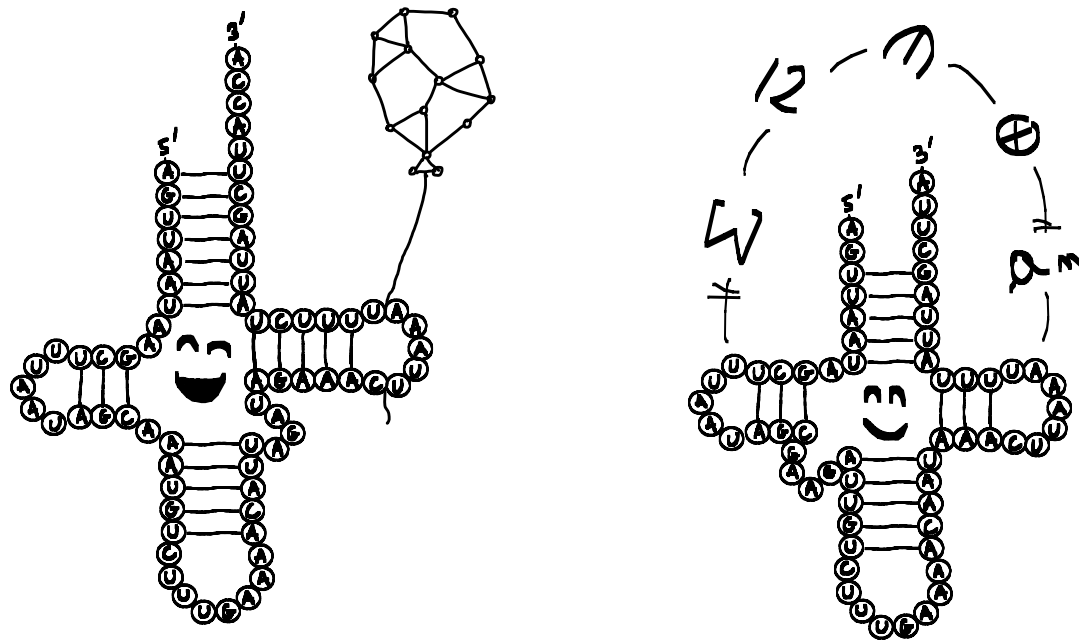# GENERATING TREE ALIGNMENTS
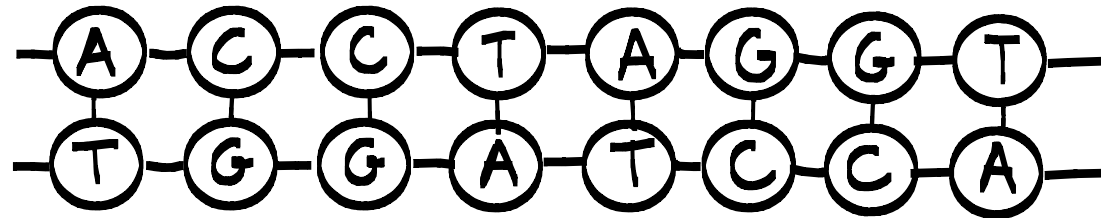## HOW COMBINATORICS CAN HELP BIOINFORMATICS

Julien COURTIEL (PiMS/Univ. of British Columbia, Vancouver)
2016 SFU Symposium on Mathematics and Computation

Co-authors: Cedric CHAUVE (Simon Fraser University, Vancouver)
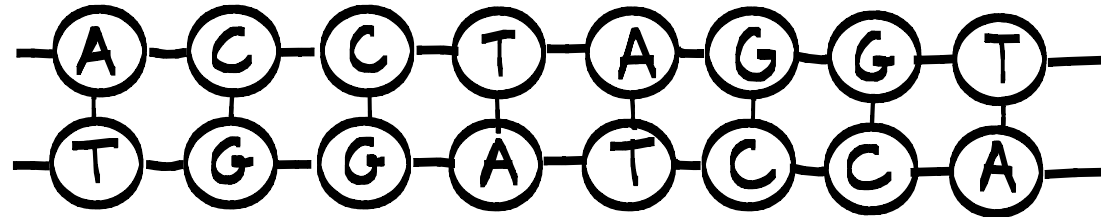Yann PONTY (CNRS/LIX, Ecole Polytechnique, Inria Saclay)
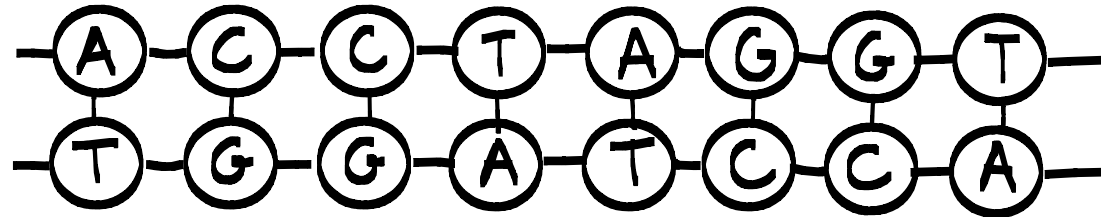
# WHAT IS RNA?

**DNA**
**the code**



Top strand: A C C T A G G T

Bottom strand: T G G A T C C A

DNA
the code

$\downarrow$

RNA

A - C - C - T - A - G - G - T
|   |   |   |   |   |   |   |
T - G - G - A - T - C - C - A

# WHAT IS RNA?

BETTER CALL POL

DNA
the code



A - C - C - T - A - G - G - T
T - G - G - A - T - C - C - A

↓

RNA

Pol

# WHAT IS RNA?

DNA
the code

↓

RNA

A - C - C - T - A - G - G - T -

T - G - G - A - T - C - C - A

Pol

- U -

# WHAT IS RNA?

DNA
the code

$\downarrow$

RNA

A — C — C — T — A — G — G — T

T — G — G — A — T — C — C — A

Pol

U — G

# WHAT IS RNA?

DNA
the code

↓

RNA

A—C—C—T—A—G—G—T

T—G—G—A—T—C—C—A

*Pol*

U—G—G

# WHAT IS RNA?

DNA
the code

DNA → RNA

A–C–C–T–A–G–G–T—

T–G–G–A–T–C–C–A

U–G–G–A–U–C–C...

Pol

# WHAT IS RNA?

DNA
the code

A - C - C - T - A - G - G - T

T - G - G - A - T - C - C - A

↓

RNA

U - G - G - A - U - C - C ...

↓

proteins

# WHAT IS RNA?

DNA
the code

A - C - C - T - A - G - G - T

T - G - G - A - T - C - C - A

↓

RNA

U - G - G - A - U - C - C ...

↓

proteins

Ribo

# WHAT IS RNA?

DNA
the code

A-C-C-T-A-G-G-T

T-G-G-A-T-C-C-A

↓

RNA

U-G-G-A-U-C-C...

Ribo

Trp

↓

proteins

# WHAT IS RNA?

DNA
the code

A—C—C—T—A—G—G—T

T—G—G—A—T—C—C—A

↓

RNA

U—G—G—A—U—C—C...

Ribo

↓

proteins

Trp—Gly

# WHAT IS RNA?

DNA
the code

A - C - C - T - A - G - G - T -
T - G - G - A - T - C - C - A -

↓

RNA

U - G - G - A - U - C - C ...

Ribo

↓

proteins

Trp - Gly - Asp - Ile - Ser

WHAT IS RNA?

DNA
the code

A - C - C - T - A - G - G - T

T - G - G - A - T - C - C - A

RNA

U - G - G - A - U - C - C ...

proteins

Trp - Gly - Asp - Ile - Ser - ...

# Classic dogma

## WHAT IS RNA?

**DNA**
the code

A - C - C - T - A - G - G - T

T - G - G - A - T - C - C - A

↓

**RNA**
the messenger

U - G - G - A - U - C - C ...

↓

**proteins**
the machine

Trp - Gly - Asp - Ile - Ser ...

# Classic dogma

## WHAT IS RNA?

**DNA**
the code

A - C - C - T - A - G - G - T
T - G - G - A - T - C - C - A

↓

**RNA**
the messenger?

U - G - G - A - U - C - C ...

BOOO RING

**proteins**
the machine

Trp - Gly - Asp - Ile - Ser ...

Dogma_2.0                    WHAT IS RNA?

DNA
the code

A—C—C—T—A—G—G—T—
T—G—G—A—T—C—C—A

↓

RNA
the messenger

U—G—G—A—U—C—C—...

actually also...        translator,      enzyme,
         ↑              regulator,       catalyst...
      recent
   breakthrough

↓

proteins
the machine

Trp—Gly—Asp—Ile—Ser—...

# WHAT IS RNA?

RNA
the messenger


U–G–G–A–U–C–C ...

actually also...   translator,   enzyme,
regulator,   catalyst...

So what is RNA?

RNA is

a single-stranded
molecule
(chain of nucleotides)...

U — G — G — A — U — C — C — ...

So what is RNA?

RNA is

a single-stranded
molecule
(chain of nucleotides)...

... stabilized by
hydrogen bonds..

# So what is RNA?

RNA is

a single-stranded molecule (chain of nucleotides)...

...stabilized by hydrogen bonds...

...which folds onto itself.

# So what is RNA?

RNA is

a single-stranded molecule (chain of nucleotides)...

... stabilized by hydrogen bonds..

...which folds onto itself.



primary structure



secondary structure



tertiary structure

# RNA COMPARISON

Interesting problem: evaluating similarity between two RNAs.

# RNA COMPARISON

Interesting problem: evaluating similarity between two RNAs.



— Why?

# RNA COMPARISON

Interesting problem: evaluating similarity between two RNAs.



O - Why?

Typical situation:

New
RNA →



ACAGUACC...

large database

AUCCAG ...          UUAGACC...

CCAGC ..           AAAGU...

# RNA COMPARISON

**Interesting problem:** evaluating similarity between two RNAs.



— Why?

**Typical situation:**

New RNA → *(RNA structure)* finding similar RNAs → large database

ACAGUACC...

AUCCAG...
CCAGC..
UUAGACC...
AAAGU...

# 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 the primary structures.

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

# MOTIVATION: RNA COMPARISON

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



First idea: compare the primary structures.
→ sequence alignment

RNA 1: AUUCGAUUA...

RNA 2: ACCAUGAUUA...

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

# MOTIVATION: RNA COMPARISON

Question: how to measure similarity between two RNAs?



Second idea: compare secondary structures.



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

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

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 $S_2$

# SEQUENCE ALIGNMENT

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

$\begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} - \\ C \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} C \\ C \end{pmatrix} \begin{pmatrix} - \\ A \end{pmatrix} \begin{pmatrix} - \\ U \end{pmatrix} \begin{pmatrix} G \\ G \end{pm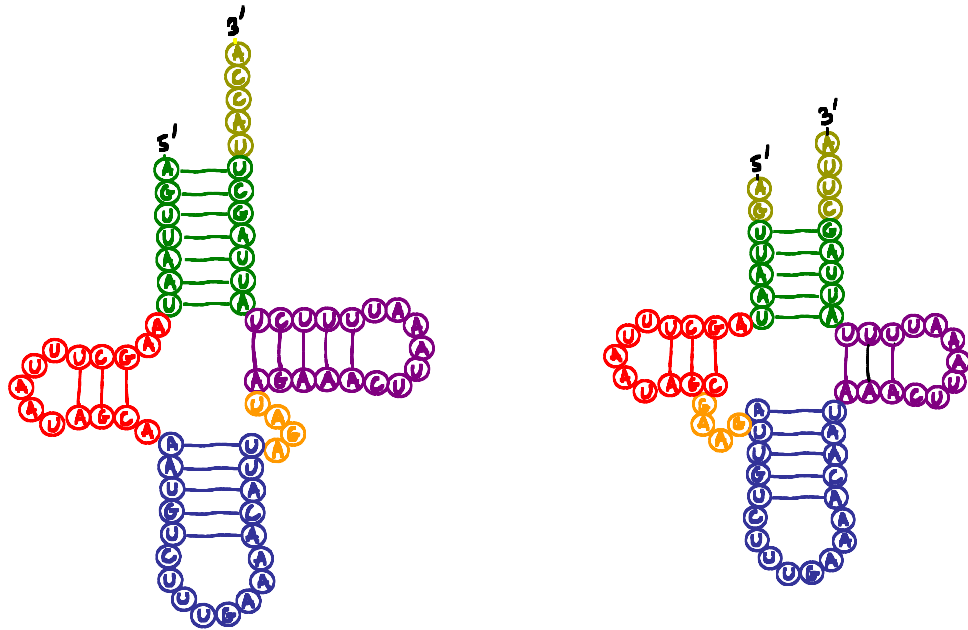atrix} \begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ U \end{pmatrix} \begin{pmatrix} A \\ U \end{pmatrix} \begin{pmatrix} C \\ A \end{pmatrix}$ $\longrightarrow$ AUUCGAUAC

$\longrightarrow$ ACCAUGAUUA

match    insertion    deletion    mismatch    projections

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

# OPTIMAL ALIGNMENT

<u>Classical problem</u>: Given $S_1$ and $S_2$, find one optimal alignment between $S_1$ and $S_2$.

<u>Solvable by Dynamic Programming</u>:

→ Needleman - Wunsch algorithm

→ Smith - Waterman algorithm

→ BLAST (heuristic)

Worst case and average time : $O(n^2)$

# 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



A

C    C    U

keep left letters

# TREE ALIGNMENTS

first projection



second projection

Keep left letters

Keep right letters

# TREE ALIGNMENTS

first projection



Keep left letters

second projection

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

# CONNECTION WITH SEQUENCE ALIGNMENTS

Tree alignments generalize sequence alignments.

**SEQUENCE**

AUUCGAUUA...     ACCAUGAUUA...

alignment :

$\left(\begin{array}{c}A\\A\end{array}\right)\left(\begin{array}{c}U\\-\end{array}\right)\left(\begin{array}{c}-\\C\end{array}\right)\left(\begin{array}{c}U\\-\end{array}\right)\left(\begin{array}{c}C\\C\end{array}\right)\left(\begin{array}{c}-\\A\end{array}\right)\left(\begin{array}{c}-\\U\end{array}\right)\left(\begin{array}{c}G\\G\end{array}\right)\left(\begin{array}{c}A\\A\end{array}\right)\left(\begin{array}{c}U\\U\end{array}\right)\left(\begin{array}{c}U\\U\end{array}\right)\left(\begin{array}{c}A\\A\end{array}\right)$...

**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



and



is the most likely?



co-optimal alignments

# SPACE OF ALIGNMENTS

Why finding one optimal alignment may be inadequate:

▶ Co-optimal alignments can be very different.

▶ 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.

?? probability of (an) alignment A

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

(Gibbs-Boltzmann distribution)



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

Legend: kT=0.1, kT=1, kT=10

Y-axis: Probability (0 to 0.5)
X-axis: Score (1 to 16)

$\underline{K=0}$ : 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,

A|A
C|- -|G
-|U C|- U|A

and

A|A
-|G
C|- C|- U|A
-|U

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

Insertions before Deletions.

# PROBLEM RAISED BY THE AMBIGUITY

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

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

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

Insertions before Deletions.

For trees, it is much more complicated!

Strategy: COMBINATORICS!

# PROBLEM RAISED BY THE AMBIGUITY

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

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

Insertions before Deletions.

For trees, it is much more complicated!

Strategy: COMBINATORICS ! ၛ☺ၟ

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

# GRAMMARS FOR SEQUENCE ALIGNMENTS

## Ambiguous grammar:

$$\mathcal{S} \leftarrow \binom{X}{Y} \boxed{\mathcal{S}} \oplus \binom{X}{-} \boxed{\mathcal{S}} \oplus \binom{-}{Y} \boxed{\mathcal{S}} \oplus \varepsilon$$

## Non-ambiguous grammar:

$$\mathcal{S} \leftarrow \binom{X}{Y} \boxed{\mathcal{S}} \oplus \binom{X}{-} \boxed{\mathcal{S}} \oplus \binom{-}{Y} \boxed{\mathcal{S}^{D}} \oplus \varepsilon$$

$$\mathcal{S}^{D} \leftarrow \binom{X}{Y} \boxed{\mathcal{S}} \oplus \binom{-}{Y} \boxed{\mathcal{S}^{D}} \oplus \varepsilon$$

# 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:

$$\mathcal{A} \leftarrow \mathcal{V}^\phi \mid \mathcal{C}_I \mid \mathcal{C}_D \mid \boxed{\mathcal{F}_I}\,\boxed{\mathcal{C}_D}$$

$$\mathcal{V}^\phi \leftarrow \mathcal{V}^\uparrow \mid \boxed{\mathcal{V}\mathcal{H}}$$

$$\mathcal{C}_I \leftarrow \boxed{\mathcal{F}_I}$$

$$\mathcal{F}_I \leftarrow \varepsilon \mid \boxed{\mathcal{F}_I}\,\boxed{\mathcal{F}_I}$$

$$\mathcal{C}_D \leftarrow \boxed{\mathcal{F}_D}$$

$$\mathcal{F}_D \leftarrow \varepsilon \mid \boxed{\mathcal{F}_D}\,\boxed{\mathcal{F}_D}$$

$$\mathcal{V}^\uparrow \leftarrow \boxed{\mathcal{H}_{ID,\phi,\phi}} \mid \boxed{\mathcal{F}_D}\,\boxed{\mathcal{V}^\uparrow}\,\boxed{\mathcal{F}_D}$$

$$\mathcal{V}\mathcal{H} \leftarrow \boxed{\mathcal{F}_I}\,\boxed{\mathcal{V}\mathcal{H}} \mid \boxed{\mathcal{V}^\phi}\,\boxed{\mathcal{F}_I} \mid \boxed{\mathcal{H}_{ID,\leftrightarrow,\phi}}\,\boxed{\mathcal{F}_I}$$

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

$$\mathcal{H}_{\mathcal{V},M,M'} \leftarrow \varepsilon \mid \boxed{\mathcal{C}_I}\,\boxed{\mathcal{H}_{\mathcal{V},M,M'}} \mid \boxed{\mathcal{C}_D}\,\boxed{\mathcal{H}_{D,M,M'}} \mid \boxed{\mathcal{V}^\phi}\,\boxed{\overline{\mathcal{H}}_{M,M'}^{i,i'}} \mid \boxed{\overline{\mathcal{H}}_{M,M'}^{i',+}}\,\boxed{\mathcal{H}_{ID,\phi,\leftrightarrow}} \mid \boxed{\overline{\mathcal{H}}_{M,M'}^{+,i'}}\,\boxed{\mathcal{H}_{ID,\leftrightarrow,\phi}}$$

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

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

only if $M' \neq \leftrightarrow$

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

# APPLICATION 1.  COUNTING.

**Theorem**  There are on average

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

between two random trees of cumulative size $n$

where $C = 0.299\ldots$

**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 $O(n_1 n_2)$.

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

## Downside
- Complicated DP scheme -

# CONCLUSION

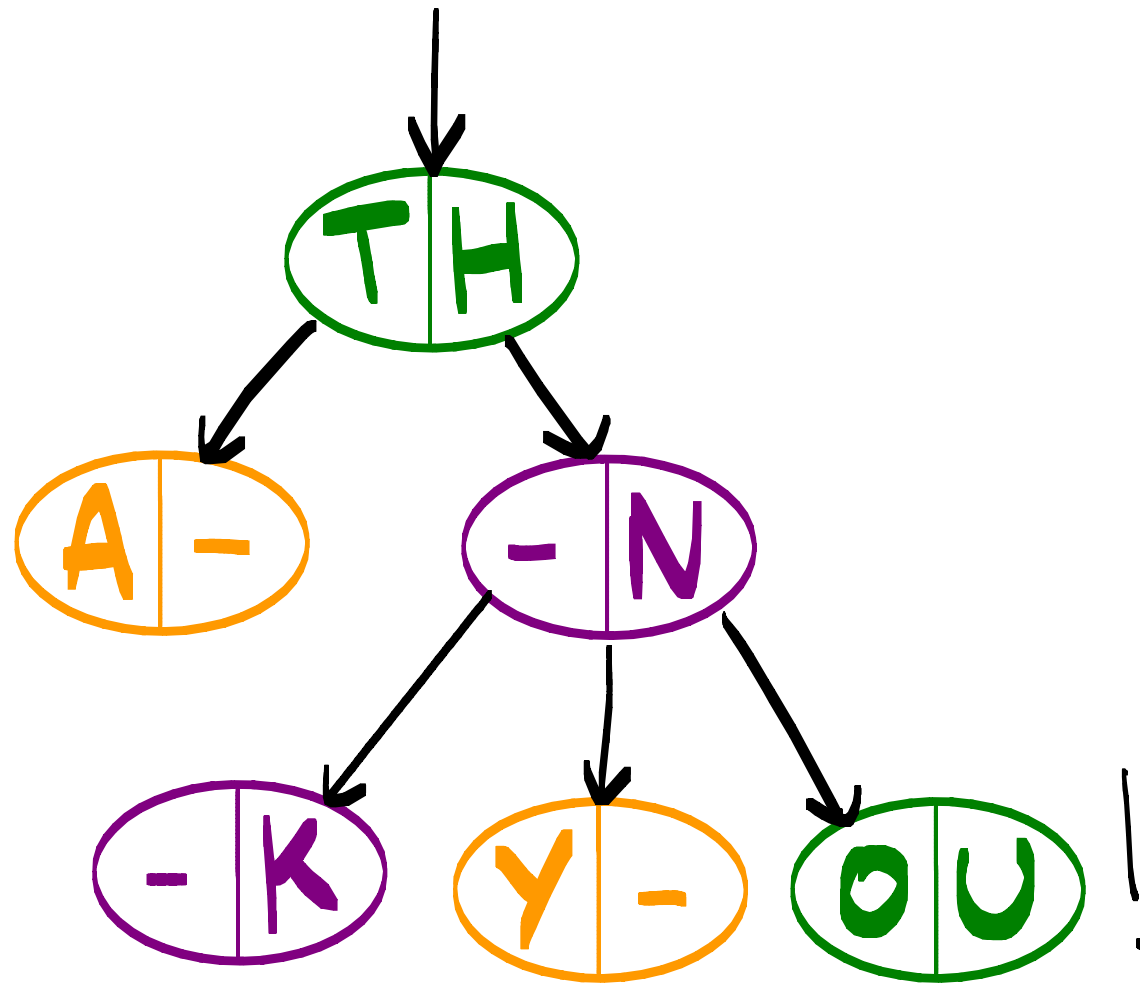"Combinatorics is a powerful tool to solve algorithmic problems."

# CONCLUSION

"Combinatorics is a powerful tool to solve algorithmic problems."

## Open questions:

→ Existence of easier decompositions?

→ Alignment problem for arc-annoted sequences?

$$\left(\begin{matrix}A\\A\end{matrix}\right)\left(\begin{matrix}U\\-\end{matrix}\right)\left(\begin{matrix}-\\C\end{matrix}\right)\left(\begin{matrix}U\\-\end{matrix}\right)\left(\begin{matrix}C\\C\end{matrix}\right)\left(\begin{matrix}-\\A\end{matrix}\right)\left(\begin{matrix}-\\U\end{matrix}\right)\left(\begin{matrix}G\\G\end{matrix}\right)\left(\begin{matrix}A\\A\end{matrix}\right)\left(\begin{matrix}U\\U\end{matrix}\right)\left(\begin{matrix}A\\U\end{matrix}\right)\left(\begin{matrix}C\\A\end{matrix}\right)\left(\begin{matrix}-\\A\end{matrix}\right)\left(\begin{matrix}U\\U\end{matrix}\right)\left(\begin{matrix}U\\U\end{matrix}\right)$$

# APPLICATION 1: COUNTING.

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

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

where

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

and

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