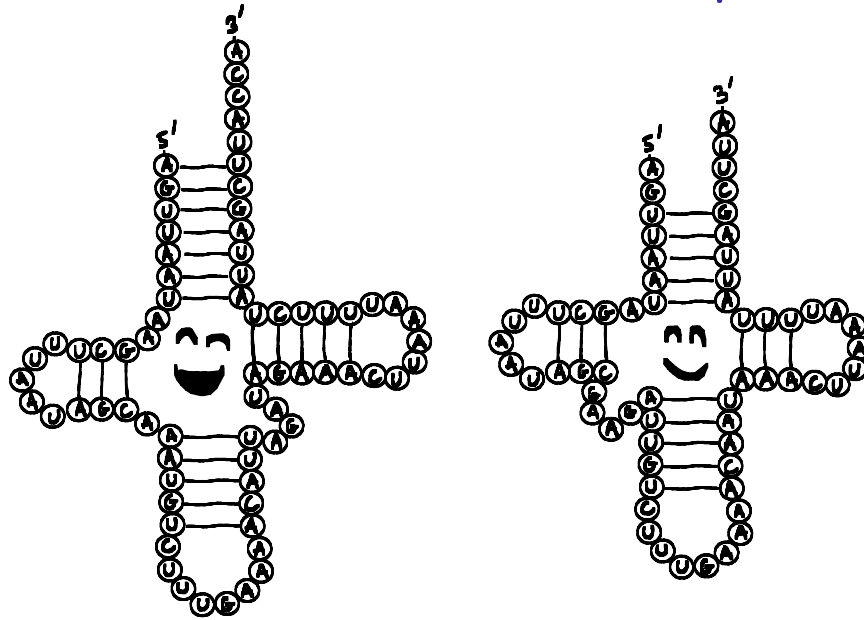# COUNTING, GENERATING AND SAMPLING TREE ALIGNMENTS
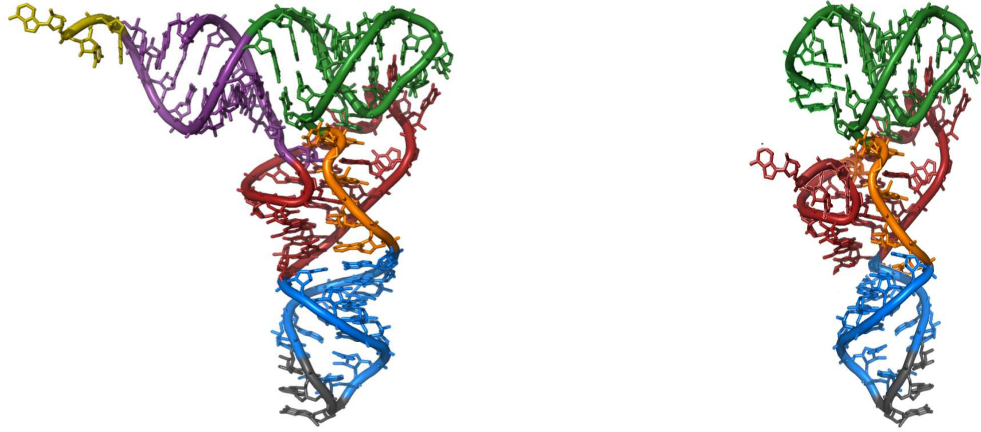
Cedric CHAUVE (Simon Fraser University, Vancouver)
Julien COURTIEL (PIMS/Univ. of British Columbia, Vancouver)
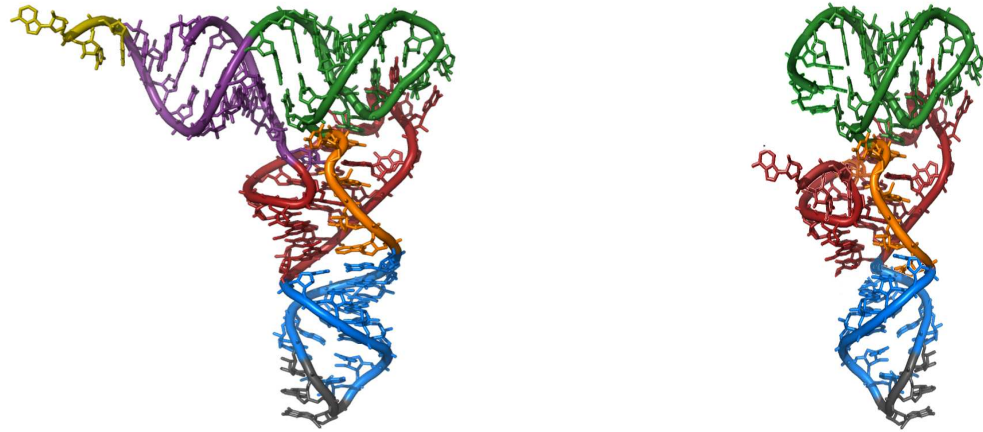Yann PONTY (CNRS/LIX, Ecole Polytechnique, Inria Saclay)

AICoB 2016

# MOTIVATION: RNA COMPARISON

Question: how to measure similarity between two RNAs?

# MOTIVATION: RNA COMPARISON

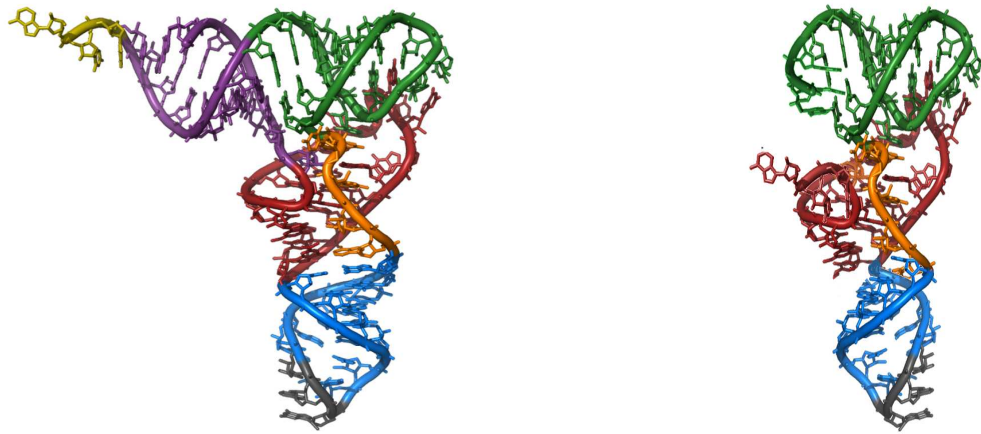Question: how to measure similarity between two RNAs?



First idea: compare nucleic acid sequences

RNA 1:   AUUCGAUUA...
RNA 2:   ACCAUGAUUA...

# MOTIVATION: RNA COMPARISON

Question: how to measure similarity between two RNAs?



First idea: compare nucleic acid sequences
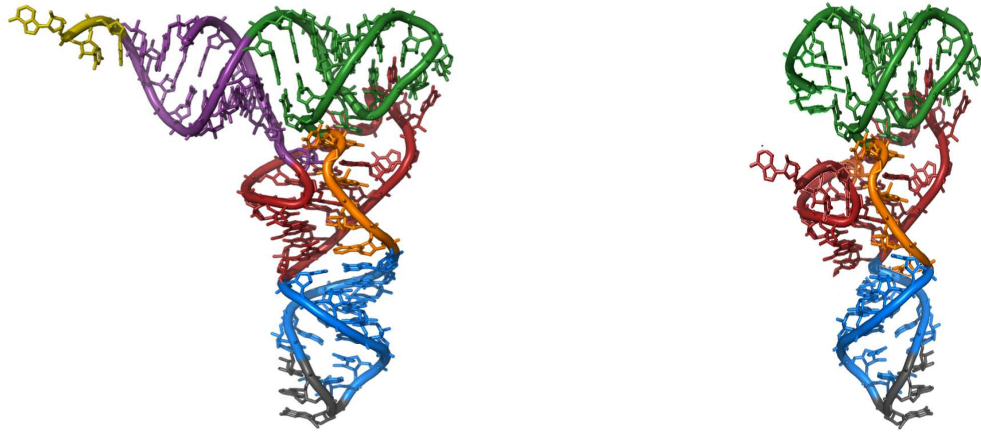→ sequence alignment

RNA 1:  AUUCGAUUA...

RNA 2:  ACCAUGAUUA...

alignment: $\begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} - \\ C \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} C \\ C \end{pmatrix} \begin{pmatrix} - \\ A \end{pmatrix} \begin{pmatrix} - \\ U \end{pmatrix} \begin{pmatrix} G \\ G \end{pmatr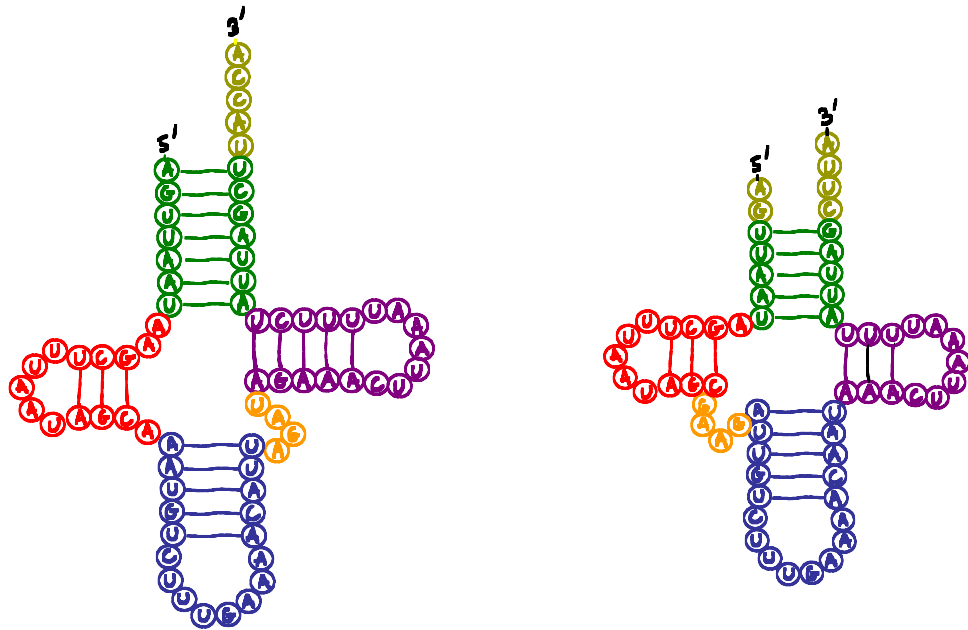ix} \begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ U \end{pmatrix} \begin{pmatrix} U \\ U \end{pmatrix} \begin{pmatrix} A \\ A \end{pmatrix}$...

# MOTIVATION: RNA COMPARISON

**Question**: how to measure similarity between two RNAs?



Second idea: compare secondary structures.



→ notion of tree alignment [Jiang, Wang, Zhang]

# OUR CONTRIBUTION

Our main result:

> An unambiguous and complete Dynamic Programming scheme for tree alignments
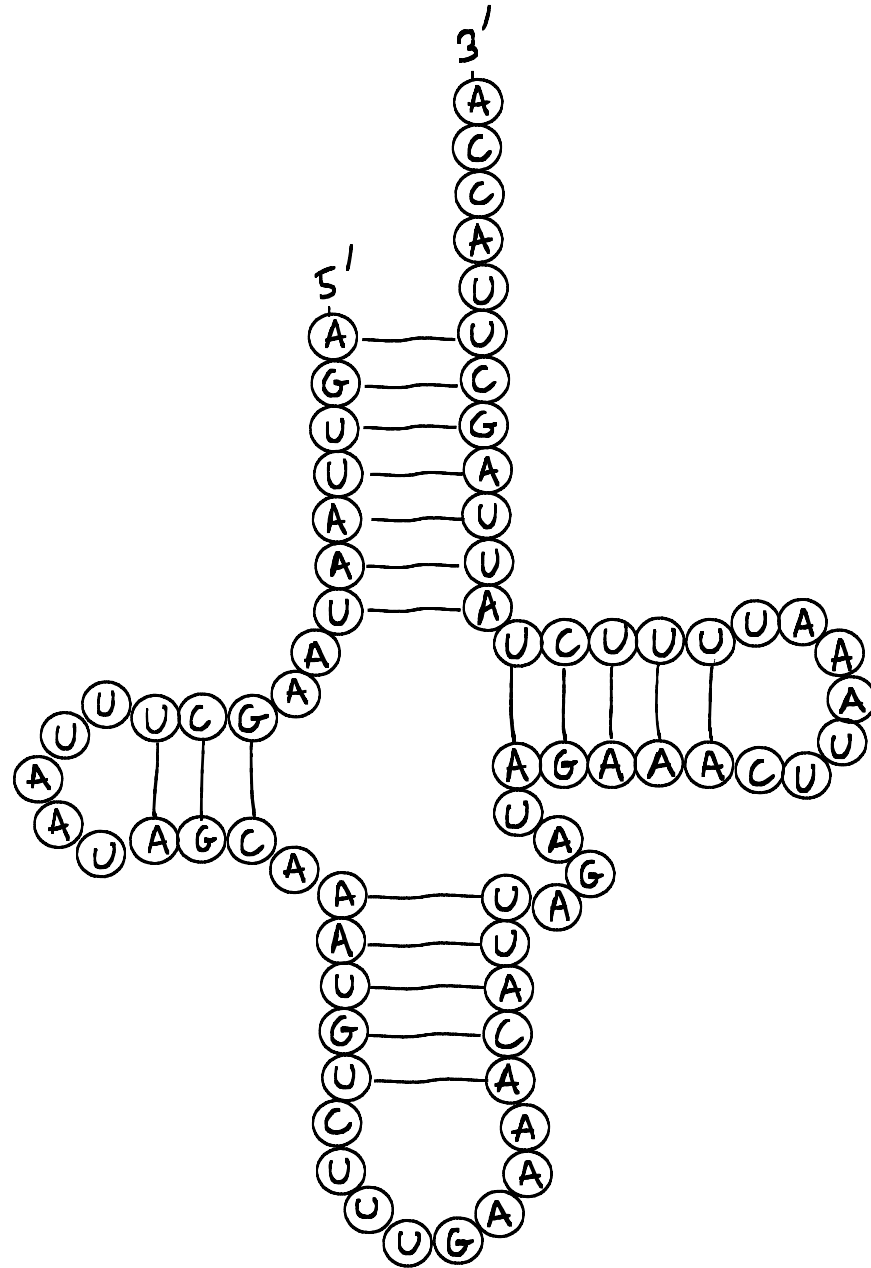
Side-products:

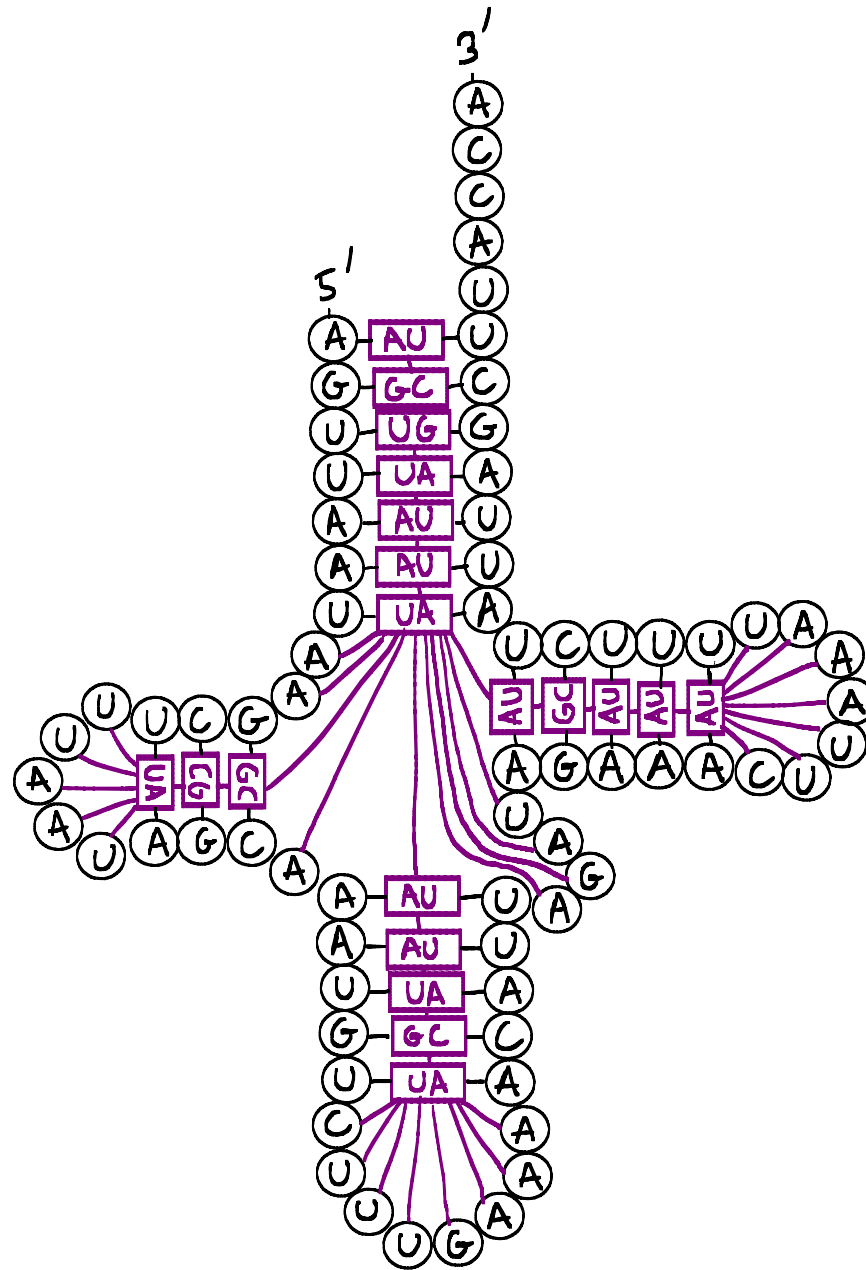- Boltzmann sampling
- Enumerative results
- Average case complexity (revisited)
- 3D alignment and motif search (upcoming)

FROM SECONDARY STRUCTURES TO TREES

# FROM SECONDARY STRUCTURES TO TREES

# FROM SECONDARY STRUCTURES TO TREES

Objective: Align trees coming from RNA $2^{ary}$ structures

# FROM SECONDARY STRUCTURES TO TREES

Objective: Align trees coming from RNA 2<sup>ary</sup> structures

# SEQUENCE ALIGNMENT

super sequence = word on $\Sigma \times \Sigma \oplus \Sigma \times \{-\} \oplus \{-\} \times \Sigma$

$$\binom{A}{A}\binom{U}{-}\binom{-}{C}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{A}{U}\binom{C}{A}$$

match     insertion     deletion     mismatch

# SEQUENCE ALIGNMENT

super sequence = word on $\Sigma \times \Sigma \oplus \Sigma \times \{-\} \oplus \{-\} \times \Sigma$

$$\begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} - \\ C \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} C \\ C \end{pmatrix} \begin{pmatrix} - \\ A \end{pmatrix} \begin{pmatrix} - \\ U \end{pmatrix} \begin{pmatrix} G \\ G \end{pmatrix} \begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ U \end{pmatrix} \begin{pmatrix} A \\ U \end{pmatrix} \begin{pmatrix} C \\ A \end{pmatrix}$$

$\longrightarrow$ AUUCGAUAC

$\longrightarrow$ ACCAUGAUUA

match    insertion    deletion    mismatch    projections

# SEQUENCE ALIGNMENT

super sequence = word on $\Sigma \times \Sigma \oplus \Sigma \times \{-\} \oplus \{-\} \times \Sigma$

$\binom{A}{A}\binom{U}{-}\binom{-}{C}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{A}{U}\binom{C}{A}$ $\longrightarrow$ AUUCGAUAC

$\longrightarrow$ ACCAUGAUUA

match    insertion    deletion    mismatch    projections
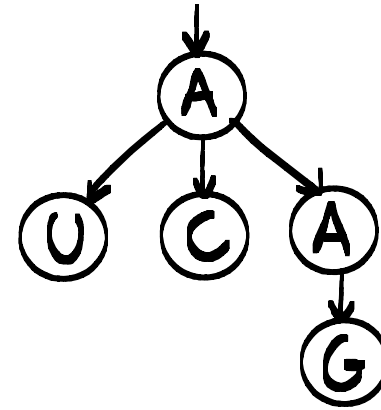
Given two sequences $S_1$ and $S_2$,

alignment between $S_1$ and $S_2$ = supersequence with projections $S_1$ and $S_2$

cost of an alignment = nb of insertions + deletions + mismatches

# SEQUENCE ALIGNMENT

super sequence = word on $\Sigma \times \Sigma \oplus \Sigma \times \{-\} \oplus \{-\} \times \Sigma$

$$\binom{A}{A}\binom{U}{-}\binom{-}{C}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{A}{U}\binom{C}{A} \longrightarrow AUUCGAUAC$$
$$\longrightarrow ACCAUGAUUA$$

match  insertion  deletion  mismatch  projections

Given two sequences $S_1$ and $S_2$,

alignment between $S_1$ and $S_2$ = supersequence with projections $S_1$ and $S_2$

cost of an alignment = nb of insertions + deletions + mismatches

Remark: sequence alignment = sequence edition

# TREES AND SUPERTREES

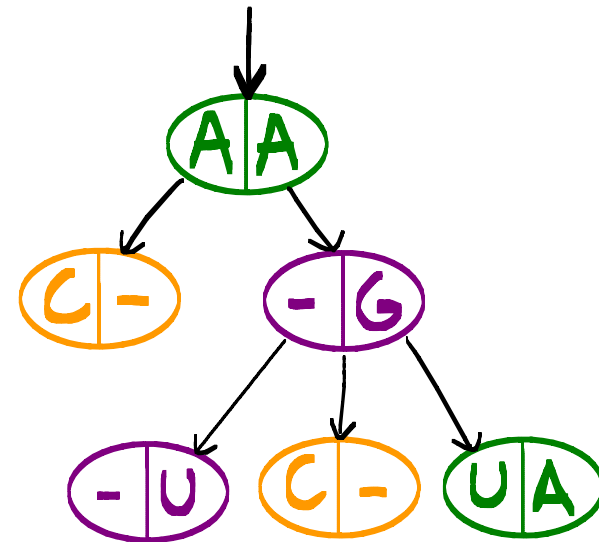Trees are plane, rooted, and vertices are labeled by an alphabet $\Sigma$.
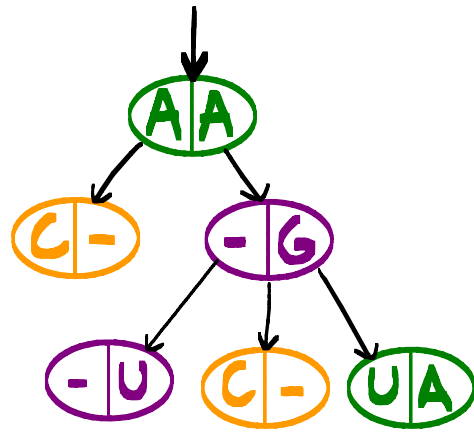
Supertree = tree with 3 types of vertices:
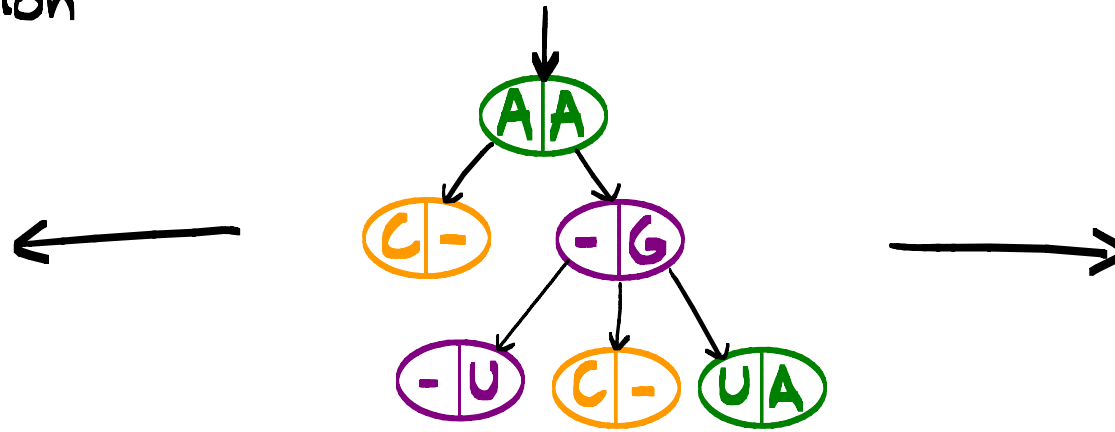
$X|Y$  (mis)match

$X|-$  insertion
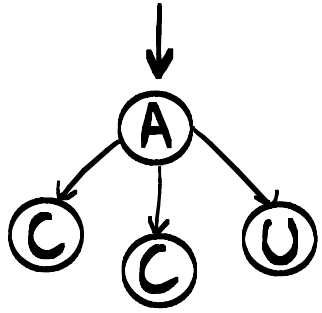
$-|Y$  deletion

# TREE ALIGNMENTS

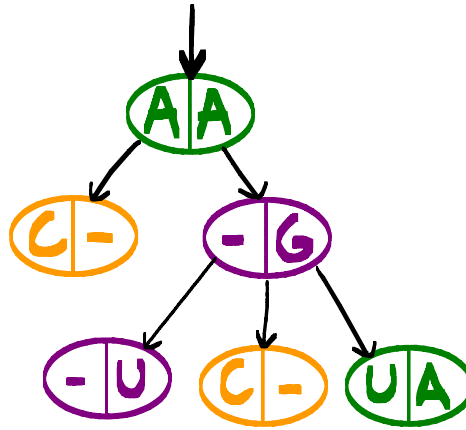# TREE ALIGNMENTS

first projection

second projection

# TREE ALIGNMENTS

**first projection**

**second projection**



keep left letters

# TREE ALIGNMENTS

first projection

second projection



keep left letters
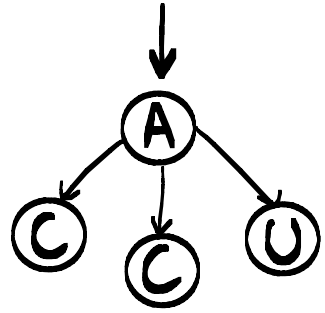
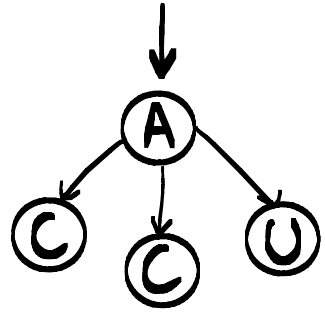# TREE ALIGNMENTS



first projection

keep left letters

second projection

keep right letters

# TREE ALIGNMENTS



first projection

second projection

keep left letters

keep right letters

# TREE ALIGNMENTS

first projection



S

second projection

T

Given two trees S and T,
alignment between S and T = supertree whose projections
are S and T.

# TREE ALIGNMENTS



first projection

second projection

S

T

Given two trees S and T,
alignment between S and T = supertree whose projections
are S and T.

cost of an alignment = nb of insertions + deletions + mismatches

Remark: tree alignment $\neq$ tree edition

# CONNECTION WITH SEQUENCE ALIGNMENTS

Tree alignments generalize sequence alignments.

**SEQUENCE**

AUUCGAUUA...          ACCAUGAUUA...

alignment:

$\binom{A}{A}\binom{U}{-}\binom{-}{C}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{U}{U}\binom{A}{A}...$

---

**TREE**



alignment:

# OPTIMAL ALIGNMENT

Classical problem: Given S and T, find one optimal alignment between S and T.

Solvable by Dynamic Programming:

Worst case time
$$O(n^4)$$
[Jiang, Wang, Zhang]

Average time
$$O(n^2)$$
[Herrbach, Denise, Dulucq]

Which alignment between  and 
is the most likely?

# SPACE OF ALIGNMENTS

Which alignment between is the most likely?

and



co-optimal alignments

# SPACE OF ALIGNMENTS

Why finding one optimal alignment may be inadequate:

▶ Co-optimal alignments can be very different. (see for instance [Vingron, Argos, 1990])

▶ Exploring the space of alignments enables the detection of high probability features.

# SPACE OF ALIGNMENTS

Objective: Sampling alignments under the Gibbs-Boltzmann probability distribution.

probability of an alignment A

$$\propto e^{-\frac{cost(A)}{K}}$$

(Gibbs-Boltzmann distribution)



co-optimal alignments

# SPACE OF ALIGNMENTS

Objective: Sampling alignments under the Gibbs-Boltzmann probability distribution.

probability of an alignment A
$$\propto e^{-\frac{cost(A)}{K}}$$
(Gibbs-Boltzmann distribution)

**Score vs Boltzmann probability (Density of states)**



$\underline{K=0}$ : Uniform distribution over optimal alignments.

$\underline{K=+\infty}$ : Uniform distribution over all alignments.

# SPACE OF ALIGNMENTS

Objective: Sampling alignments under the Gibbs-Boltzmann probability distribution.

**Score vs Boltzmann probability (Density of states)**

probability of ?? (an) alignment A

$$\propto e^{-\frac{cost(A)}{K}}$$

(Gibbs-Boltzmann distribution)

$\underline{K=0}$ : Uniform distribution over optimal alignments.

$\underline{K=+\infty}$ : Uniform distribution over all alignments.

# AMBIGUITY OF ALIGNMENTS

For sequences,

$$\binom{A}{A}\binom{U}{-}\binom{-}{C}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{U}{U}\binom{A}{A}$$

is the same alignment as

$$\binom{A}{A}\binom{-}{C}\binom{U}{-}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{U}{U}\binom{A}{A}$$

# AMBIGUITY OF ALIGNMENTS

For trees,



and

induce the same alignment between

# AMBIGUITY OF ALIGNMENTS

The two supertrees



do not induce the same alignment between the trees

# PROBLEM RAISED BY THE AMBIGUITY

For sequences, we can deal with the ambiguity by defining canonical alignments.

Ex:
$$\begin{pmatrix} A \\ A \end{pmatrix}\begin{pmatrix} U \\ - \end{pmatrix}\begin{pmatrix} - \\ C \end{pmatrix}\begin{pmatrix} U \\ - \end{pmatrix}\begin{pmatrix} C \\ C \end{pmatrix}\begin{pmatrix} - \\ A \end{pmatrix}\begin{pmatrix} - \\ U \end{pmatrix}\begin{pmatrix} G \\ G \end{pmatrix}\begin{pmatrix} A \\ A \end{pmatrix}\begin{pmatrix} U \\ U \end{pmatrix}\begin{pmatrix} U \\ U \end{pmatrix}\begin{pmatrix} A \\ A \end{pmatrix}$$

# PROBLEM RAISED BY THE AMBIGUITY

For sequences, we can deal with the ambiguity by defining canonical alignments.

Ex :
$$\binom{A}{A}\binom{U}{-}\binom{U}{-}\binom{-}{C}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{U}{U}\binom{A}{A}$$

Insertions before Deletions.

# PROBLEM RAISED BY THE AMBIGUITY

For sequences, we can deal with the ambiguity
by defining canonical alignments.

Ex:
$$\binom{A}{A}\binom{U}{-}\binom{U}{-}\binom{-}{C}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{U}{U}\binom{A}{A}$$

Insertions before Deletions.

For trees, it is much more complicated!

# PROBLEM RAISED BY THE AMBIGUITY

For sequences, we can deal with the ambiguity by defining canonical alignments.

Ex: $\begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} - \\ C \end{pmatrix} \begin{pmatrix} C \\ C \end{pmatrix} \begin{pmatrix} - \\ A \end{pmatrix} \begin{pmatrix} - \\ U \end{pmatrix} \begin{pmatrix} G \\ G \end{pmatrix} \begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ U \end{pmatrix} \begin{pmatrix} U \\ A \end{pmatrix}$

Insertions before Deletions.

For trees, it is much more complicated!

Strategy: Build a context-free grammar that generates every alignment exactly once

# GRAMMARS FOR SEQUENCE ALIGNMENTS

## Ambiguous grammar:

$$\mathcal{S} \leftarrow \begin{pmatrix} X \\ Y \end{pmatrix} \boxed{\mathcal{S}} \oplus \begin{pmatrix} X \\ - \end{pmatrix} \boxed{\mathcal{S}} \oplus \begin{pmatrix} - \\ Y \end{pmatrix} \boxed{\mathcal{S}}$$

## Non-ambiguous grammar:

$$\mathcal{S} \leftarrow \begin{pmatrix} X \\ Y \end{pmatrix} \boxed{\mathcal{S}} \oplus \begin{pmatrix} X \\ - \end{pmatrix} \boxed{\mathcal{S}} \oplus \begin{pmatrix} - \\ Y \end{pmatrix} \boxed{\mathcal{S}^D}$$

$$\mathcal{S}^D \leftarrow \begin{pmatrix} X \\ Y \end{pmatrix} \boxed{\mathcal{S}} \oplus \begin{pmatrix} - \\ Y \end{pmatrix} \boxed{\mathcal{S}^D}$$

# A GRAMMAR FOR ALIGNMENTS

For trees, an ambiguous grammar can be derived from [Jiang, Wang, Zhang].

Our result:

Theorem: The set 𝔸 generated by the following grammar contains every tree alignment exactly once.

Our (complicated) non-ambiguous grammar:

$$\mathcal{A} \leftarrow \mathcal{V}^\phi \mid \mathcal{C}_I \mid \mathcal{C}_D \mid \; [\text{x}|-] \to \boxed{\mathcal{F}_I}\,\box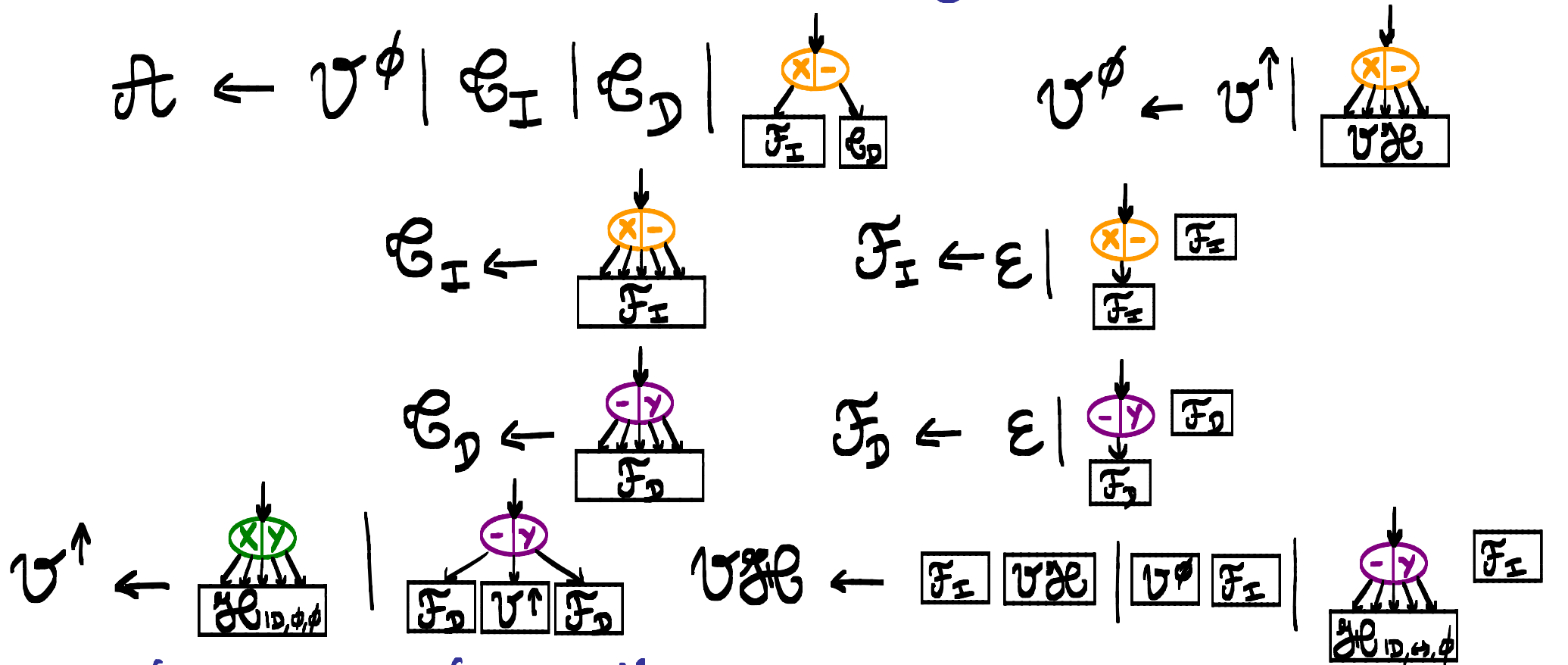ed{\mathcal{C}_D} \qquad\qquad \mathcal{V}^\phi \leftarrow \mathcal{V}^\uparrow \mid [\text{x}|-] \to \boxed{\mathcal{VH}}$$

$$\mathcal{C}_I \leftarrow [\text{x}|-] \to \boxed{\mathcal{F}_I} \qquad\qquad \mathcal{F}_I \leftarrow \varepsilon \mid [\text{x}|-] \to \boxed{\mathcal{F}_I}\;\boxed{\mathcal{F}_I}$$

$$\mathcal{C}_D \leftarrow [-|\text{y}] \to \boxed{\mathcal{F}_D} \qquad\qquad \mathcal{F}_D \leftarrow \varepsilon \mid [-|\text{y}] \to \boxed{\mathcal{F}_D}\;\boxed{\mathcal{F}_D}$$

$$\mathcal{V}^\uparrow \leftarrow [\text{x}|\text{y}] \to \boxed{\mathcal{H}_{1D,\phi,\phi}} \;\mid\; [-|\text{y}] \to \boxed{\mathcal{F}_D}\,\boxed{\mathcal{V}^\uparrow}\,\boxed{\mathcal{F}_D}$$

$$\mathcal{VH} \leftarrow \boxed{\mathcal{F}_I}\,\boxed{\mathcal{VH}} \;\mid\; \boxed{\mathcal{V}^\phi}\,\boxed{\mathcal{F}_I} \;\mid\; [-|\text{y}] \to \boxed{\mathcal{H}_{1D,\leftrightarrow,\phi}}\;\boxed{\mathcal{F}_I}$$

For $\mathcal{V} \in \{D, \mathcal{D}\}$, $(M,M') \in \{\phi, \rightarrow, \leftrightarrow\}^2$:

$$\mathcal{H}_{V,M,M'} \leftarrow \varepsilon \;\mid\; \boxed{\mathcal{C}_I}\,\boxed{\mathcal{H}_{V,M,M'}} \;\mid\; \boxed{\mathcal{C}_D}\,\boxed{\mathcal{H}_{D,M,M'}} \;\mid\; \boxed{\mathcal{V}^\phi}\,\boxed{\mathcal{H}_{M,M'}^{1'}} \;\mid\; [\text{x}|-] \to \boxed{\mathcal{H}_{1D,\phi,\leftrightarrow}}\;\boxed{\mathcal{H}_{M,M'}^{1',+}} \;\mid\; [-|\text{y}] \to \boxed{\mathcal{H}_{1D,\leftrightarrow,\phi}}\;\boxed{\mathcal{H}_{M,M'}^{+,1'}}$$

only if $(M,M')=(\phi,\phi)$

only if $\mathcal{V}\neq D$ and $M\neq\leftrightarrow$

only if $M'\neq\leftrightarrow$

no room for $\overline{\mathcal{H}}^{i,o}_{M,M'}$...

# APPLICATION 1. COUNTING.

**Theorem** There are on average

$$C \times 1.5^n \text{ alignments}$$

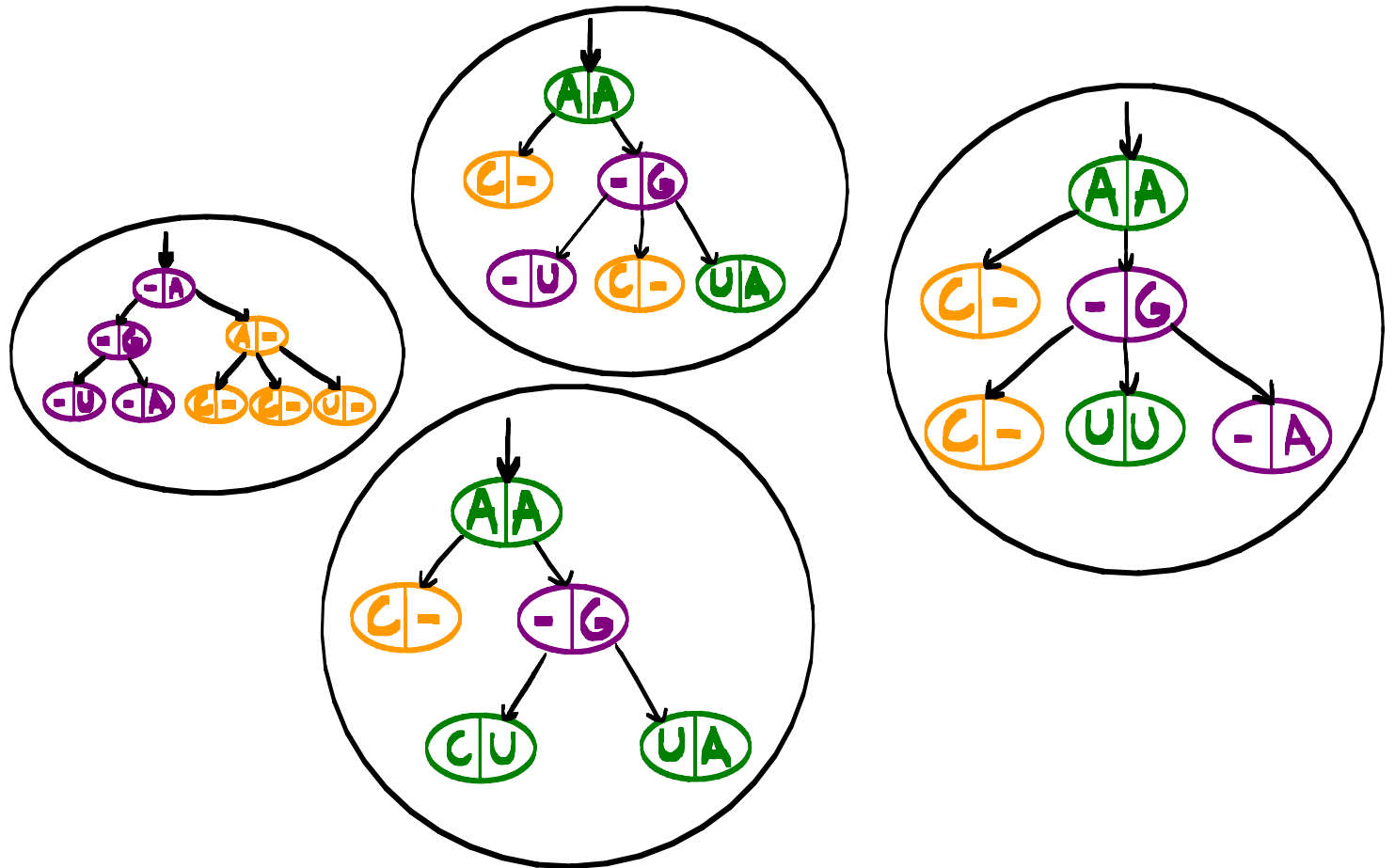between two random trees of cumulative size $n$

where $C = 0.299...$

**Corollary:** A same alignment was repeated

$$\sim 0.875 \times 1.412^n \text{ times on}$$

average in Jiang et al.'s ambiguous grammar.

# APPLICATION 2. SAMPLING

Objective: Sampling alignments under the Gibbs-Boltzmann probability distribution.

Strategy:

→ Filter the grammar to obtain a new grammar that only generates alignments between two fixed trees S and T

→ Use dynamic programming.

# SAMPLING

**Theorem** Let $S$ and $T$ be two trees of size $n_1$ and $n_2$. Sampling alignments between $S$ and $T$ under the Gibbs-Boltzmann distribution can be done with worst-case time and space complexities $O(n_1 n_2 (n_1 + n_2)^2)$ and with average-case time and space complexities $O(n_1 n_2)$.

# SAMPLING

**Theorem** Let $S$ and $T$ be two trees of size $n_1$ and $n_2$. Sampling alignments between $S$ and $T$ under the Gibbs-Boltzmann distribution can be done with worst-case time and space complexities $O(n_1 n_2 (n_1 + n_2)^2)$ and with average-case time and space complexities $O(n_1 n_2)$.

Proof inspired by

[Herrbach, Denise, Dulucq]

# SAMPLING

**Theorem** Let $S$ and $T$ be two trees of size $n_1$ and $n_2$. Sampling alignments between $S$ and $T$ under the Gibbs-Boltzmann distribution can be done with worst-case time and space complexities $O(n_1 n_2 (n_1 + n_2)^2)$ and with average-case time and space complexities $O(n_1 n_2)$.

<u>Upsides</u>:
- No additional complexity cost (except constants, moderate)
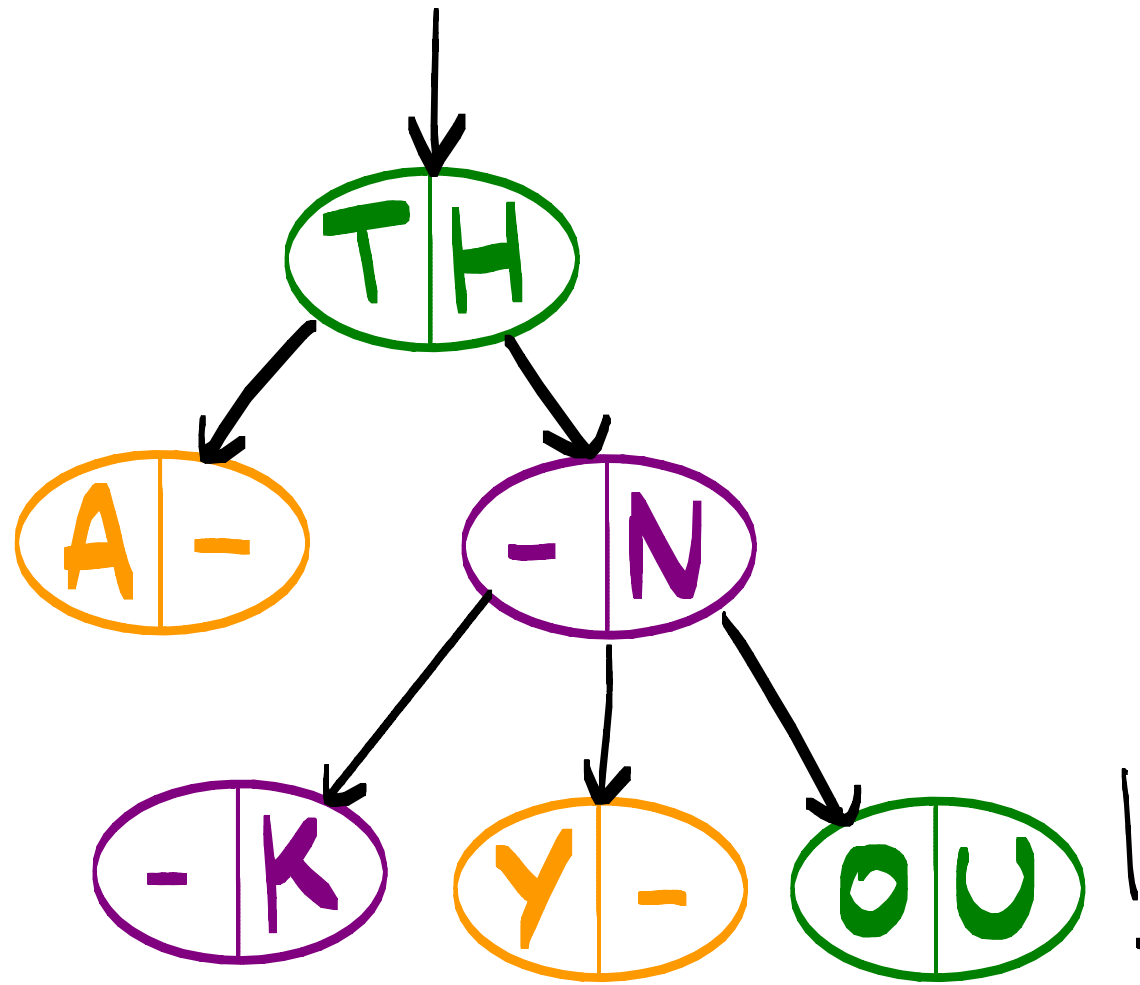- Flexibility of the sampling algorithm.
- Already implemented.

<u>Downside</u>
- Complicated DP scheme -

# CONCLUSION

→ Algorithm that generates the neighborhood of a tree?

→ Existence of easier decompositions?

→ Alignment problem for arc-annoted sequences?

$$\binom{A}{A}\binom{U}{-}\binom{-}{C}\binom{U}{-}\binom{C}{C}\binom{-}{A}\binom{-}{U}\binom{G}{G}\binom{A}{A}\binom{U}{U}\binom{A}{U}\binom{C}{A}\binom{-}{A}\binom{U}{U}\binom{U}{U}$$

# APPLICATION 1: COUNTING.

Theorem: The generating function $A(z, u)$ of tree alignments satisfies

$$A(z, u) = \left(z^2 + z - u z^2 + \frac{z}{\sqrt{1-4z}}\right) \times B(z, u)$$

where

$$\left(u z \, C(z)^2 - z^2 C(z)^2 + 2z\right) B(z, u)^2 + \left(z^2 C^4(z) - 2z C(z)^2 - 1\right) B(z, u) + C^2(z) = 0$$

and

$$C(z) = \frac{1 - \sqrt{1-4z}}{2z} \qquad \text{Catalan generating function}$$