.k.a. seaweed braids Can be used either for dynamic programming, or for divide-and-conquer Divide-and-conquer is better for:

- comparing strings dynamically (truncation, concatenation)
- comparing compressed srtrings (e.g. LZ-compression)
- comparing strings in parallel

The "conquer" step is non-trivial — but not so much that it couldn't have been discovered much earlier!

- 4 週 ト - 4 三 ト - 4 三 ト -

The "squared paper" notation: Integers {... - 2, -1, 0, 1, 2, ...} $x^- = x - \frac{1}{2}$ $x^+ = x + \frac{1}{2}$ Half-integers {... - $\frac{3}{2}$, $-\frac{1}{2}$, $\frac{1}{2}$, $\frac{3}{2}$, $\frac{5}{2}$, ...} = {... (-2)⁺, (-1)⁺, 0⁺, 1⁺, 2⁺} Planar dominance:

- $(i, j) \ll (i', j')$ iff i < i' and j < j' (the "above-left" relation)
- $(i,j) \ge (i',j')$ iff i > i' and j < j' (the "below-left" relation)

where "above/below" follow matrix convention (not the Cartesian one!)

イロト イ理ト イヨト イヨト 二日

A permutation matrix is a $0/1\ matrix$ with exactly one nonzero per row and per column



イロト イポト イヨト イヨト

Given matrix *D*, its distribution matrix is made up of \geq -dominance sums: $D^{\Sigma}(i,j) = \sum_{i>i,j<j} D(i,j)$

3

Given matrix D, its distribution matrix is made up of \geq -dominance sums: $D^{\Sigma}(i,j) = \sum_{i>i,j < j} D(i,j)$

Given matrix *E*, its density matrix is made up of quadrangle differences: $E^{\Box}(\hat{\imath}, \hat{\jmath}) = E(\hat{\imath}^{-}, \hat{\jmath}^{+}) - E(\hat{\imath}^{-}, \hat{\jmath}^{-}) - E(\hat{\imath}^{+}, \hat{\jmath}^{+}) + E(\hat{\imath}^{+}, \hat{\jmath}^{-})$

where D^{Σ} , *E* over integers; *D*, E^{\Box} over half-integers

Given matrix D, its distribution matrix is made up of \geq -dominance sums: $D^{\Sigma}(i,j) = \sum_{i>i,j < j} D(i,j)$

Given matrix *E*, its density matrix is made up of quadrangle differences: $E^{\Box}(\hat{\imath}, \hat{\jmath}) = E(\hat{\imath}^{-}, \hat{\jmath}^{+}) - E(\hat{\imath}^{-}, \hat{\jmath}^{-}) - E(\hat{\imath}^{+}, \hat{\jmath}^{+}) + E(\hat{\imath}^{+}, \hat{\jmath}^{-})$

where D^{Σ} , E over integers; D, E^{\Box} over half-integers

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{\Sigma} = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \qquad \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}^{\sqcup} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

イロト 不得下 イヨト イヨト 二日

Given matrix D, its distribution matrix is made up of \geq -dominance sums: $D^{\Sigma}(i,j) = \sum_{i>i,j < j} D(i,j)$

Given matrix E, its density matrix is made up of quadrangle differences: $E^{\Box}(\hat{\imath}, \hat{\jmath}) = E(\hat{\imath}^{-}, \hat{\jmath}^{+}) - E(\hat{\imath}^{-}, \hat{\jmath}^{-}) - E(\hat{\imath}^{+}, \hat{\jmath}^{+}) + E(\hat{\imath}^{+}, \hat{\jmath}^{-})$

where D^{Σ} , E over integers; D, E^{\Box} over half-integers

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{\Sigma} = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \qquad \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}^{\Box} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

 $(D^{\Sigma})^{\Box} = D$ for all D

Matrix *E* is simple, if $(E^{\Box})^{\Sigma} = E$; equivalently, if it has all zeros in the left column and bottom row

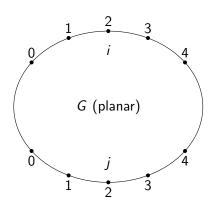
イロト 不得下 イヨト イヨト 二日

Introduction Terminology and notation

Matrix *E* is Monge, if E^{\Box} is nonnegative

Intuition: boundary-to-boundary distances in a (weighted) planar graph





G. Monge (1746-1818)

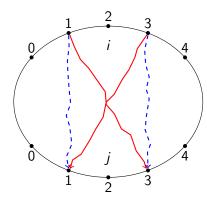
Alexander Tiskin (Warwick)

Introduction Terminology and notation

Matrix *E* is Monge, if E^{\Box} is nonnegative

Intuition: boundary-to-boundary distances in a (weighted) planar graph





G. Monge (1746-1818)

 $blue + blue \leq red + red$

Alexander Tiskin (Warwick)

Matrix *E* is unit-Monge, if E^{\Box} is a permutation matrix Intuition: boundary-to-boundary distances in a grid-like graph

< A[™] →

-∢ ∃ ▶

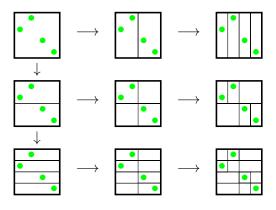
Matrix *E* is unit-Monge, if E^{\Box} is a permutation matrix Intuition: boundary-to-boundary distances in a grid-like graph Simple unit-Monge matrix: P^{Σ} , where *P* is a permutation matrix P^{Σ} sometimes called the rank function of *P*

Seaweed matrix: P used as an implicit representation of P^{Σ}

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{\Sigma} = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Efficient P^{Σ} queries: range tree on nonzeros of P [Bentley: 1980]

- binary search tree by *i*-coordinate
- under every node, binary search tree by *j*-coordinate



Efficient P^{Σ} queries: (contd.)

Every node of the range tree represents a canonical range (rectangular region), and stores its nonzero count

Overall, $\leq n \log n$ canonical ranges are non-empty

A P^{Σ} query is equivalent to \geq -dominance counting: how many nonzeros are \geq -dominated by query point?

Answer: sum up nonzero counts in $\leq \log^2 n$ disjoint canonical ranges Total size $O(n \log n)$, query time $O(\log^2 n)$

Efficient P^{Σ} queries: (contd.)

Every node of the range tree represents a canonical range (rectangular region), and stores its nonzero count

Overall, $\leq n \log n$ canonical ranges are non-empty

A P^{Σ} query is equivalent to \geq -dominance counting: how many nonzeros are \geq -dominated by query point?

Answer: sum up nonzero counts in $\leq \log^2 n$ disjoint canonical ranges

Total size $O(n \log n)$, query time $O(\log^2 n)$

There are asymptotically more efficient (but less practical) data structures

Total size O(n), query time $O(\frac{\log n}{\log \log n})$ [JáJá+: 2004] [Chan, Pǎtrascu: 2010]

・ロト ・四ト ・ヨト ・ヨト ・ヨ



milloudection

- 2 Matrix distance multiplication
- 3 Semi-local string comparison
- 4 Compressed string comparison
- 5 Conclusions and future work

Alexander Tiskin (Warwick)

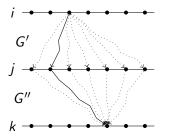
Distance algebra (a.k.a (min, +) or tropical algebra):

 \bullet addition \oplus given by min, multiplication \odot given by +

 $\mathsf{Matrix} \odot \text{-multiplication}$

 $A \odot B = C$ $C(i,k) = \bigoplus_{j} (A(i,j) \odot B(j,k)) = \min_{j} (A(i,j) + B(j,k))$

Intuition: gluing distances in a union of graphs



Alexander Tiskin (Warwick)

Matrix classes closed under \odot -multiplication (for given *n*):

- ullet general (integer, real) matrices \sim general weighted graphs
- ullet Monge matrices \sim planar weighted graphs
- ullet simple unit-Monge matrices \sim grid-like graphs

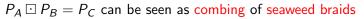
Implicit \odot -multiplication: define $P_A \boxdot P_B = P_C$ as $P_A^{\Sigma} \odot P_B^{\Sigma} = P_C^{\Sigma}$ The seaweed monoid \mathcal{T}_n :

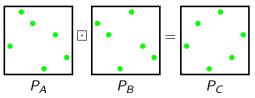
- ullet simple unit-Monge matrices under \odot
- equivalently, permutation (seaweed) matrices under 🖸

Also known as the 0-Hecke monoid of the symmetric group $H_0(S_n)$

・ロト ・ 同ト ・ ヨト ・ ヨト - -

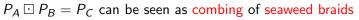
Seaweed braids

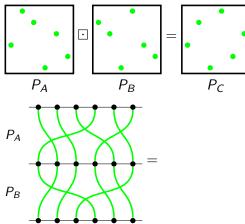




3 🕨 3

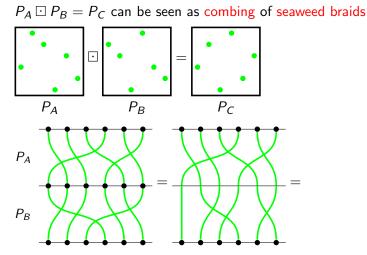
Seaweed braids





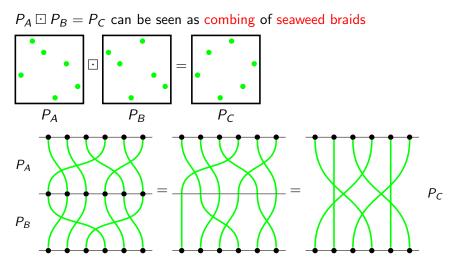
3 ▶ 3

Seaweed braids



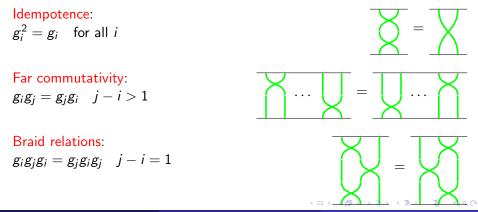
э

Seaweed braids



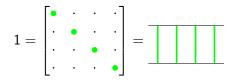
The seaweed monoid T_n :

- n! elements (permutations of size n)
- n-1 generators $g_1, g_2, \ldots, g_{n-1}$ (elementary crossings)

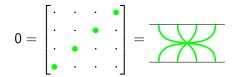


Seaweed braids

Identity: $1 \boxdot x = x$



Zero: $0 \boxdot x = 0$



Related structures:

- positive braids: far comm; braid relations
- braids: $g_i g_i^{-1} = 1$; far comm; braid relations
- Coxeter's presentation of S_n : $g_i^2 = 1$; far comm; braid relations
- locally free idempotent monoid: idem; far comm [Vershik+: 2000]

Generalisations:

- general 0-Hecke monoids [Fomin, Greene: 1998; Buch+: 2008]
 Coxeter monoids [Tsaranov: 1990; Richardson, Springer: 1990]
- \mathcal{J} -trivial monoids

Richardson, Springer: 1990 [Denton+: 2011]

Seaweed braids

Computation in the seaweed monoid: a confluent rewriting system can be obtained by software (SEMIGROUPE, GAP)

B ▶ < B ▶

Seaweed braids

Computation in the seaweed monoid: a confluent rewriting system can be obtained by software (SEMIGROUPE, GAP)

$$\mathcal{T}_3$$
: 1, $a = g_1$, $b = g_2$; ab , ba ; $aba = 0$

aa
ightarrow a bb
ightarrow b bab
ightarrow 0 aba
ightarrow 0

B ▶ < B ▶

Computation in the seaweed monoid: a confluent rewriting system can be obtained by software (SEMIGROUPE, GAP)

$$\mathcal{T}_3$$
: 1, $a = g_1, \ b = g_2; \ ab, \ ba; \ aba = 0$
 $aa o a \qquad bb o b \qquad bab o 0 \qquad aba o 0$

 \mathcal{T}_4 : 1, $a = g_1$, $b = g_2$, $c = g_3$; ab, ac, ba, bc, cb, aba, abc, acb, bac, bcb, cba, abac, abcb, acba, bacb, bcba, abacb, abcba; abacba = 0

- 4 @ > - 4 @ > - 4 @ > -

Computation in the seaweed monoid: a confluent rewriting system can be obtained by software (SEMIGROUPE, GAP)

$$\mathcal{T}_3$$
: 1, $a=g_1, \ b=g_2; \ ab, \ ba; \ aba=0$
 $aa o a bb o b bb bb bab o 0 abab o 0$

 \mathcal{T}_4 : 1, $a = g_1$, $b = g_2$, $c = g_3$; ab, ac, ba, bc, cb, aba, abc, acb, bac, bcb, cba, abac, abcb, acba, bacb, bcba, abacb, abcba; abacba = 0

Easy to use, but not an efficient algorithm

イロト イポト イヨト イヨト 二日

Seaweed matrix multiplication

The implicit unit-Monge matrix \odot -multiplication problem

Given permutation matrices P_A , P_B , compute P_C , such that $P_A^{\Sigma} \odot P_B^{\Sigma} = P_C^{\Sigma}$ (equivalently, $P_A \boxdot P_B = P_C$)

Seaweed matrix multiplication

The implicit unit-Monge matrix \odot -multiplication problem

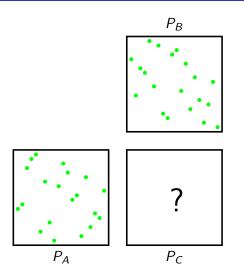
Given permutation matrices P_A , P_B , compute P_C , such that $P_A^{\Sigma} \odot P_B^{\Sigma} = P_C^{\Sigma}$ (equivalently, $P_A \boxdot P_B = P_C$)

Matrix ⊙-multiplication: running time

type	time	
general	$O(n^3)$	standard
	$O\left(\frac{n^3(\log\log n)^3}{\log^2 n}\right)$	[Chan: 2007]
Monge	$O(n^2)$	via [Aggarwal+: 1987]
implicit unit-Monge	$O(n^{1.5})$	[T: 2006]
	$O(n \log n)$	[T: 2010]

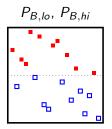
イロト イ押ト イヨト イヨト

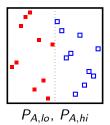
Seaweed matrix multiplication



2

Seaweed matrix multiplication

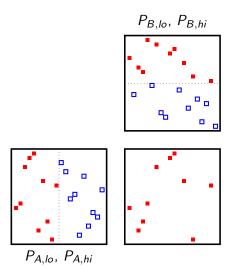




Alexander Tiskin (Warwick)

э

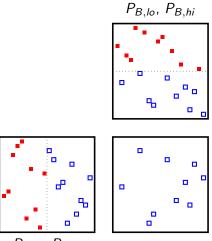
Seaweed matrix multiplication



æ

э

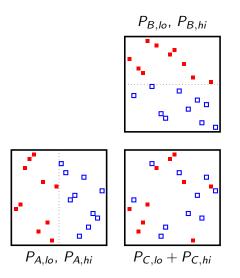
Seaweed matrix multiplication



 $P_{A,lo}, P_{A,hi}$

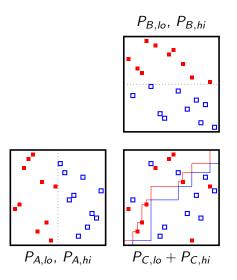
э

Seaweed matrix multiplication

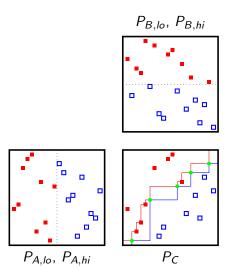


æ

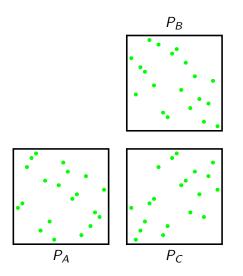
Seaweed matrix multiplication



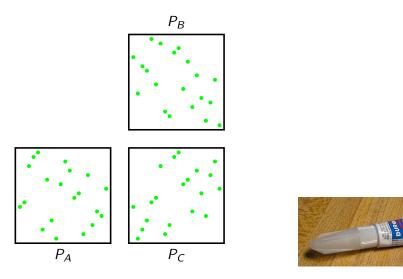
Seaweed matrix multiplication



Seaweed matrix multiplication



Seaweed matrix multiplication



Seaweed matrix multiplication

Implicit unit-Monge matrix \odot -multiplication: the algorithm

$$P_C^{\Sigma}(i,k) = \min_j \left(P_A^{\Sigma}(i,j) + P_B^{\Sigma}(j,k) \right)$$

Divide-and-conquer on the range of j: two subproblems of size n/2 $P_{A,lo}^{\Sigma} \odot P_{B,lo}^{\Sigma} = P_{C,lo}^{\Sigma}$ $P_{A,hi}^{\Sigma} \odot P_{B,hi}^{\Sigma} = P_{C,hi}^{\Sigma}$ Conquer: tracing a balanced path from bottom-left to top-right P_C Path invariant: equal number of \ll -dominated nonzeros in $P_{C,hi}$ and \ll -dominating nonzeros in $P_{C,lo}$ All such nonzeros are "errors" in the output, must be removed

Must compensate by placing nonzeros on the path, time O(n)

Overall time $O(n \log n)$

イロト イ理ト イヨト イヨト

Bruhat order

Bruhat order

Permutation A is lower ("more sorted") than permutation B in the Bruhat order $(A \leq B)$, if $B \rightsquigarrow A$ by successive pairwise sorting (equivalently, $A \rightsquigarrow B$ by anti-sorting) of arbitrary pairs

Permutation matrices: $P_A \leq P_B$, if $P_B \rightsquigarrow P_A$ by successive 2 × 2 submatrix sorting: $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

Bruhat order

Bruhat order

Permutation A is lower ("more sorted") than permutation B in the Bruhat order $(A \leq B)$, if $B \rightsquigarrow A$ by successive pairwise sorting (equivalently, $A \rightsquigarrow B$ by anti-sorting) of arbitrary pairs

Permutation matrices: $P_A \leq P_B$, if $P_B \rightsquigarrow P_A$ by successive 2 × 2 submatrix sorting: $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

Plays an important role in group theory and algebraic geometry (inclusion order of Schubert varieties)

Describes pivoting order in Gaussian elimination (LUP/Bruhat decomposition)

イロト イ理ト イヨト イヨト

Bruhat order

Bruhat comparability: running time

$O(n^2)$
$O(n \log n)$
$O\left(\frac{n\log n}{\log\log n}\right)$

[Ehresmann: 1934; Proctor: 1982; Grigoriev: 1982] [T: 2013] [Gawrychowski: NEW]

Ehresmann's criterion (dot criterion, related to tableau criterion) $P_A \preceq P_B$ iff $P_A^{\Sigma} \leq P_B^{\Sigma}$ elementwise

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{\Sigma} = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \leq \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 2 & 2 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}^{\Sigma}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{\Sigma} = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \not\leq \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{\Sigma}$$
Time $O(n^2)$

- ∢ ∃ ▶

Matrix distance multiplication Bruhat order

Seaweed criterion

 $P_A \preceq P_B$ iff $P_A^R \boxdot P_B = Id^R$, where $P^R =$ clockwise rotation of matrix P

Intuition: permutations represented by seaweed braids

 $P_A \preceq P_B$, iff

- no pair of seaweeds is crossed in P_A , while the "corresponding" pair is uncrossed in P_B
- equivalently, no pair is uncrossed in P_A^R , while the "corresponding" pair is uncrossed in P_B

• equivalently,
$$P_A^R \boxdot P_B = Id^R$$

Time $O(n \log n)$ by seaweed matrix multiplication

イロト イポト イヨト イヨト 二日

Bruhat order

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \preceq \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{R} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} = Id^{R}$$



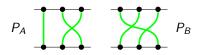
æ

イロト イヨト イヨト イヨト

Alexander Tiskin (Warwick)

Bruhat order

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \preceq \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{R} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} = Id^{R}$$



< □ > < ---->

æ

∃ ► < ∃ ►</p>

Bruhat order

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \preceq \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{R} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} = Id^{R}$$
$$P_{A} \longrightarrow P_{B} \xrightarrow{P_{A}} P_{B} \longrightarrow P_{B}$$

æ

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Bruhat order

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \preceq \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{R} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \boxdot \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} = Id^{R}$$
$$P_{A} \longrightarrow P_{B} \qquad P_{B} \qquad P_{B} \qquad P_{B} \qquad P_{B} \qquad Id^{R}$$

æ

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Bruhat order

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \not\preceq \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{R} \boxdot \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \boxdot \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \neq Id^{R}$$



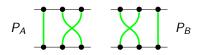
æ

イロト イヨト イヨト イヨト

Alexander Tiskin (Warwick)

Bruhat order

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \not\preceq \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}^{R} \boxdot \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \boxdot \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \neq Id^{R}$$

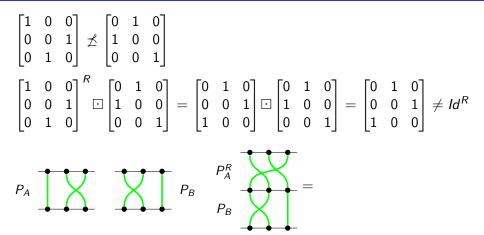


< □ > < ---->

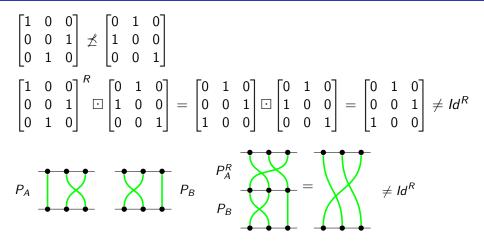
æ

∃ ► < ∃ ►</p>

Bruhat order



Bruhat order



Alternative solution: clever implementation of Ehresmann's criterion [Gawrychowski: 2013]

The online partial sums problem: maintain array X[1:n], subject to

- $update(k, \Delta)$: $X[k] \leftarrow X[k] + \Delta$
- prefixsum(k): return $\sum_{1 \le i \le k} X[i]$

Query time:

Θ(log n) in semigroup or group model
 Θ(log n) in RAM model on integers [Pătrașcu, Demaine: 2004]

Gives Bruhat comparability in time $O(\frac{n \log n}{\log \log n})$ in RAM model

Open problem: seaweed multiplication in time $O(\frac{n \log n}{\log \log n})$?

(日) (四) (王) (王) (王)



- Semi-local string comparison

Alexander Tiskin (Warwick)

Consider strings (= sequences) over an alphabet of size σ

- Contiguous substrings vs not necessarily contiguous subsequences
- Special cases of substring: prefix, suffix
- Notation: strings a, b of length m, n respectively
- Assume where necessary: $m \le n$; m, n reasonably close

Consider strings (= sequences) over an alphabet of size σ

- Contiguous substrings vs not necessarily contiguous subsequences
- Special cases of substring: prefix, suffix
- Notation: strings a, b of length m, n respectively
- Assume where necessary: $m \le n$; m, n reasonably close
- The longest common subsequence (LCS) score:
 - length of longest string that is a subsequence of both a and b
 - equivalently, alignment score, where score(match) = 1 and score(mismatch) = 0

In biological terms, "loss-free alignment" (unlike efficient but "lossy" BLAST)

The LCS problem

Alexander Tiskin (Warwick)

Give the LCS score for a vs b

Image: Image:

э

.∋...>

The LCS problem

Give the LCS score for a vs b

LCS: running time

 $O\left(\frac{mn(\log\log n)^2}{\log^2 n}\right)$

$$O(mn) \\ O\left(\frac{mn}{\log^2 n}\right) \qquad \sigma = O(1)$$

[Wagner, Fischer: 1974] [Masek, Paterson: 1980] [Crochemore+: 2003] [Paterson, Dančík: 1994] [Bille, Farach-Colton: 2008]

Running time varies depending on the RAM model version

We assume word-RAM with word size $\log n$ (where it matters)

LCS computation by dynamic programming

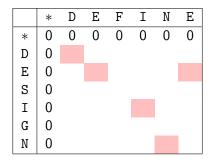
$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

Image: A matrix of the second seco

∃ ► < ∃ ►</p>

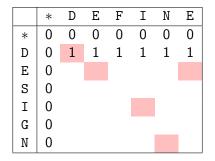
LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$



LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

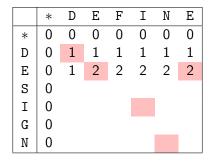


- < A

∃ >

LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

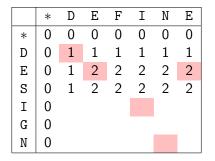


- < A

∃ >

LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

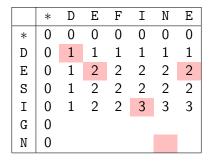


< A

3 🕨 3

LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

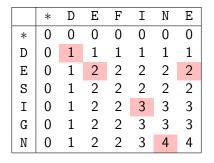


< A

3) 3

LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

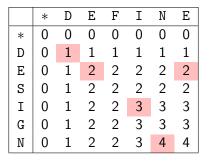


< <>></>

3) 3

LCS computation by dynamic programming

$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$

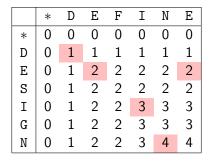


< <>></>

3) 3

LCS computation by dynamic programming

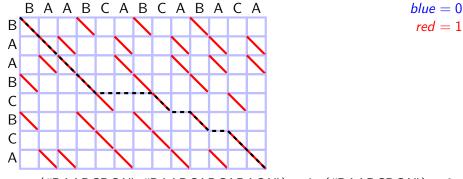
$$\begin{aligned} & lcs(a, \emptyset) = 0 \\ & lcs(\emptyset, b) = 0 \end{aligned} \qquad \qquad lcs(a\alpha, b\beta) = \begin{cases} \max(lcs(a\alpha, b), lcs(a, b\beta)) & \text{if } \alpha \neq \beta \\ & lcs(a, b) + 1 & \text{if } \alpha = \beta \end{cases} \end{aligned}$$



lcs("DEFINE", "DESIGN") = 4LCS(a, b) can be traced back through the dynamic programming table at no

extra asymptotic time cost

LCS on the alignment graph (directed, acyclic)



score("BAABCBCA", "BAABCABCABACA") = *len*("BAABCBCA") = 8

LCS = highest-score path from top-left to bottom-right

LCS: dynamic programming

Sweep cells in any \ll -compatible order

```
Cell update: time O(1)
```

Overall time O(mn)

Alexander Tiskin (Warwick)

[WF: 1974]

LCS: micro-block dynamic programming

[MP: 1980; BF: 2008]

Sweep cells in micro-blocks, in any $\ll\text{-compatible order}$

Micro-block size:

- $t = O(\log n)$ when $\sigma = O(1)$
- $t = O(\frac{\log n}{\log \log n})$ otherwise

Micro-block interface:

- O(t) characters, each $O(\log \sigma)$ bits, can be reduced to $O(\log t)$ bits
- O(t) small integers, each O(1) bits

Micro-block update: time O(1), by precomputing all possible interfaces

Overall time $O(\frac{mn}{\log^2 n})$ when $\sigma = O(1)$, $O(\frac{mn(\log \log n)^2}{\log^2 n})$ otherwise



'Begin at the beginning,' the King said gravely, 'and go on till you come to the end: then stop.'

L. Carroll, Alice in Wonderland



'Begin at the beginning,' the King said gravely, 'and go on till you come to the end: then stop.'

L. Carroll, Alice in Wonderland

Dynamic programming: begins at empty strings, proceeds by appending characters, then stops

What about:

• prepending/deleting characters (dynamic LCS)

< < p>< < p>

- concatenating strings (LCS on compressed strings; parallel LCS)
- taking substrings (= local alignment)



Dynamic programming from both ends: better by $\times 2$, but still not good enough

Is dynamic programming strictly necessary to solve sequence alignment problems?

Eppstein+, *Efficient algorithms for sequence analysis*, 1991

The semi-local LCS problem

Give the (implicit) matrix of $O((m + n)^2)$ LCS scores:

- string-substring LCS: string *a* vs every substring of *b*
- prefix-suffix LCS: every prefix of a vs every suffix of b
- suffix-prefix LCS: every suffix of a vs every prefix of b
- substring-string LCS: every substring of a vs string b

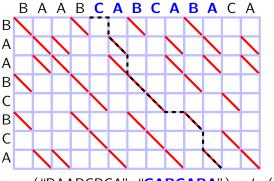
The semi-local LCS problem

Give the (implicit) matrix of $O((m + n)^2)$ LCS scores:

- string-substring LCS: string *a* vs every substring of *b*
- prefix-suffix LCS: every prefix of a vs every suffix of b
- suffix-prefix LCS: every suffix of a vs every prefix of b
- substring-string LCS: every substring of a vs string b



Semi-local LCS on the alignment graph



blue = 0red = 1

score("BAABCBCA", "CABCABA") = len("ABCBA") = 5

String-substring LCS: all highest-score top-to-bottom paths Semi-local LCS: all highest-score boundary-to-boundary paths

Semi-local string comparison Score matrices and seaweed matrices

The score matrix H

0	1	2	3	4	5	6	6	7	8	8	8	8	8
-1	0	1	2	3	4	5	5	6	7	7	7	7	7
-2	-1	0	1	2	3	4	4	5	6	6	6	6	7
-3	-2	-1	0	1	2	3	3	4	5	5	6	6	7
-4	-3	-2	-1	0	1	2	2	3	4	4	5	5	6
-5	-4	-3	-2	-1	0	1	2	3	4	4	5	5	6
-6	-5	-4	-3	-2	-1	0	1	2	3	3	4	4	5
-7	-6	-5	-4	-3	-2	-1	0	1	2	2	3	3	4
-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	3	4
-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4
-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3
-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2
-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1
-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0

- a = "BAABCBCA"
- b = "BAAB**CABCABA**CA"
- $H(i,j) = score(a, b\langle i : j \rangle)$

$$H(4, 11) = 5$$

$$H(i,j) = j - i \text{ if } i > j$$

< 🗇

э

- ∢ ∃ ▶

Score matrices and seaweed matrices

Semi-local LCS: output representation and running time

size	query time			
$\overline{O(n^2)}$	<i>O</i> (1)			trivial
$O(m^{1/2}n)$	$O(\log n)$	string-substring		[Alves+: 2003]
O(n)	<i>O</i> (<i>n</i>)	string-substring		[Alves+: 2005]
$O(n \log n)$	$O(\log^2 n)$			[T: 2006]
or any	2D orthogona	al range counting	data structi	ure
running tim	e			
$O(mn^2)$				naive
O(mn)	string-	substring	[Schmidt:	1998; Alves+: 2005]
O(mn)				[T: 2006]
$O\left(\frac{mn}{\log^{0.5}n}\right)$				[T: 2006]
$O\Big(\frac{mn(\log\log \log n)}{\log^2 n}\Big)$	$\frac{(n)^2}{(n-1)^2}$			[T: 2007]

The score matrix H and the seaweed matrix P

H(i,j): the number of matched characters for *a* vs substring $b\langle i : j \rangle$ j - i - H(i,j): the number of unmatched characters Properties of matrix j - i - H(i,j):

- simple unit-Monge
- therefore, $= P^{\Sigma}$, where $P = -H^{\Box}$ is a permutation matrix

P is the seaweed matrix, giving an implicit representation of *H* Range tree for *P*: memory $O(n \log n)$, query time $O(\log^2 n)$

- ・ 伺 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Semi-local string comparison Score matrices and seaweed matrices

The score matrix H and the seaweed matrix P

0	1	2	3	4	5	6	6	7	8	8	8	8	8
-1	0	1	2	3	4	5	5	6	7	7	7	7	7
-2	-1	0	1	2	3	4	4	5	6	6	6	6	7
-3	-2	-1	0	1	2	3	3	4	5	5	6	6	7
-4	-3	-2	-1	0	1	2	2	3	4	4	5	5	6
-5	-4	-3	-2	-1	0	1	2	3	4	4	5	5	6
-6	-5	-4	-3	-2	-1	0	1	2	3	3	4	4	5
-7	-6	-5	-4	-3	-2	-1	0	1	2	2	3	3	4
-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	3	4
-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4
-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3
-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2
-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1
-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0

- a = "BAABCBCA"
- *b* = "BAAB**CABCABA**CA"

$$H(i,j) = score(a, b\langle i : j \rangle)$$

$$H(4, 11) = 5$$

$$H(i,j) = j - i \text{ if } i > j$$

The score matrix H and the seaweed matrix P

0	1	2	3	4	5	6	6	7	8	8	8	8	8
-1	0	1	2	3	4	5	5	6	7	7	7	7	7
-2	-1	0	1	2	3	4	4	5	6	6	6	6	7
-3	-2	-1	0	1	2	3	3	4	5	5	6	6	7
-4	-3	-2	-1	0	1	2	2	3	4	4	5	5	6
-5	-4	-3	-2	-1	0	1	2	3	4	4	5	5	6
-6	-5	-4	-3	-2	-1	0	1	2	3	3	4	4	5
-7	-6	-5	-4	-3	-2	-1	0	1	2	2	3	3	4
-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	3	4
-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4
-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3
-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2
-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1
-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0

b = "BAAB**CABCABA**CA"

$$H(i,j) = score(a, b\langle i : j \rangle)$$

$$H(4, 11) = 5$$

$$H(i,j) = j - i \text{ if } i > j$$

blue: difference in H is 0

red: difference in H is 1

The score matrix H and the seaweed matrix P

0	1	2	3	4	5	6	6	7	8	8	8	8	8
-1	0	1	2	3	4	5	5	6	7	7	7	7	7
-2	-1	0	1	2	3	4	4	5	6	6	6	6	7
-3	-2	-1	0	1	2	3	3	4	5	5	6	6	7
-4	-3	-2	-1	0	1	2	2	3	4	4	5	5	6
-5	-4	-3	-2	-1	0	1	2	3	4	4	5	5	6
-6	-5	-4	-3	-2	-1	0	1	2	3	3	4	4	5
-7	-6	-5	-4	-3	-2	-1	0	1	2	2	3	3	4
-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	3	4
-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4
-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3
-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2
-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1
-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	0

$$a =$$
 "BAABCBCA"

$$b =$$
 "BAAB**CABCABA**CA"

$$H(i,j) = score(a, b\langle i : j \rangle)$$

$$H(4, 11) = 5$$

$$H(i,j) = j - i \text{ if } i > j$$

blue: difference in H is 0

red: difference in H is 1

green: P(i,j) = 1

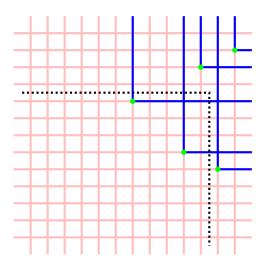
Image: Image:

$$H(i,j) = j - i - P^{\Sigma}(i,j)$$

- ∢ ∃ →

Score matrices and seaweed matrices

The score matrix H and the seaweed matrix P



a = "BAABCBCA"

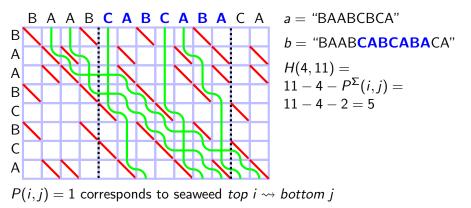
$$H(4, 11) =$$

 $11 - 4 - P^{\Sigma}(i, j) =$
 $11 - 4 - 2 = 5$

Image: Image:

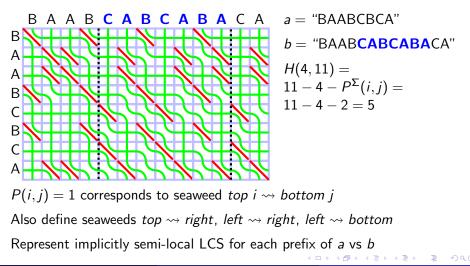
э

The (combed) seaweed braid in the alignment graph



Semi-local string comparison Score matrices and seaweed matrices

The (combed) seaweed braid in the alignment graph



Score matrices and seaweed matrices



Seaweed braid: a highly symmetric object (element of $H_0(S_n)$) Can be built by assembling subbraids: divide-and-conquer Flexible approach to local alignment, compressed approximate matching, parallel computation... The LCS problem is a special case of the weighted alignment score problem with weighted matches (w_M) , mismatches (w_X) and gaps (w_G)

- LCS score: $w_M = 1$, $w_X = w_G = 0$
- Levenshtein score: $w_M = 2$, $w_X = 1$, $w_G = 0$

The LCS problem is a special case of the weighted alignment score problem with weighted matches (w_M) , mismatches (w_X) and gaps (w_G)

• LCS score: $w_M = 1$, $w_X = w_G = 0$

Alexander Tiskin (Warwick)

• Levenshtein score: $w_M = 2$, $w_X = 1$, $w_G = 0$

Alignment score is rational, if w_M , w_X , w_G are rational numbers

The LCS problem is a special case of the weighted alignment score problem with weighted matches (w_M) , mismatches (w_X) and gaps (w_G)

- LCS score: $w_M = 1$, $w_X = w_G = 0$
- Levenshtein score: $w_M = 2$, $w_X = 1$, $w_G = 0$

Alignment score is rational, if w_M , w_X , w_G are rational numbers

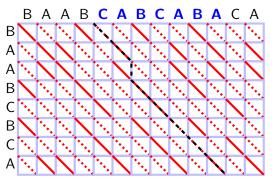
Edit distance: minimum cost to transform a into b by weighted character edits (insertion, deletion, substitution)

Corresponds to weighted alignment score with $w_M = 0$:

- insertion/deletion cost $-w_G$
- substitution weight $-w_X$

Semi-local string comparison Weighted alignment

Weighted alignment graph for a, b

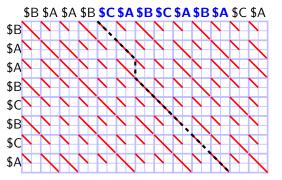


blue = 0red (solid) = 2 red (dotted) = 1

Levenshtein("BAABCBCA", "CABCABA") = 11

Semi-local string comparison Weighted alignment

Reduction: ordinary alignment graph for blown-up a, b

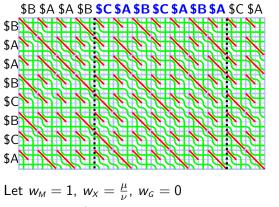


blue = 0red = 1 or 2

Levenshtein("BAABCBCA", "**CABCABA**") = *lcs*("\$B\$A\$A\$B\$C\$B\$C\$A", "**\$C\$A\$B\$C\$A\$B\$A**") = 11

Weighted alignment

Reduction: rational-weighted semi-local alignment to semi-local LCS



Increase $\times \nu^2$ in complexity (can be reduced to ν)



- 4 Compressed string comparison

Notation: pattern p of length m; text t of length n

A GC-string (grammar-compressed string) t is a straight-line program (context-free grammar) generating $t = t_{\bar{n}}$ by \bar{n} assignments of the form

• $t_k = \alpha$, where α is an alphabet character

• $t_k = uv$, where each of u, v is an alphabet character, or t_i for i < k

In general, $n = O(2^{\bar{n}})$

Example: Fibonacci string "ABAABABAABAAB"

 $t_1 = A$ $t_2 = t_1 B$ $t_3 = t_2 t_1$ $t_4 = t_3 t_2$ $t_5 = t_4 t_3$ $t_6 = t_5 t_4$

Grammar-compression covers various compression types, e.g. LZ78, LZW (not LZ77 directly)

Simplifying assumption: arithmetic up to n runs in O(1)

This assumption can be removed by careful index remapping

Compressed string comparison Extended substring-string LCS on GC-strings

LCS: running time ($r=m+n, \ ar{r}=ar{m}+ar{n})$										
p	t									
plain	plain	<i>O</i> (<i>mn</i>)		[Wagner, Fischer: 1974]						
		$O\left(\frac{mn}{\log^2 m}\right)$		[Masek, Paterson: 1980]						
		5		[Crochemore+: 2003]						
plain	GC	$O(m^3\bar{n}+\ldots)$	gen. CFG	[Myers: 1995]						
		$O(m^{1.5}\bar{n})$	ext subs-s	[T: 2008]						
		$O(m \log m \cdot \bar{n})$	ext subs-s	[T: 2010]						
GC	GC	NP-hard		[Lifshits: 2005]						
		$O(r^{1.2}\bar{r}^{1.4})$	R weights	[Hermelin+: 2009]						
		$O(r \log r \cdot \overline{r})$		[T: 2010]						
		$O(r \log(r/\overline{r}) \cdot \overline{r})$		[Hermelin+: 2010]						
		$O(r \log^{1/2}(r/\overline{r}) \cdot \overline{r})$		[Gawrychowski: 2012]						

Extended substring-string LCS (plain pattern, GC text): the algorithm

For every k, compute by recursion the appropriate part of seaweed matrix P_{p,t_k} , using matrix \boxdot -multiplication: time $O(m \log m \cdot \bar{n})$

Overall time $O(m \log m \cdot \bar{n})$

Compressed string comparison Subsequence recognition on GC-strings

The global subsequence recognition problem

Does pattern p appear in text t as a subsequence?

Global subsequence recognition: running time

р	t		
plain	plain	O(n)	greedy
plain	GC	$O(m\bar{n})$	greedy
GC	GC	NP-hard	[Lifshits: 2005]

The local subsequence recognition problem

Find all minimally matching substrings of t with respect to p

Substring of t is matching, if p is a subsequence of t

Matching substring of t is minimally matching, if none of its proper substrings are matching

Compressed string comparison

Subsequence recognition on GC-strings

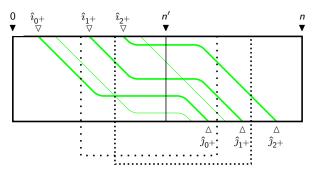
Local subsequence recognition: running time (+ output)

р	t		
plain	plain	<i>O</i> (<i>mn</i>)	[Mannila+: 1995]
		$O\left(\frac{mn}{\log m}\right)$	[Das+: 1997]
		$O(c^{m}+n)$	[Boasson+: 2001]
		$O(m + n\sigma)$	[Troniček: 2001]
plain	GC	$O(m^2 \log m\bar{n})$	[Cégielski+: 2006]
		$O(m^{1.5}\bar{n})$	[T: 2008]
		$O(m \log m \cdot \bar{n})$	[T: 2010]
		$O(m \cdot \bar{n})$	[Yamamoto+: 2011]
GC	GC	NP-hard	[Lifshits: 2005]

イロト イ押ト イヨト イヨト

Compressed string comparison

Subsequence recognition on GC-strings



 $b\langle i:j \rangle$ matching iff box [i:j] not pierced left-to-right

Determined by \ll -chain of \gtrless -maximal seaweeds

 $b\langle i:j\rangle$ minimally matching iff (i,j) is in the interleaved skeleton \ll -chain

Compressed string comparison

Subsequence recognition on GC-strings

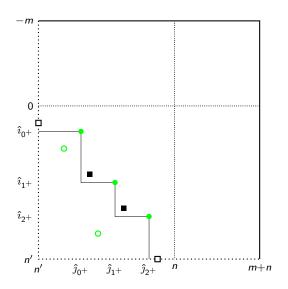


Image: Image:

66 / 75

.∋...>

Local subsequence recognition (plain pattern, GC text): the algorithm

For every k, compute by recursion the appropriate part of seaweed matrix P_{p,t_k} , using matrix \boxdot -multiplication: time $O(m \log m \cdot \bar{n})$

Local subsequence recognition (plain pattern, GC text): the algorithm

For every k, compute by recursion the appropriate part of seaweed matrix P_{p,t_k} , using matrix \boxdot -multiplication: time $O(m \log m \cdot \bar{n})$

Given an assignment t = t't'', find by recursion

- minimally matching substrings in t'
- minimally matching substrings in t''

Local subsequence recognition (plain pattern, GC text): the algorithm

For every k, compute by recursion the appropriate part of seaweed matrix P_{p,t_k} , using matrix \boxdot -multiplication: time $O(m \log m \cdot \bar{n})$

Given an assignment t = t't'', find by recursion

- minimally matching substrings in t'
- minimally matching substrings in t''

Then, find \ll -chain of \gtrless -maximal seaweeds in time $\bar{n} \cdot O(m) = O(m\bar{n})$ Its skeleton \ll -chain: minimally matching substrings in t overlapping t', t''Overall time $O(m \log m \cdot \bar{n}) + O(m\bar{n}) = O(m \log m \cdot \bar{n})$

- 本間 と 本語 と 本語 と

Threshold approximate matching

The threshold approximate matching problem

Find all matching substrings of t with respect to p, according to a threshold k

Substring of t is matching, if the edit distance for p vs t is at most k

Compressed string comparison

Threshold approximate matching

Threshold approximate matching: running time (+output)

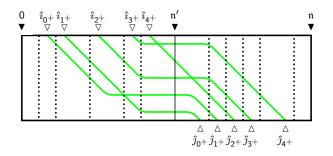
р	t		
plain	plain	O(mn)	[Sellers: 1980]
		O(mk)	[Landau, Vishkin: 1989]
		$O(m+n+\frac{nk^4}{m})$	[Cole, Hariharan: 2002]
plain	GC	$O(m\bar{n}k^2)$	[Kärkkäinen+: 2003]
		$O(m\bar{n}k+\bar{n}\log n)$	[LV: 1989] via [Bille+: 2010]
		$O(m\bar{n}+\bar{n}k^4+\bar{n}\log n)$	[CH: 2002] via [Bille+: 2010]
		$O(m \log m \cdot \bar{n})$	[T: NEW]
GC	GC	NP-hard	[Lifshits: 2005]

(Also many specialised variants for LZ compression)

- ∢ ศ⊒ ▶

Compressed string comparison

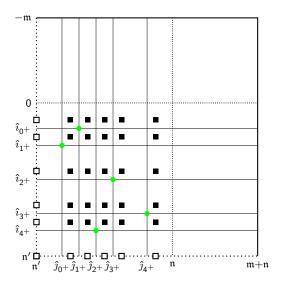
Threshold approximate matching



Blow up: weighted alignment on strings p, t of size m, n equivalent to LCS on strings \mathfrak{p} , \mathfrak{t} of size $\mathfrak{m} = \nu m$, $\mathfrak{n} = \nu n$

Compressed string comparison

Threshold approximate matching



æ

イロト イポト イヨト イヨト

Threshold approximate matching

Threshold approx matching (plain pattern, GC text): the algorithm

Algorithm structure similar to local subsequence recognition

«-chains replaced by $m \times m$ submatrices

Extra ingredients:

- the blow-up technique: reduction of edit distances to LCS scores
- implicit matrix searching, replaces ≪-chain interleaving

Monge row minima:"SMAWK" O(m)[Aggarwal+: 1987]Implicit unit-Monge row minima:

 O(m log log m)
 [T: 2012]

 O(m)
 [Gawrychowski: 2012]

Overall time $O(m \log m \cdot \bar{n}) + O(m\bar{n}) = O(m \log m \cdot \bar{n})$

1 Introduction

- 2 Matrix distance multiplication
- 3 Semi-local string comparison
- 4 Compressed string comparison
- 5 Conclusions and future work

Conclusions and future work

A powerful alternative to dynamic programming Implicit unit-Monge matrices:

• the seaweed monoid, multiplication in time $O(n \log n)$

Conclusions and future work

A powerful alternative to dynamic programming Implicit unit-Monge matrices:

• the seaweed monoid, multiplication in time $O(n \log n)$

Semi-local LCS problem:

• isomorphic to seaweed monoid, generalises to rational scores

A powerful alternative to dynamic programming Implicit unit-Monge matrices:

• the seaweed monoid, multiplication in time $O(n \log n)$

Semi-local LCS problem:

• isomorphic to seaweed monoid, generalises to rational scores Approximate matching in GC-text in time $O(m \log m \cdot \bar{n})$ Other applications:

- maximum clique in a circle graph in time $O(n \log^2 n)$
- parallel LCS in time $O(\frac{mn}{p})$, comm $O(\frac{m+n}{p^{1/2}})$ per processor
- identification of evolutionary-conserved regions in DNA

Thank you



2