

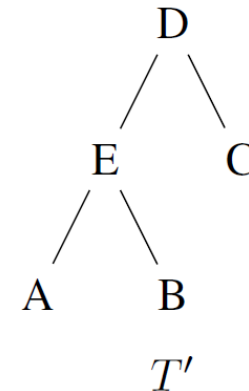
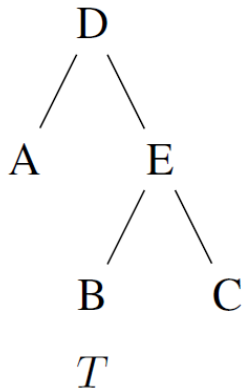
Counting tree alignments and algorithmic consequences

Cédric CHAUVÉ
Simon Fraser Univ.

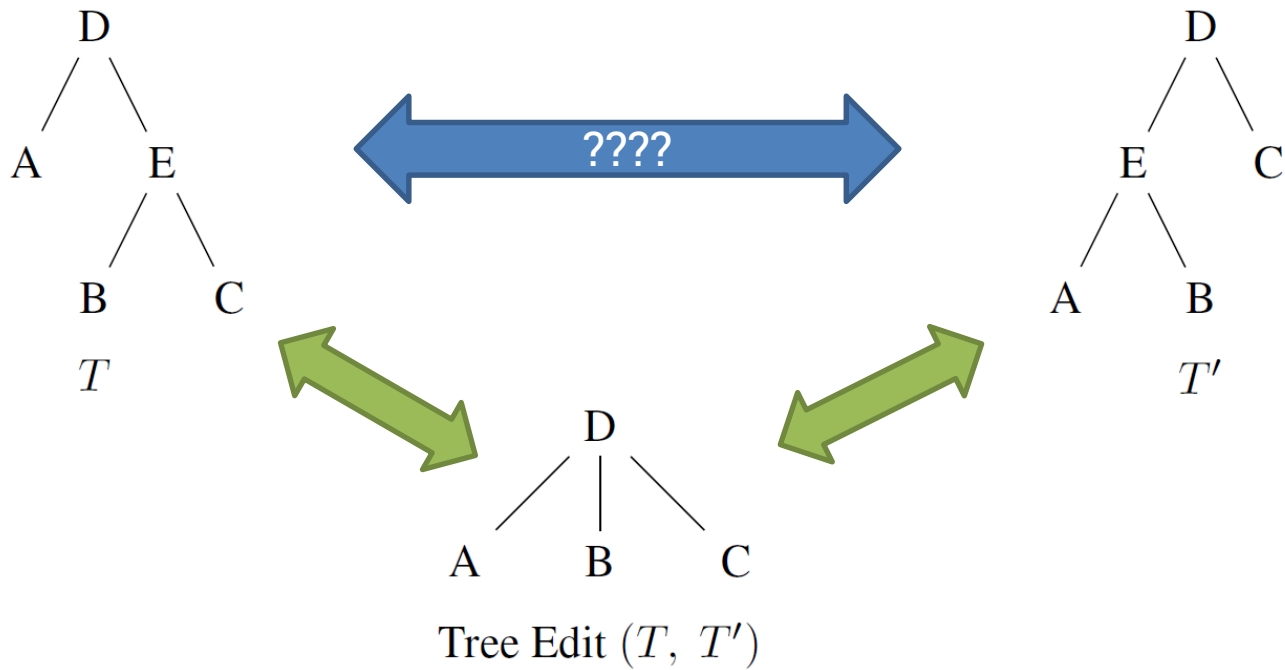
Julien COURTIÉL
Univ. de Caen

Yann PONTY
Ecole Polytechnique

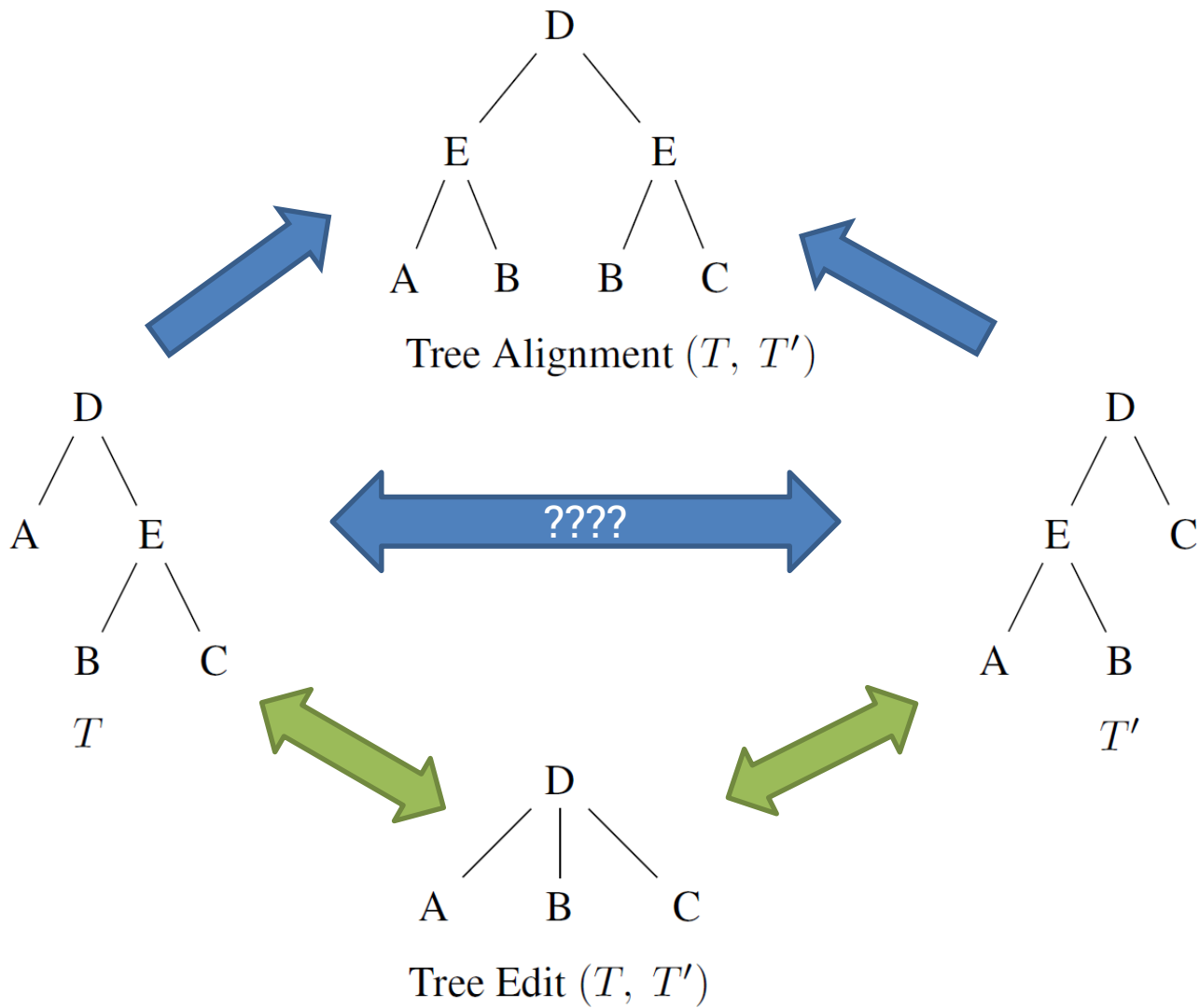
Alignment distance vs edit distance



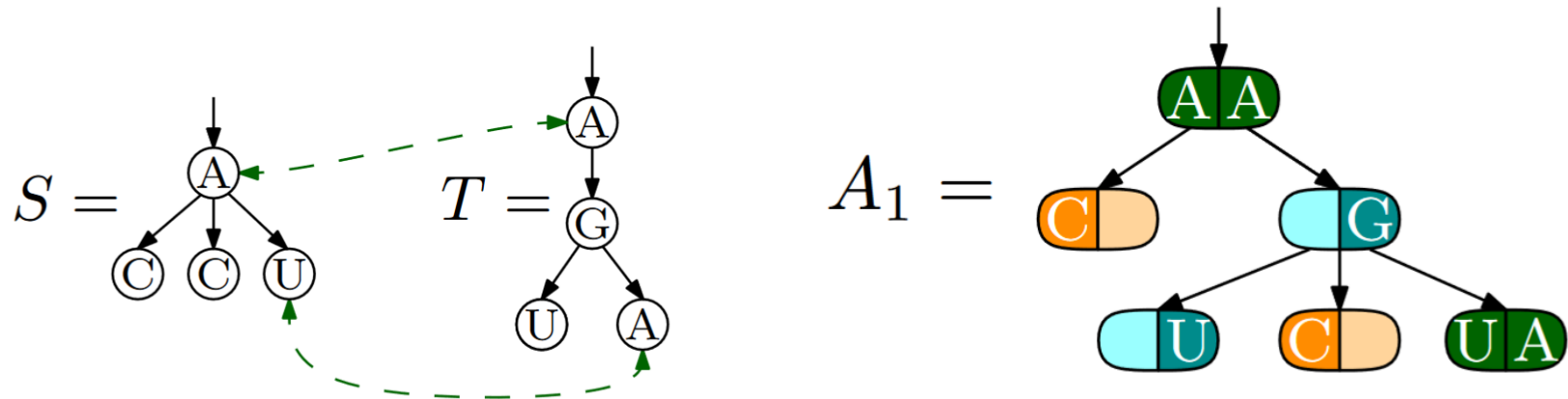
Alignment distance vs edit distance



Alignment distance vs edit distance



Tree alignments and supertrees



- ▶ **Alignment** = Set of correspondences, *aka* (mis)matches (\rightarrow Set of indels)
 + **Validity**: Consistency with ancestries in S and T
- ▶ Alt. **Alignment** = **Supertree** S such that S and T can be recovered from S

Tree alignment problem (Jiang-Wang-Zhang 1995)

Input Two trees S ($|S| = n_1$) and T ($|T| = n_2$)

Output Supertree maximizing weighted sum over (mis)matches

Jiang, Wang and Zhang (JWZ) DP algorithm

$$\begin{array}{l}
 \text{Align} \left(\begin{array}{c} \text{Tree} \\ \text{vs Tree} \end{array} \right) = \min \left\{ \begin{array}{l}
 \text{Align} \left(\begin{array}{c} \text{Tree} \\ \text{vs Tree} \end{array} \right) + \text{Del}(\bullet) \\
 \text{Align} \left(\begin{array}{c} \text{Tree} \\ \text{vs Tree} \end{array} \right) + \text{Ins}(\bullet) \\
 \text{Align} \left(\begin{array}{c} \text{Tree} \\ \text{vs Tree} \end{array} \right) + \text{Subst}(\bullet, \bullet)
 \end{array} \right. \\
 \\
 \text{Align} \left(\begin{array}{c} \text{Forest} \\ \text{vs Forest} \end{array} \right) = \min \left\{ \begin{array}{l}
 \min_{\text{Forest} = \text{Tree} + \text{Forest}} \left\{ \text{Align}(\text{Tree}) + \text{Align}(\text{Forest}) + \text{Del}(\bullet) \right\} \\
 \min_{\text{Forest} = \text{Forest} + \text{Tree}} \left\{ \text{Align}(\text{Forest}) + \text{Align}(\text{Tree}) + \text{Ins}(\bullet) \right\} \\
 \text{Align}(\text{Tree}) + \text{Align}(\text{Tree})
 \end{array} \right.
 \end{array}$$

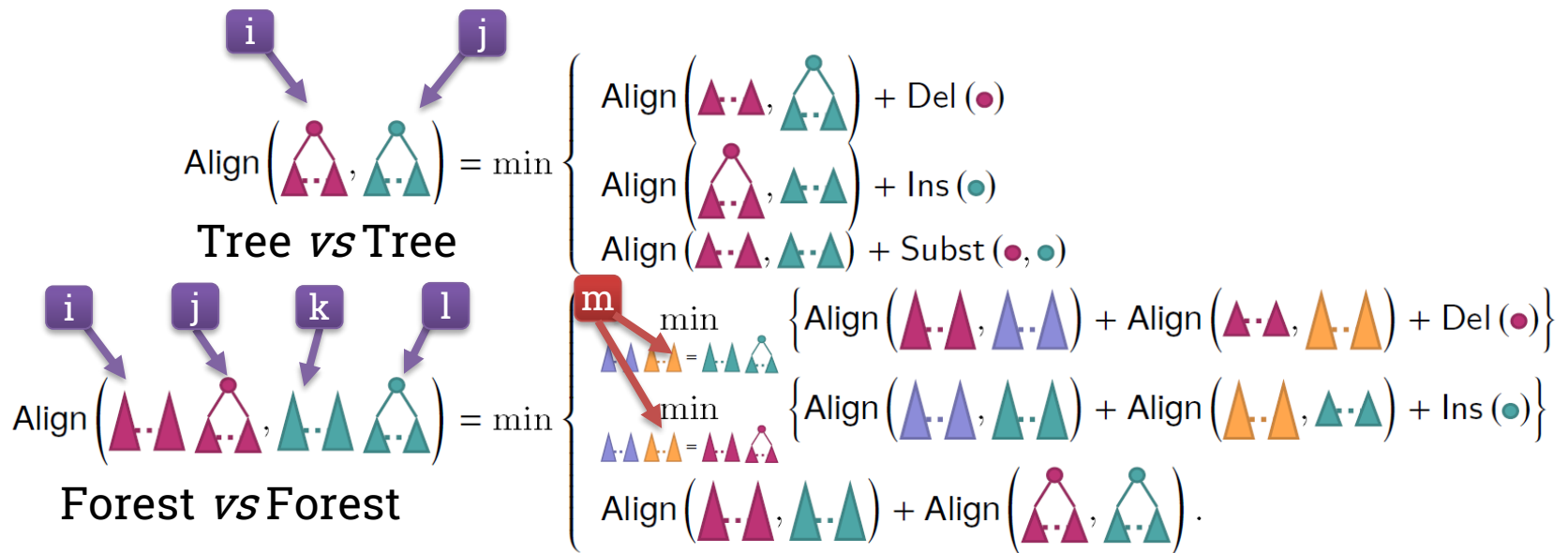
Complexities

▶ **Worst-case** $\rightarrow O(n_1^2 \cdot n_2^2)$ **space**, $O(n_1^2 \cdot n_2^2 \cdot \max(n_1, n_2))$ **time**

But at least one of (i, j, k, l) is first/last of its siblings

$\rightarrow \theta(n_1 \cdot n_2 \cdot \max(n_1, n_2))$ **space**, $\theta(n_1 \cdot n_2 \cdot \max(n_1, n_2)^2)$ **time**

Jiang, Wang and Zhang (JWZ) DP algorithm



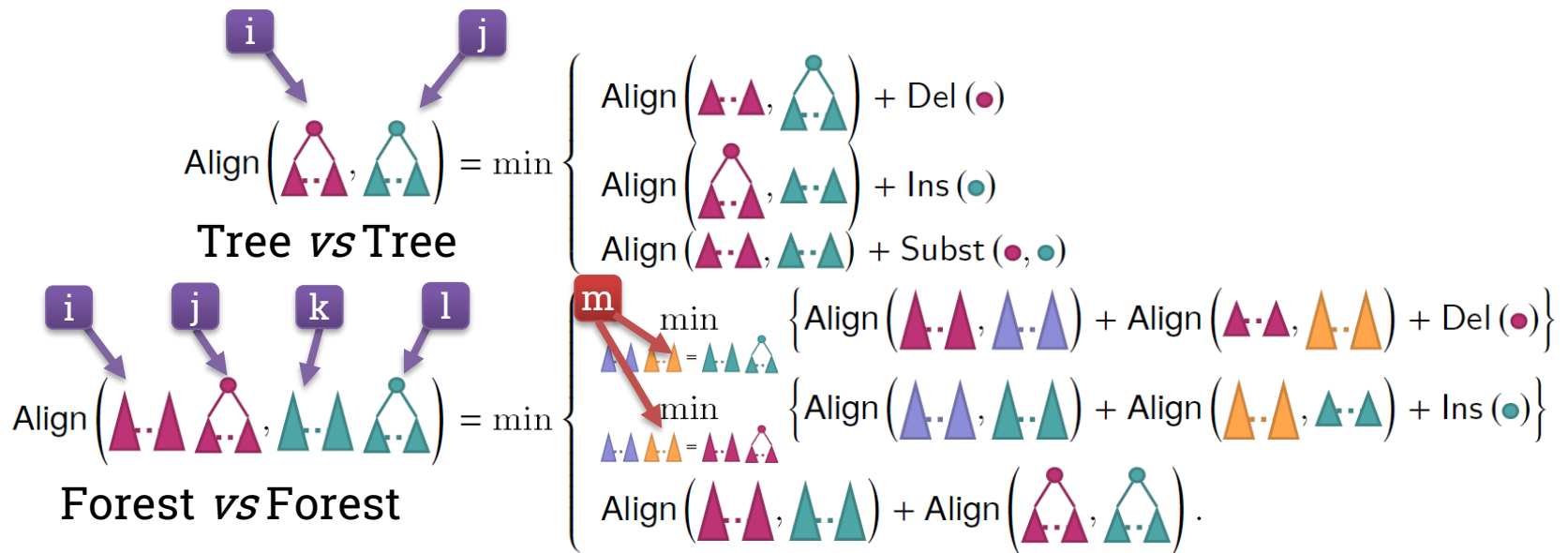
Complexities

► **Worst-case** $\rightarrow O(n_1^2 \cdot n_2^2)$ **space**, $O(n_1^2 \cdot n_2^2 \cdot \max(n_1, n_2))$ **time**

But at least one of (i, j, k, l) is first/last of its siblings

$\rightarrow \theta(n_1 \cdot n_2 \cdot \max(n_1, n_2))$ **space**, $\theta(n_1 \cdot n_2 \cdot \max(n_1, n_2)^2)$ **time**

Jiang, Wang and Zhang (JWZ) DP algorithm



Complexities

▶ **Worst-case** $\rightarrow O(n_1^2 \cdot n_2^2)$ space, $O(n_1^2 \cdot n_2^2 \cdot \max(n_1, n_2))$ time

But at least one of (i, j, k, l) is first/last of its siblings

$\rightarrow \theta(n_1 \cdot n_2 \cdot \max(n_1, n_2))$ space, $\theta(n_1 \cdot n_2 \cdot \max(n_1, n_2)^2)$ time

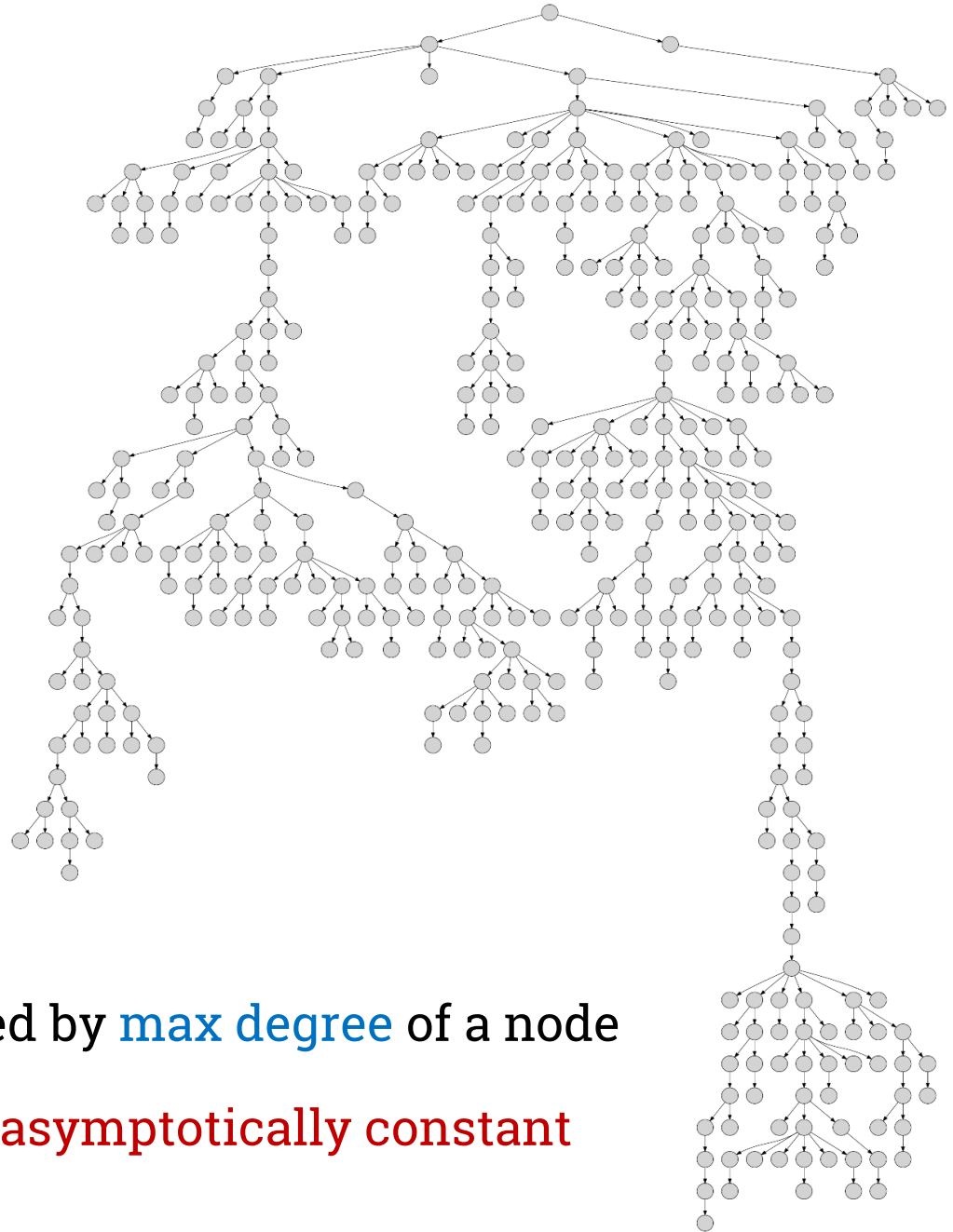
▶ **Average-case:** Random, uniformly distributed, trees of length n_1 and n_2

$\rightarrow \theta(n_1 \cdot n_2)$ space, $\theta(n_1 \cdot n_2)$ time

[Herrbach, Dulucq, Denise, TCS 2010]

Remark: Holds for Boltzmann-distributed RNA 2D struct. (homopolymer model)

Intuition



Wow



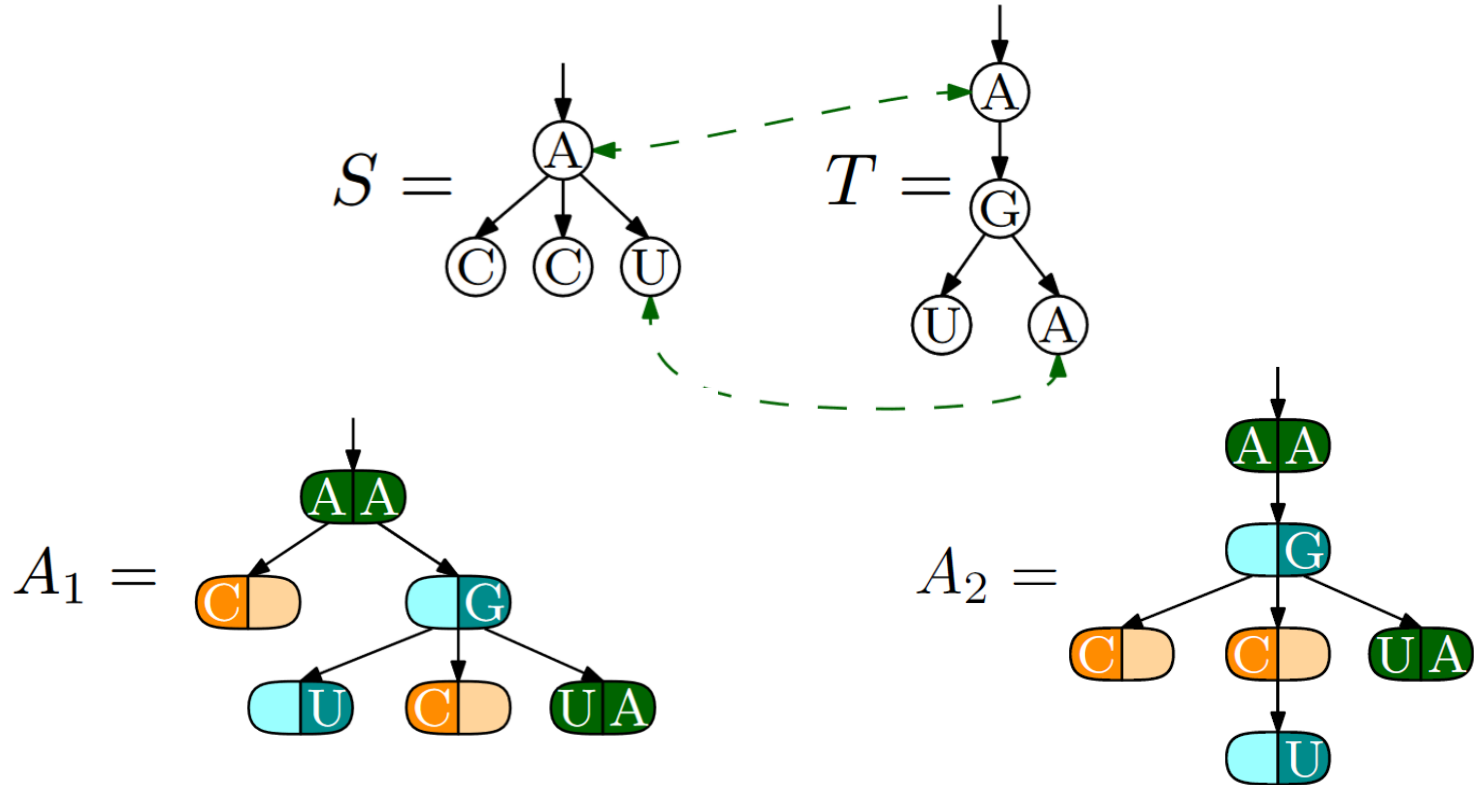
Time complexity dictated by **max degree** of a node

Max. degree on average **asymptotically constant**

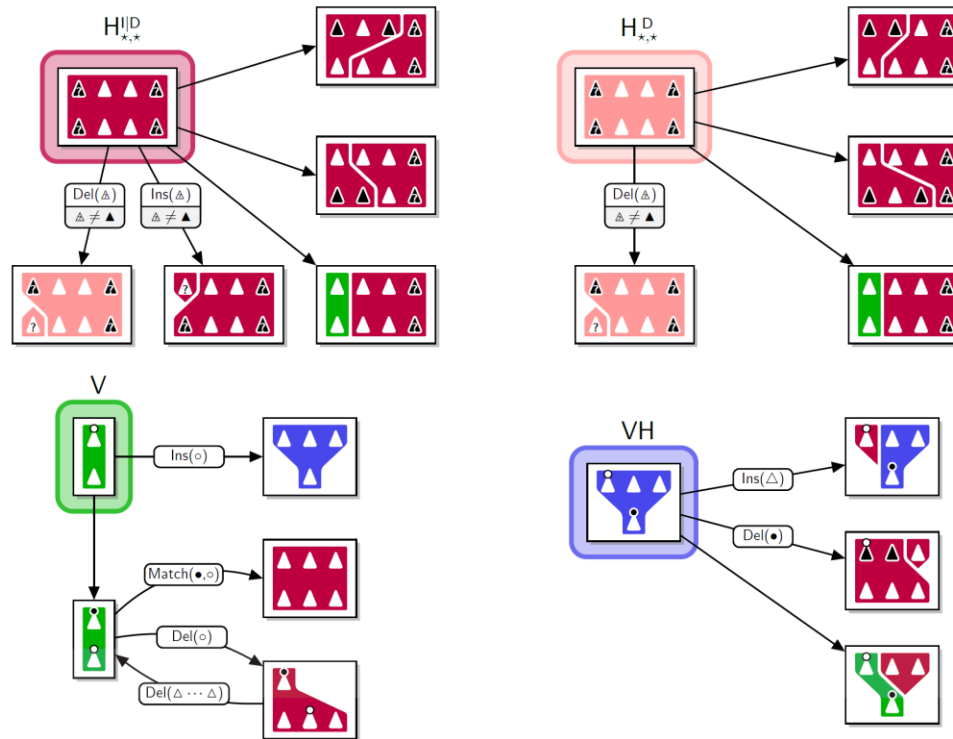
Counting alignments

Motivation: Ensemble analyses (*e.g.* MEA alignment, Bolz. Prob...)

Problem: DP scheme of Jiang-Wang-Zhang is **ambiguous**



Unambiguous decomposition/DP scheme



Theorem: DP scheme **complete**, **unambiguous** + complexities of JWZ

Proof: by intimidation, mainly!

(induction, seriously...)

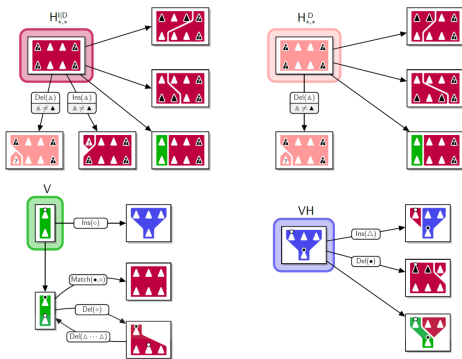
Alternatively: Über-simple unambiguous decomposition (yet mildly incomplete)

[Berkemer, Höhner zu Siederdisen, Stadler, Algorithms 2017]

Counting tree alignments

[Chauve Courtiel P, IJFCS 2018]

Scary DP scheme



$$\Delta \approx \nu^{\emptyset} \oplus \text{orange triangle} \oplus \text{cyan triangle} \oplus \text{rooted tree with orange and cyan triangles}$$

$$\nu^{\emptyset} \approx \nu^{\dagger} \oplus \text{rooted tree with orange triangles and } \nu\mathcal{H}$$

$$\nu^{\dagger} \approx \text{rooted tree with orange triangles and } \mathcal{H}_{ID, \emptyset, \emptyset} \oplus \text{rooted tree with cyan triangles and } \nu^{\dagger}$$

$$\nu\mathcal{H} \approx \text{orange triangle } \nu\mathcal{H} \oplus \nu^{\emptyset} \text{ orange triangles} \oplus \text{rooted tree with cyan triangles and } \mathcal{H}_{D, LR, \emptyset}$$

$$\mathcal{H}_{\nu, M, M'} \approx \varepsilon \oplus \text{orange triangle } \mathcal{H}_{\nu, M, M'} \oplus \text{cyan triangle } \mathcal{H}_{D, M, M'} \oplus \nu^{\emptyset} \mathcal{H} \oplus \text{rooted tree with orange triangles and } \mathcal{H}_{ID, \emptyset, LR} \oplus \text{rooted tree with cyan triangles and } \mathcal{H}_{D, LR, \emptyset}$$

if $(M, M') = (\emptyset, \emptyset)$ if $\nu \neq D$ and $M \neq LR$ if $M' \neq LR$

Scarier combinatorial specification/grammar

System of (algebraic) functional equations → Generating functions

→ Singularity analysis → Asymptotic properties of tree alignments

Asymptotic properties of tree alignments

- ▶ **#Tree alignments**, over a total of n nodes, equivalent to

$$\frac{\sqrt{2}(3 - \sqrt{3})}{24\sqrt{\pi}} \frac{6^n}{n\sqrt{n}}$$

→ $\approx 1.5^n$ tree alignments per pair of tree

- ▶ **#(Mis)matches** in random tree alignment:

- ▶ **Expectation** $\sim n/6$
- ▶ **Variance** $\sim n/6$

- ▶ Avg #supertrees per tree alignment **exponential** on length **yet**, for all n , **unique supertree** for certain alignments
 - Exponential bias induced by JWZ decomposition

Trees easier to align than sequences?!

Input: Pair of random uniform **sequences** of length n_1 and n_2

▶ **Theorem:** Needleman-Wunsch runs in $\theta(n_1 \cdot n_2)$ **expected time**

Input: Pair of random uniform **trees** of length n_1 and n_2

▶ **Reminder:** JWZ algorithm (+ ours) run in $\theta(n_1 \cdot n_2)$ **expected time**

Trees easier to align than sequences?!

Input: Pair of random uniform **sequences** of length n_1 and n_2

▶ **Theorem:** Needleman-Wunsch runs in $\theta(n_1 \cdot n_2)$ **expected time**

Input: Pair of random uniform **trees** of length n_1 and n_2

▶ **Reminder:** JWZ algorithm (+ ours) run in $\theta(n_1 \cdot n_2)$ **expected time**

Input: Random uniform **pair of sequences** of **cumulated length** n

▶ **Theorem:** Needleman-Wunsch runs in $\theta(n^2)$ **expected time**

Input: Random uniform **pair of trees** of **cumulated length** n

Trees easier to align than sequences?!

Input: Pair of random uniform **sequences** of length n_1 and n_2

▶ **Theorem:** Needleman-Wunsch runs in $\theta(n_1 \cdot n_2)$ **expected time**

Input: Pair of random uniform **trees** of length n_1 and n_2

▶ **Reminder:** JWZ algorithm (+ ours) run in $\theta(n_1 \cdot n_2)$ **expected time**

Input: Random uniform **pair of sequences** of **cumulated length** n

▶ **Theorem:** Needleman-Wunsch runs in $\theta(n^2)$ **expected time**

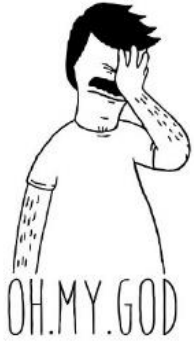
Input: Random uniform **pair of trees** of **cumulated length** n

▶ **Theorem:** JWZ algorithm (+ ours) run in $\theta(n \sqrt{n})$ **expected time**



Aligning **trees** easier than aligning **sequences**?!

Of course, this is cheating...



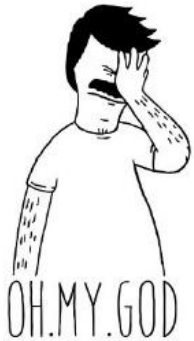
n = cumulated length of a pair

▶ For sequences:

▶ Both have expected length $\frac{n}{2}$

▶ For trees:

Of course, this is cheating...



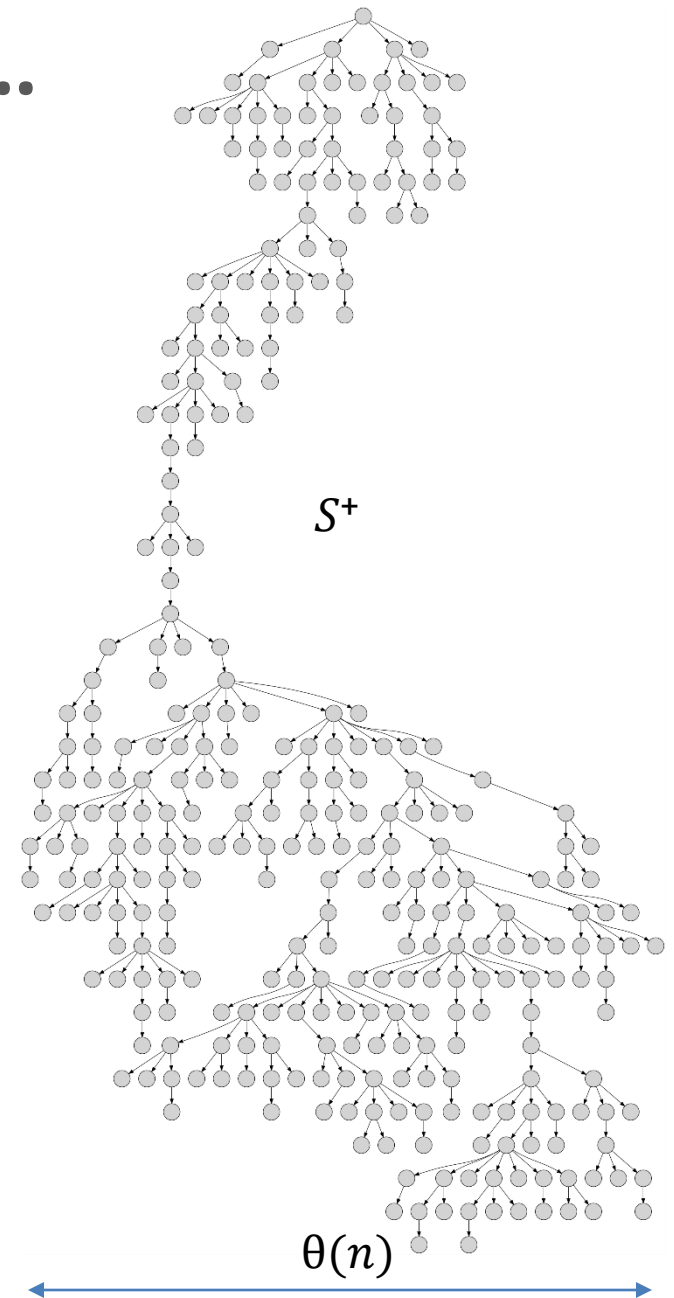
n = cumulated length of a pair

▶ For sequences:

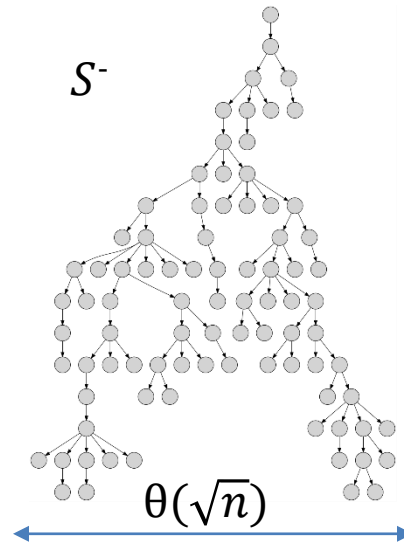
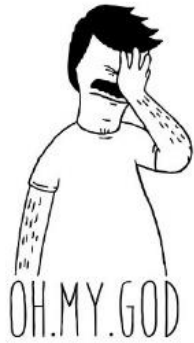
▶ Both have expected length $\frac{n}{2}$

▶ For trees:

▶ Largest tree S^+ has expected length in $\theta(n)$



Of course, this is cheating...



n = cumulated length of a pair

► For sequences:

► Both have expected length $\frac{n}{2}$

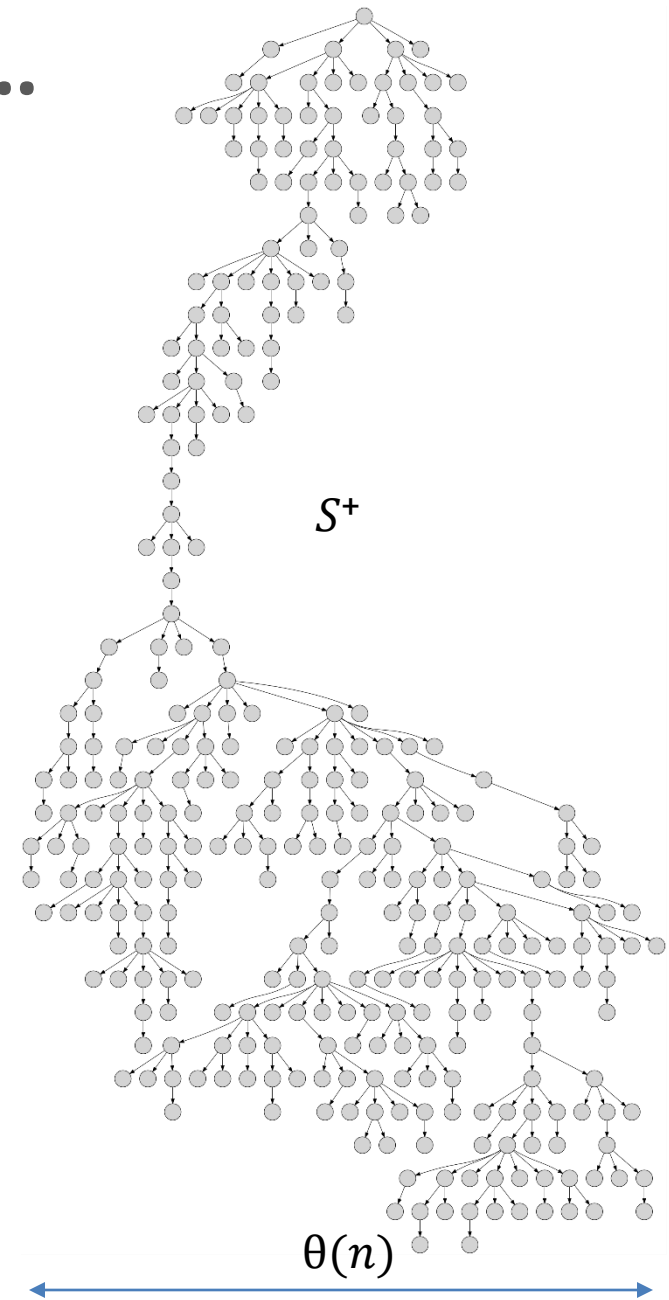
► For trees:

► Largest tree S^+ has expected length in $\theta(n)$

► Smallest tree S^- has expected length in $\theta(\sqrt{n})$

The subquadratic complexity of JWZ is

only an artifact of the length distribution!



Conclusions/thanks

- ▶ Don't be fooled by **combinatorics** people! (including this one...)
- ▶ Tree alignment amenable to **ensemble analyses**
- ▶ Who would like to **co-implement** our ~~monster~~ DP beauty?

Thanks to Bled Organizers + YOU



Julien Courtiel



Cédric Chauve

