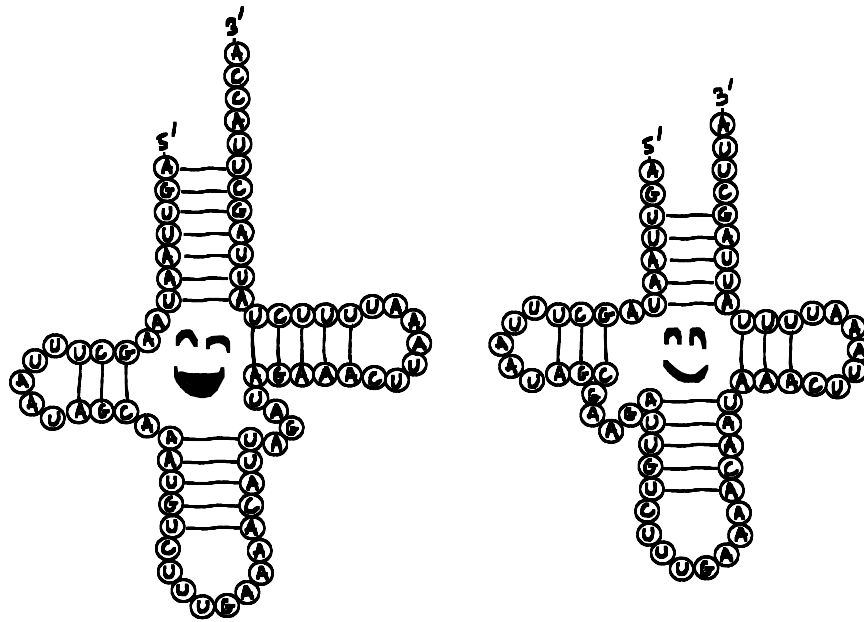

COUNTING, GENERATING AND SAMPLING TREE ALIGNMENTS

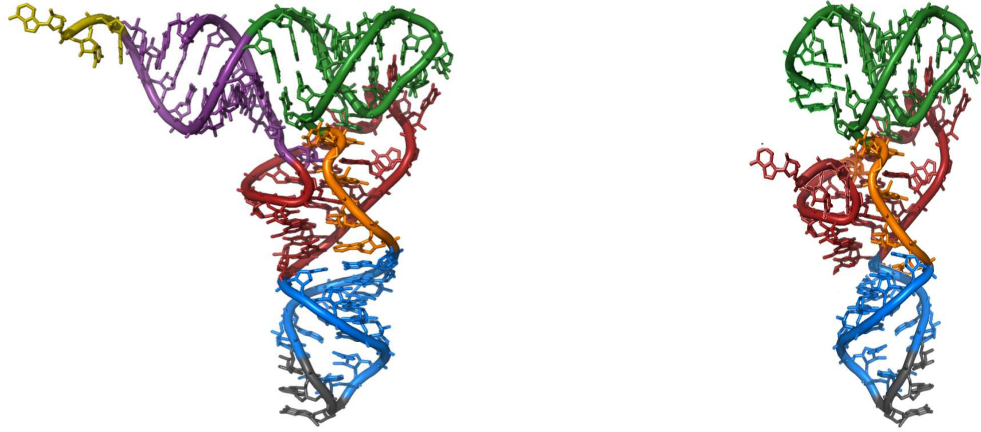
Cedric CHAUVÉ (Simon Fraser University, Vancouver)
Julien COURTIEL (PIMS/Univ. of British Columbia, Vancouver)
Yann PONTY (CNRS/LIX)



SeqBio 2015

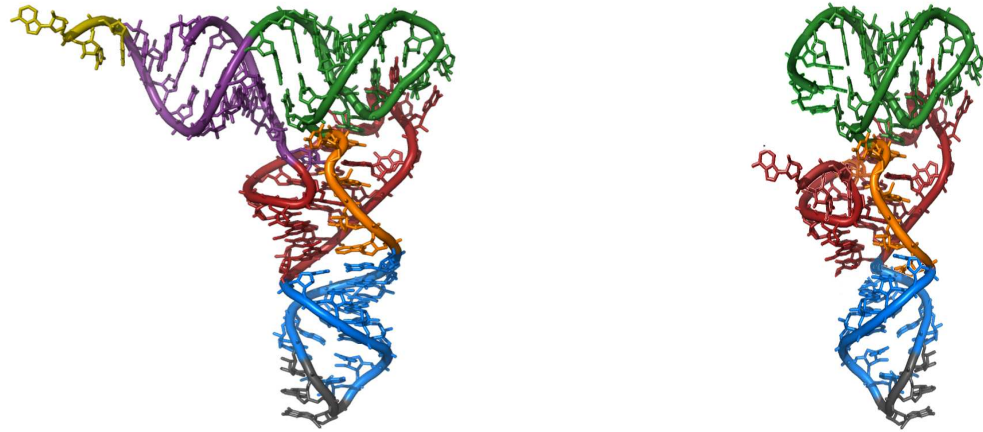
MOTIVATION: RNA COMPARISON

Question: how to measure similarity between two RNAs?



MOTIVATION: RNA COMPARISON

Question: how to measure similarity between two RNAs?



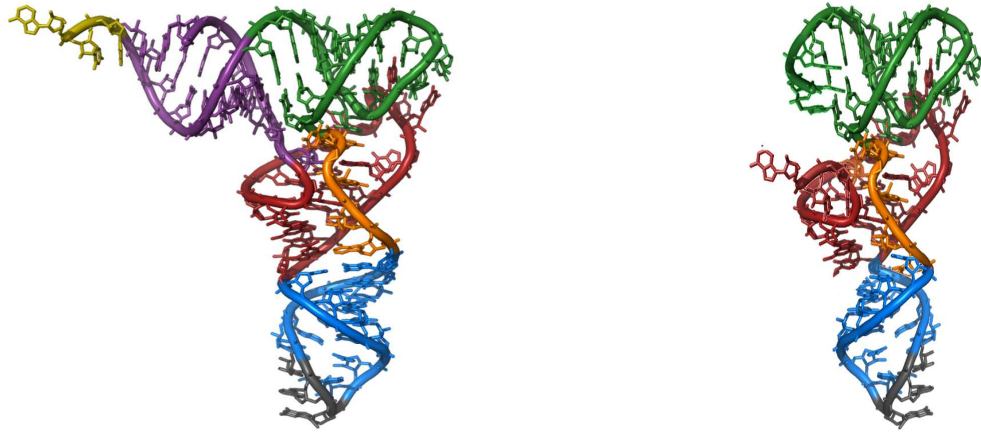
First idea: compare nucleic acid sequences.

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

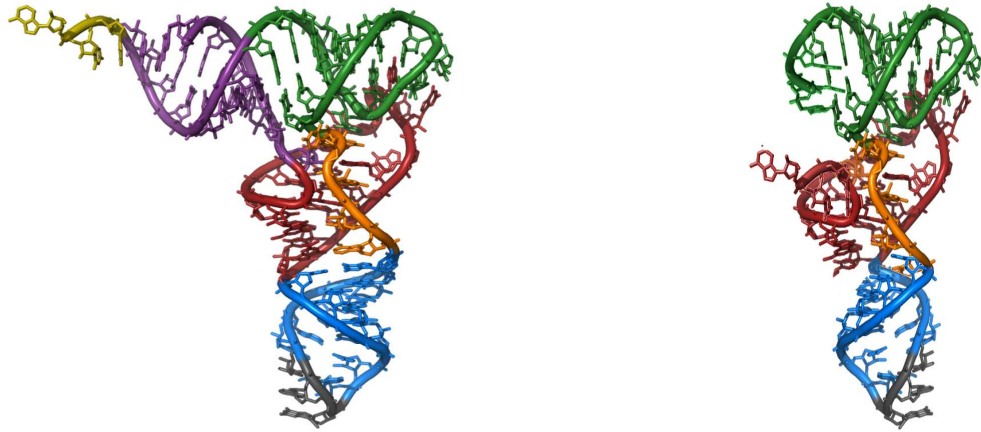
RNA 1: AUUCGAUUA...

RNA 2: ACCAUGAUUA...

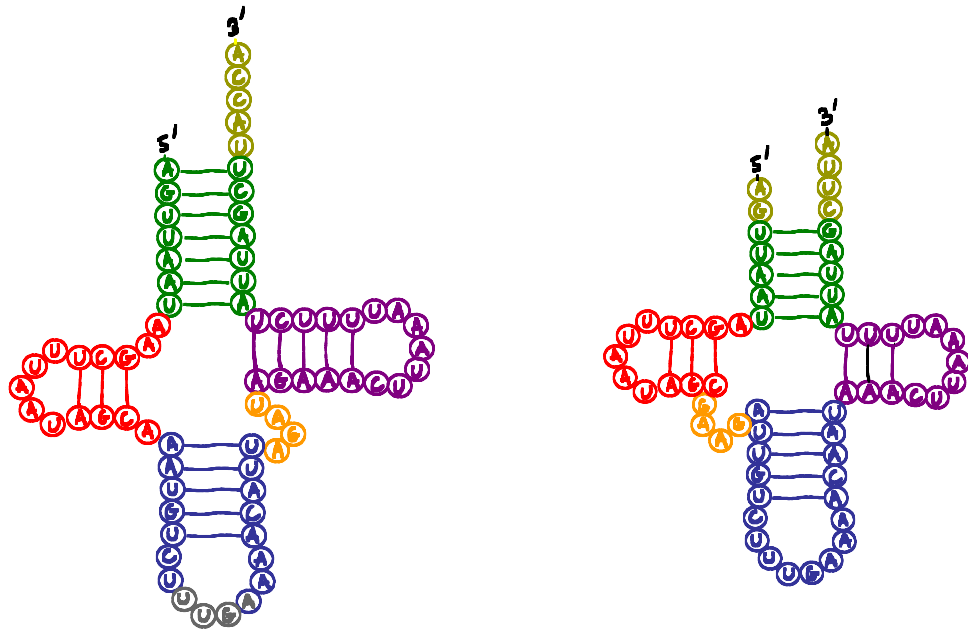
alignment: $\begin{pmatrix} A \\ A \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} - \\ C \end{pmatrix} \begin{pmatrix} U \\ - \end{pmatrix} \begin{pmatrix} C \\ C \end{pmatrix} \begin{pmatrix} - \\ A \end{pmatrix} \begin{pmatrix} - \\ U \end{pmatrix} \begin{pmatrix} G \\ G \end{pmatrix}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

MOTIVATION: RNA COMPARISON

Question: how to measure similarity between two RNAs?

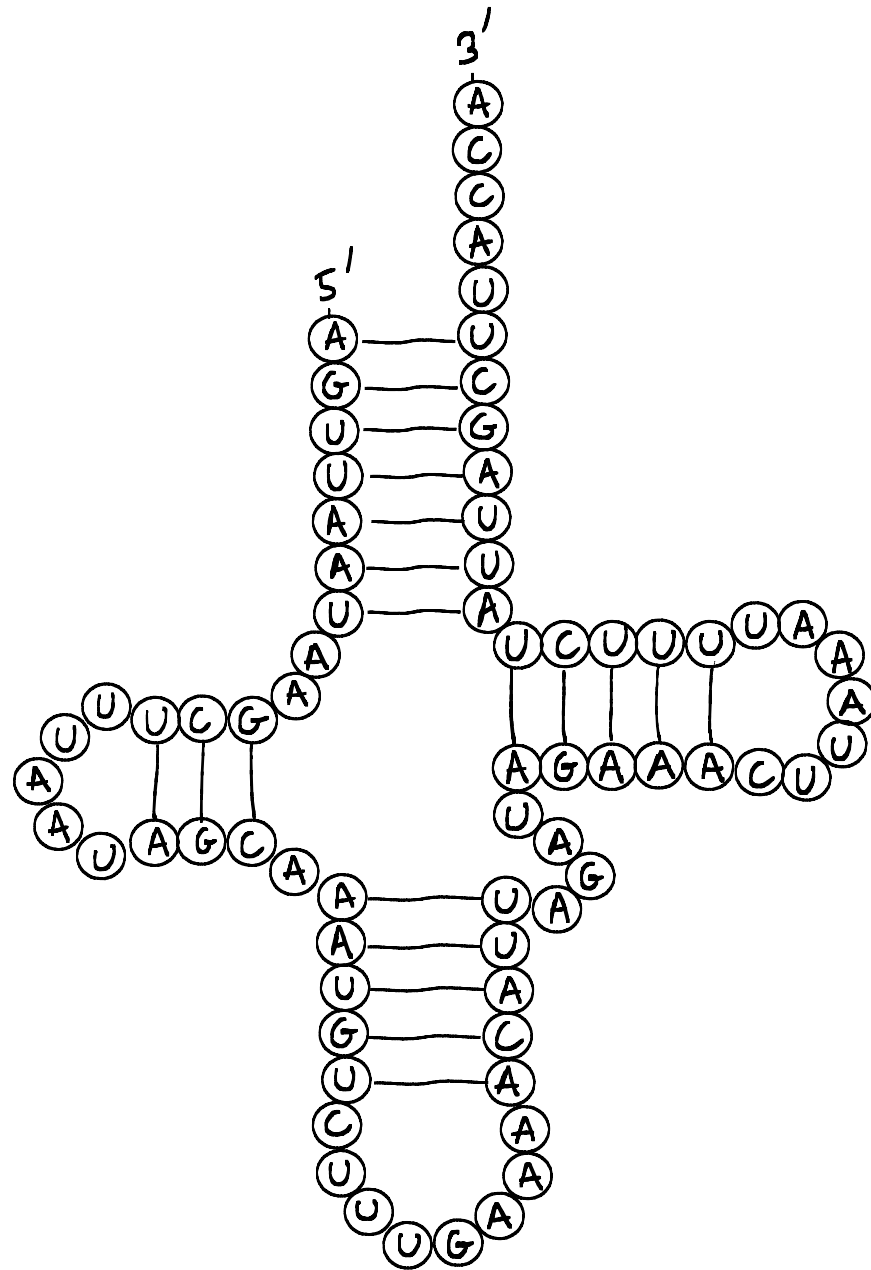


Second idea: compare secondary structures.

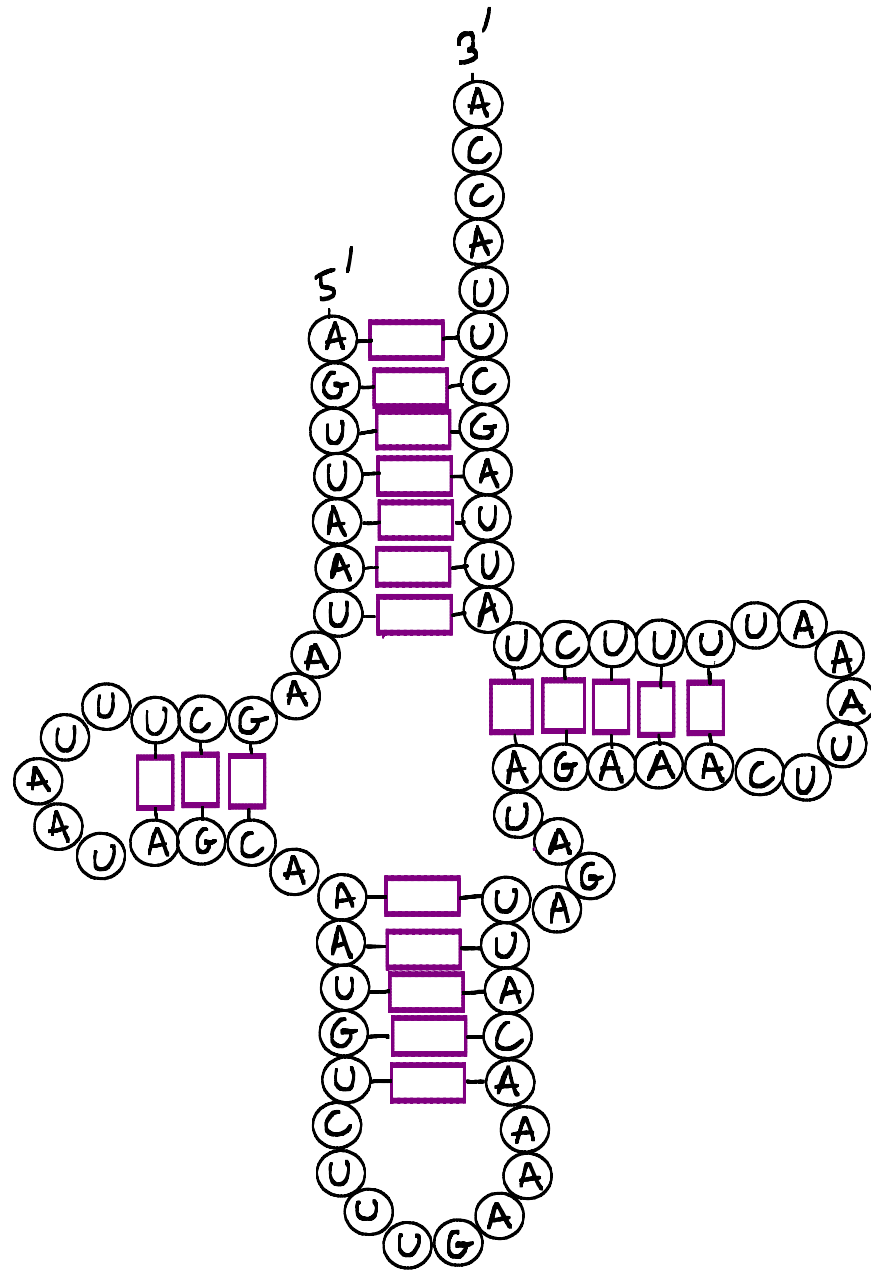


→ notion of
tree alignment
[Jiang, Wang,
Zhang]

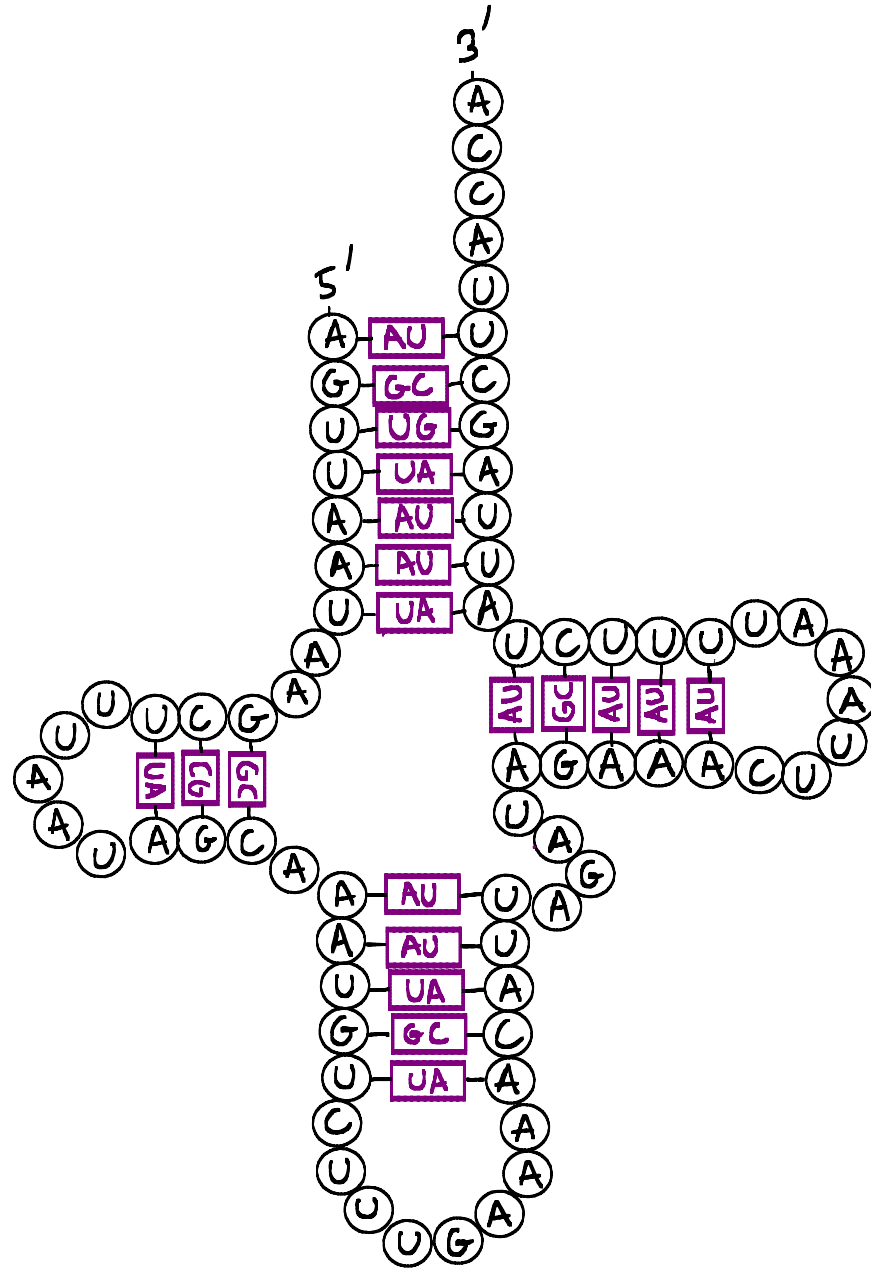
FROM SECONDARY STRUCTURES TO TREES



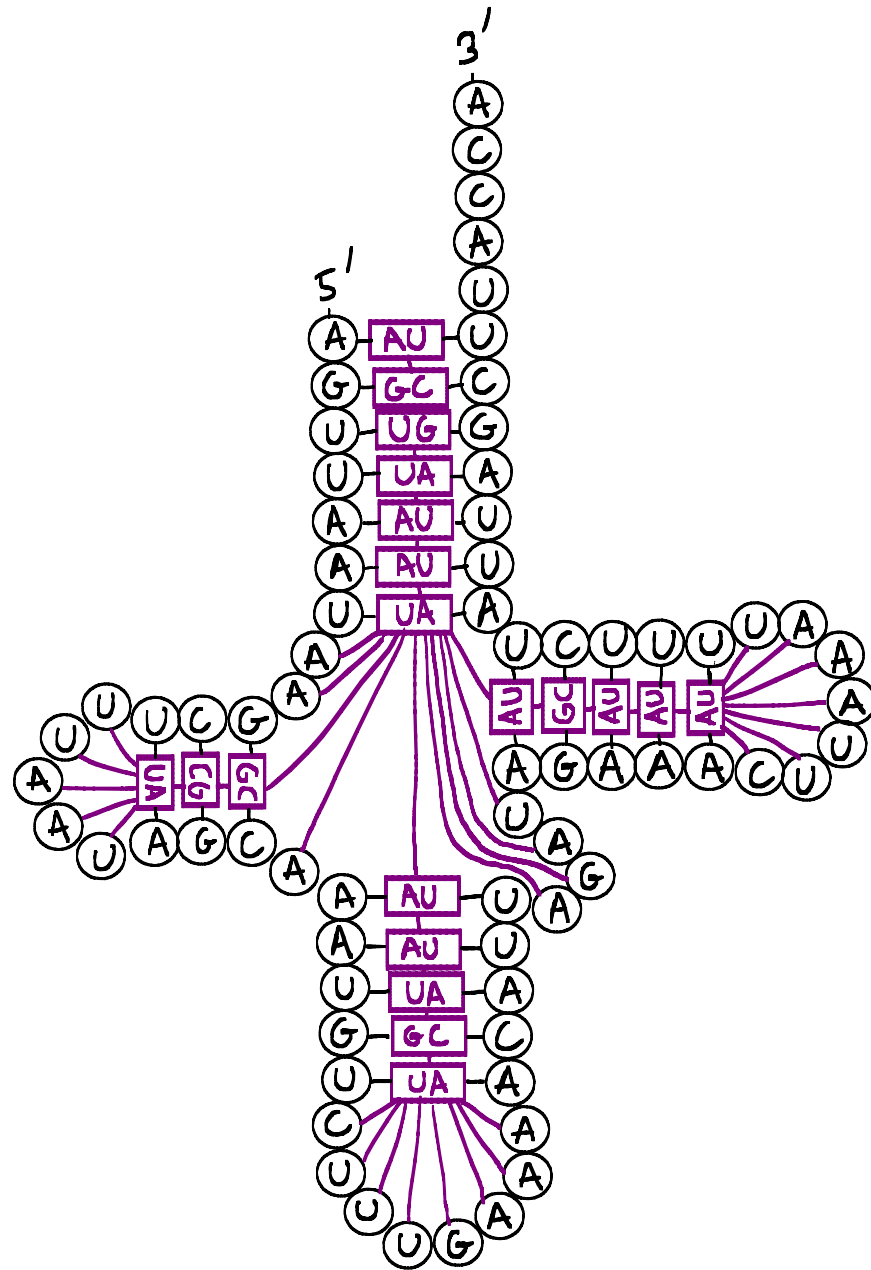
FROM SECONDARY STRUCTURES TO TREES



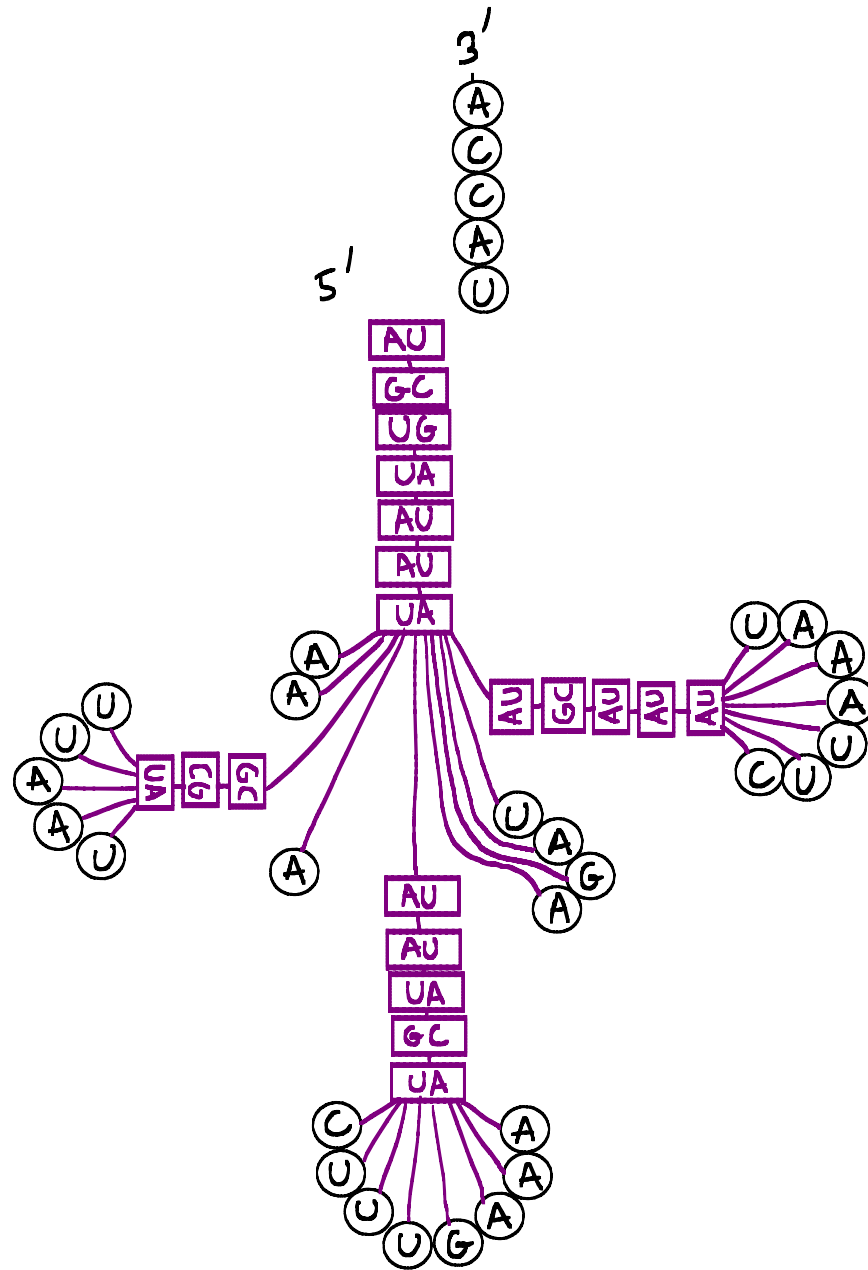
FROM SECONDARY STRUCTURES TO TREES



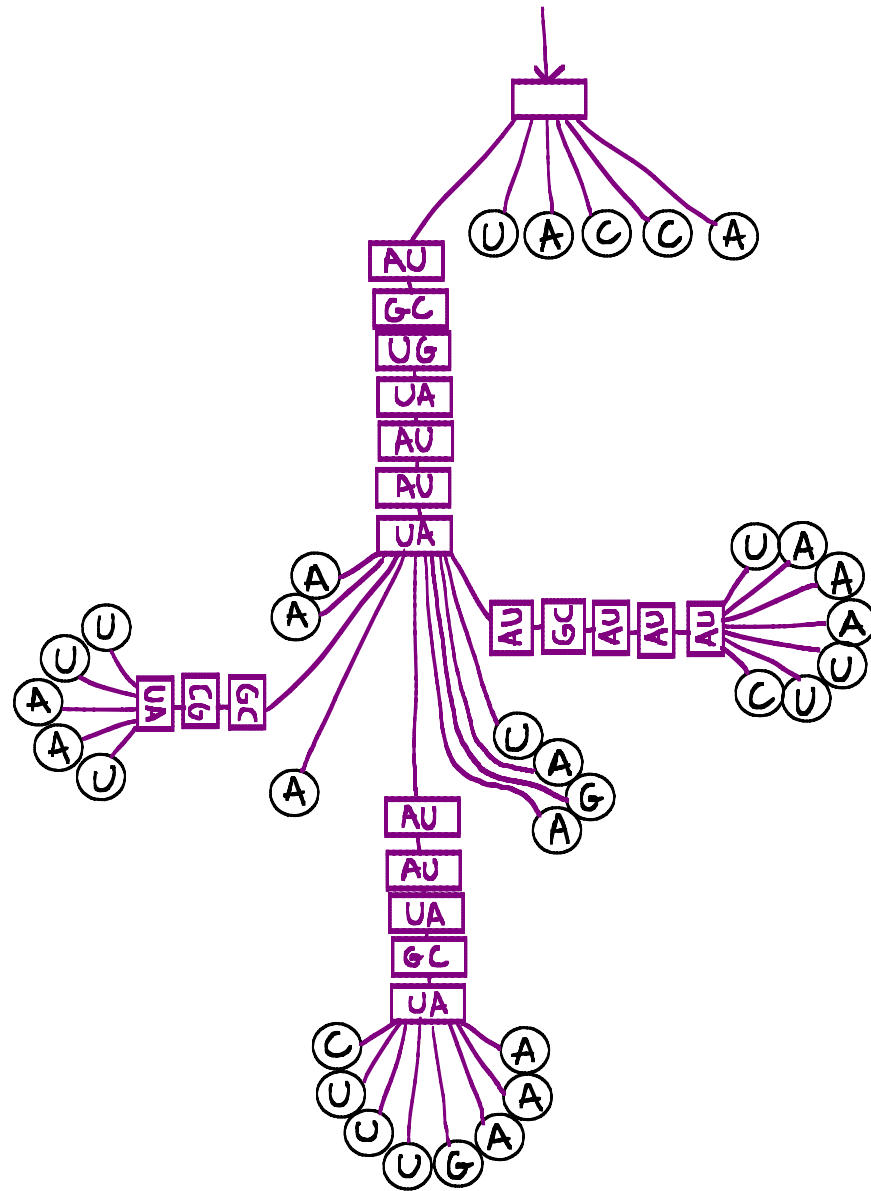
FROM SECONDARY STRUCTURES TO TREES



FROM SECONDARY STRUCTURES TO TREES

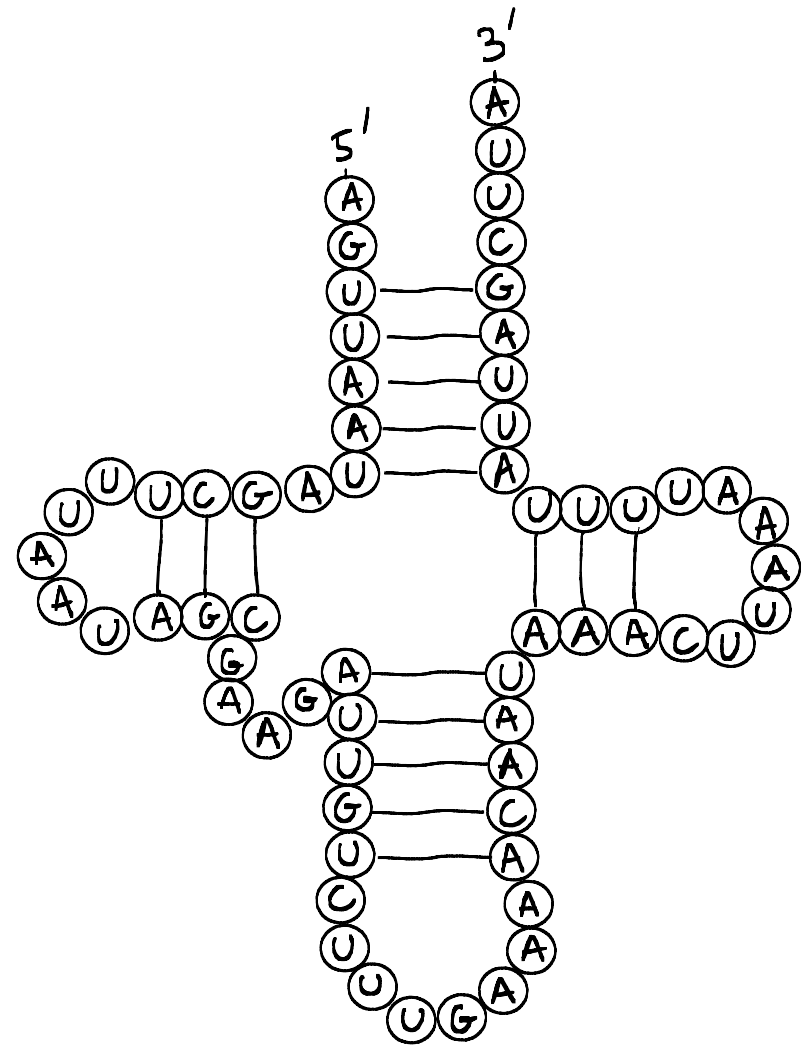
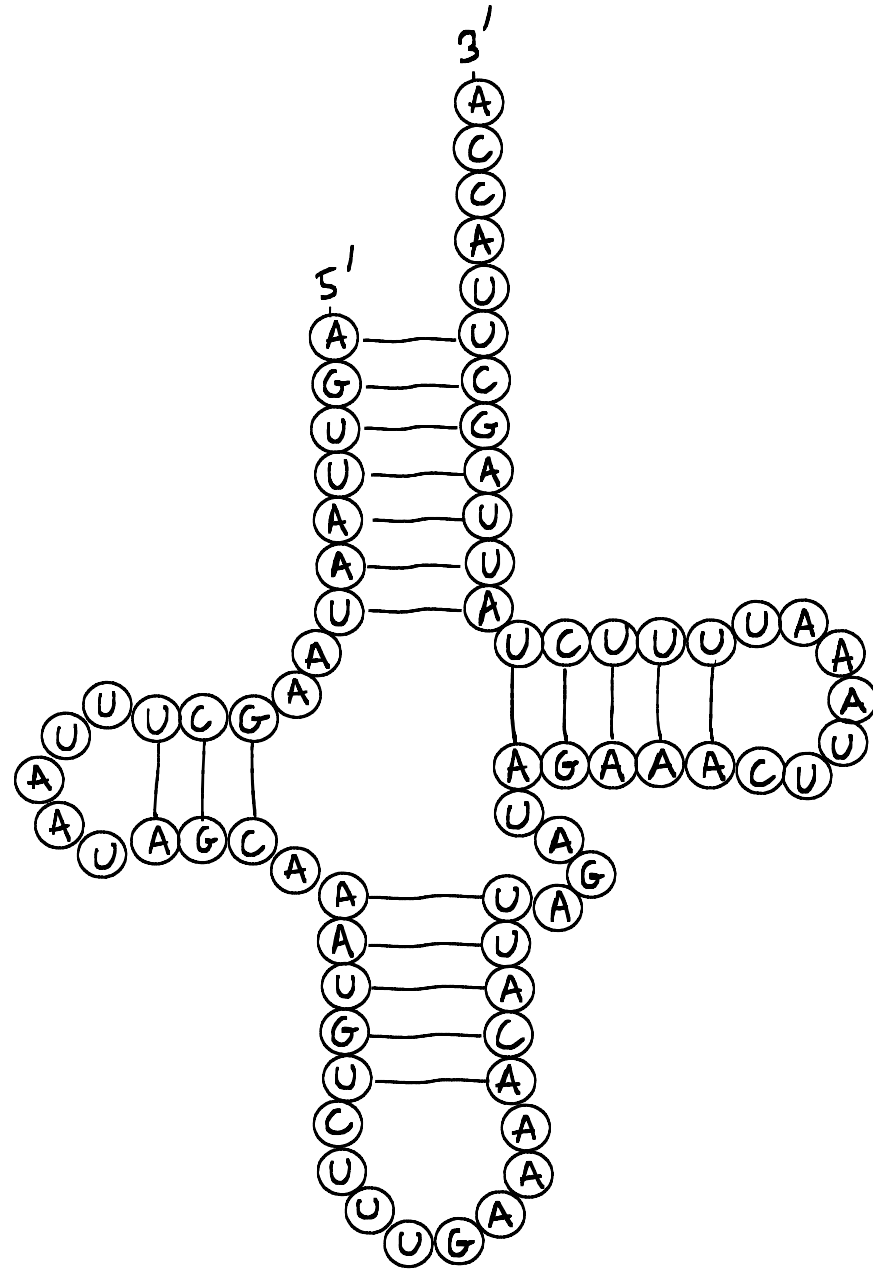


FROM SECONDARY STRUCTURES TO TREES



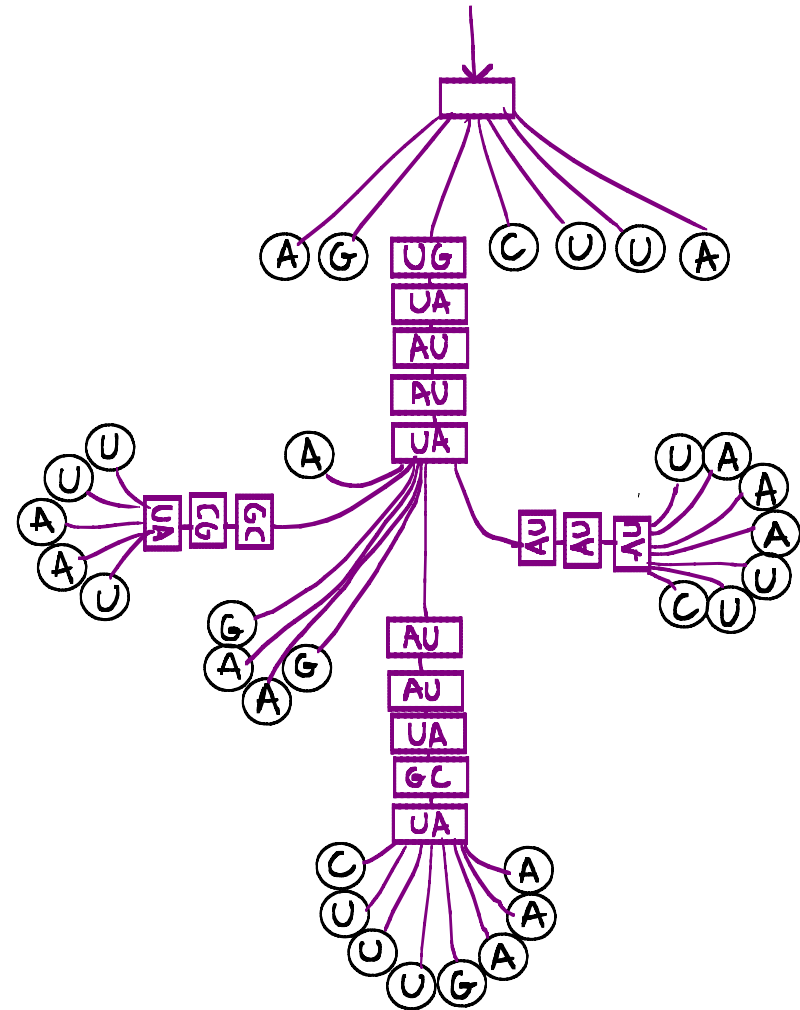
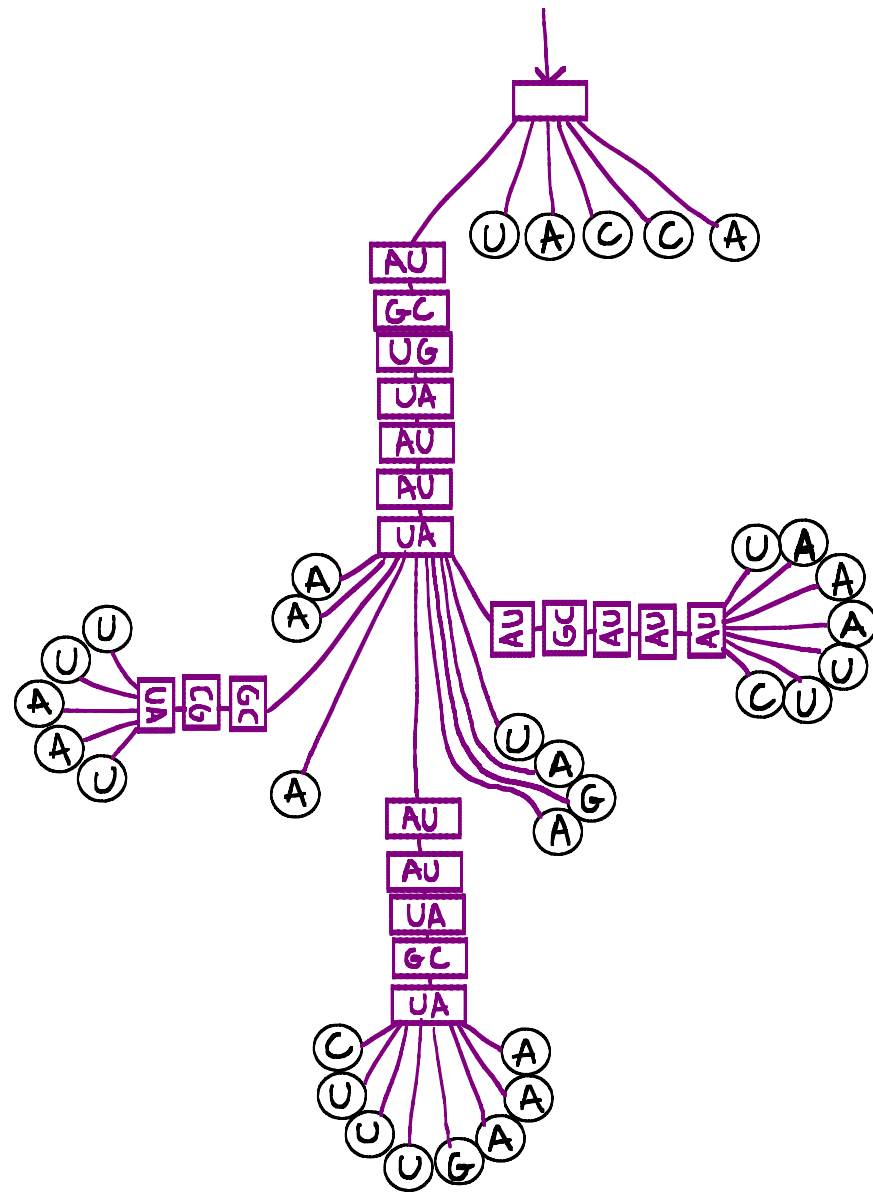
FROM SECONDARY STRUCTURES TO TREES

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



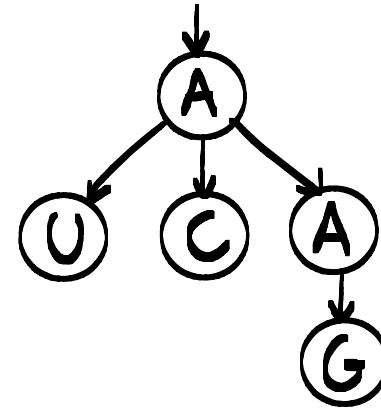
FROM SECONDARY STRUCTURES TO TREES

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



TREES AND SUPERTREES

Trees are plane, rooted, and vertices are labeled by an alphabet Σ .

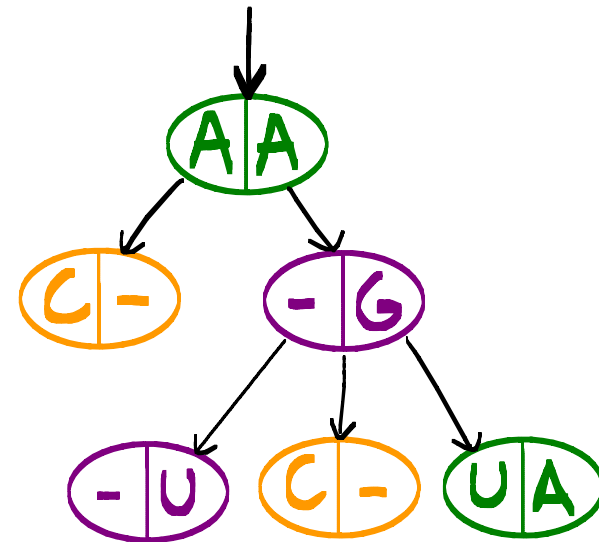


Supertree = tree with 3 types of vertices:

$(X|Y)$ (mis)match

$(X|-)$ insertion

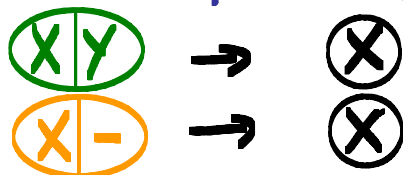
$(-|Y)$ deletion



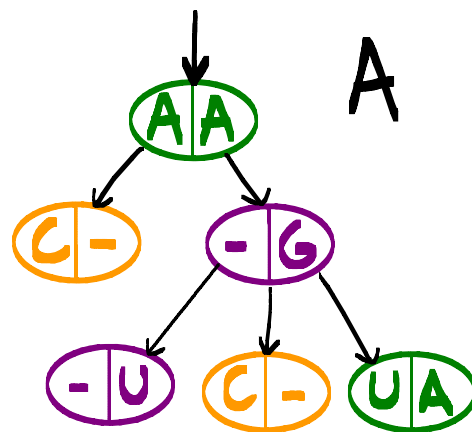
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

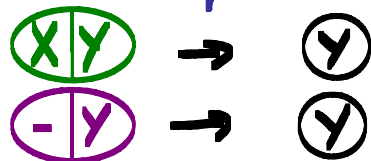
$\Pi_1(A)$ = tree
obtained by changing



and removing -|Y .



$\Pi_2(A)$ = tree
obtained by changing



and removing X|- .

SUPERTREES INDUCE TREE ALIGNMENTS

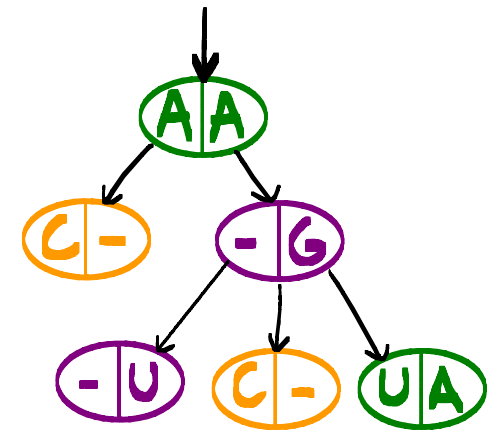
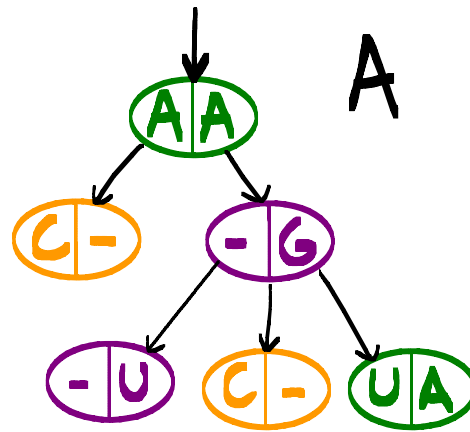
Let A be a supertree,

$\Pi_1(A)$ = tree
obtained by changing

$\begin{array}{|c|c|} \hline X & Y \\ \hline \end{array} \rightarrow \begin{array}{|c|} \hline X \\ \hline \end{array}$

$\begin{array}{|c|c|} \hline X & - \\ \hline \end{array} \rightarrow \begin{array}{|c|} \hline X \\ \hline \end{array}$

and removing $\begin{array}{|c|c|} \hline - & Y \\ \hline \end{array}$.



$\Pi_2(A)$ = tree
obtained by changing

$\begin{array}{|c|c|} \hline X & Y \\ \hline \end{array} \rightarrow \begin{array}{|c|} \hline Y \\ \hline \end{array}$

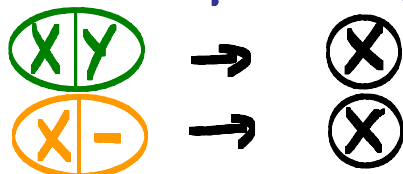
$\begin{array}{|c|c|} \hline - & Y \\ \hline \end{array} \rightarrow \begin{array}{|c|} \hline Y \\ \hline \end{array}$

and removing $\begin{array}{|c|c|} \hline X & - \\ \hline \end{array}$.

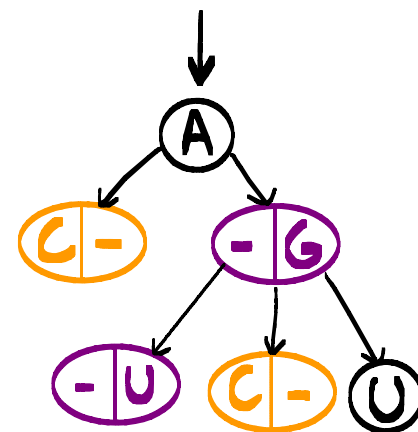
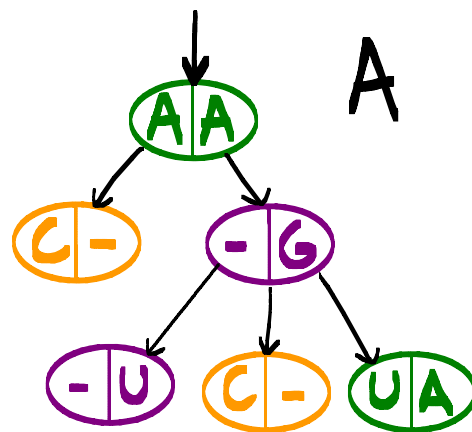
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

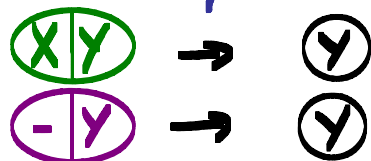
$\Pi_1(A)$ = tree
obtained by changing



and removing -|Y .



$\Pi_2(A)$ = tree
obtained by changing

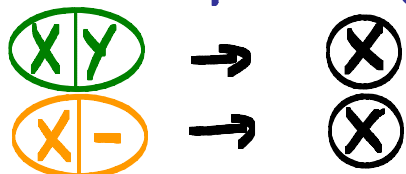


and removing X|- .

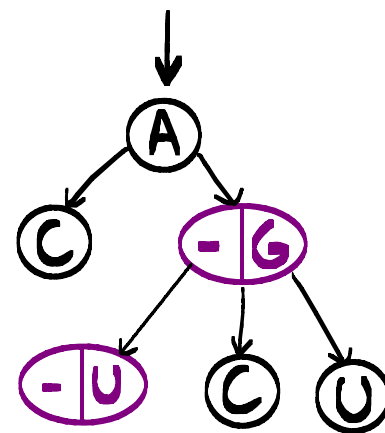
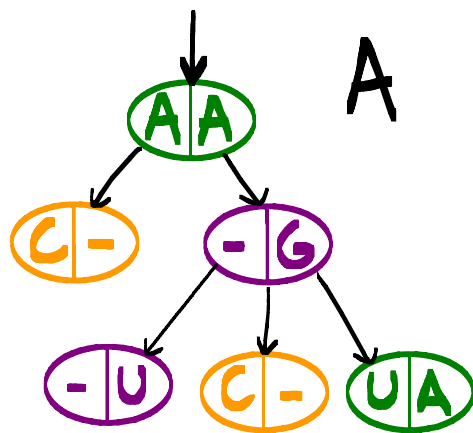
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

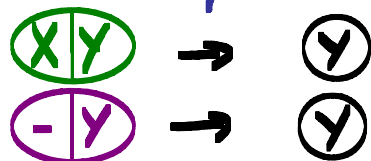
$\Pi_1(A)$ = tree
obtained by changing



and removing -|Y .



$\Pi_2(A)$ = tree
obtained by changing

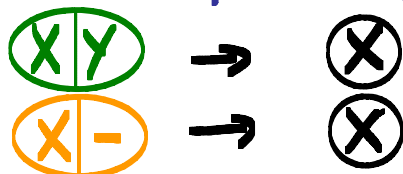


and removing X|- .

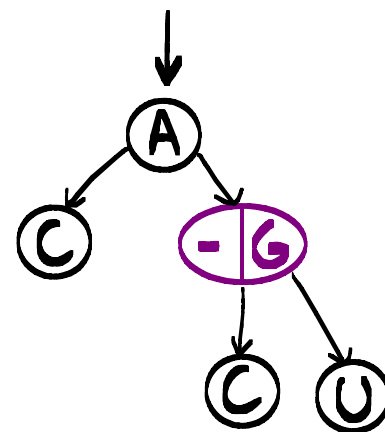
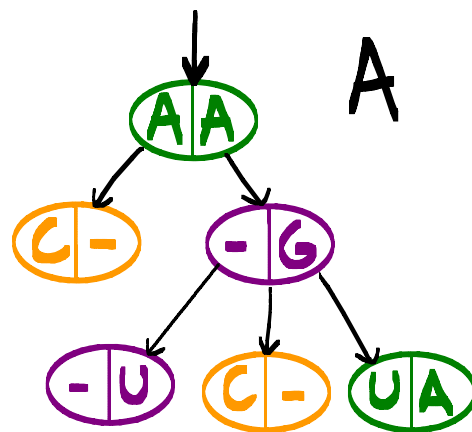
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

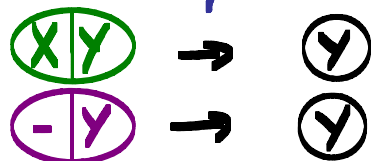
$\Pi_1(A)$ = tree
obtained by changing



and removing -|Y .



$\Pi_2(A)$ = tree
obtained by changing

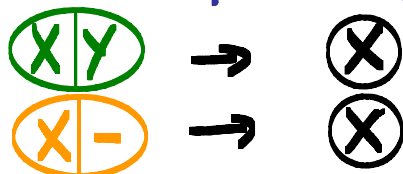


and removing X|- .

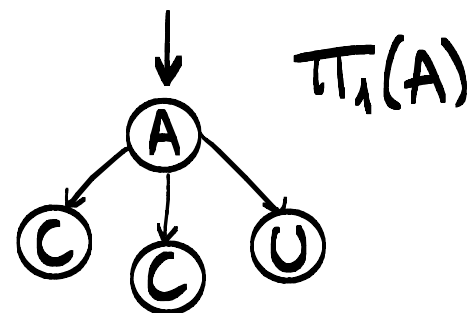
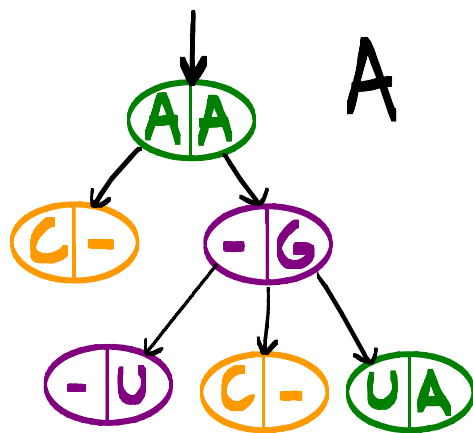
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

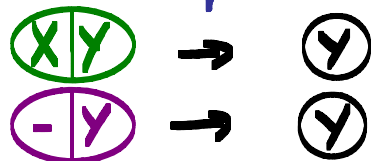
$\pi_1(A)$ = tree
obtained by changing



and removing $\textcircled{-|Y}$.



$\pi_2(A)$ = tree
obtained by changing

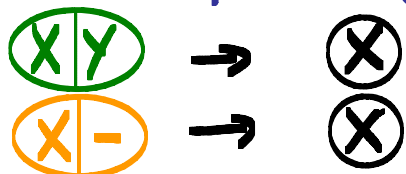


and removing $\textcircled{X|-}$.

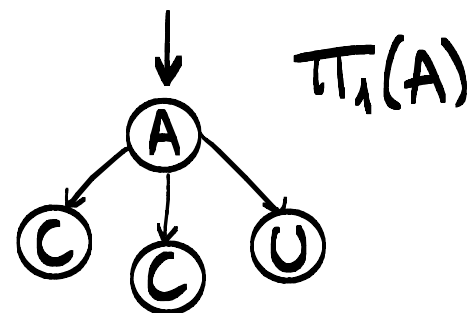
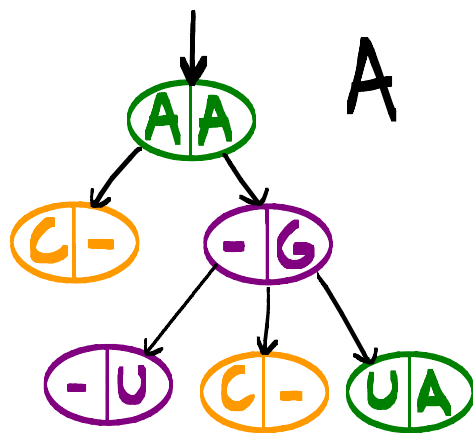
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

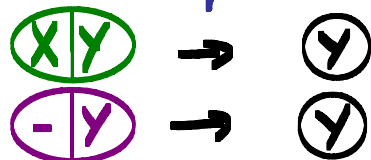
$\pi_1(A)$ = tree
obtained by changing



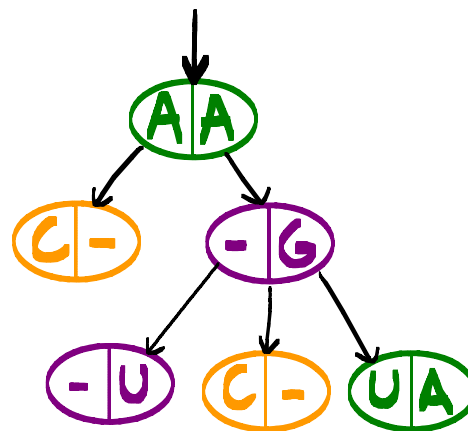
and removing -|Y .



$\pi_2(A)$ = tree
obtained by changing



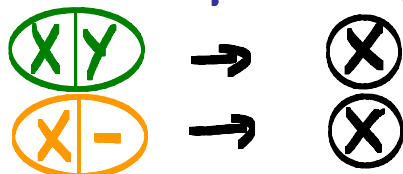
and removing X|- .



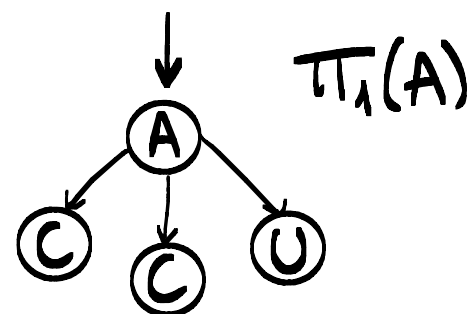
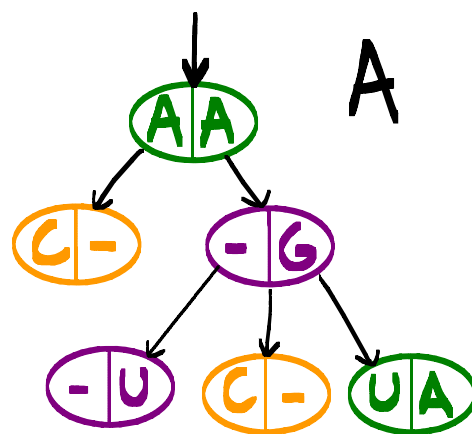
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

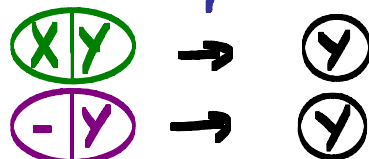
$\pi_1(A)$ = tree
obtained by changing



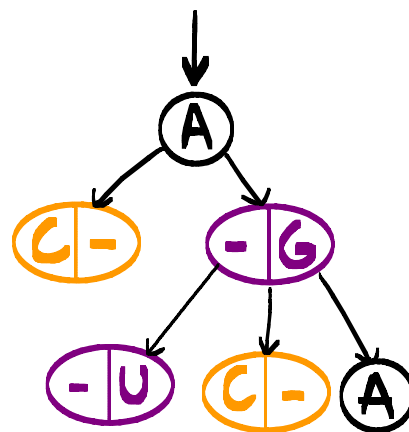
and removing -|Y .



$\pi_2(A)$ = tree
obtained by changing



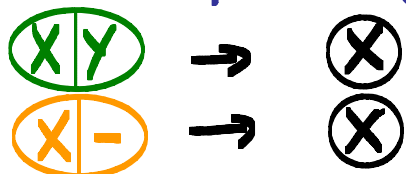
and removing X|- .



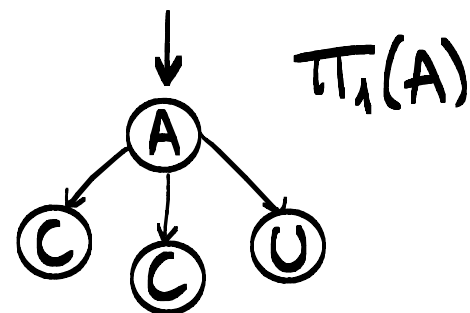
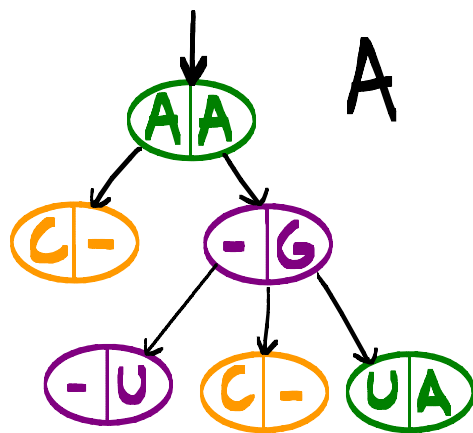
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

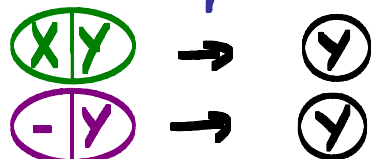
$\pi_1(A)$ = tree
obtained by changing



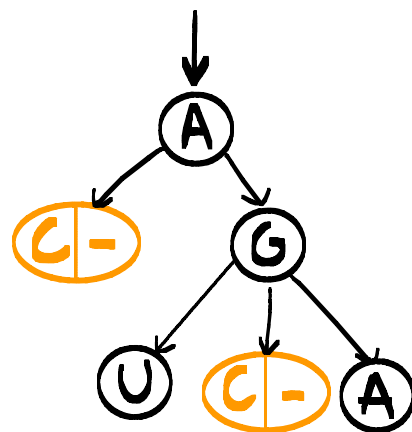
and removing -Y .



$\pi_2(A)$ = tree
obtained by changing



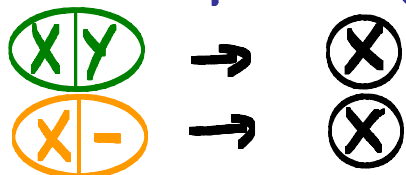
and removing X- .



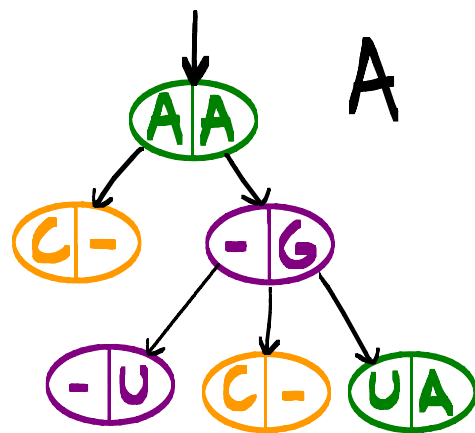
SUPERTREES INDUCE TREE ALIGNMENTS

Let A be a supertree,

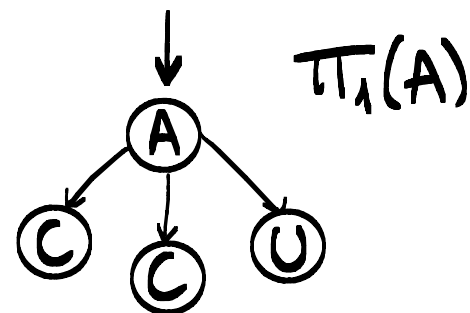
$\pi_1(A)$ = tree
obtained by changing



and removing -|Y .

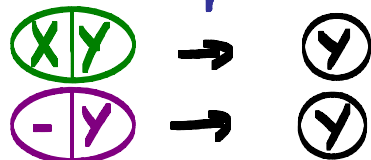


A

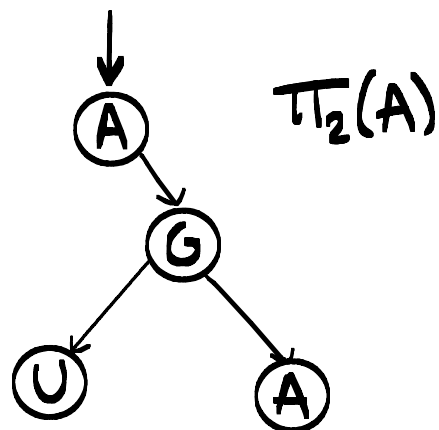


$\pi_1(A)$

$\pi_2(A)$ = tree
obtained by changing

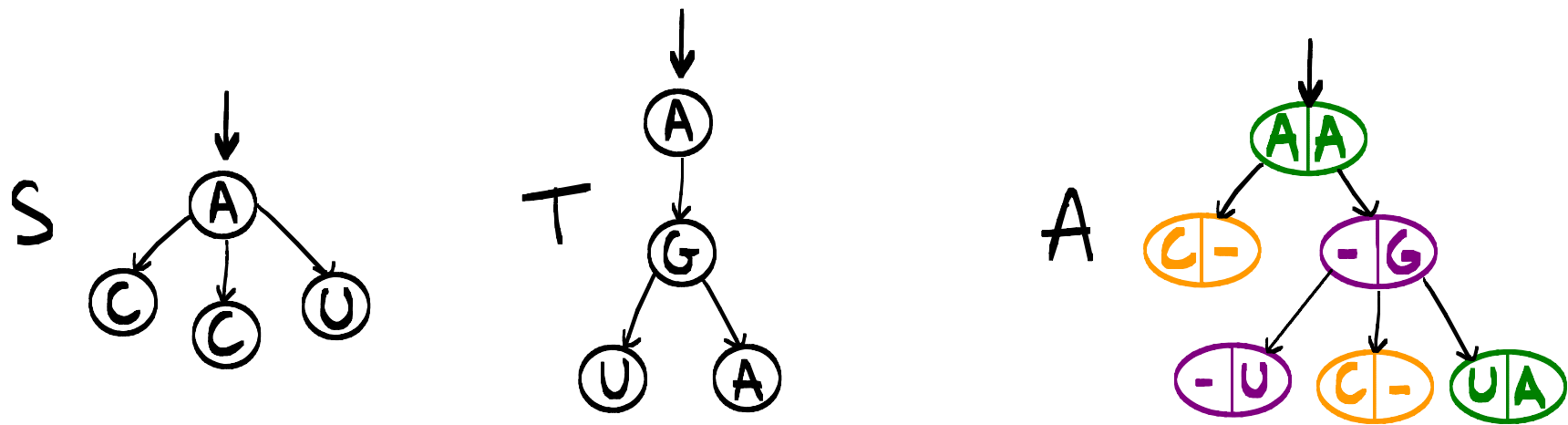


and removing X|- .



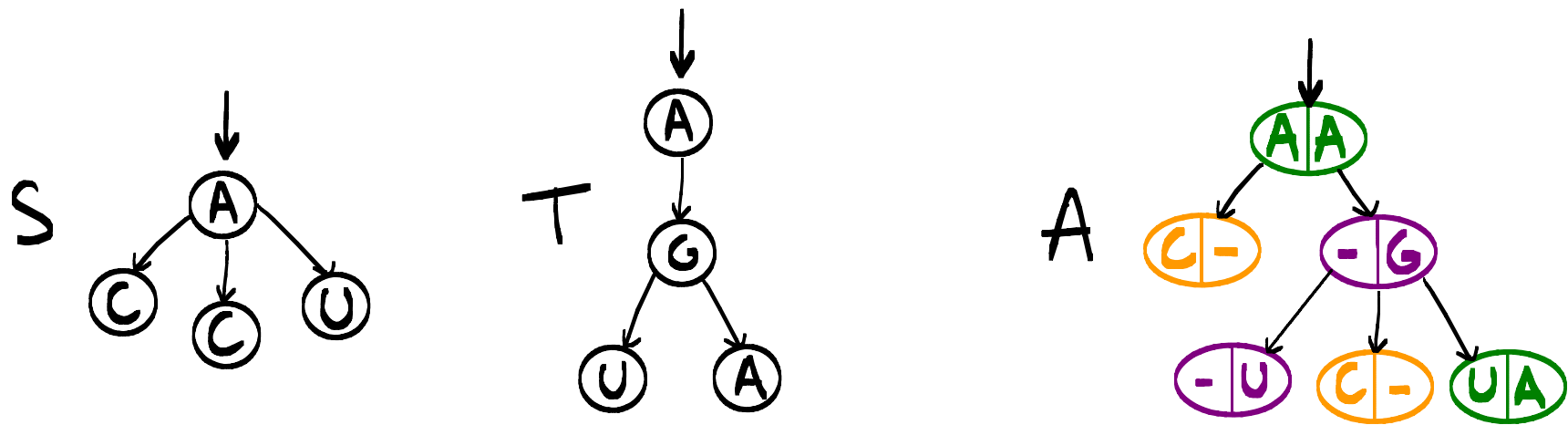
$\pi_2(A)$

SUPERTREES INDUCE TREE ALIGNMENTS



Given two trees S and T ,
a supertree A defines an alignment between S and T
if $\pi_1(A) = S$ and $\pi_2(A) = T$.

SUPERTREES INDUCE TREE ALIGNMENTS



Given two trees S and T ,
a supertree A defines an alignment between S and T
if $\pi_1(A) = S$ and $\pi_2(A) = T$.

$\text{cost}(A) = \text{nb of insertions} + \text{deletions} + \text{mismatches}$
(can be changed more complicated models)

CONNECTION WITH SEQUENCE ALIGNMENTS

Tree alignments generalize sequence alignments.

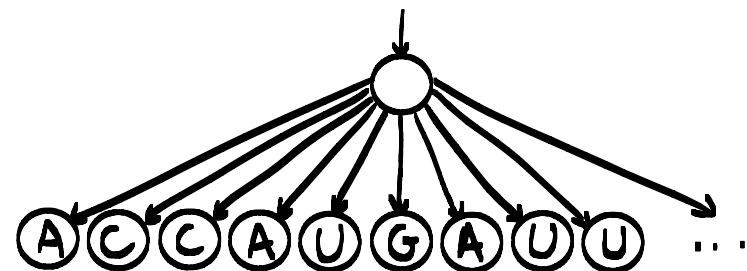
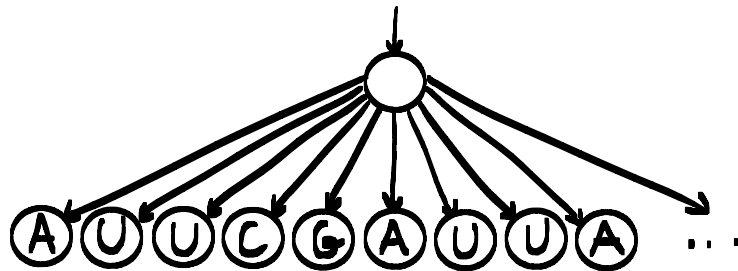
MCZWCQMS

AUUCGAUUA...

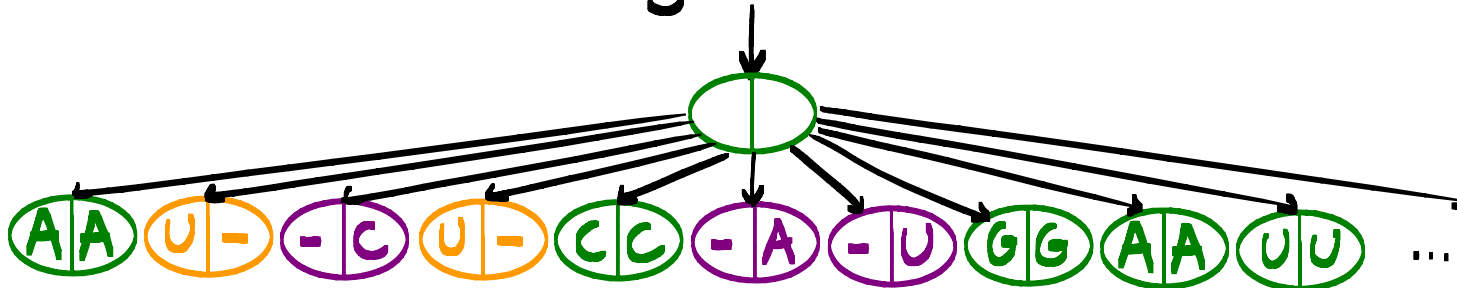
ACCAUGAUUA...

alignment:

(A)(U)(-)(U)(C)(-)(-)(G)(A)(U)(U)(A) ...
(A)(-)(C)(-)(C)(A)(U)(G)(A)(U)(U)(A) ...



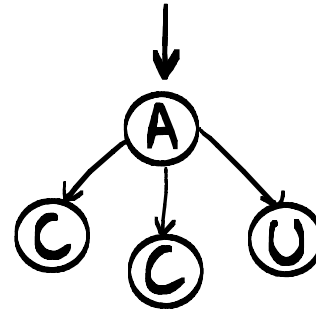
alignment:



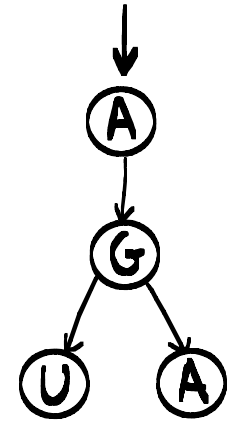
MMR-I

SPACE OF ALIGNMENTS

Which alignment between
is the most likely?

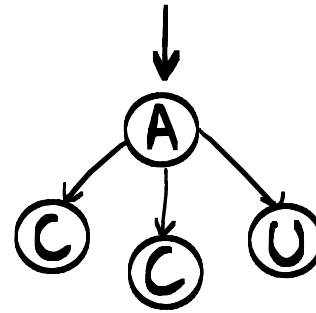


and

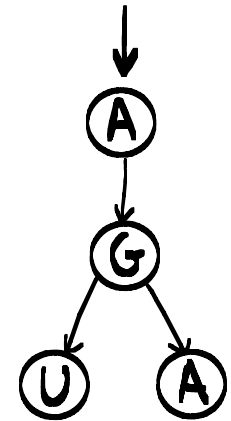


SPACE OF ALIGNMENTS

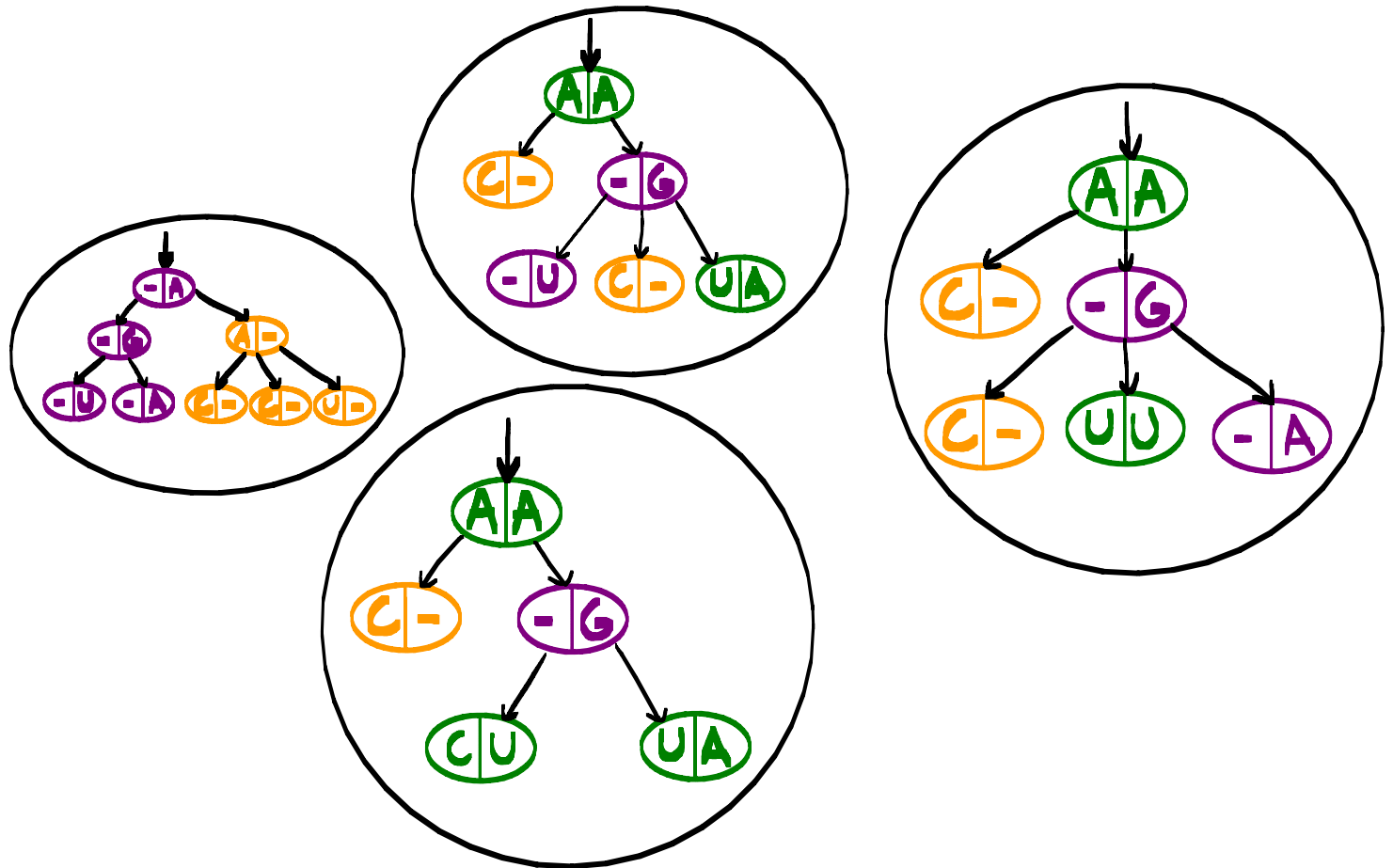
Which alignment between
is the most likely?



and

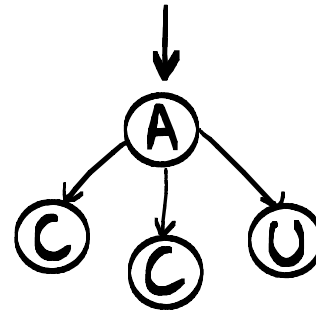


probability of
an alignment A
 $\propto e^{-\frac{\text{cost}(A)}{k}}$
(Gibbs-Boltzmann
distribution)

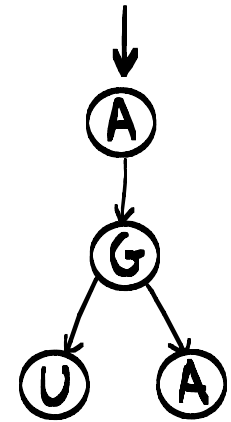


SPACE OF ALIGNMENTS

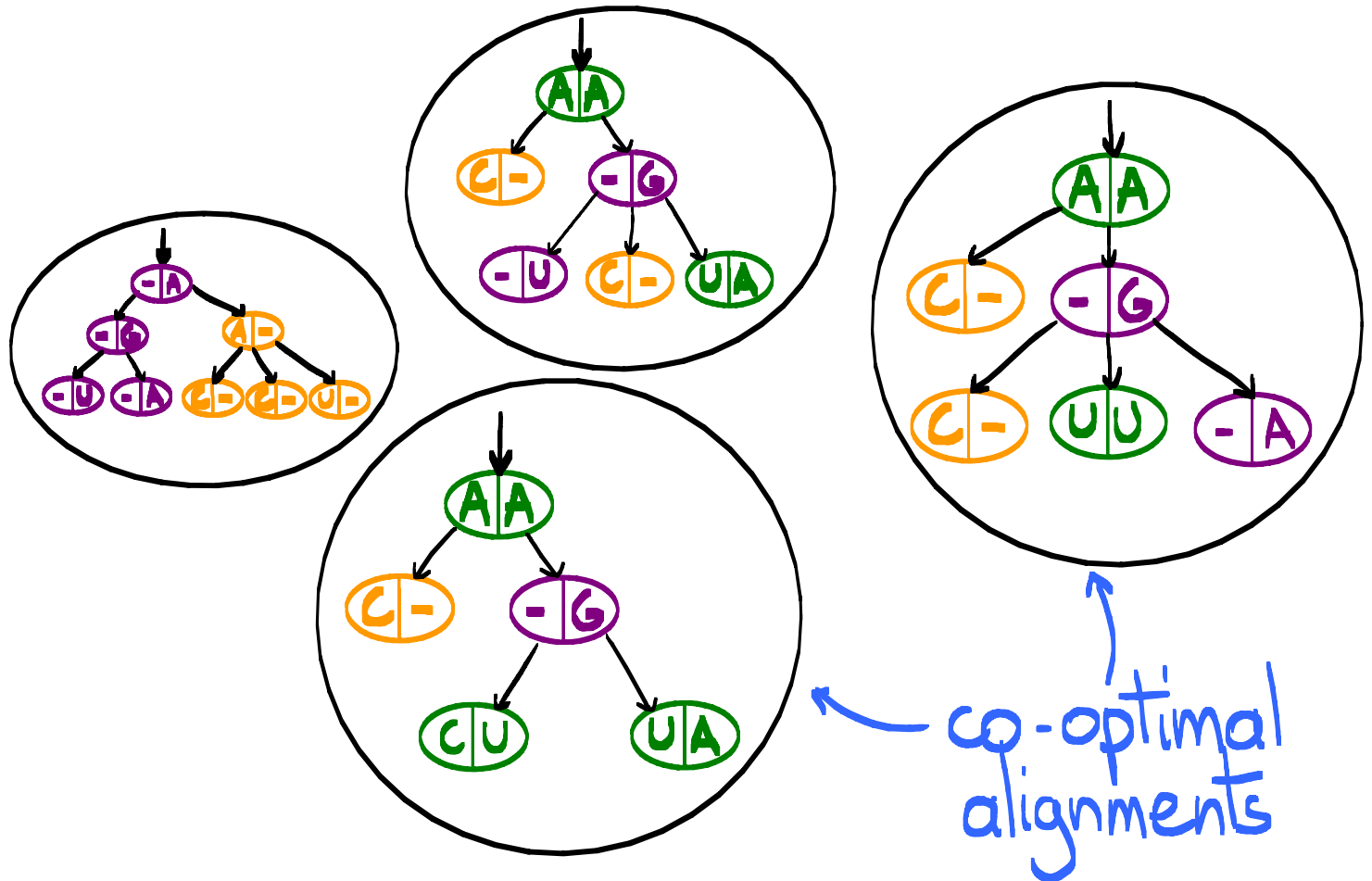
Which alignment between
is the most likely?



and



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



SPACE OF ALIGNMENTS

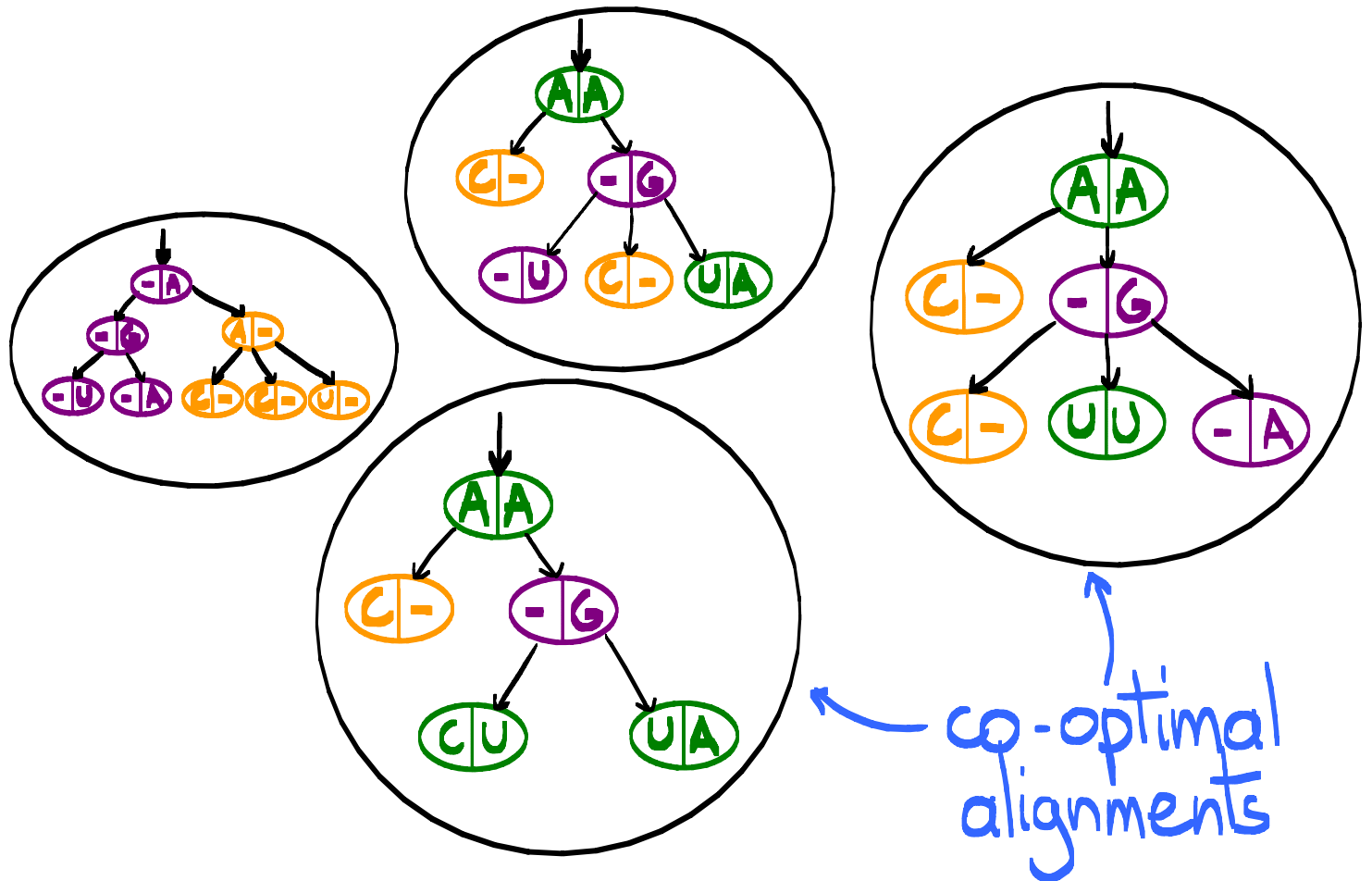
Why finding one optimal alignment may be inadequate:

- ▶ 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.

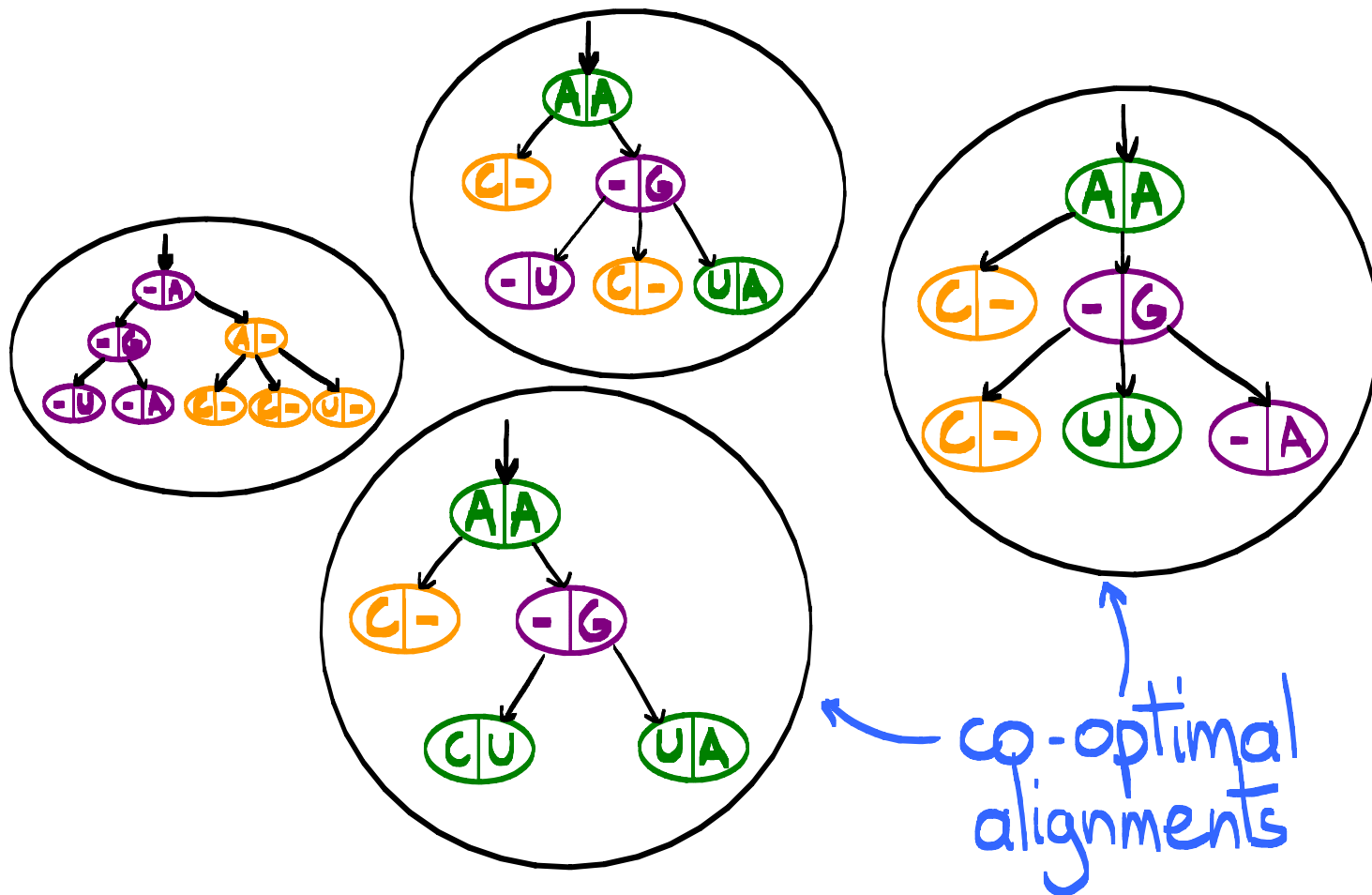
probability of an alignment A
 $\propto e^{-\frac{\text{cost}(A)}{k}}$
(Gibbs-Boltzmann distribution)



SPACE OF ALIGNMENTS

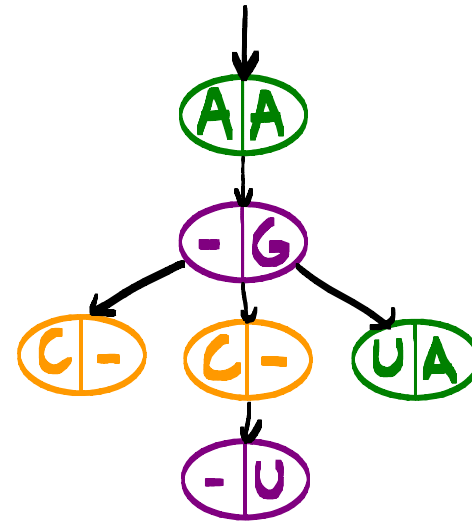
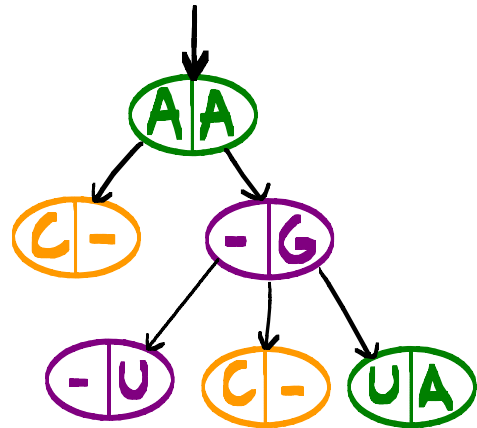
Objective: Sampling alignments under the Gibbs-Boltzmann probability distribution.

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

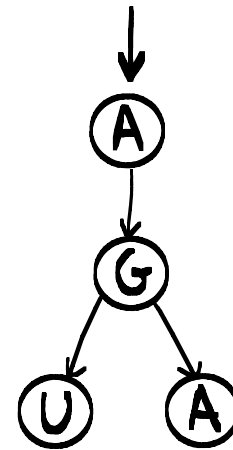
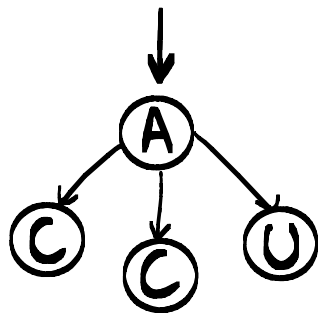


AMBIGUITY OF ALIGNMENTS

The two supertrees

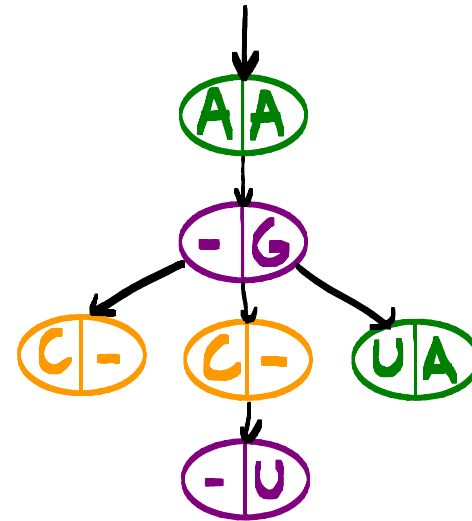
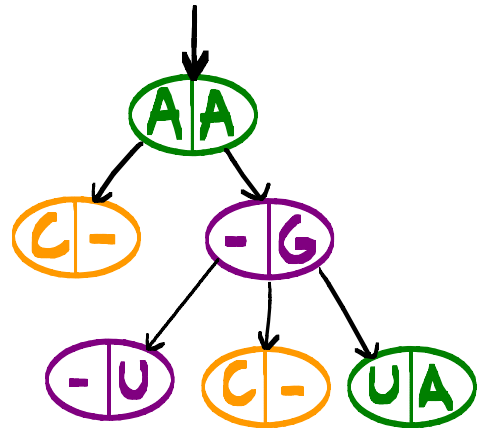


induce the same alignment between the trees

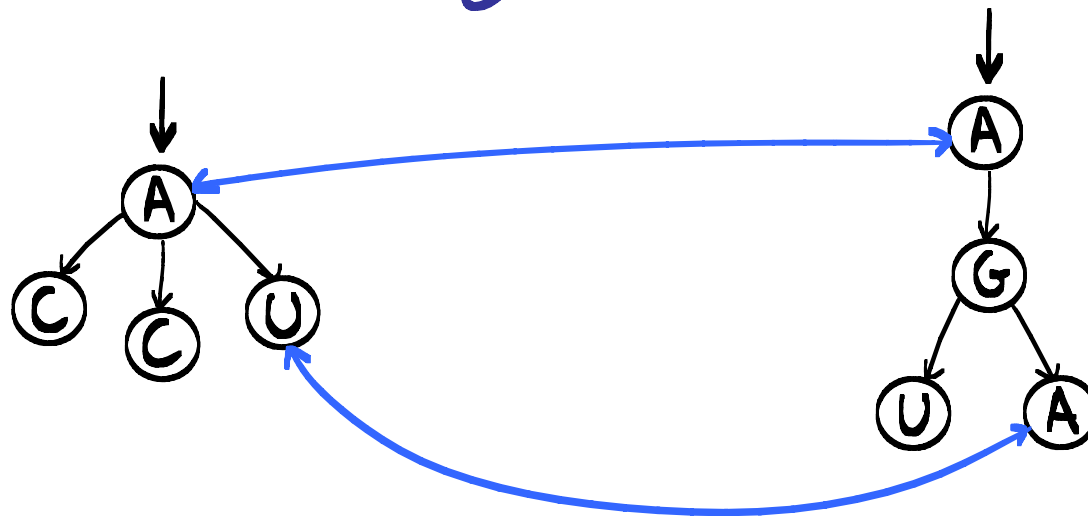


AMBIGUITY OF ALIGNMENTS

The two supertrees

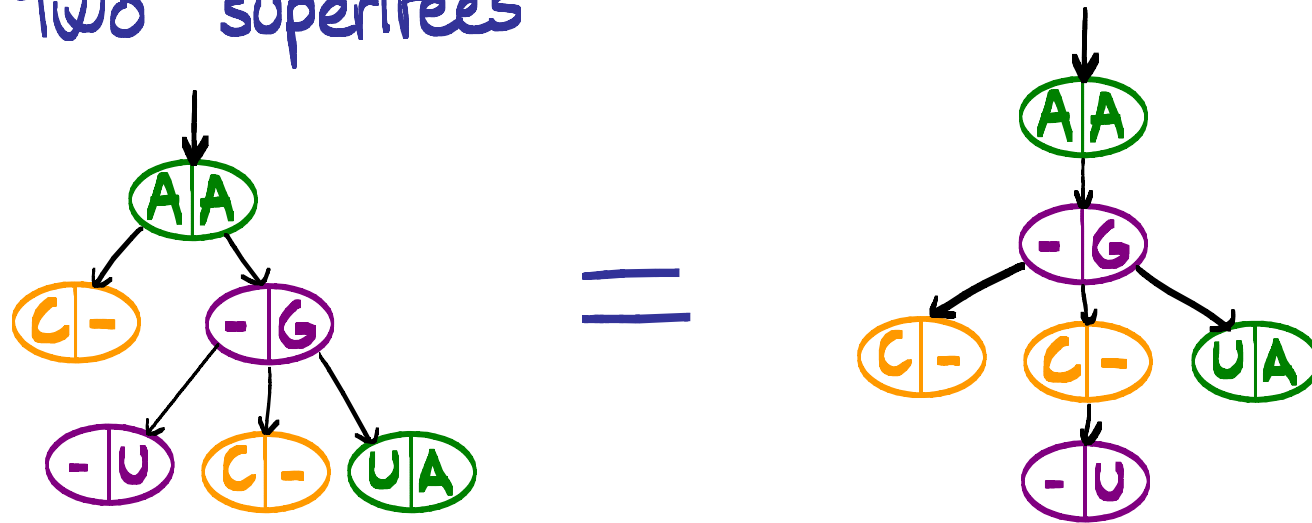


induce the same alignment between the trees

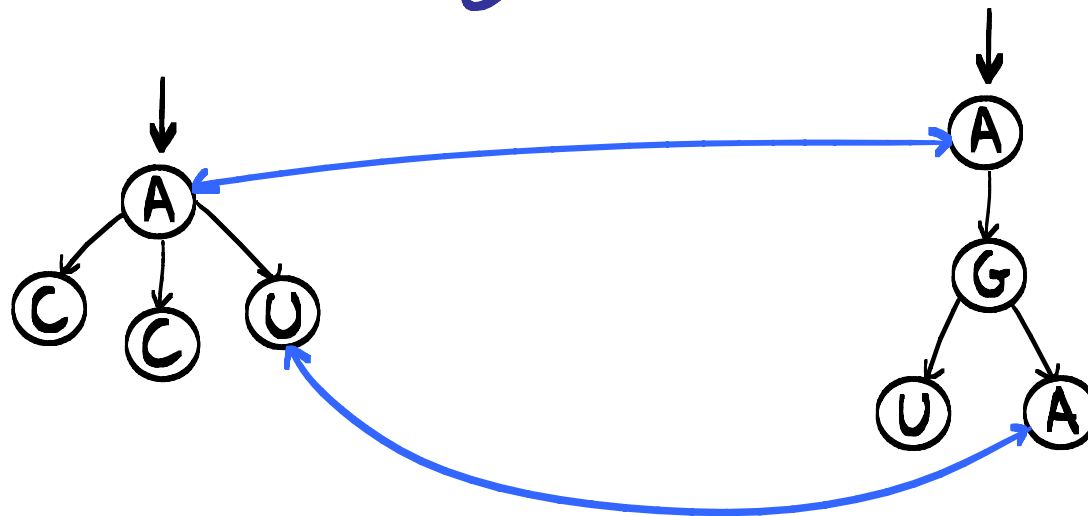


AMBIGUITY OF ALIGNMENTS

The two supertrees



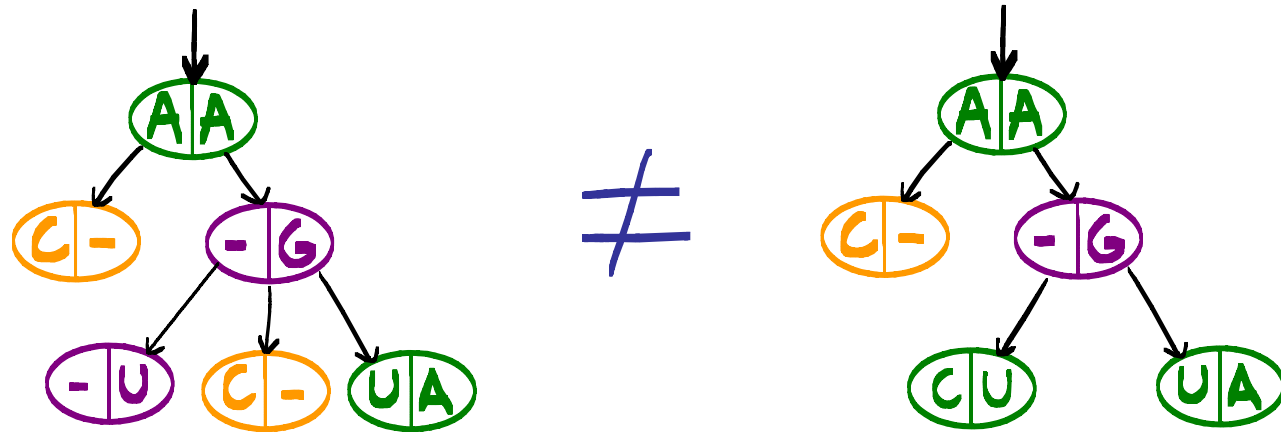
induce the same alignment between the trees



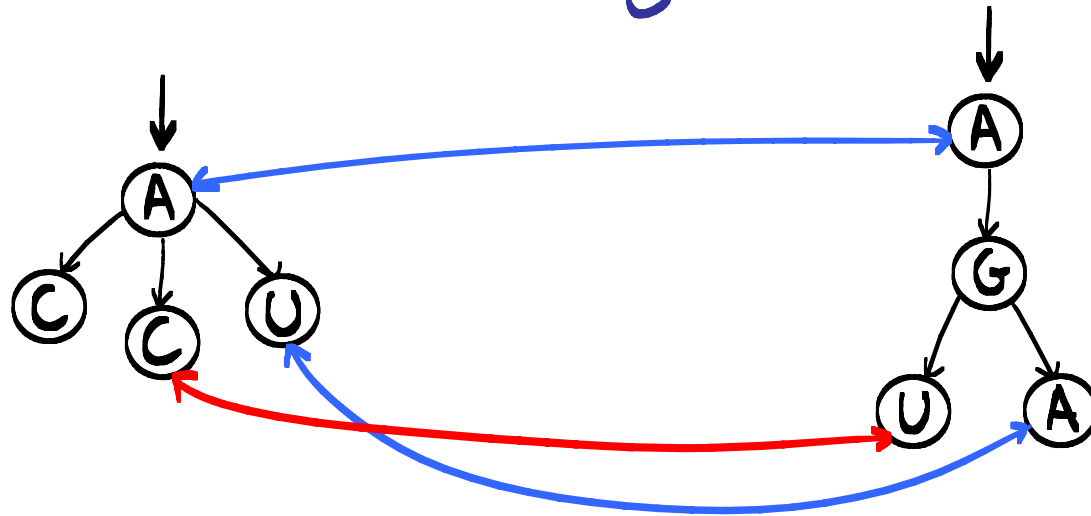
They are the same!

AMBIGUITY OF ALIGNMENTS

The two supertrees



do not induce the same alignment between the trees



A GRAMMAR FOR ALIGNMENTS

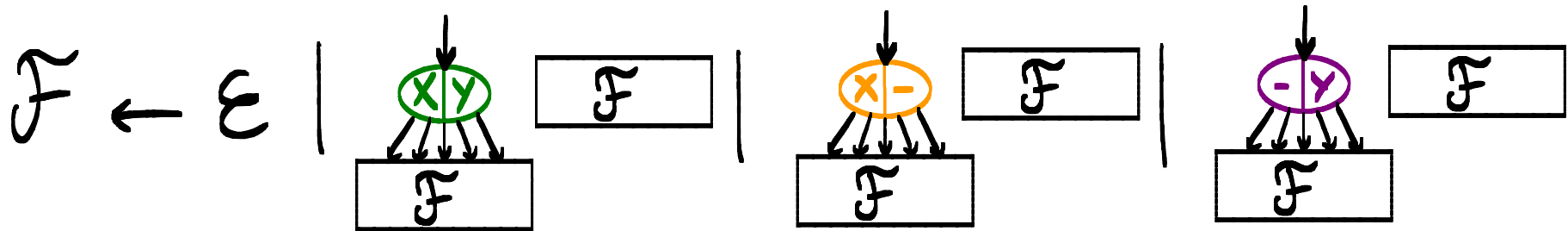
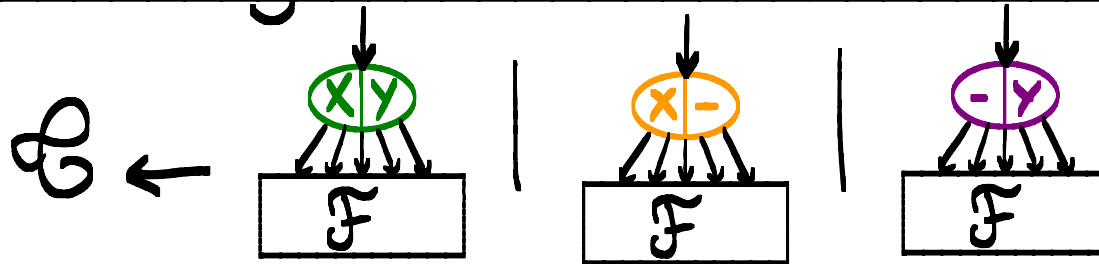
Strategy: Build a context-free grammar that generates every alignment exactly once

A GRAMMAR FOR ALIGNMENTS

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

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]

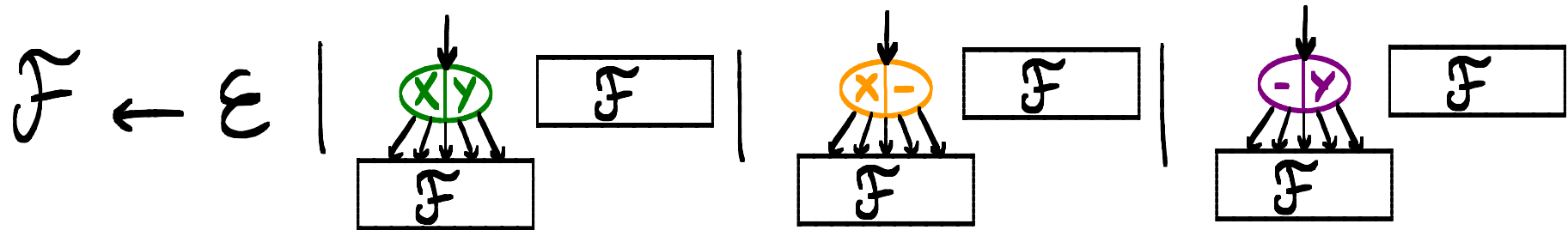
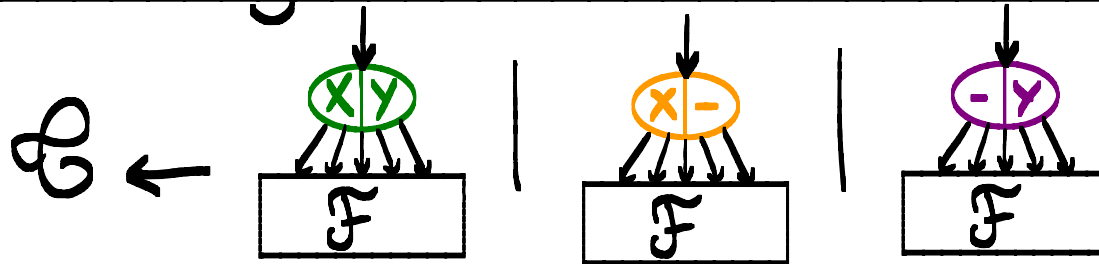


A GRAMMAR FOR ALIGNMENTS

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

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

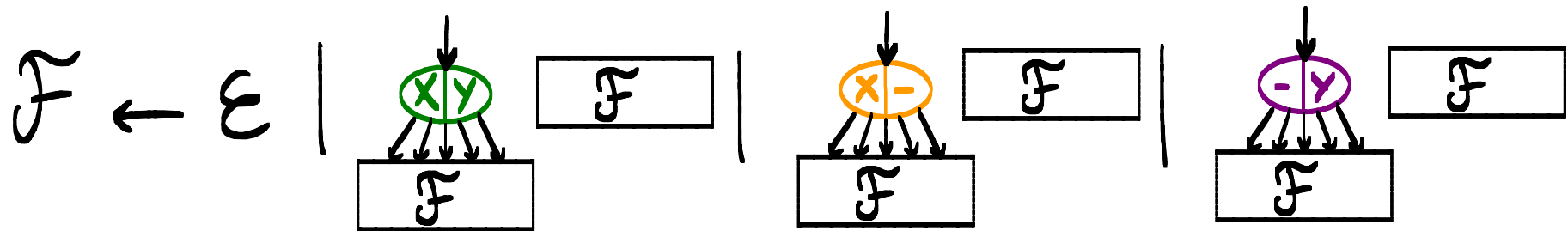
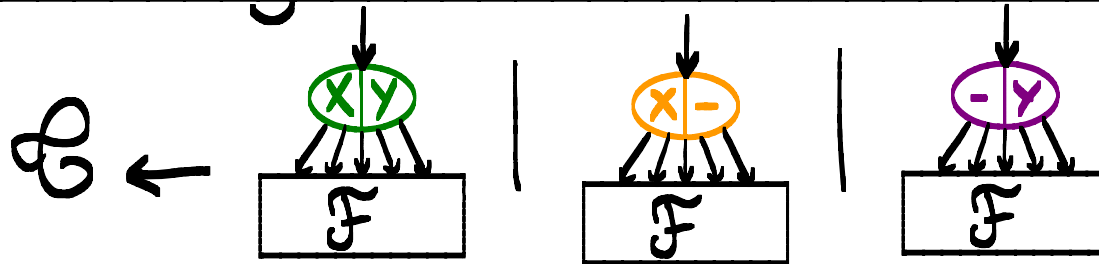


A GRAMMAR FOR ALIGNMENTS

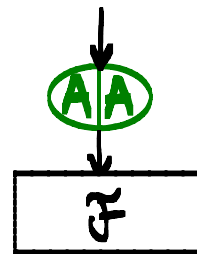
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

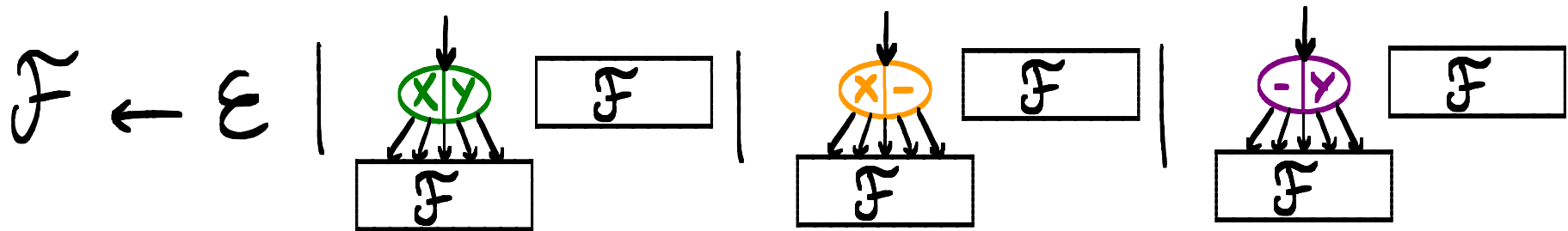
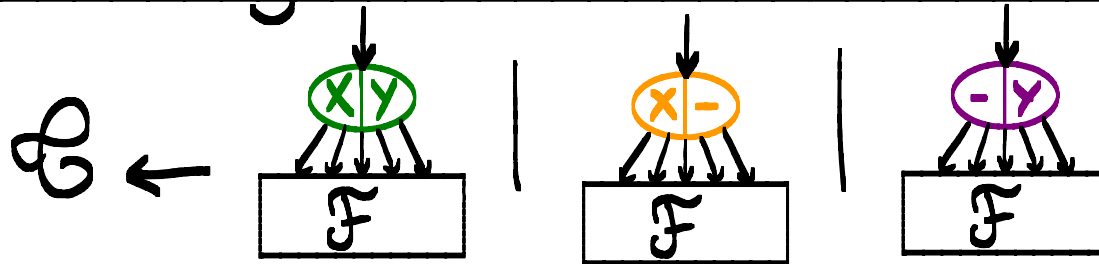


A GRAMMAR FOR ALIGNMENTS

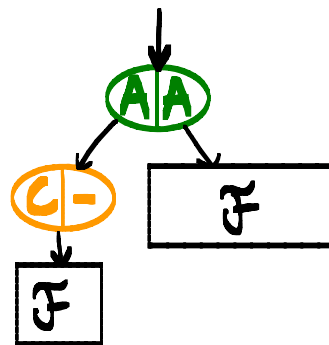
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

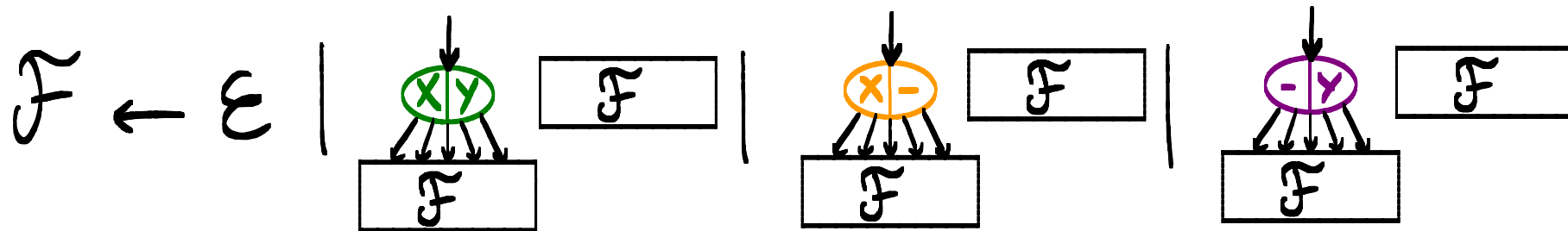
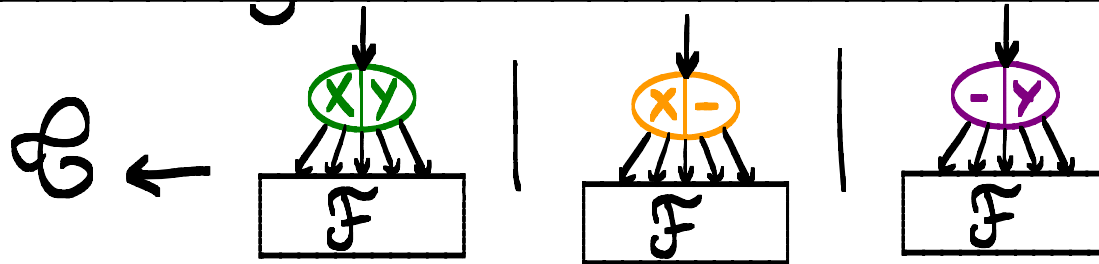


A GRAMMAR FOR ALIGNMENTS

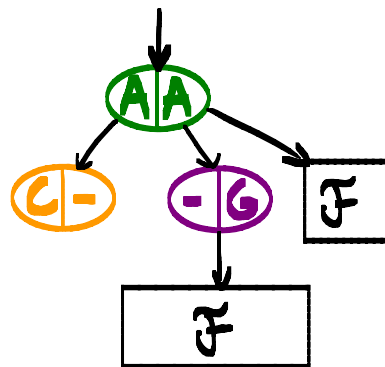
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

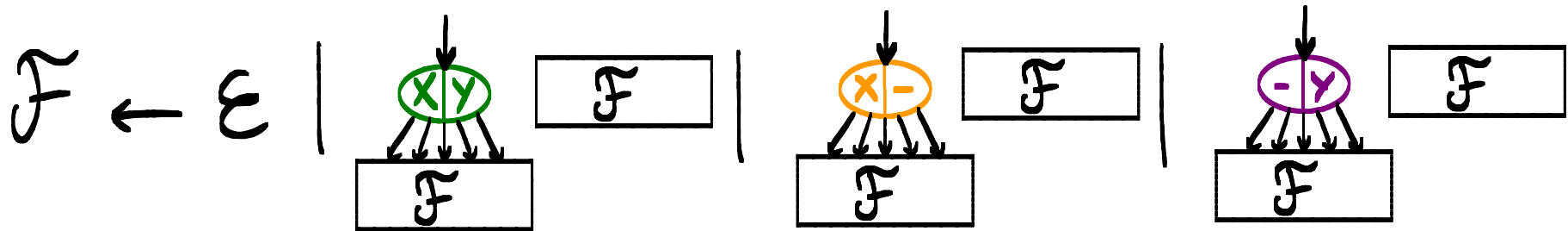
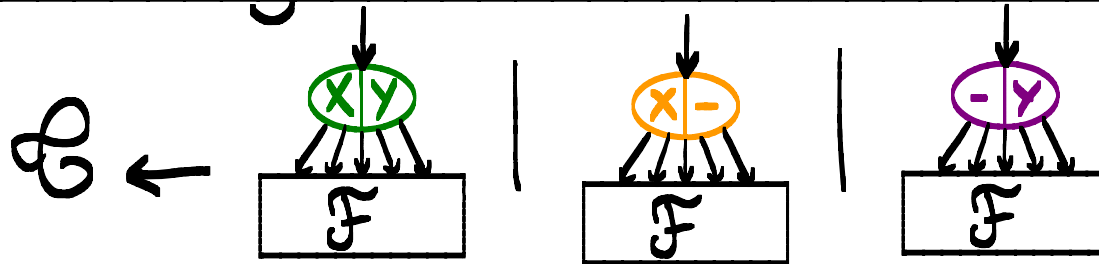


A GRAMMAR FOR ALIGNMENTS

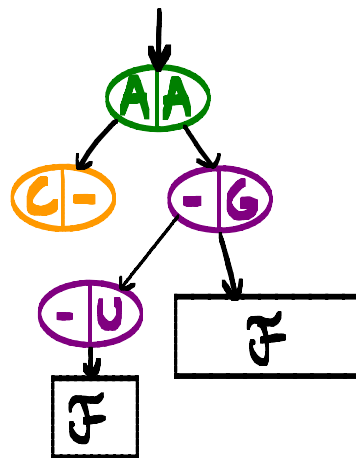
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

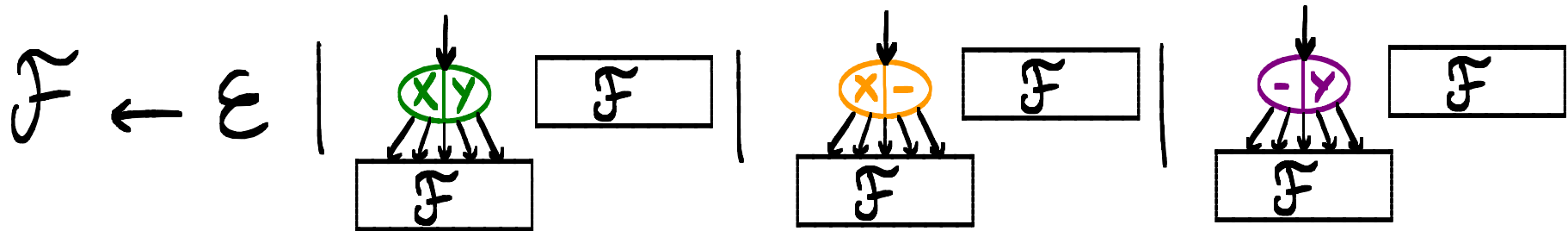
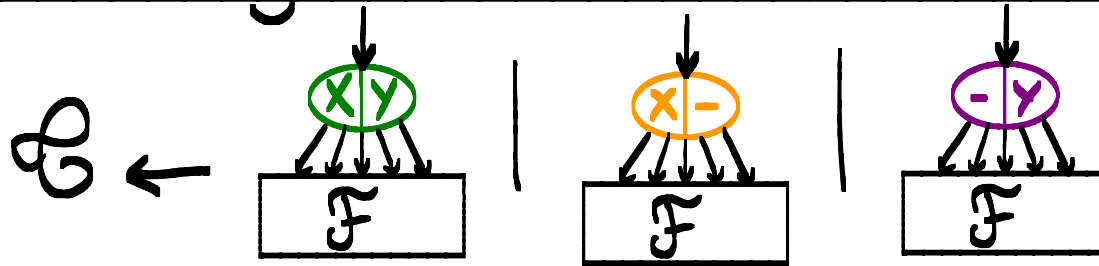


A GRAMMAR FOR ALIGNMENTS

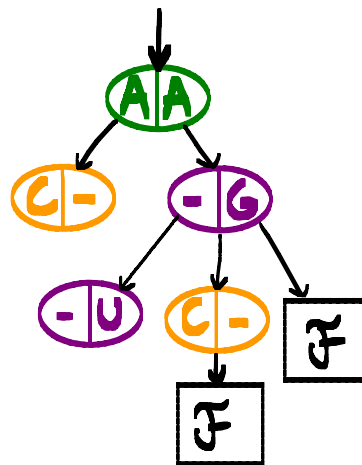
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

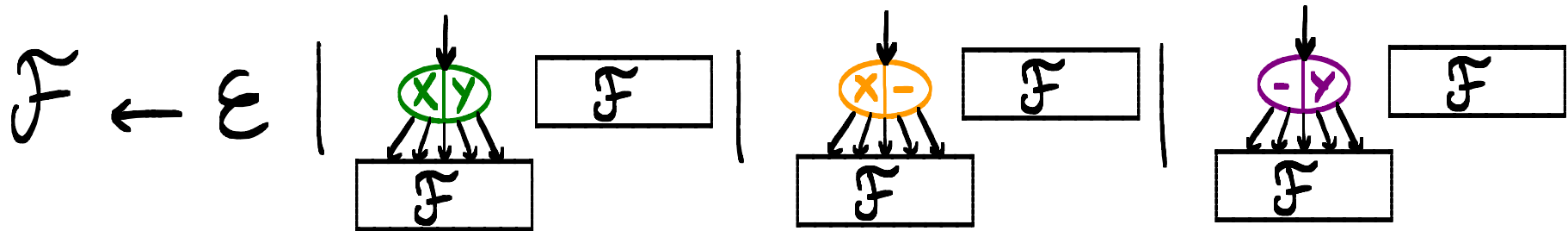
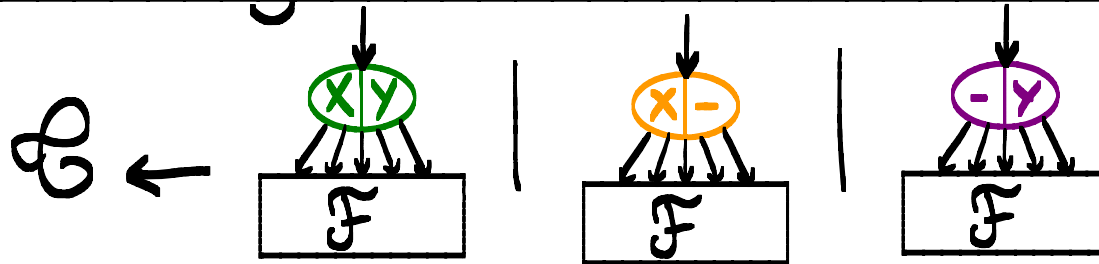


A GRAMMAR FOR ALIGNMENTS

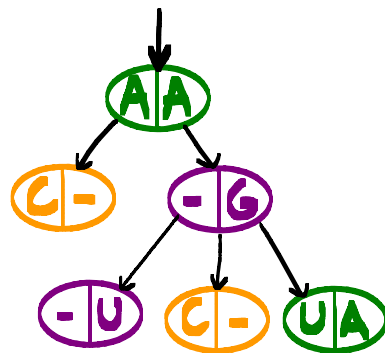
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

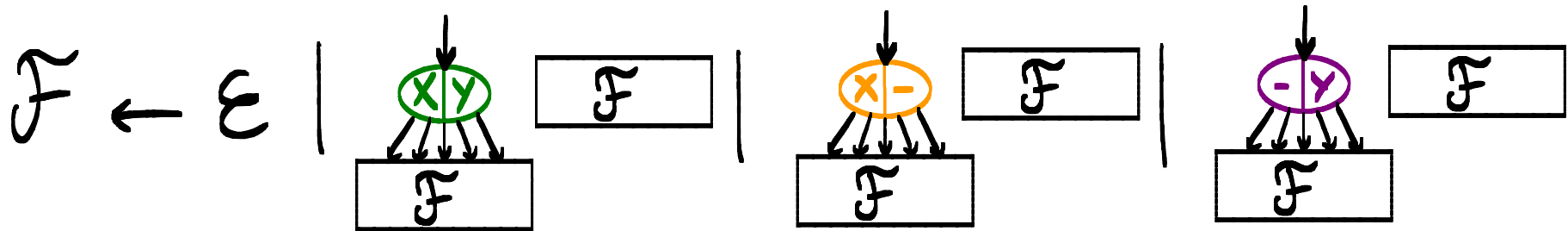
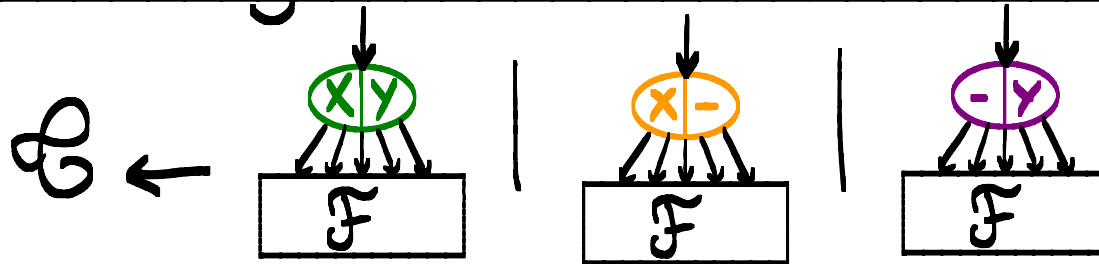


A GRAMMAR FOR ALIGNMENTS

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

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

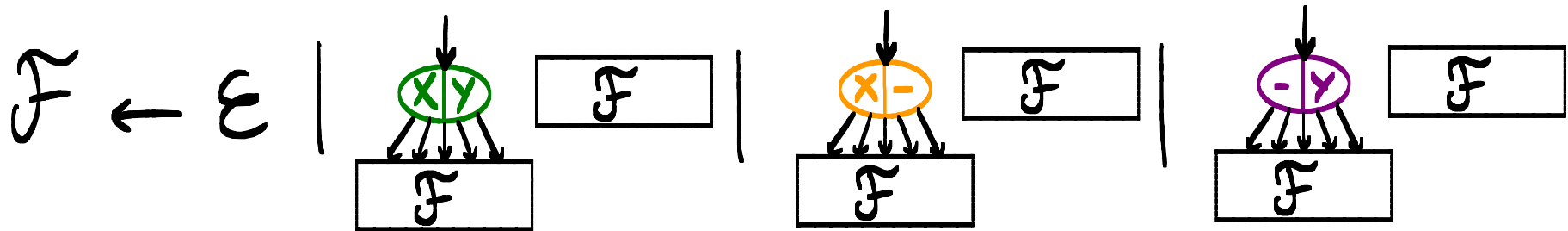
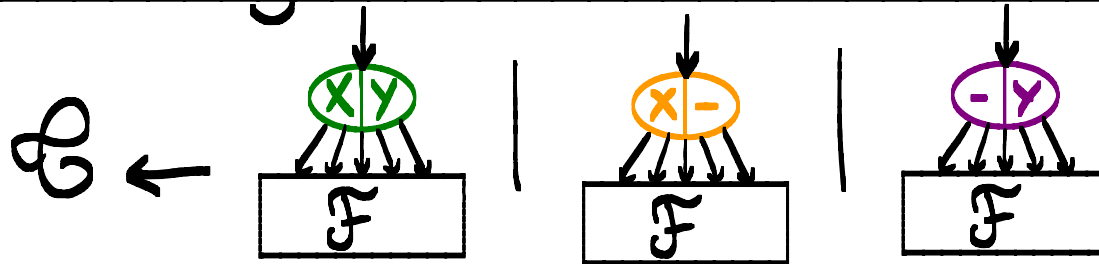


A GRAMMAR FOR ALIGNMENTS

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

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

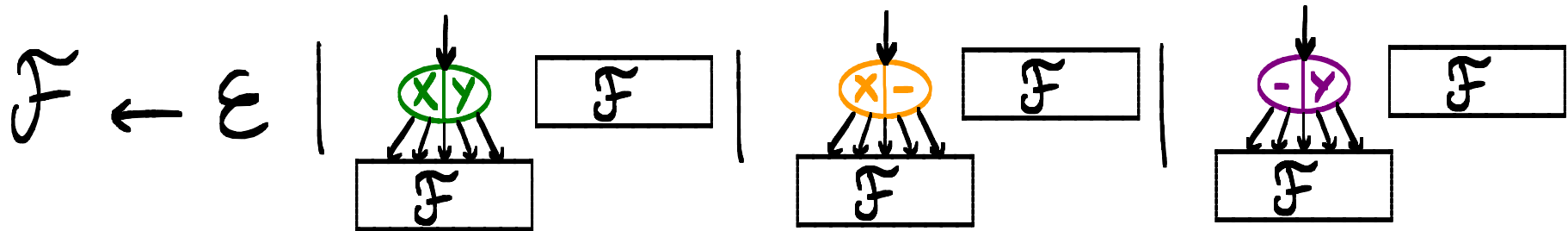
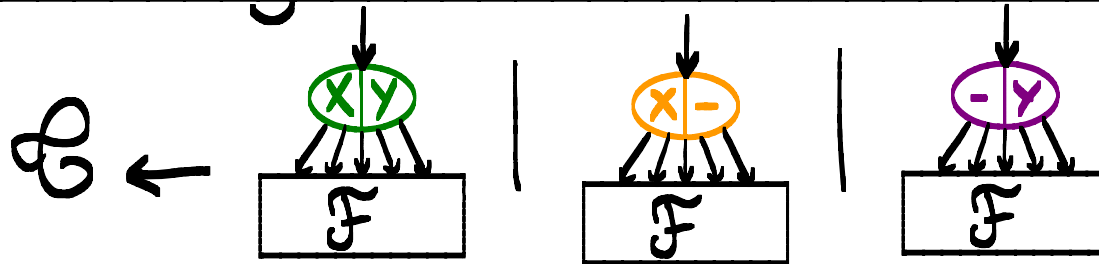


A GRAMMAR FOR ALIGNMENTS

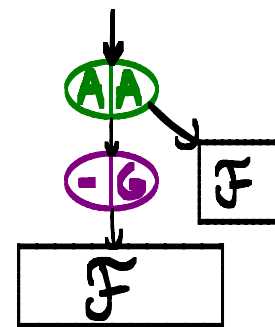
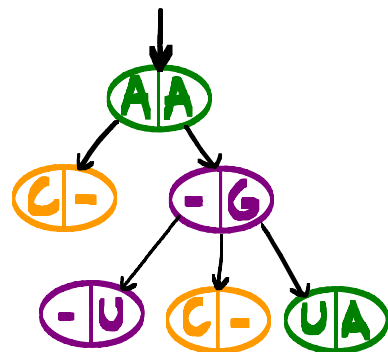
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

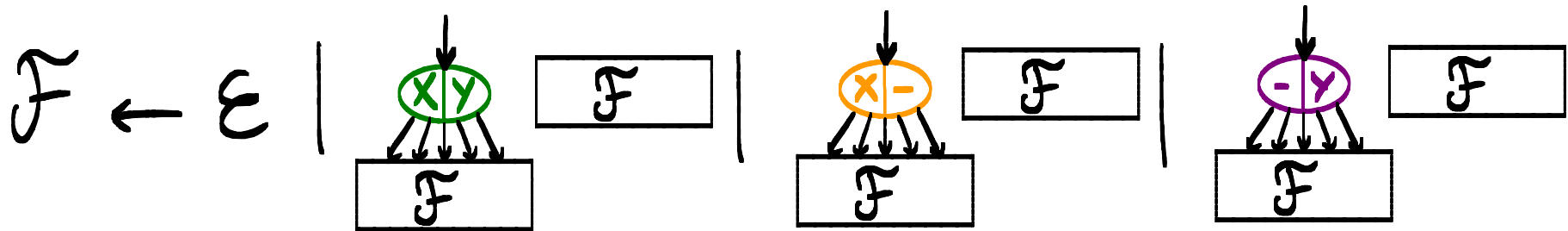
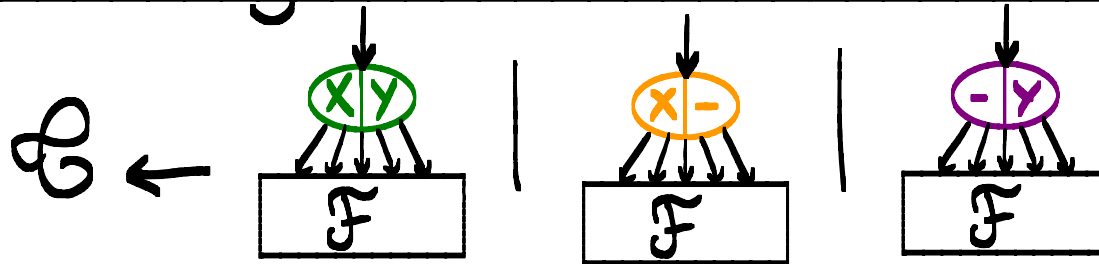


A GRAMMAR FOR ALIGNMENTS

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

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

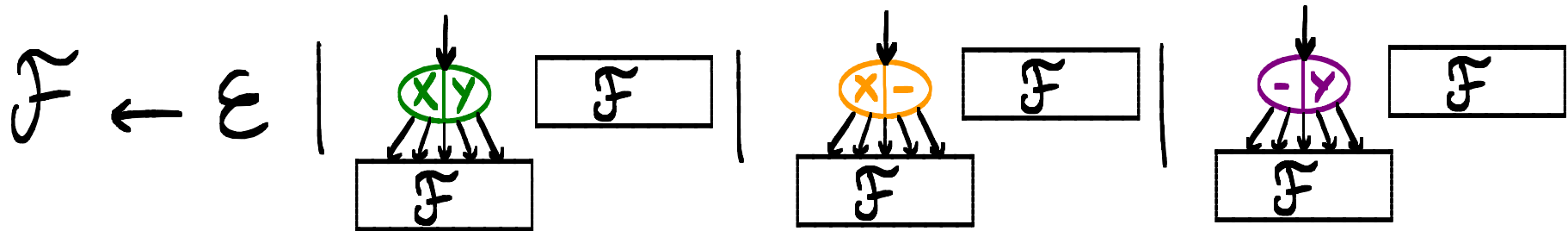
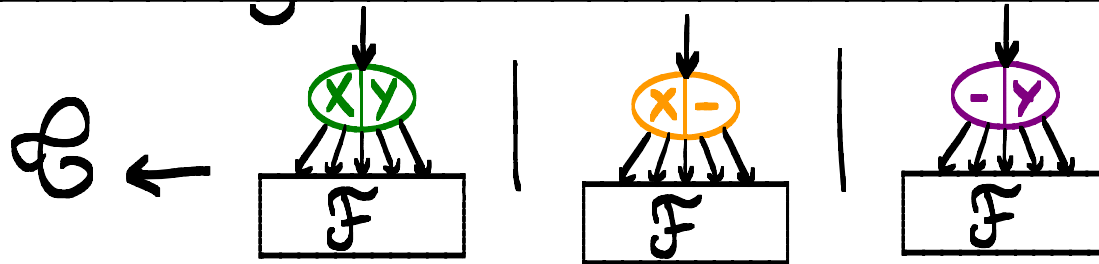


A GRAMMAR FOR ALIGNMENTS

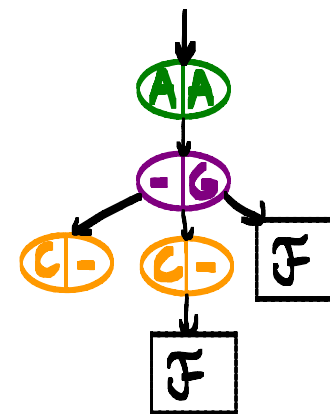
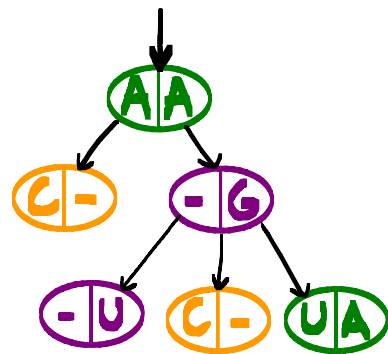
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

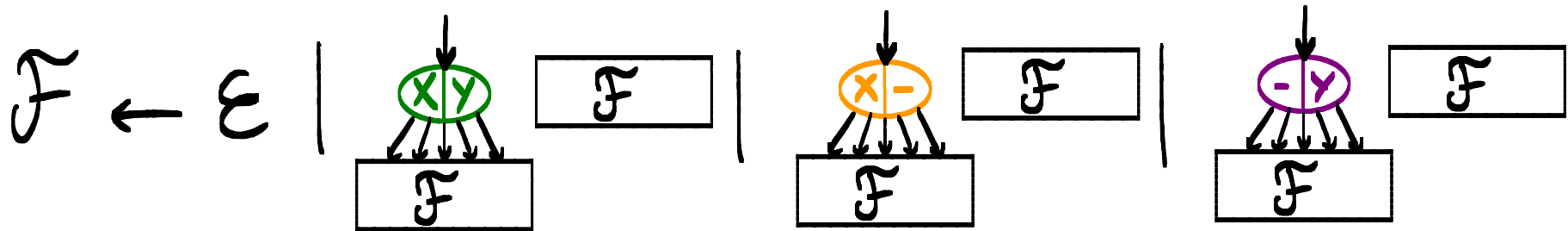
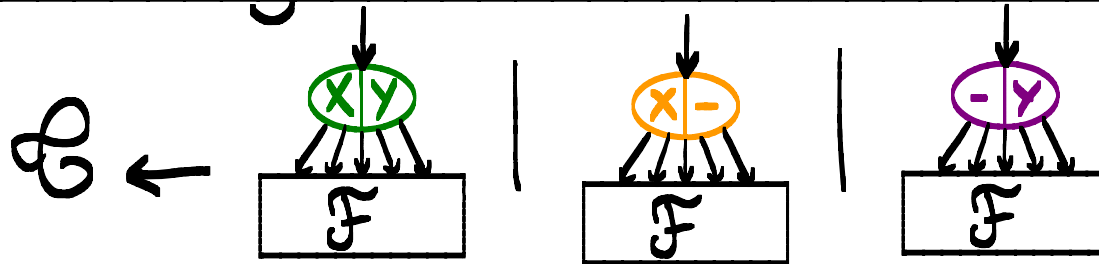


A GRAMMAR FOR ALIGNMENTS

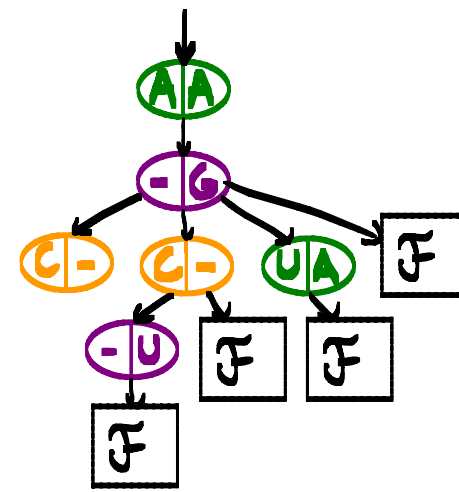
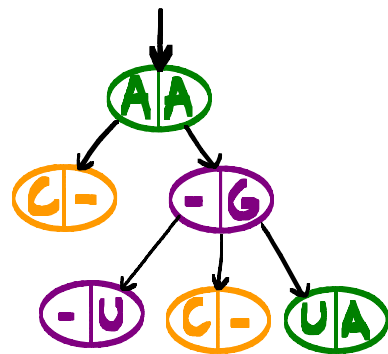
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

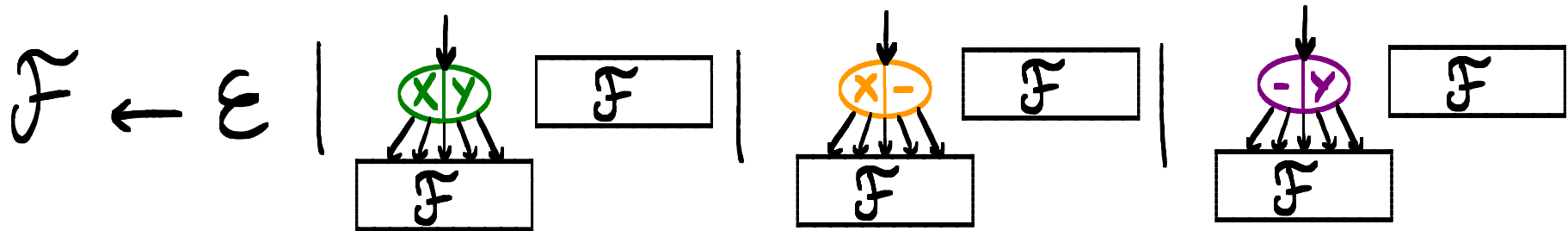
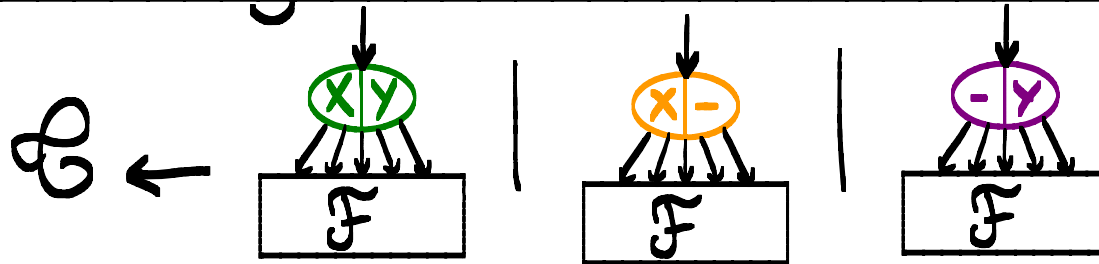


A GRAMMAR FOR ALIGNMENTS

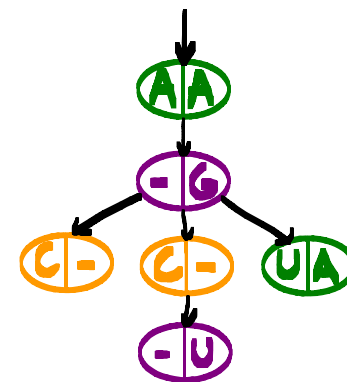
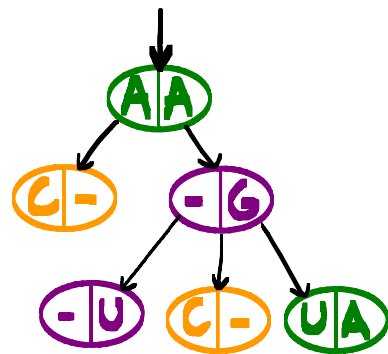
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:

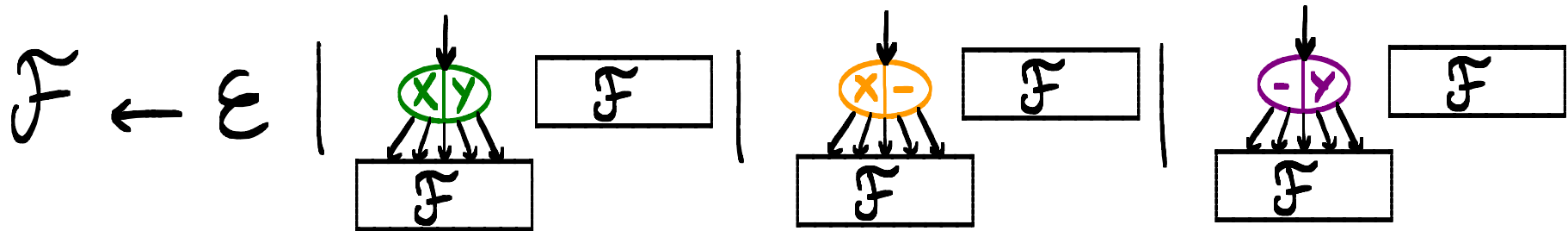
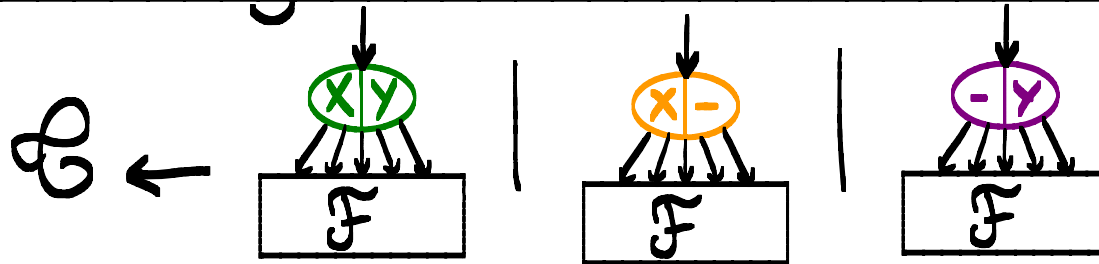


A GRAMMAR FOR ALIGNMENTS

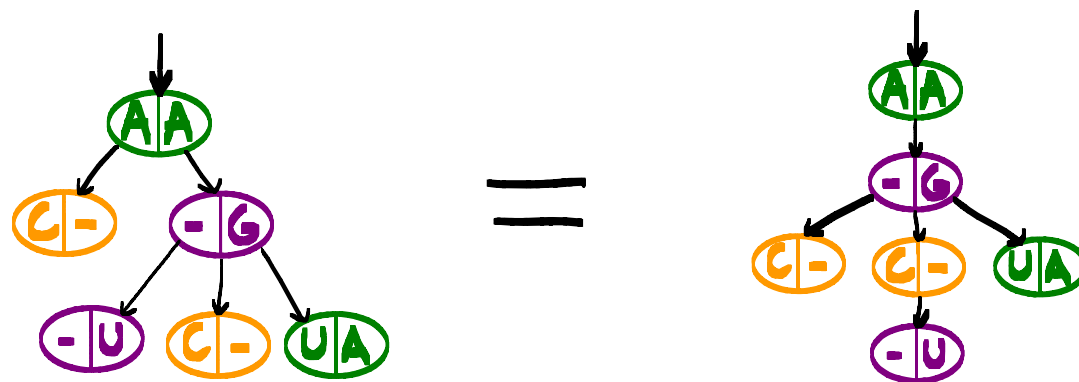
Strategy: Build a context-free grammar that generates every alignment exactly once

An example of grammar that does not work:

[Jiang,
Wang,
Zhang]



Ex:



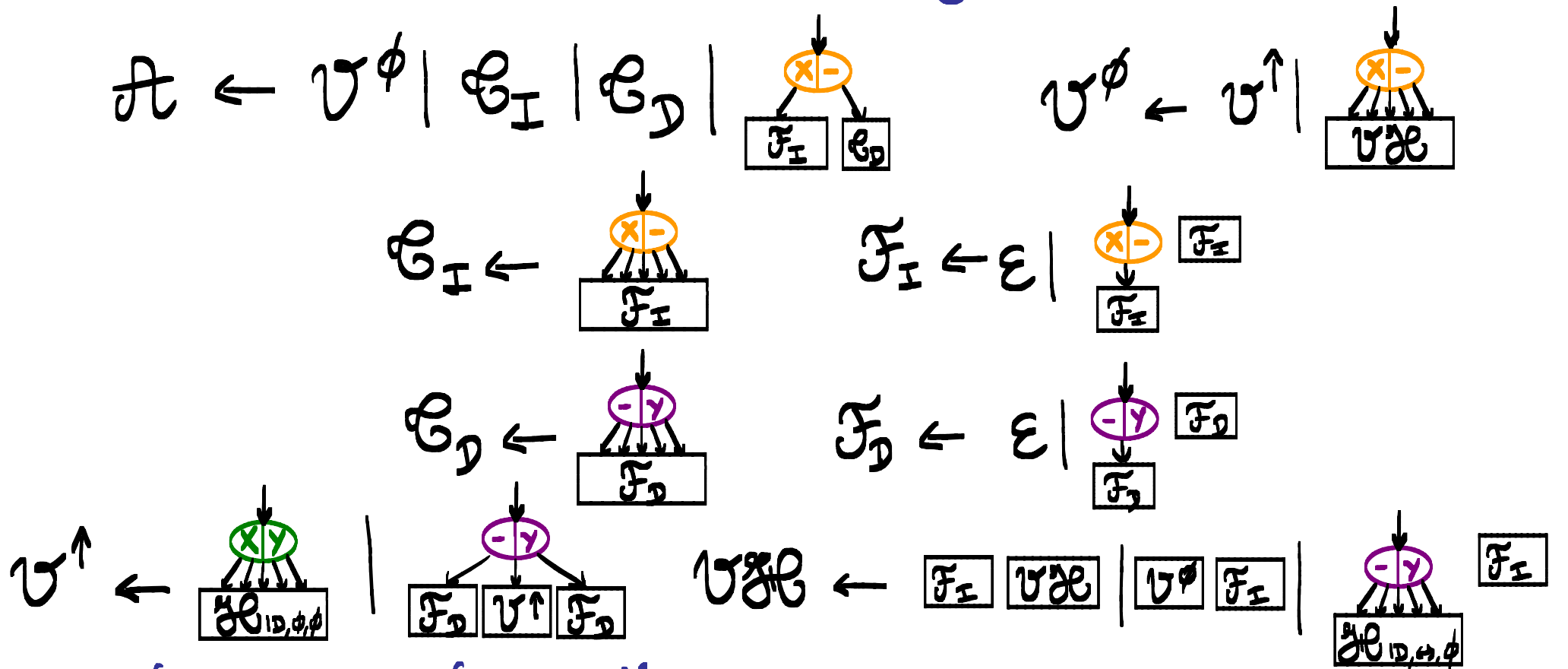
ambiguous
grammar

A GRAMMAR FOR ALIGNMENTS

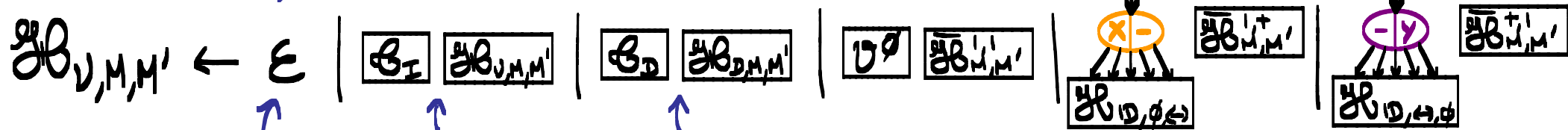
Theorem : The set \mathcal{A} generated by the following grammar contains every tree alignment exactly once.

A GRAMMAR FOR ALIGNMENTS

Our (complicated) non-ambiguous grammar:



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



only if $(M, M') = (\phi, \phi)$ only if $V \neq D$ and $M \neq \Leftrightarrow$ only if $M' \neq \Leftrightarrow$

no room for $B_{M',M'}^{i,j}$...

SOME STATISTICAL PROPERTIES

Theorem There are on average
 $C \times 1.5^n$ alignments
between two random trees of cumulative size n
where $C = 0.299\dots$

Corollary: A same alignment was repeated
 $\sim 0.875 \times 1.412^n$ times on
average in the previous
ambiguous grammar.

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 $O(n_1 n_2)$.

Strategy:

- Filter the grammar to obtain a new grammar that only generates alignments between 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 $O(n_1 n_2)$.

Proof inspired by
[Herrbach, Denise, Dulucq]

CONCLUSION

→ We are using our grammar and adapted dynamic programming algorithms to revisit the 3D alignments of RNA structures.

→ more general method?

new way to design
dynamic programming algorithms?

