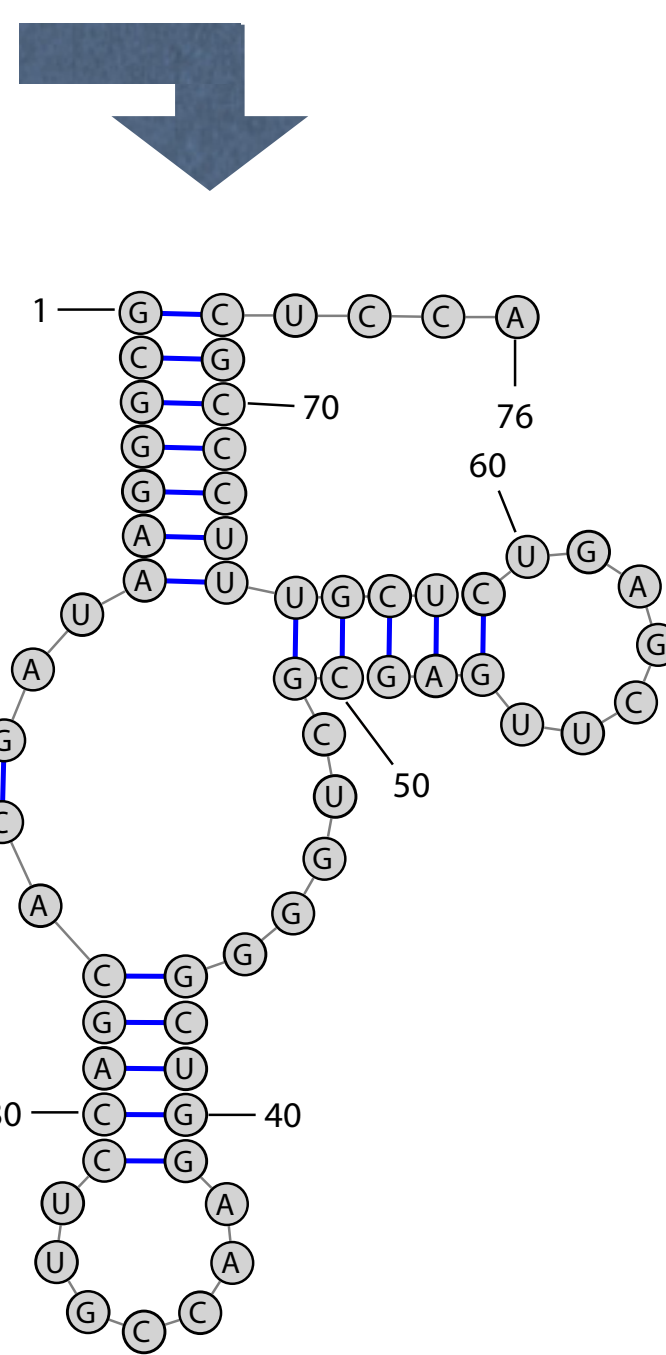


# LinearFold: Linear-Time Approximate RNA Folding by 5'-to-3' dynamic programming and beam search

$x$  GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA

$y$  (((((((((.....))))).((((.....))))).(((.....)))))).....((((.....)))))).....



**Liang Huang \***

Oregon State University & Baidu Research USA

Joint work with He Zhang \*\*, Dezhong Deng \*\*, Kai Zhao, Kaibo Liu, David Hendrix and David Mathews



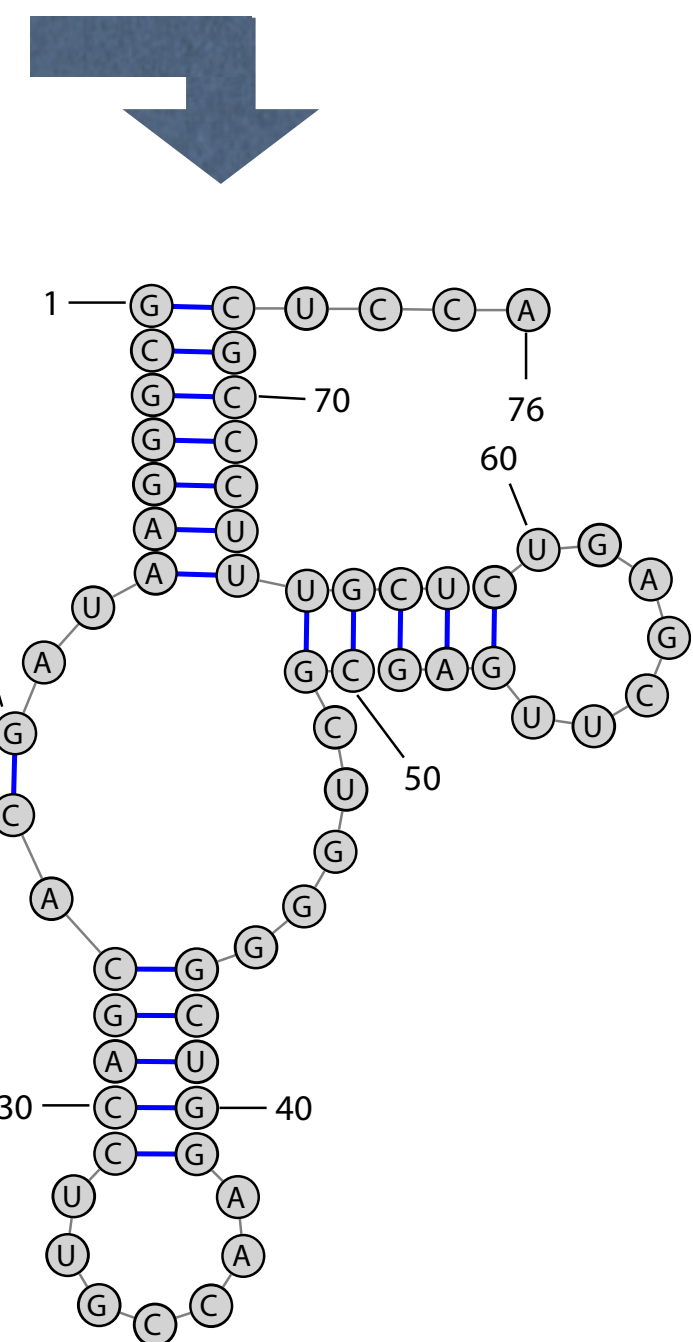
ISMB 2019 Proceedings Talk



\* corresponding author \*\* equal contribution

# LinearFold: Linear-Time Approximate RNA Folding by 5'-to-3' dynamic programming and beam search

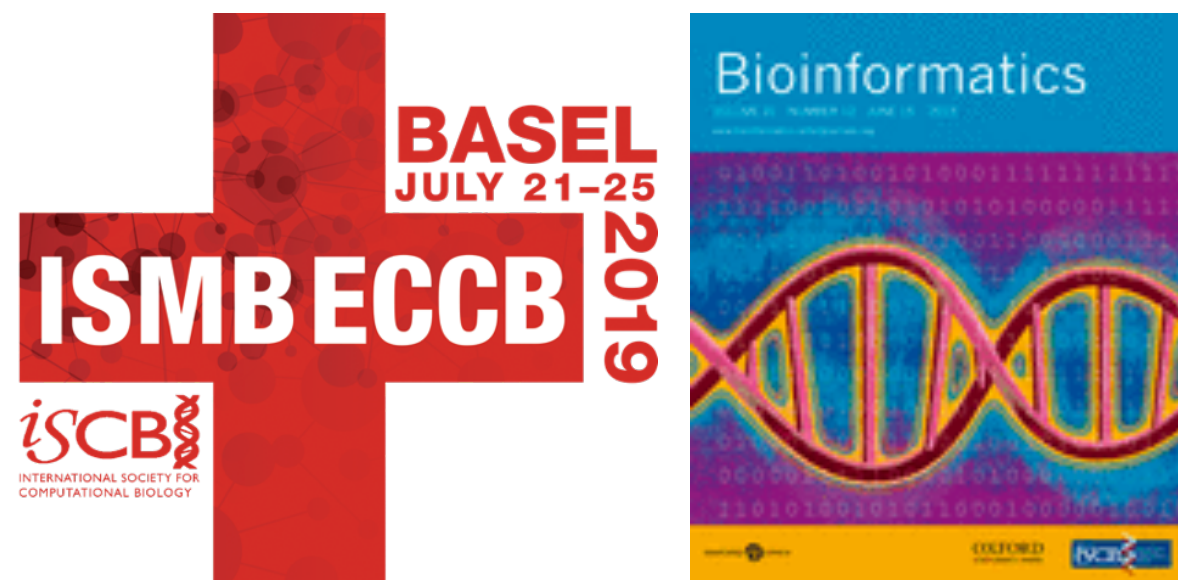
$x$  GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA  
 $y$  (((((((((..((((.....))))).((((.....))))).(((.....))))).((((.....)))))).....



**Liang Huang \***

Oregon State University & Baidu Research USA

Joint work with He Zhang \*\*, Dezhong Deng \*\*, Kai Zhao, Kaibo Liu, David Hendrix and David Mathews



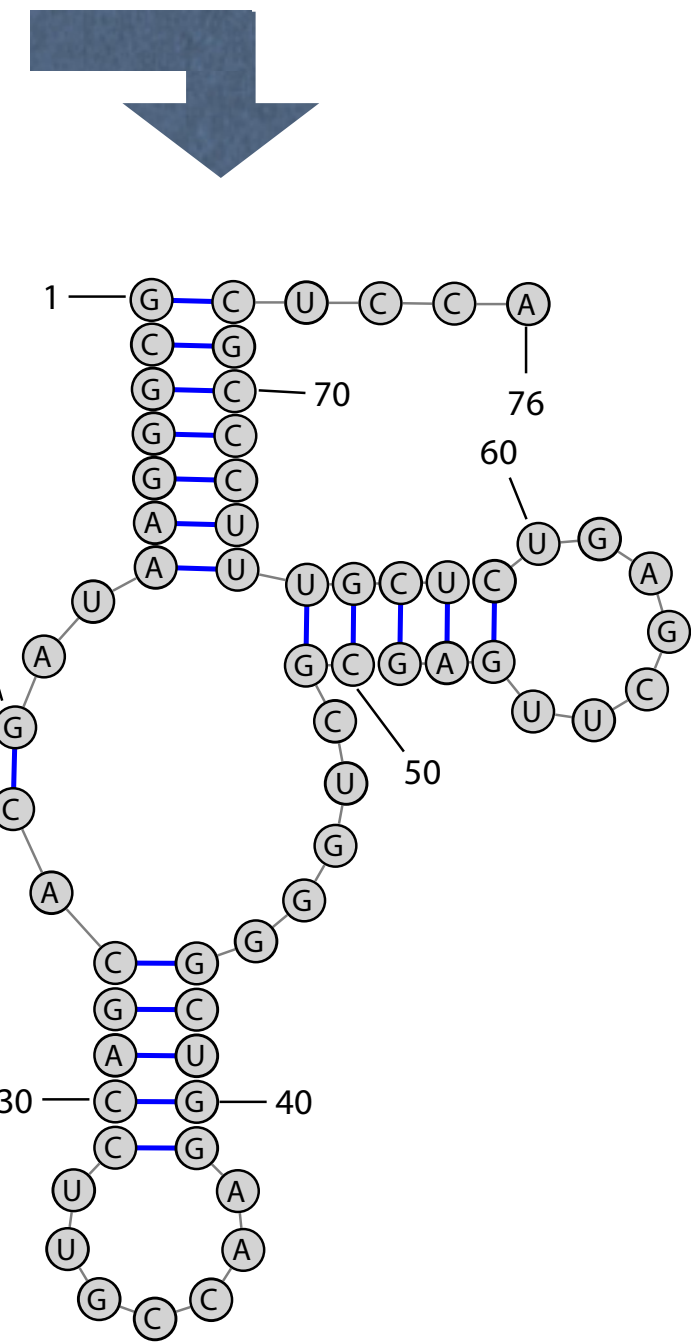
ISMB 2019 Proceedings Talk



\* corresponding author \*\* equal contribution

# LinearFold: Linear-Time Approximate RNA Folding by 5'-to-3' dynamic programming and beam search

$x$  GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUCCCGCUCCA  
 $y$  (((((((((..((((.....))))).((((.....))))).(((.....))))).(((.....)))))).....



**Liang Huang \***

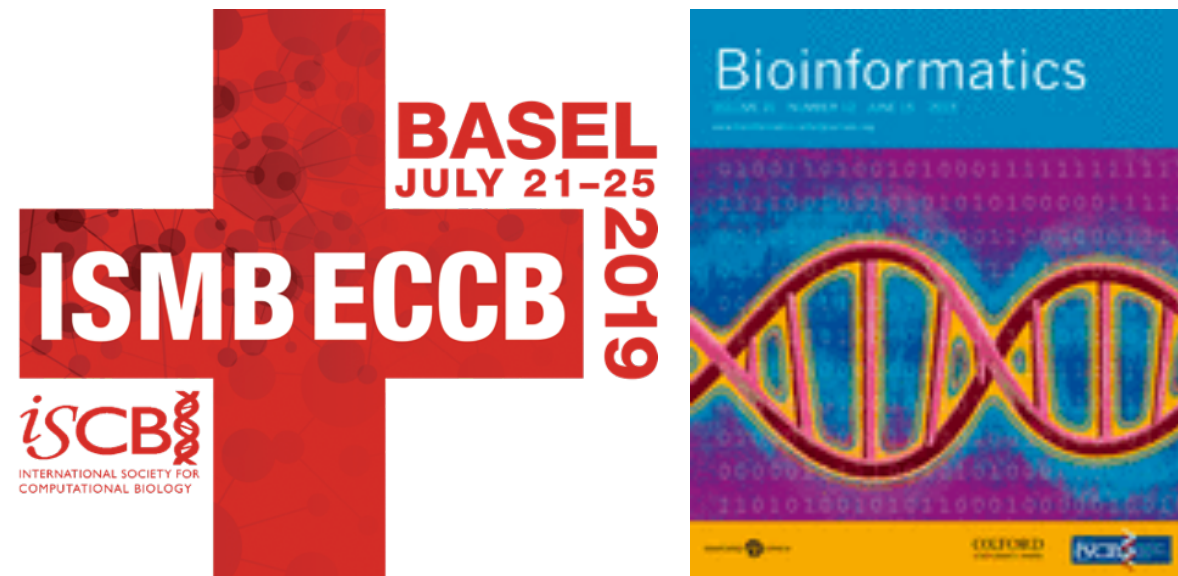
Oregon State University & Baidu Research USA

Joint work with He Zhang \*\*, Dezhong Deng \*\*, Kai Zhao, Kaibo Liu, David Hendrix and David Mathews

first  $O(n)$  (approx.) RNA folding algorithm & server ([linearfold.org](http://linearfold.org)) with even higher accuracy than  $O(n^3)$  algorithms

**Baidu Research**

\* corresponding author \*\* equal contribution



ISMB 2019 Proceedings Talk



# RNA Secondary Structure Prediction

allowed pairs: G-C A-U G-U

assume no crossing pairs

(no pseudoknots)

example: transfer RNA (tRNA)

input  $x$  GCGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

# RNA Secondary Structure Prediction

allowed pairs: G-C A-U G-U

example: transfer RNA (tRNA)

assume no crossing pairs

(no pseudoknots)

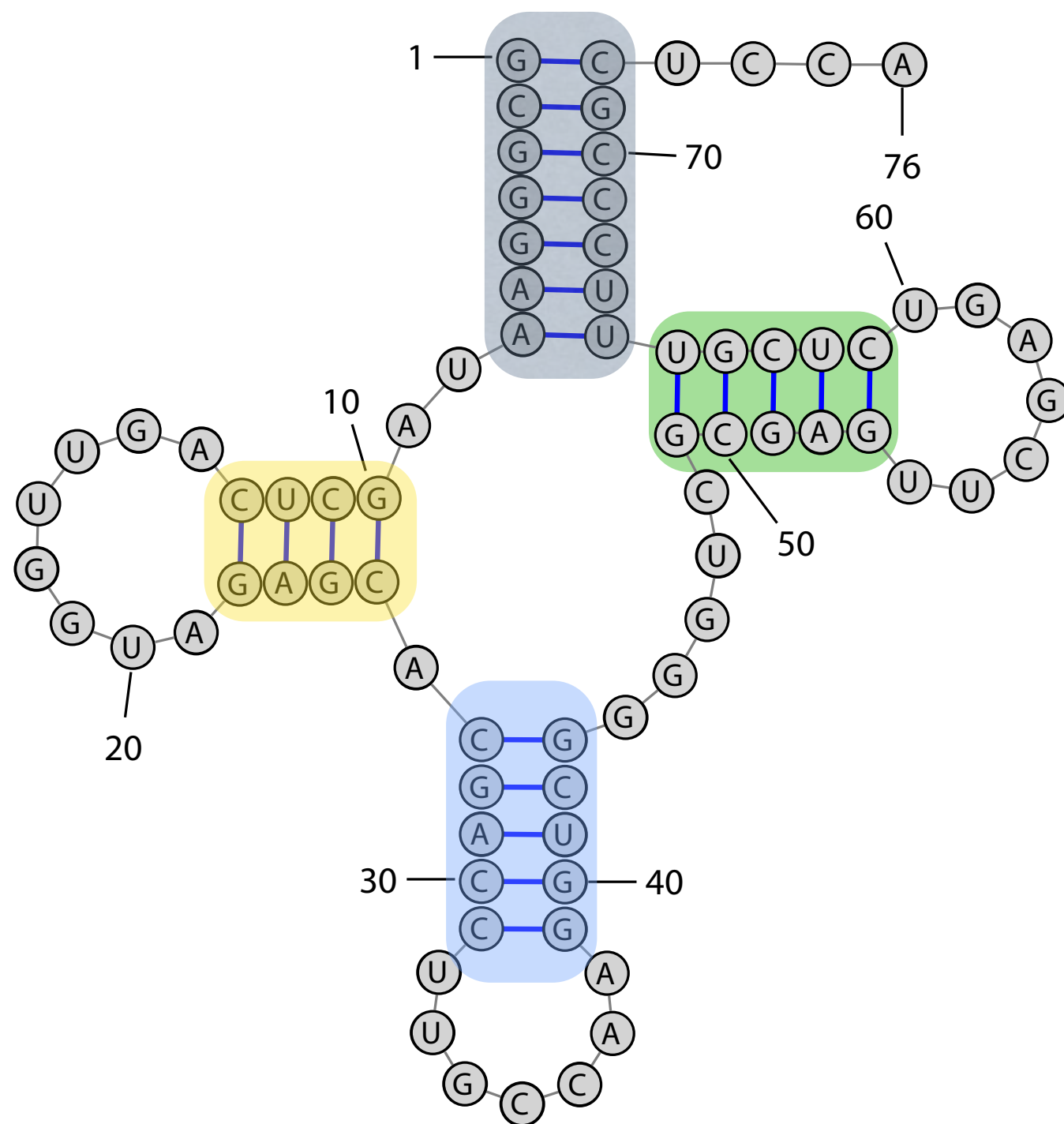
input	$x$	GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA
output	$y$	((((((((..(((.....))))).(((.....))))).(((.....))))))))).....

# RNA Secondary Structure Prediction

allowed pairs: G-C A-U G-U  
 assume no crossing pairs  
 (no pseudoknots)

example: transfer RNA (tRNA)

input  $x$  GCGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUCCCGCUCCA  
 output  $y$  (((((((...(((.....))))).(((.....))))).(((.....)))))).....

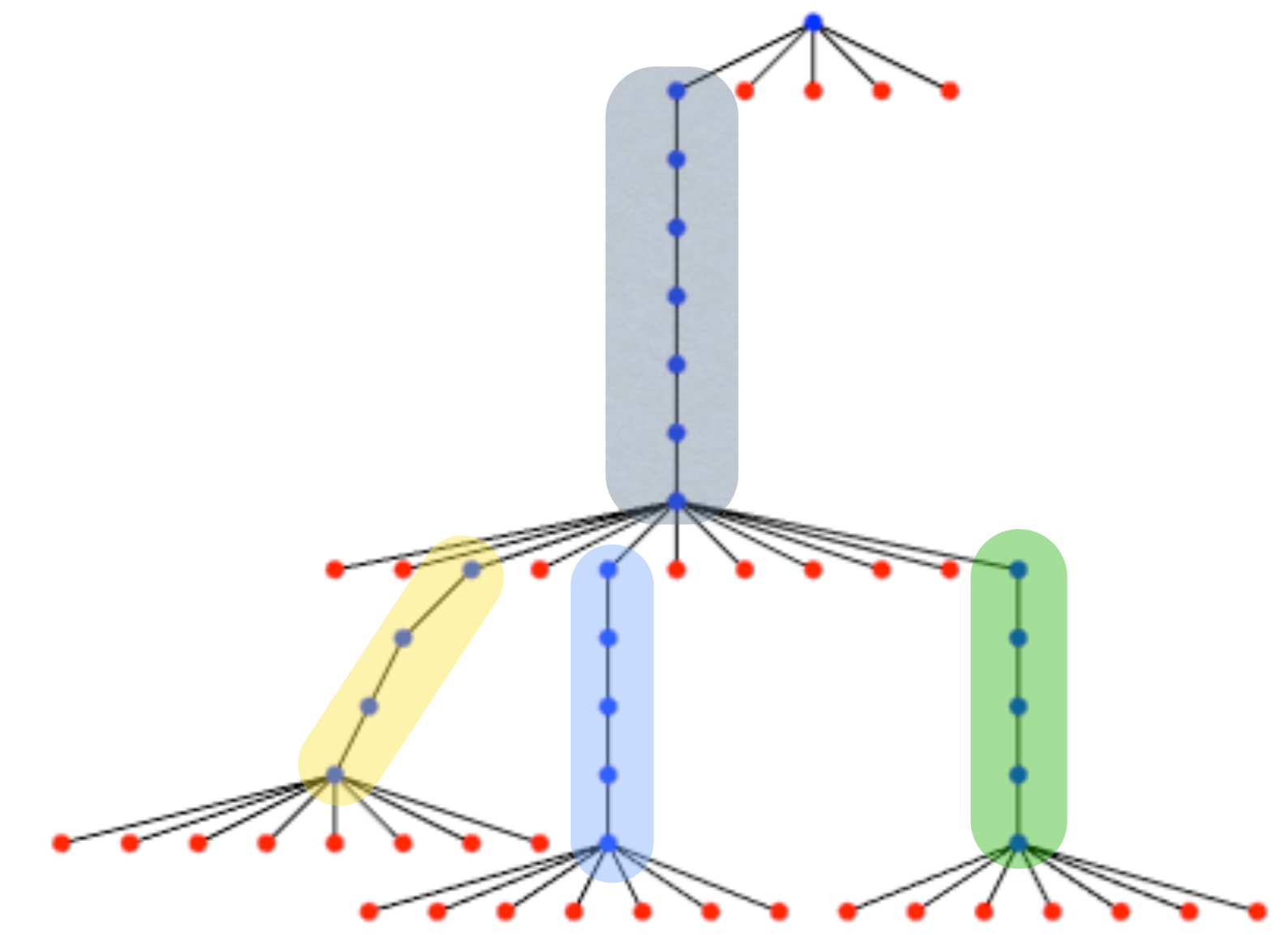
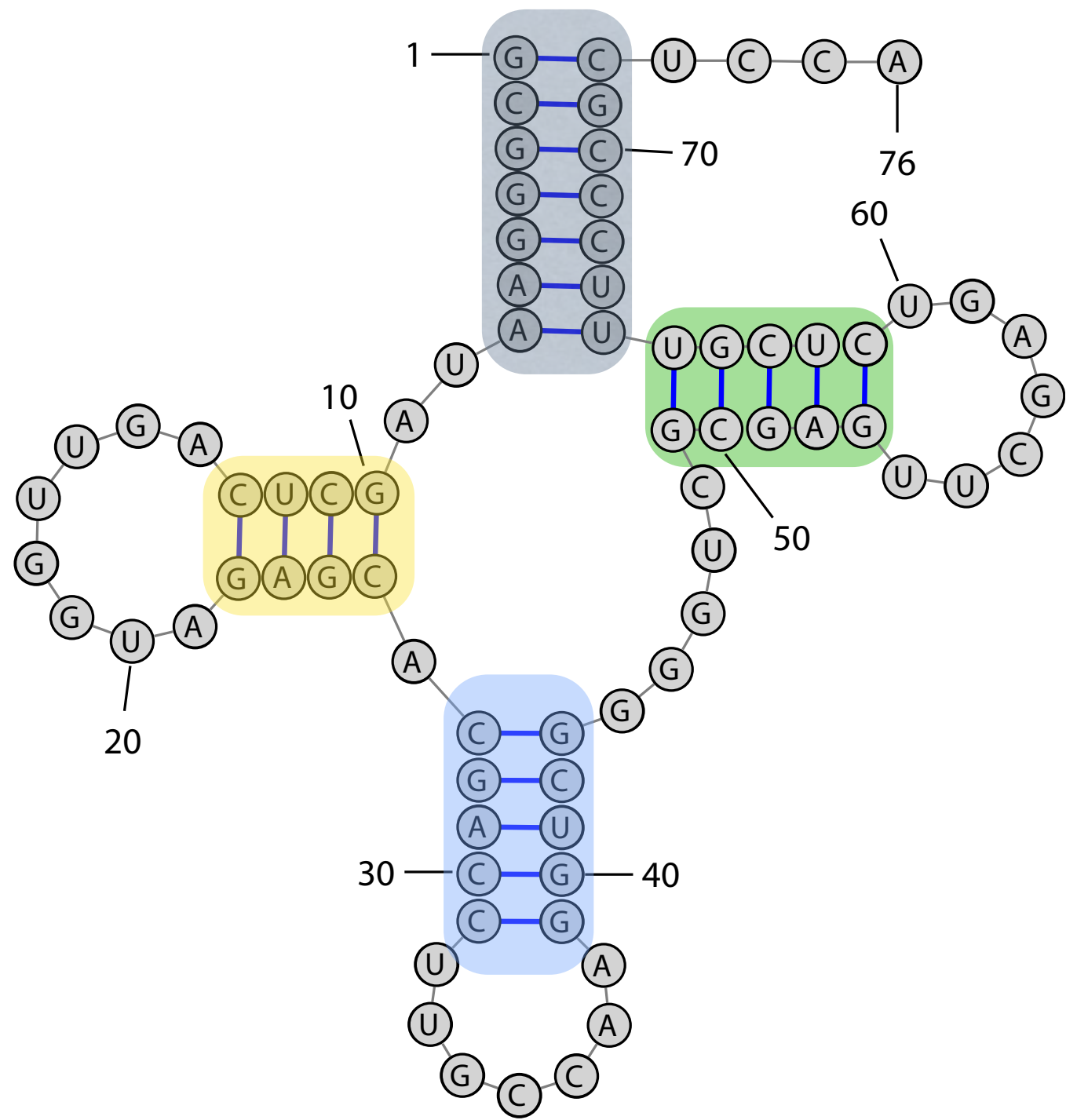


# RNA Secondary Structure Prediction

allowed pairs: G-C A-U G-U  
 assume no crossing pairs  
 (no pseudoknots)

example: transfer RNA (tRNA)

input  $x$  GCGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUCCCGCUCCA  
 output  $y$  (((((((...(((.....))))).((((.....))))).(((.....)))))).....



● Paired  
 ● Unpaired

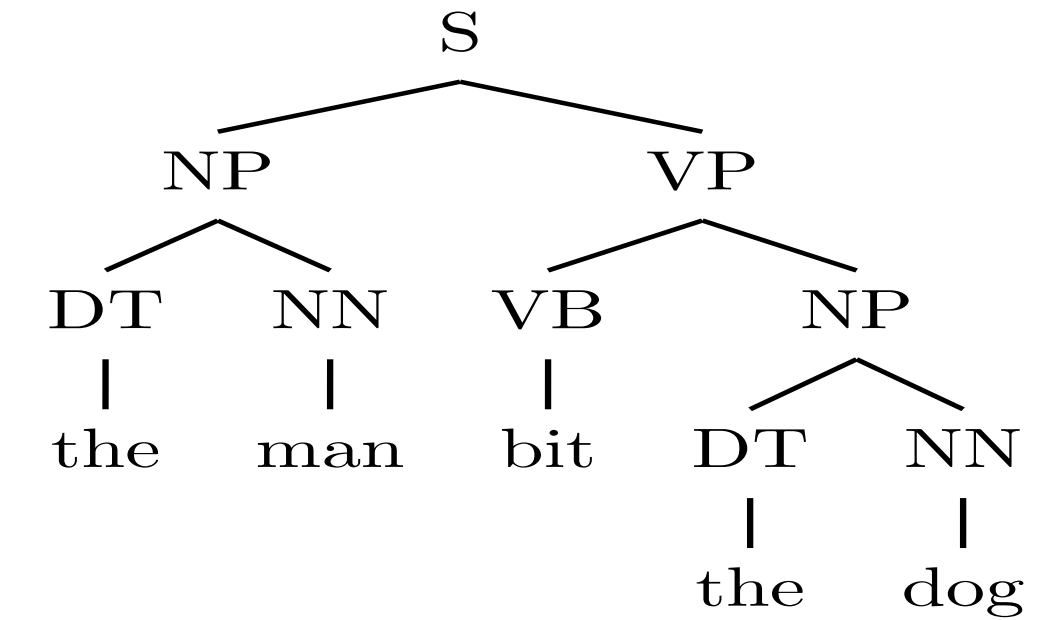
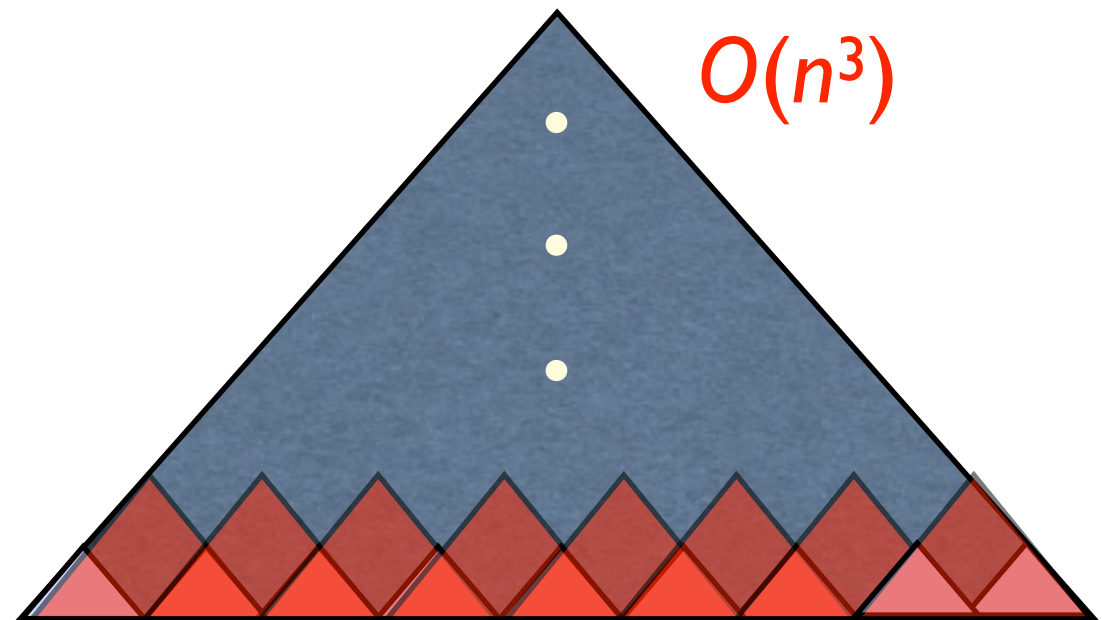
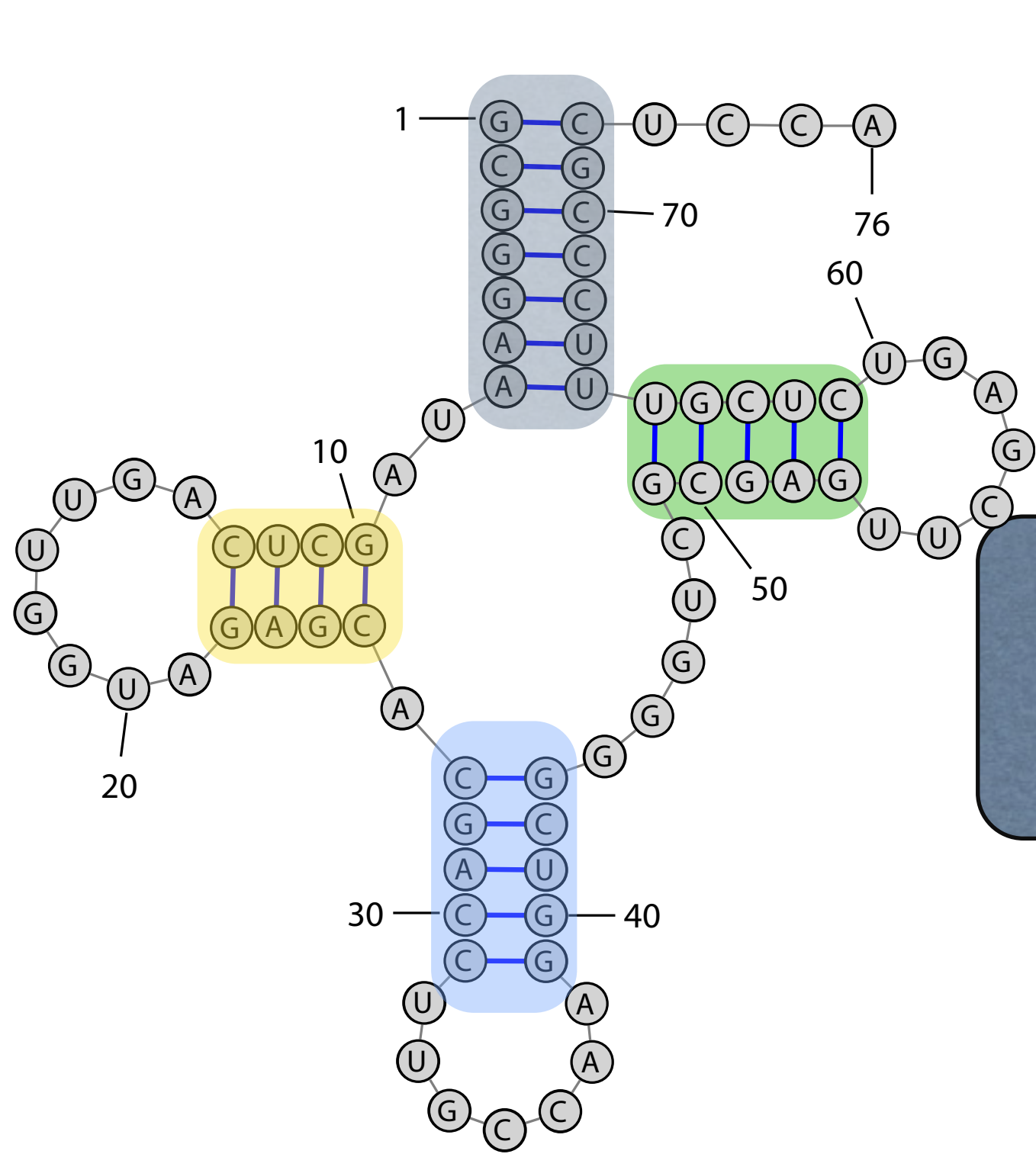
parse tree

# RNA Secondary Structure Prediction

allowed pairs: G-C A-U G-U  
 assume no crossing pairs  
 (no pseudoknots)

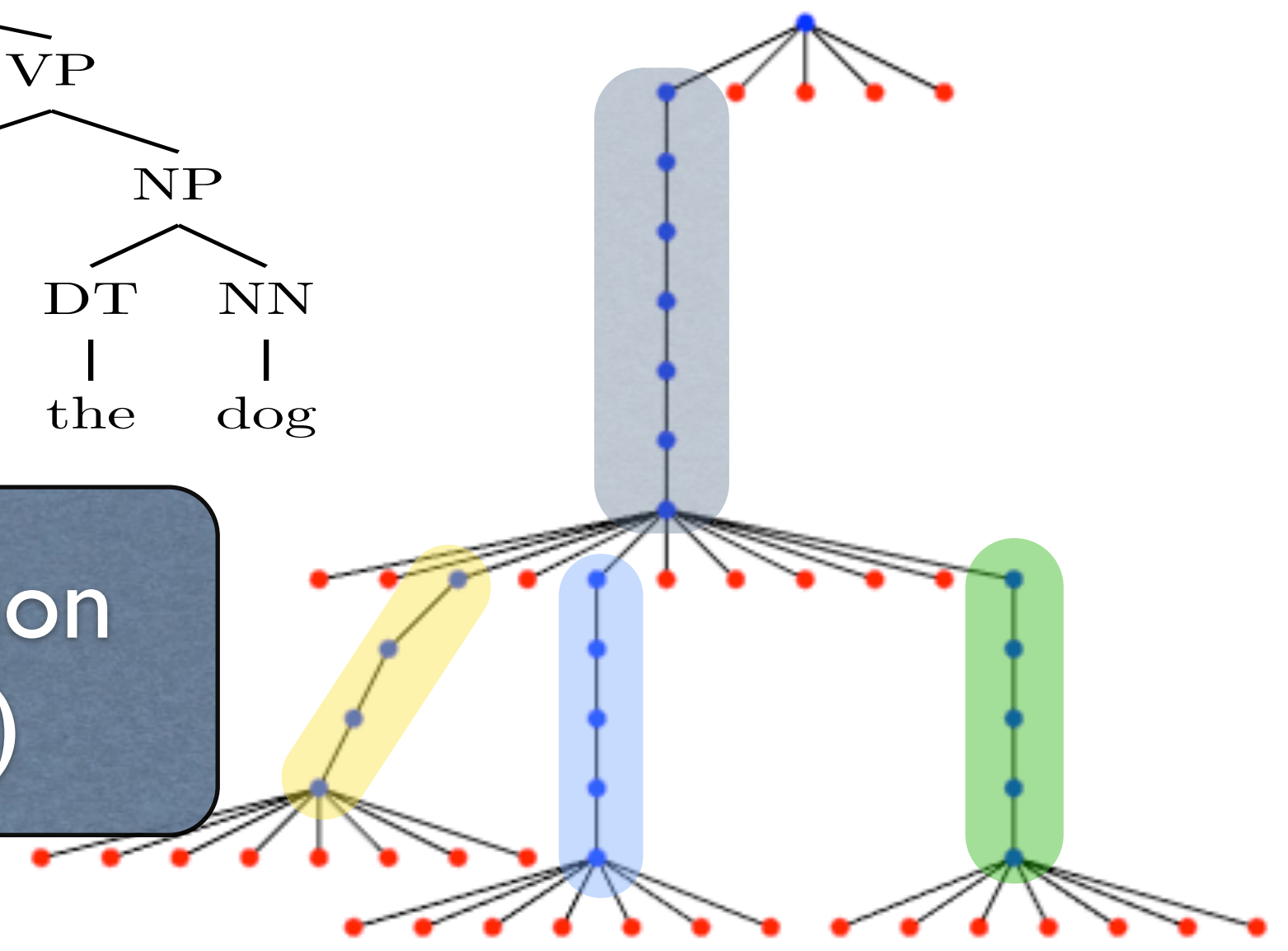
example: transfer RNA (tRNA)

input  $x$  GCGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUCCCGCUCCA  
 output  $y$  (((((((...(((.....)))))).....((((.....)))))).....)))))).....



problem: standard structure prediction algorithms are way too slow:  $O(n^3)$

● Paired  
 ● Unpaired



parse tree

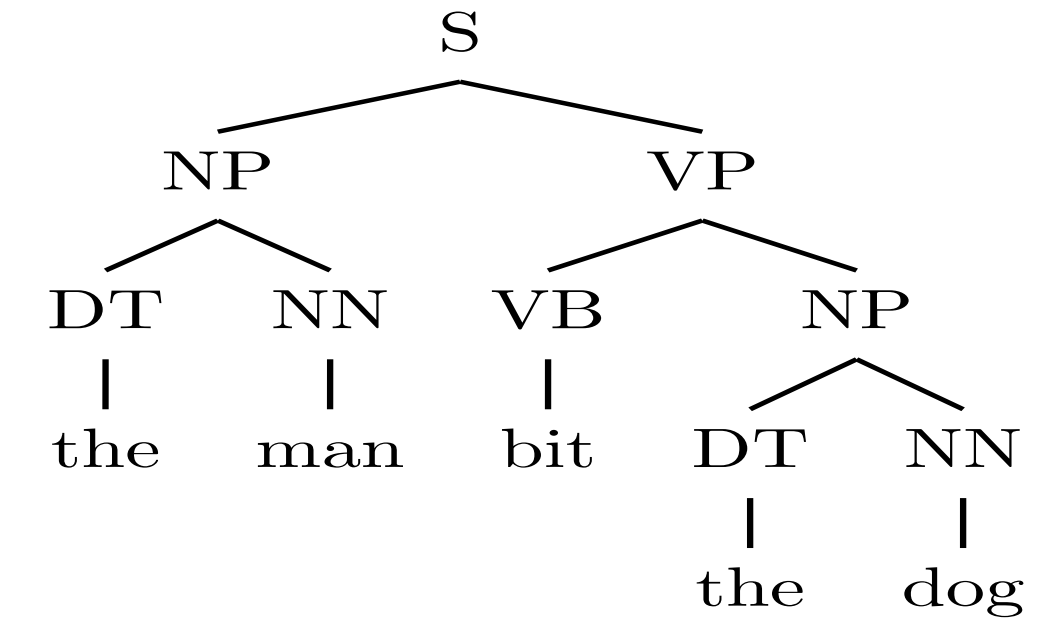
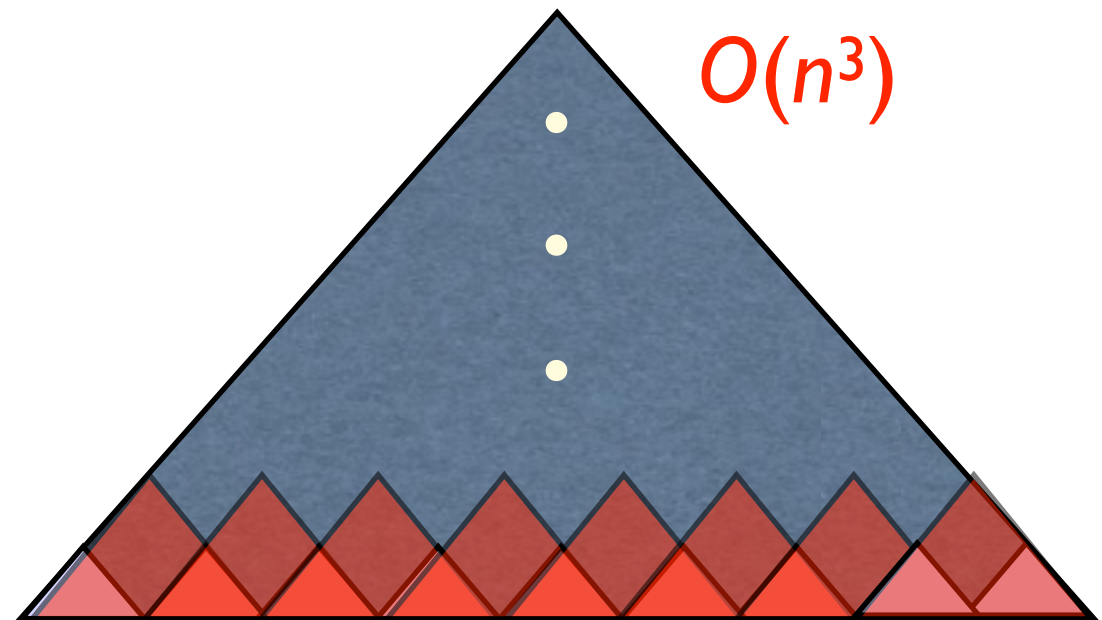
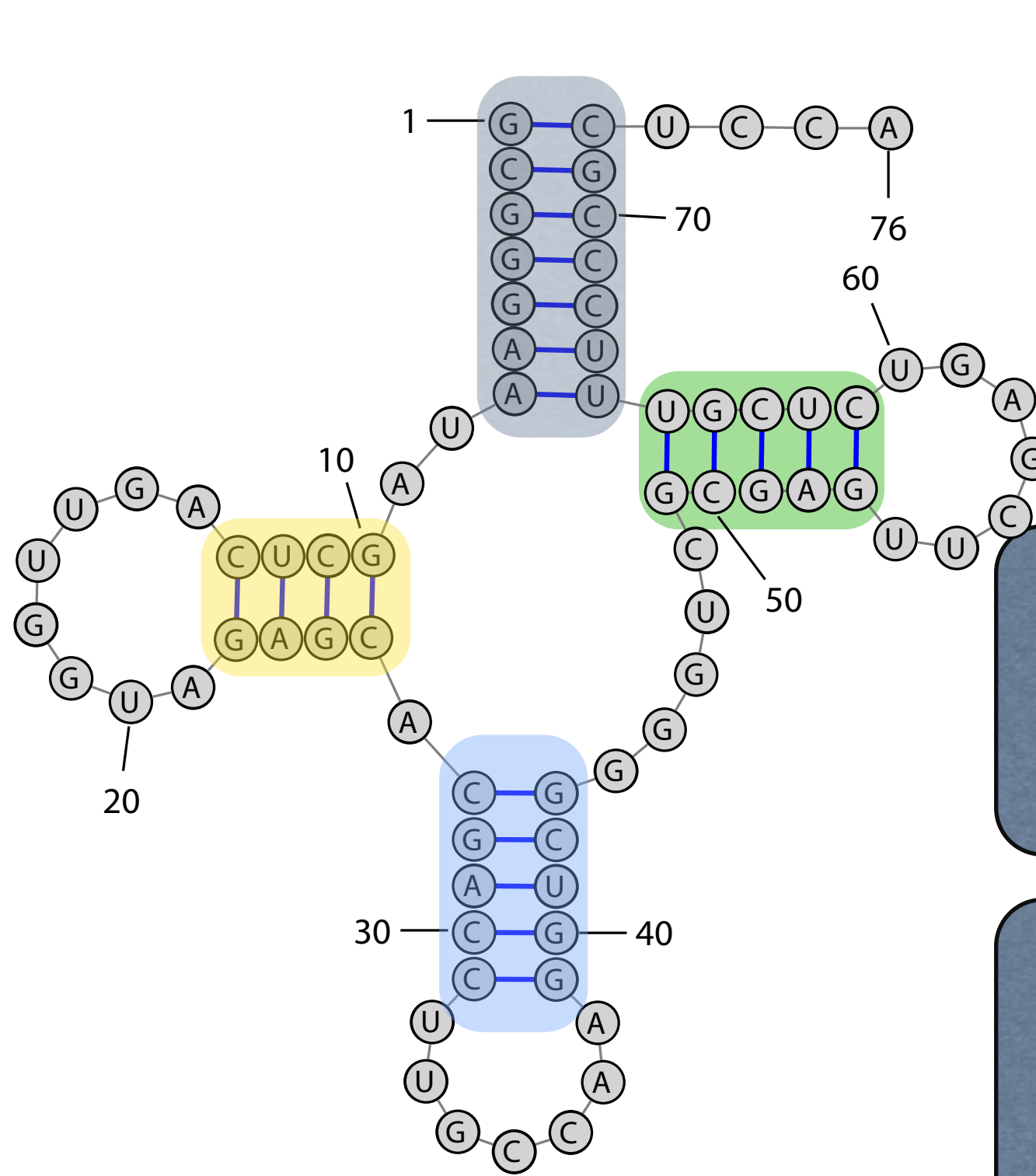


# RNA Secondary Structure Prediction

allowed pairs: G-C A-U G-U  
 assume no crossing pairs  
 (no pseudoknots)

example: transfer RNA (tRNA)

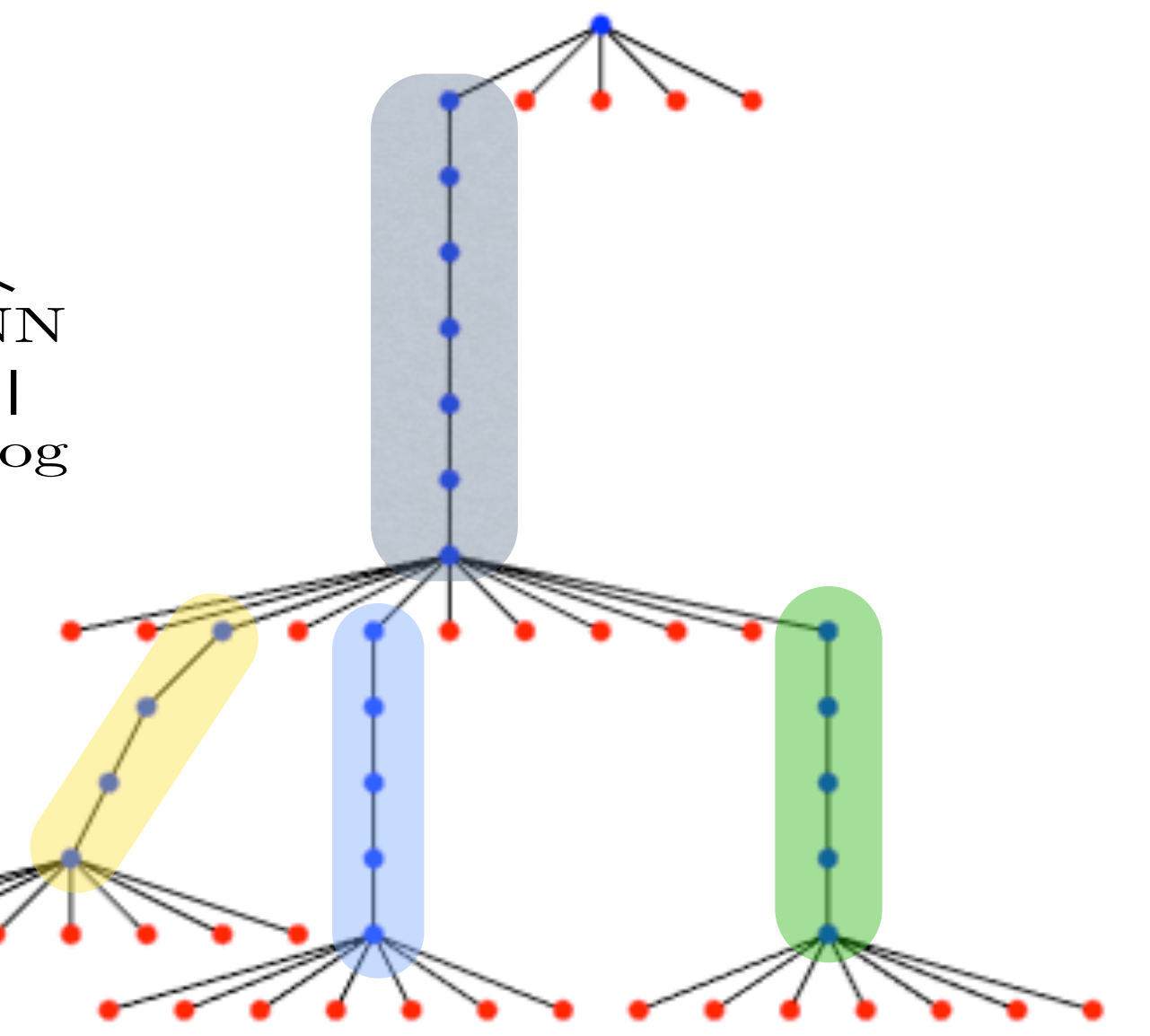
input  $x$  GCGGGAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUCCCGCUCCA  
 output  $y$  ((((((...(((...))))))...(((...))))))...))



● Paired  
 ● Unpaired

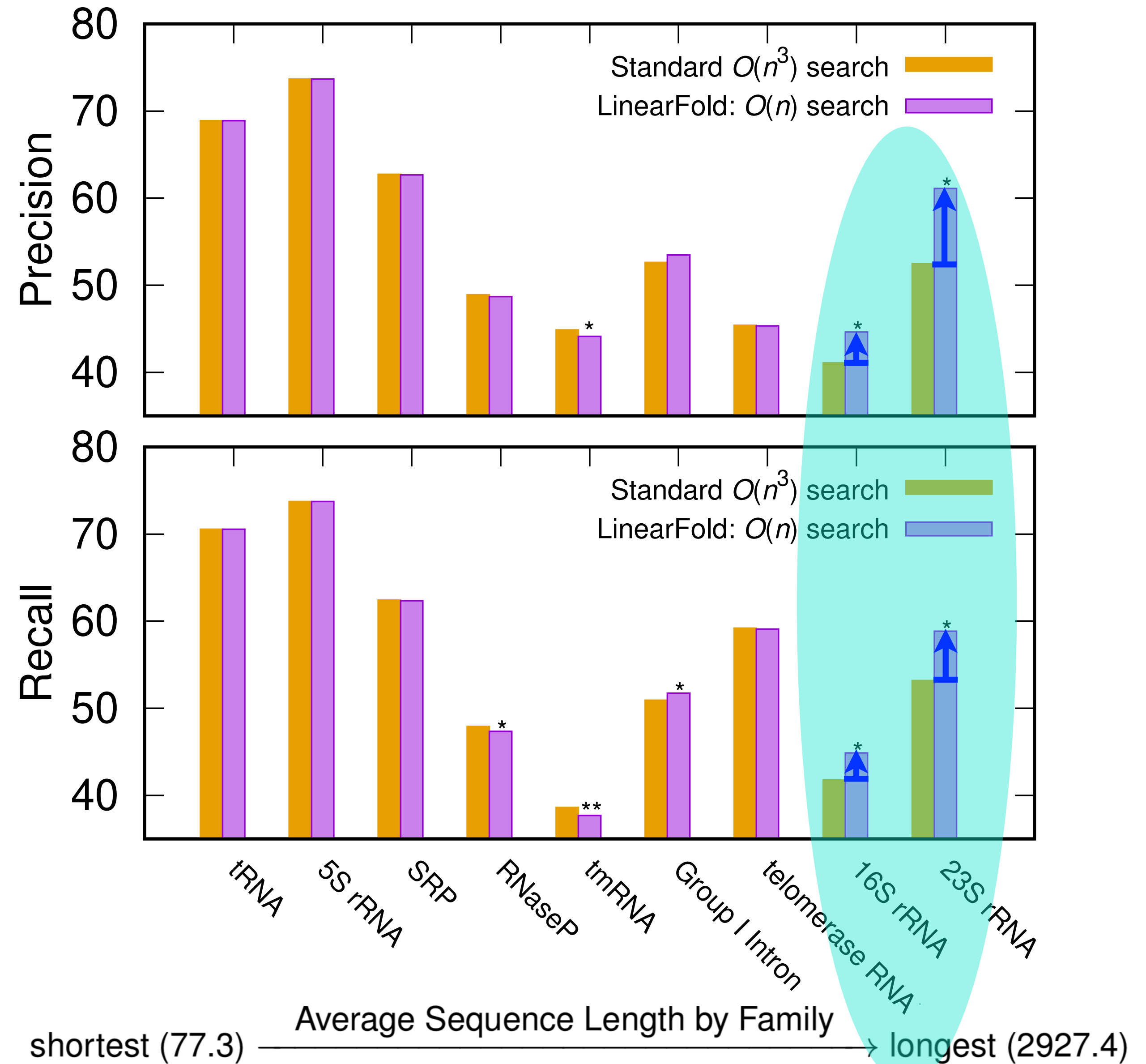
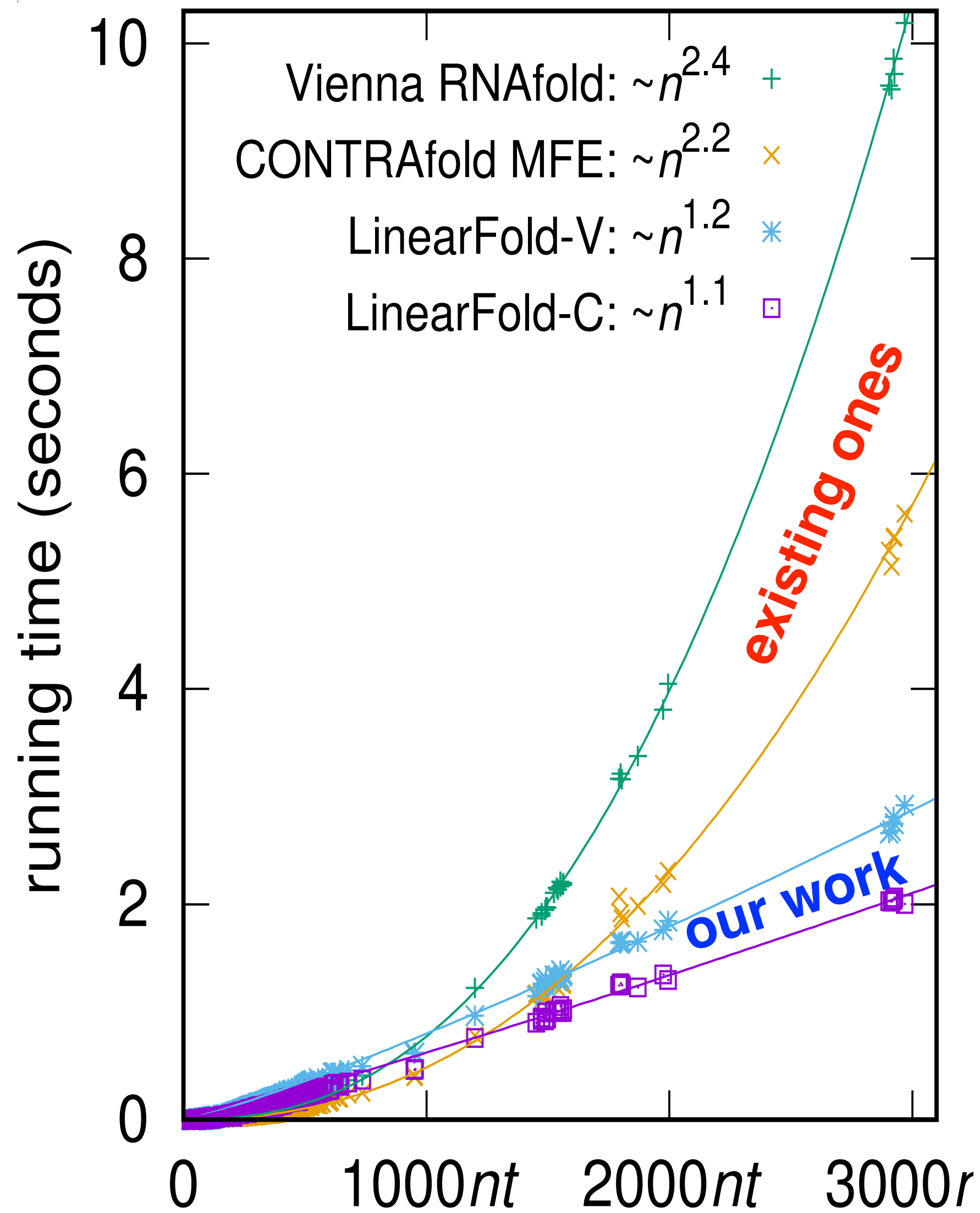
problem: standard structure prediction algorithms are way too slow:  $O(n^3)$

solution: adapt my linear-time dynamic programming algorithms from parsing



parse tree

# Results: LinearFold is Much Faster and More Accurate

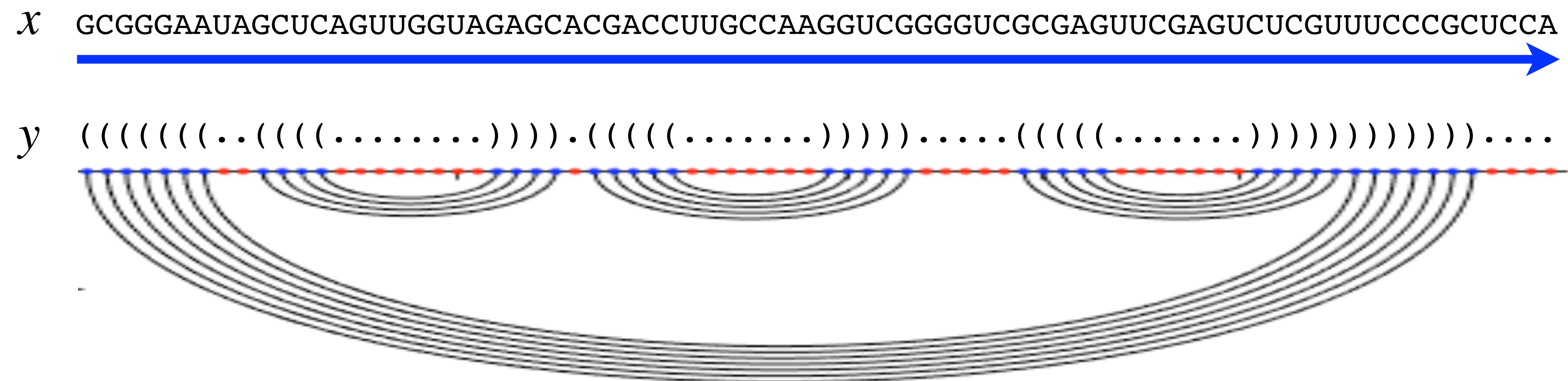


# From Linguistics to Biology

$x$  GCGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

$y$  (((((((((..((((.....))))).((((.....))))). .... (((((((.....))))))))))....

# From Linguistics to Biology



# Computational Linguistics => Computational Biology

## linguistics

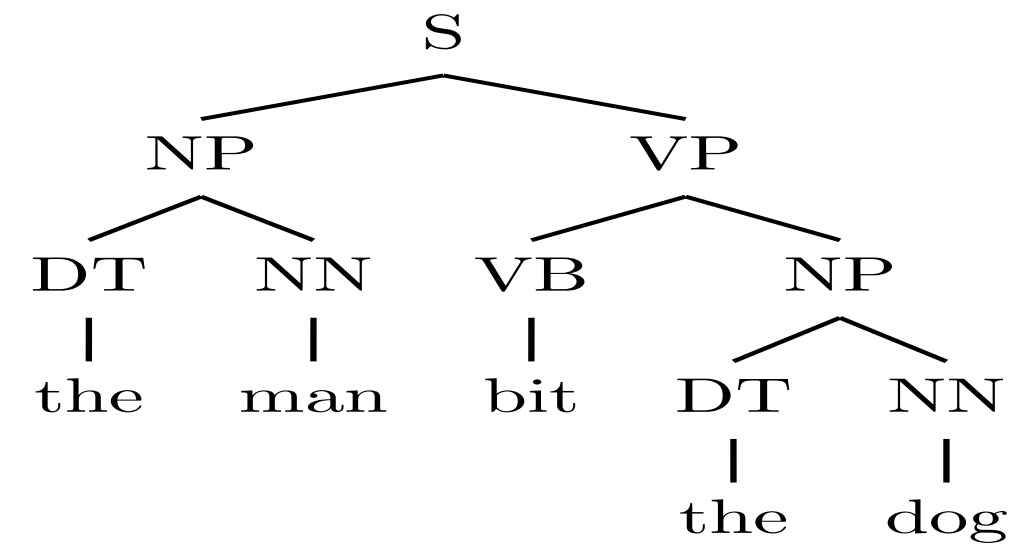
## compiler theory

## comp. linguistics

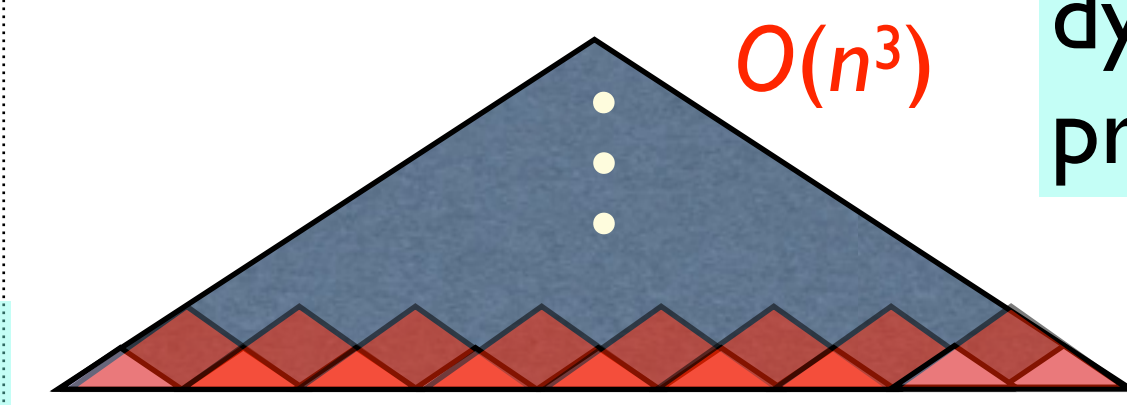
## computational biology

1955 Chomsky: context-free grammars (CFGs)

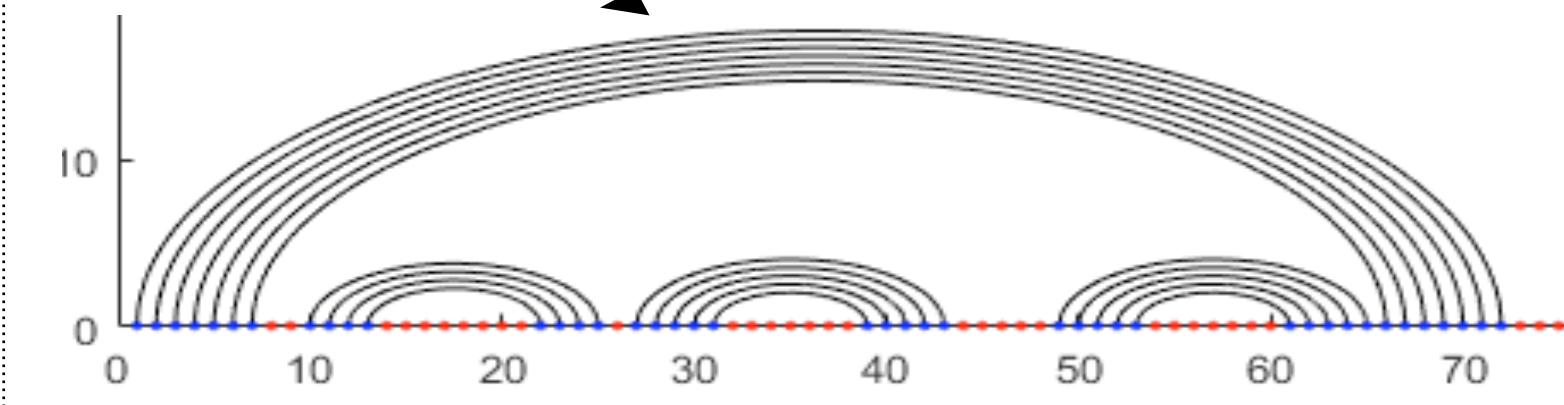
1958 Backus & Naur: CFGs for program, lang.



1964 Cocke \ bottom-up  
1965 Kasami - **CKY**  $O(n^3)$   
1967 Younger / for all CFGs



dynamic programming



1978: Nussinov  $O(n^3)$  RNA folding  
1981: Zuker & Siegler

# Computational Linguistics => Computational Biology

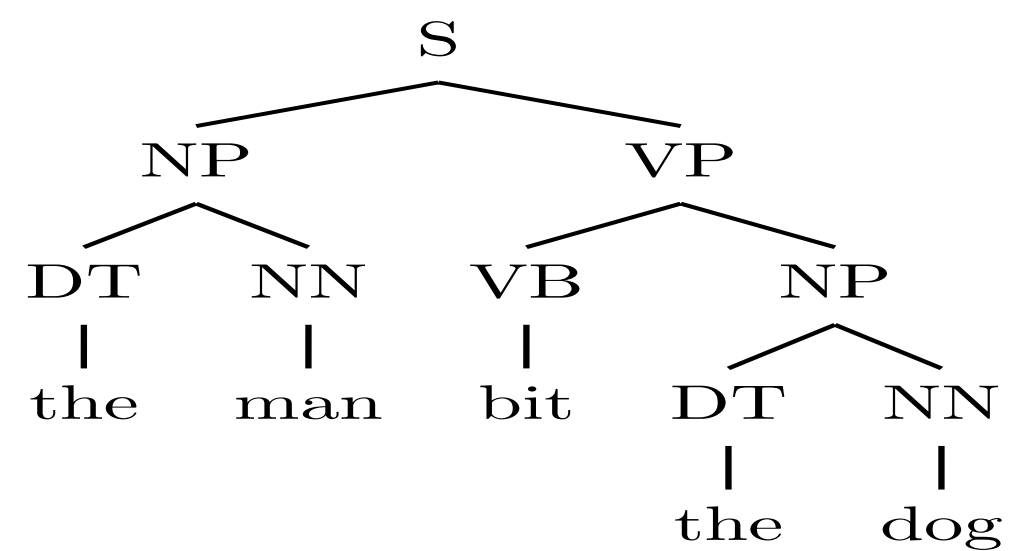
## linguistics

## compiler theory

## comp. linguistics

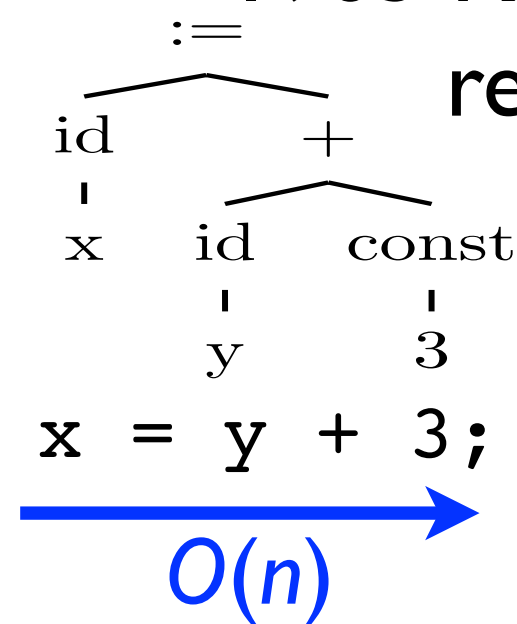
## computational biology

1955 Chomsky: context-free grammars (CFGs)

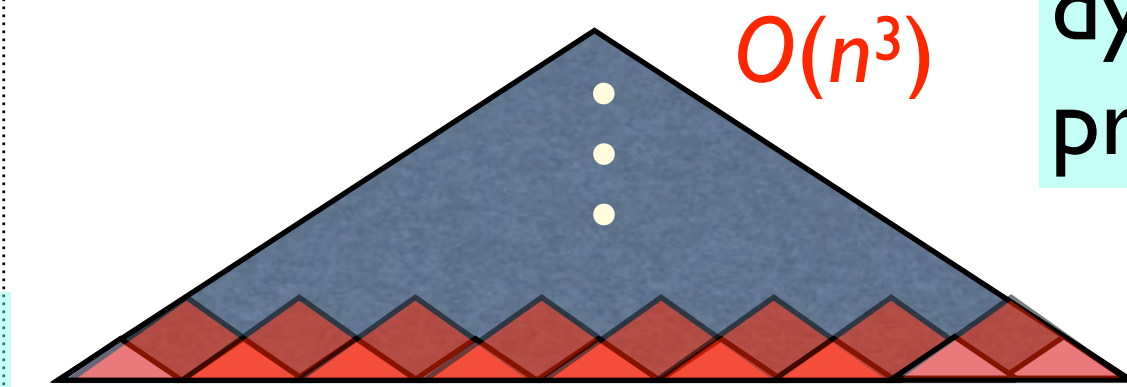


1958 Backus & Naur: CFGs for program, lang.

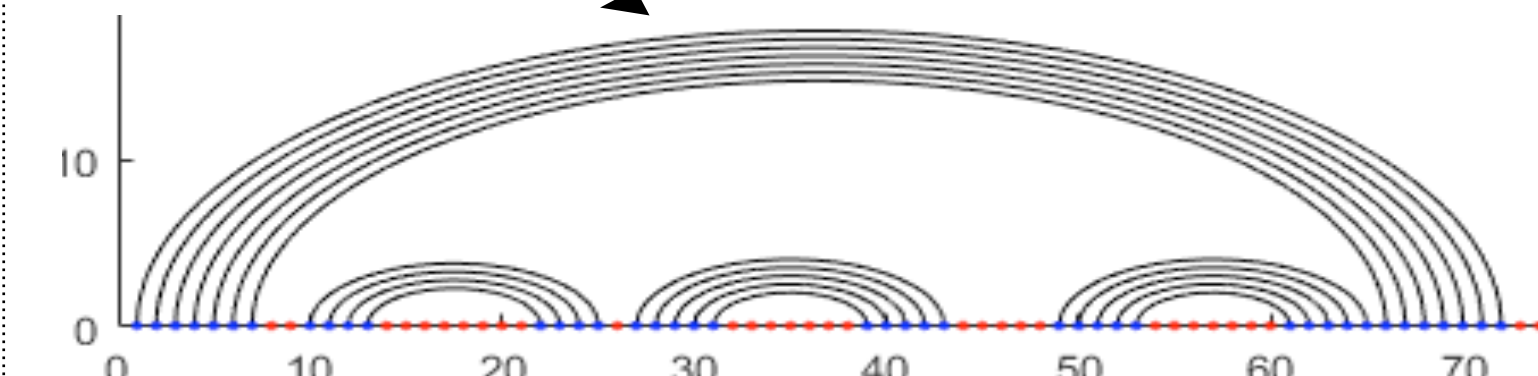
1965 Knuth: **LR** parsing for restricted CFGs:  $O(n)$



1964 Cocke \ bottom-up  
1965 Kasami - **CKY**  $O(n^3)$   
1967 Younger / for all CFGs



dynamic programming



1978: Nussinov  $O(n^3)$  RNA folding  
1981: Zuker & Siegler

# Computational Linguistics => Computational Biology

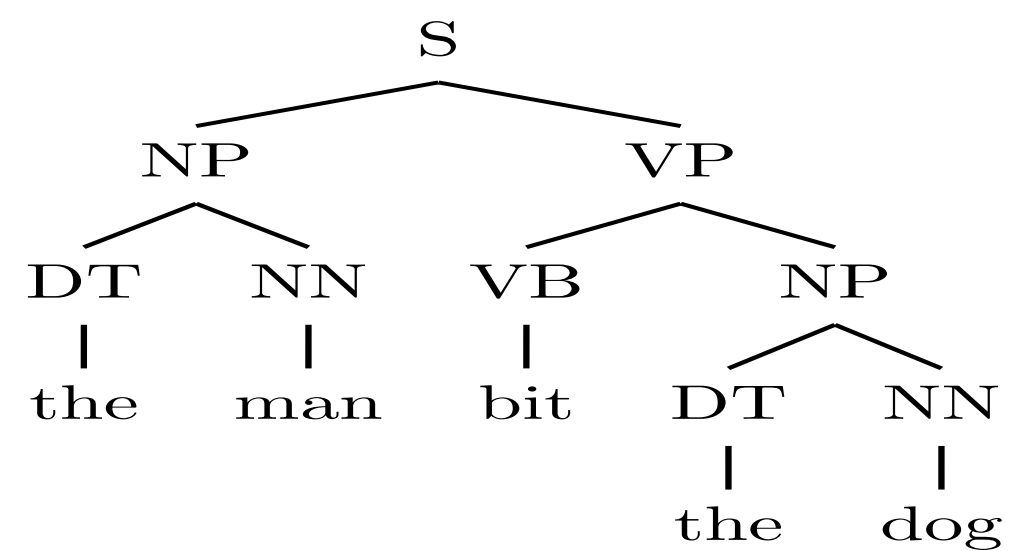
## linguistics

## compiler theory

## comp. linguistics

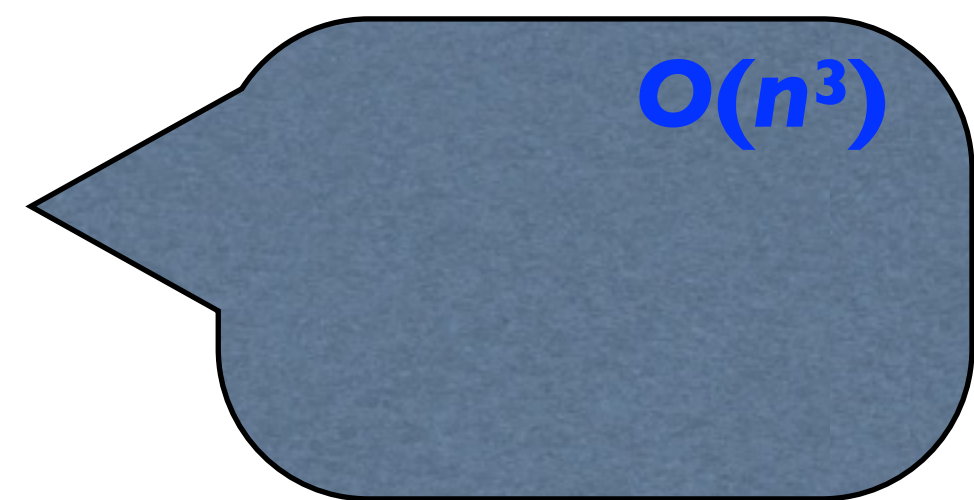
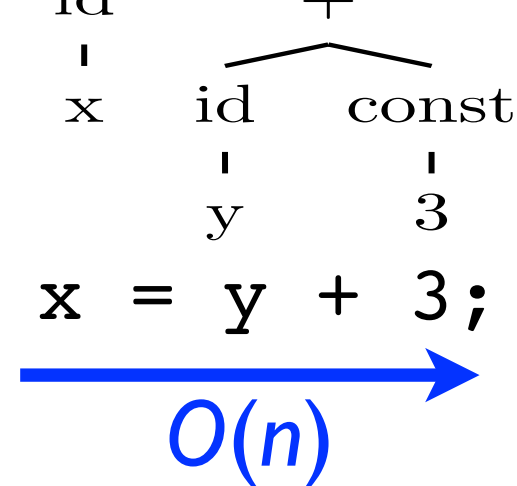
## computational biology

1955 Chomsky: context-free grammars (CFGs)



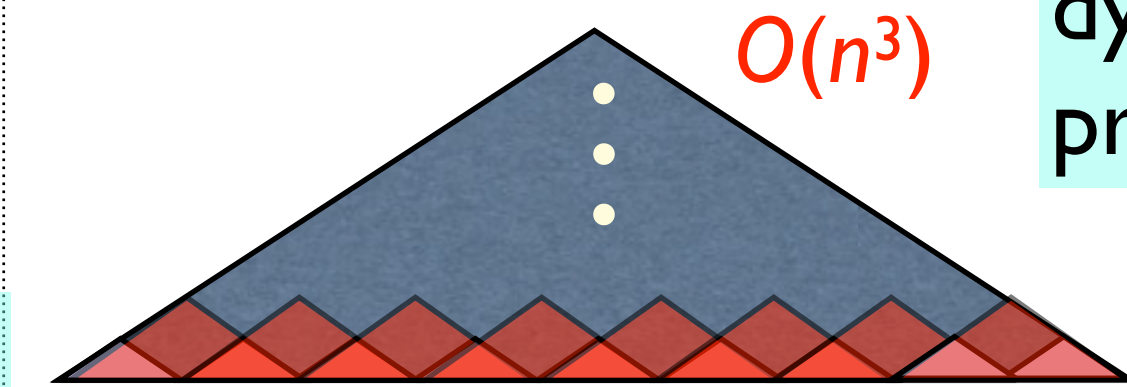
1958 Backus & Naur: CFGs for program, lang.

1965 Knuth: **LR** parsing for restricted CFGs:  $O(n)$

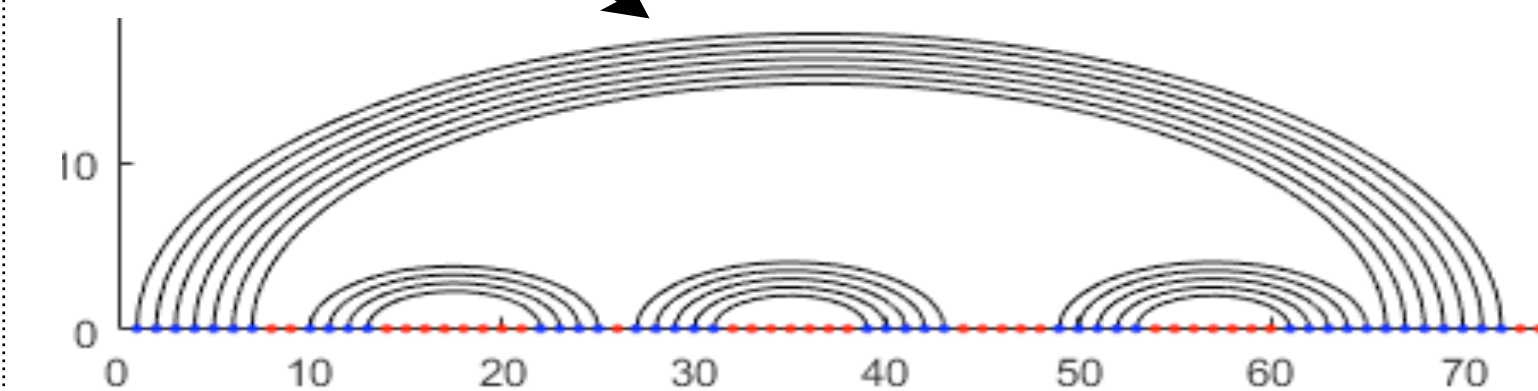


1964 Cocke \ bottom-up  
1965 Kasami - **CKY**  $O(n^3)$   
1967 Younger / for all CFGs

1986 Tomita: **Generalized LR** for all CFGs:  $O(n^3)$



dynamic programming



1978: Nussinov  $O(n^3)$  RNA folding  
1981: Zuker & Siegler

# Computational Linguistics => Computational Biology

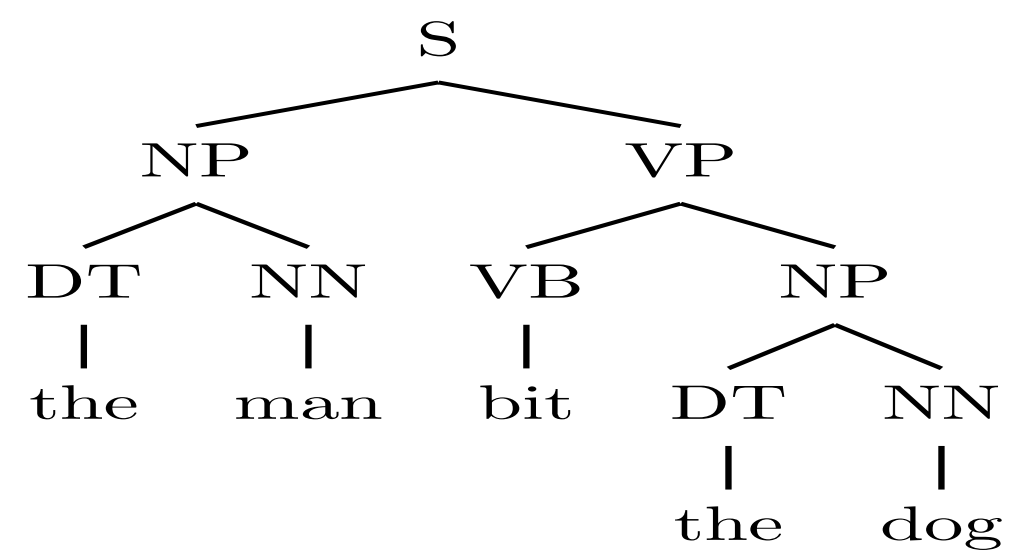
## linguistics

## compiler theory

## comp. linguistics

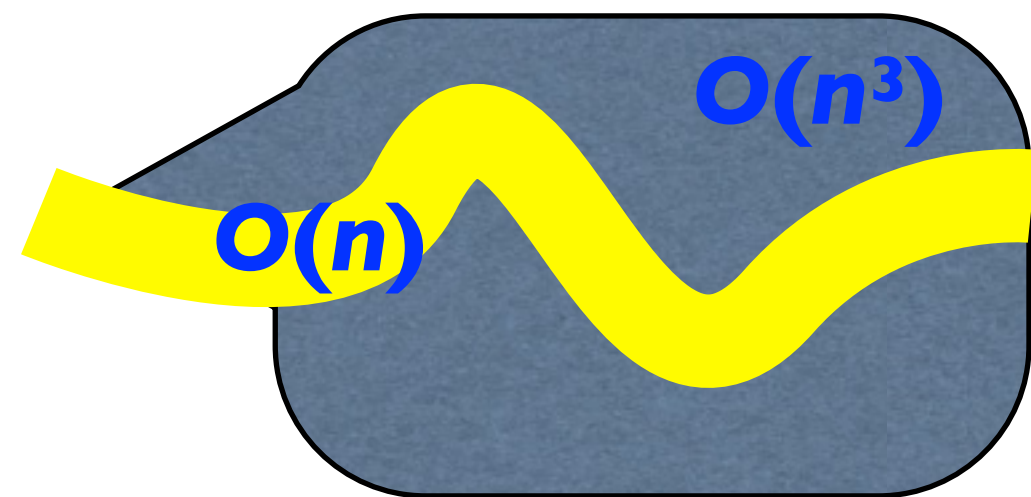
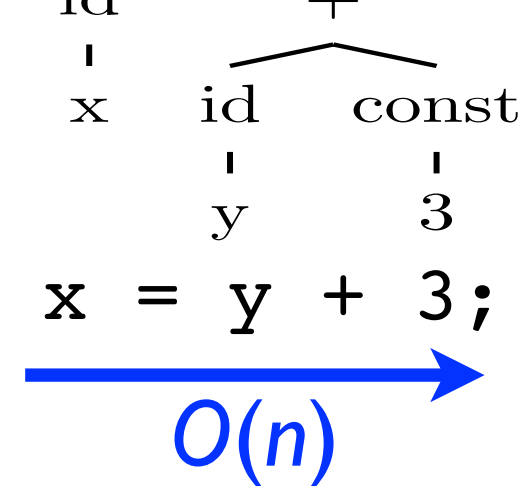
## computational biology

1955 Chomsky: context-free grammars (CFGs)



1958 Backus & Naur: CFGs for program, lang.

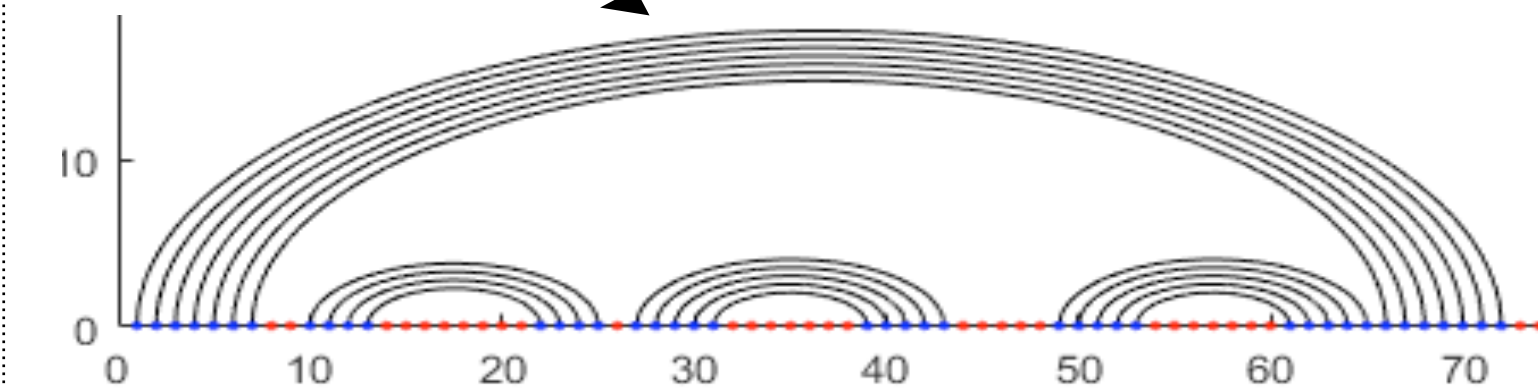
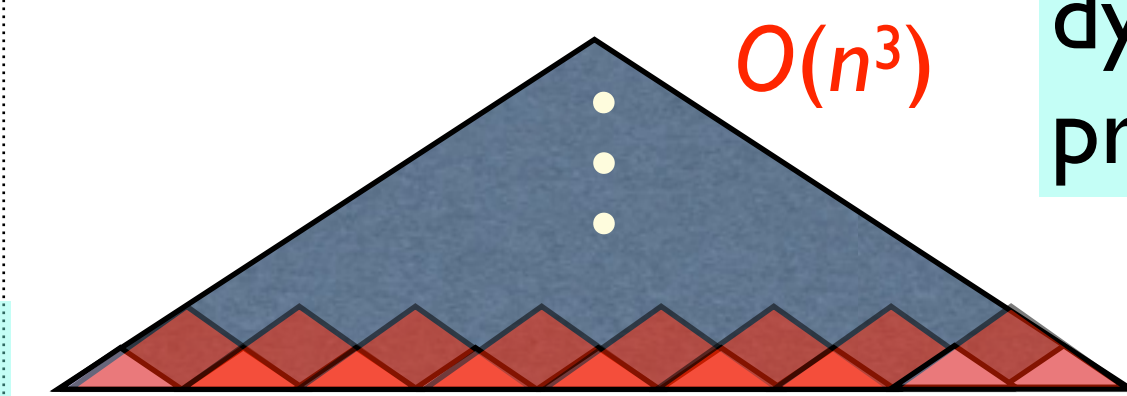
1965 Knuth: **LR** parsing for restricted CFGs:  $O(n)$



1964 Cocke \ bottom-up  
1965 Kasami - **CKY**  $O(n^3)$   
1967 Younger / for all CFGs

1986 Tomita: **Generalized LR** for all CFGs:  $O(n^3)$

2010: Huang & Sagae:  $O(n)$  (approx.) DP for all CFGs



1978: Nussinov  $O(n^3)$  RNA folding  
1981: Zuker & Siegler



# Computational Linguistics => Computational Biology

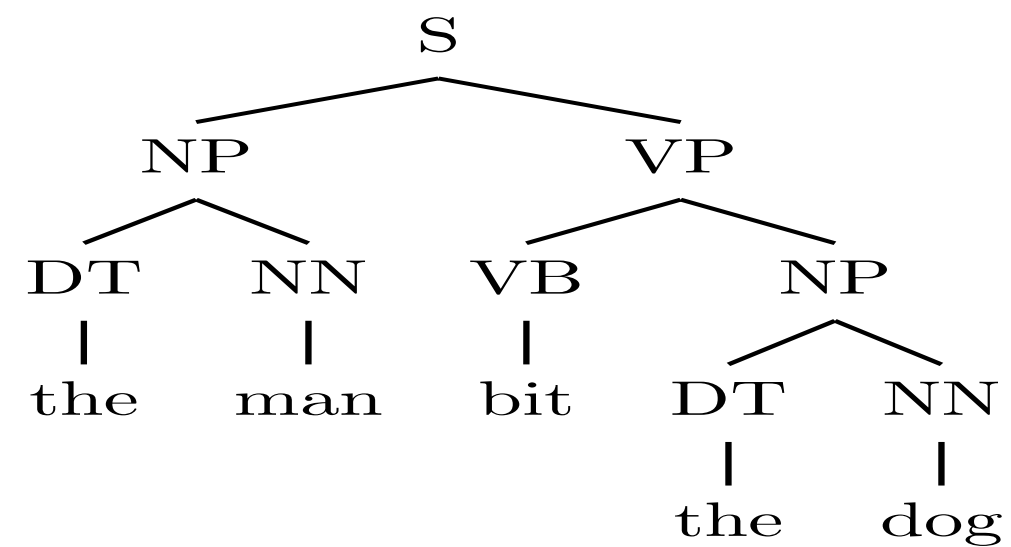
## linguistics

## compiler theory

## comp. linguistics

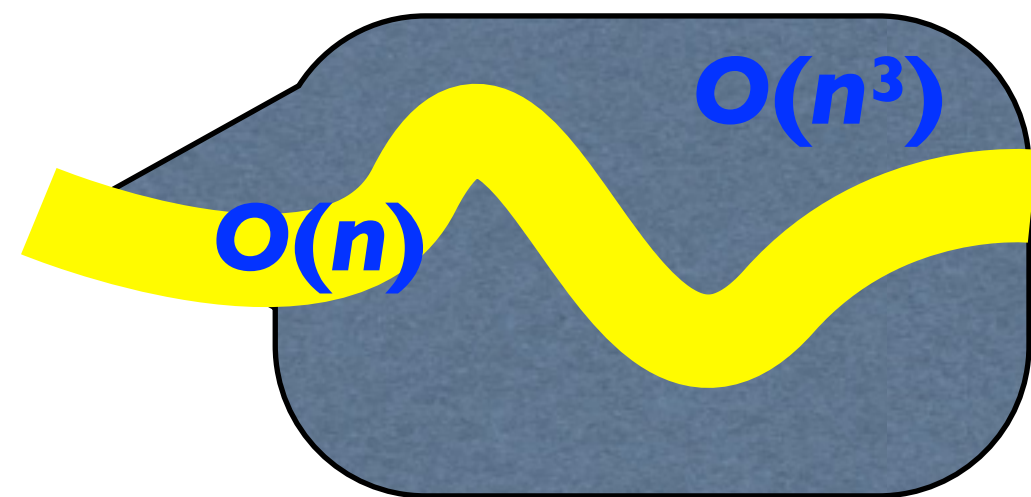
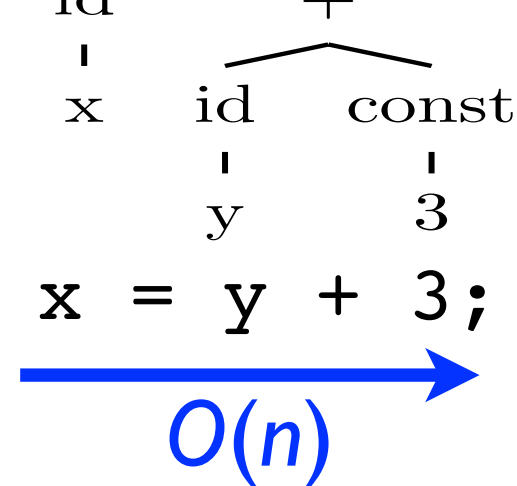
## computational biology

1955 Chomsky: context-free grammars (CFGs)



1958 Backus & Naur: CFGs for program, lang.

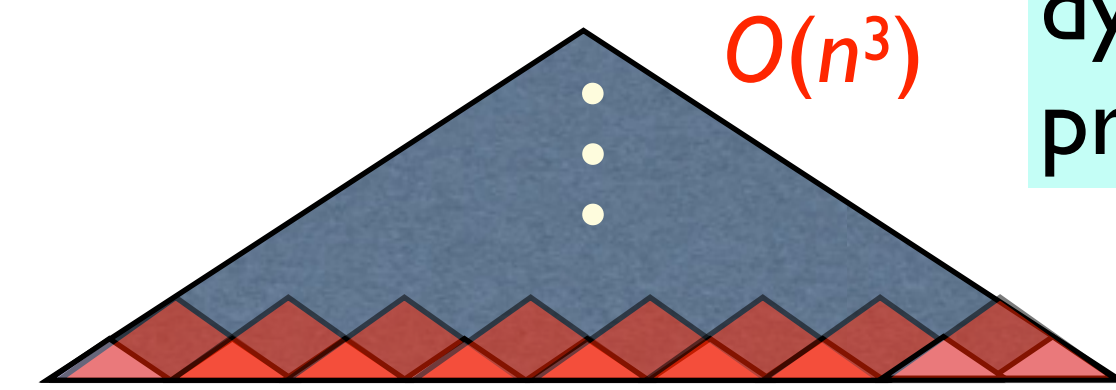
1965 Knuth: **LR** parsing for restricted CFGs:  $O(n)$



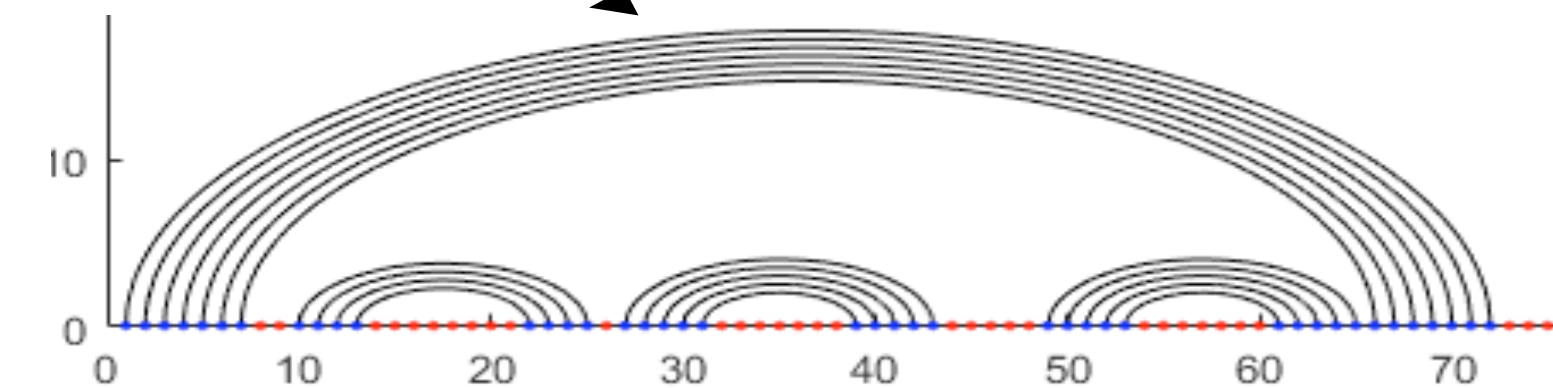
1964 Cocke \ bottom-up  
1965 Kasami - **CKY**  $O(n^3)$   
1967 Younger / for all CFGs

1986 Tomita: **Generalized LR** for all CFGs:  $O(n^3)$

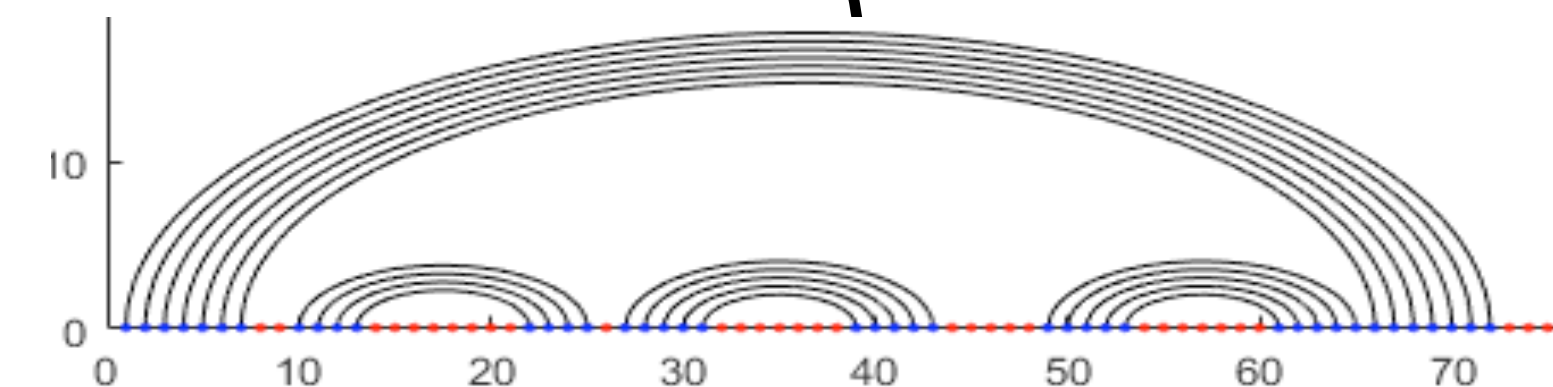
2010: Huang & Sagae:  $O(n)$  (approx.) DP for all CFGs



dynamic programming



1978: Nussinov  $O(n^3)$  RNA folding  
1981: Zuker & Sieglar

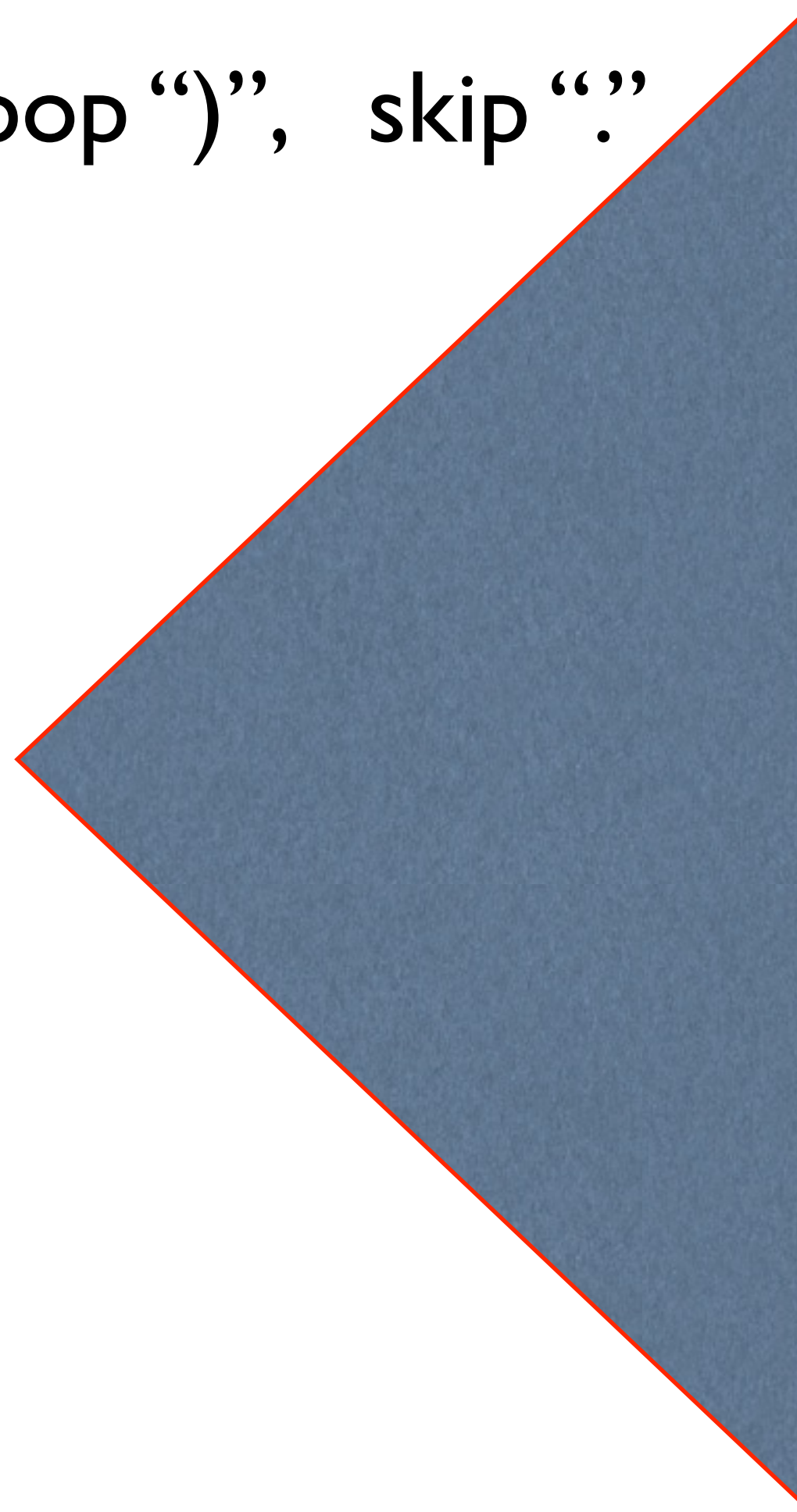
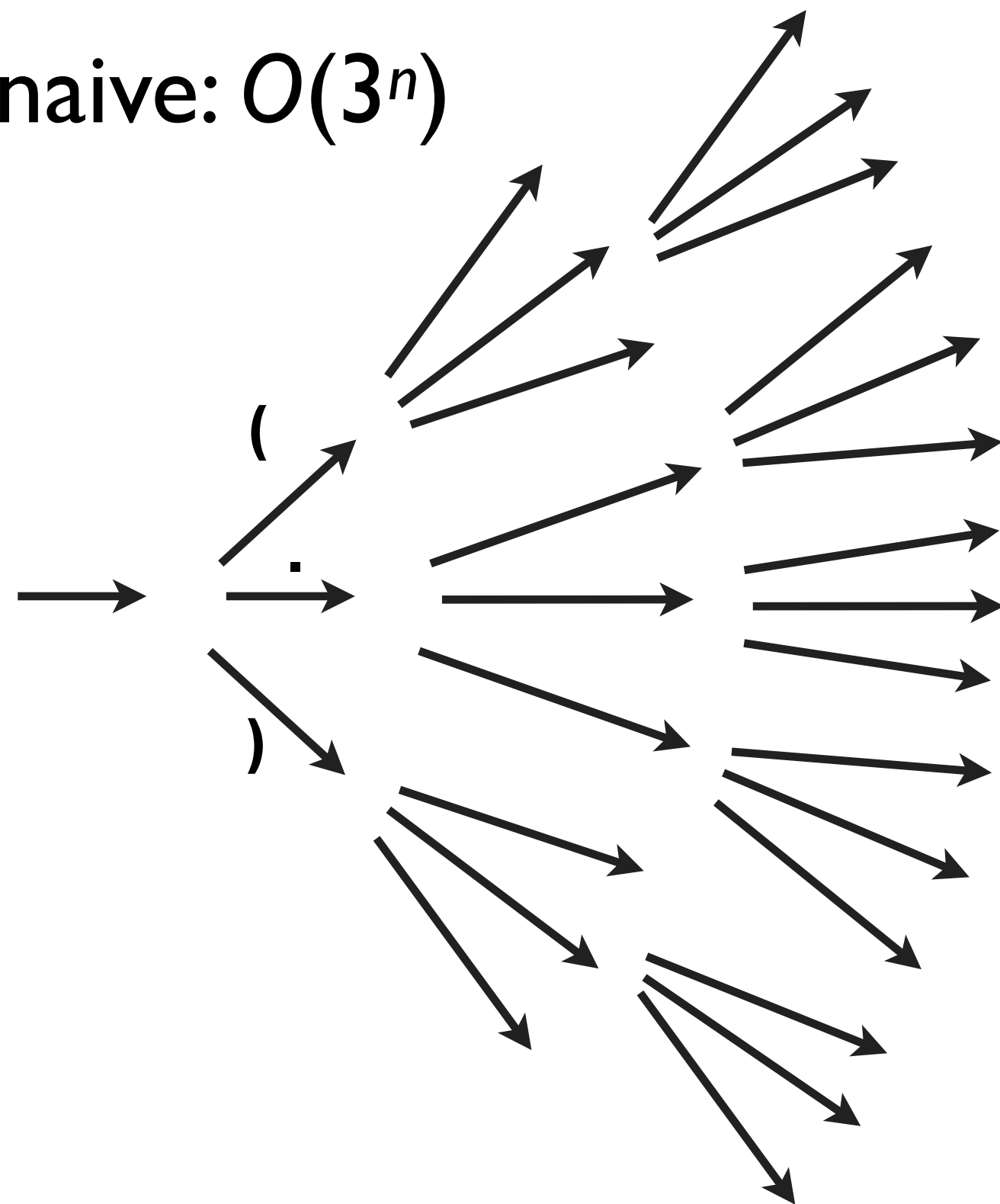


2019: **LinearFold**  $O(n)$  (approx.) RNA folding

# How to Fold RNAs in Linear-Time?


GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

- idea 0: tag each nucleotide from left to right
  - maintain a stack: push “(”, pop “)”, skip “.”
  - naive:  $O(3^n)$



# How to Fold RNAs in Linear-Time?

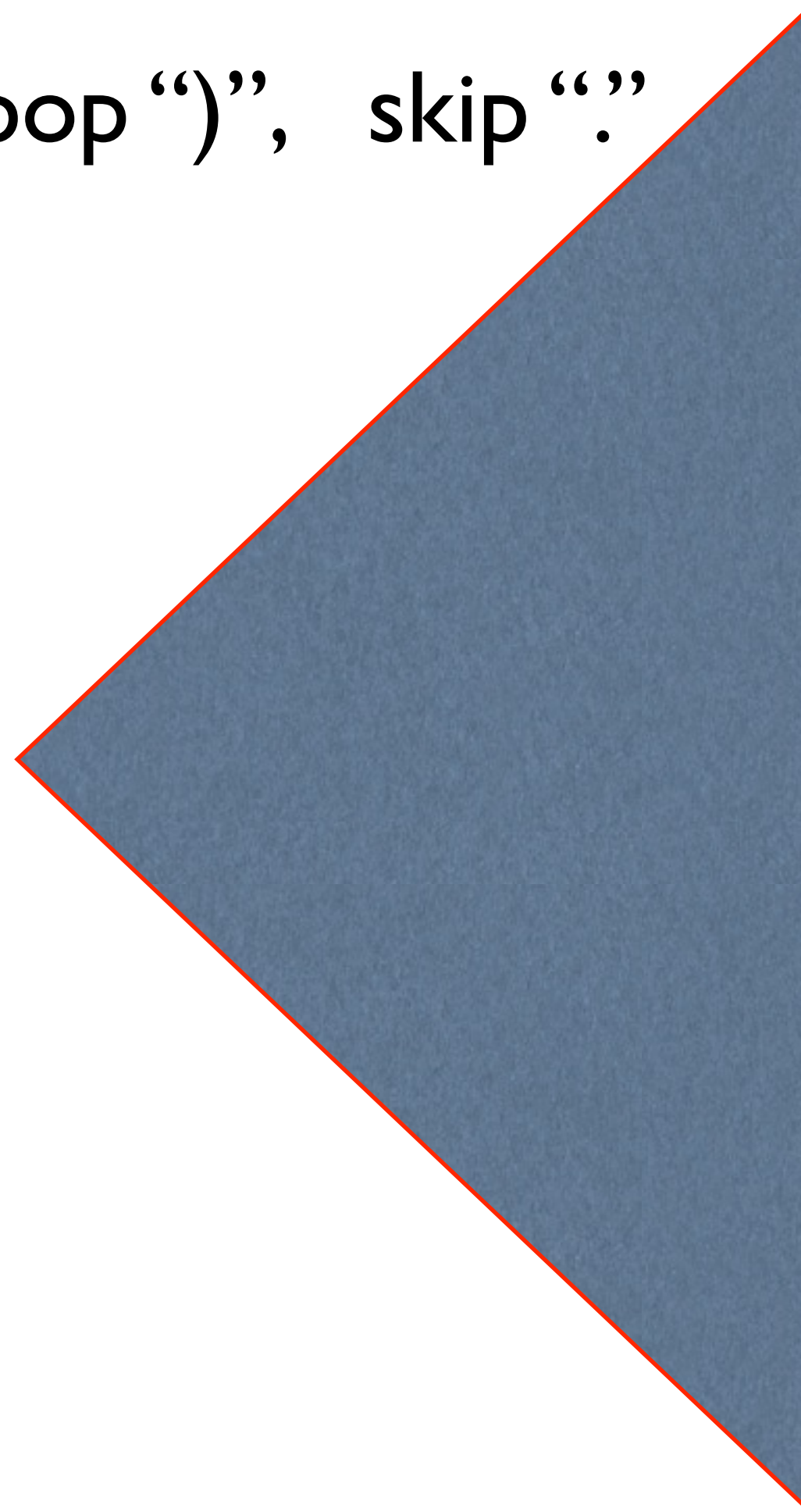
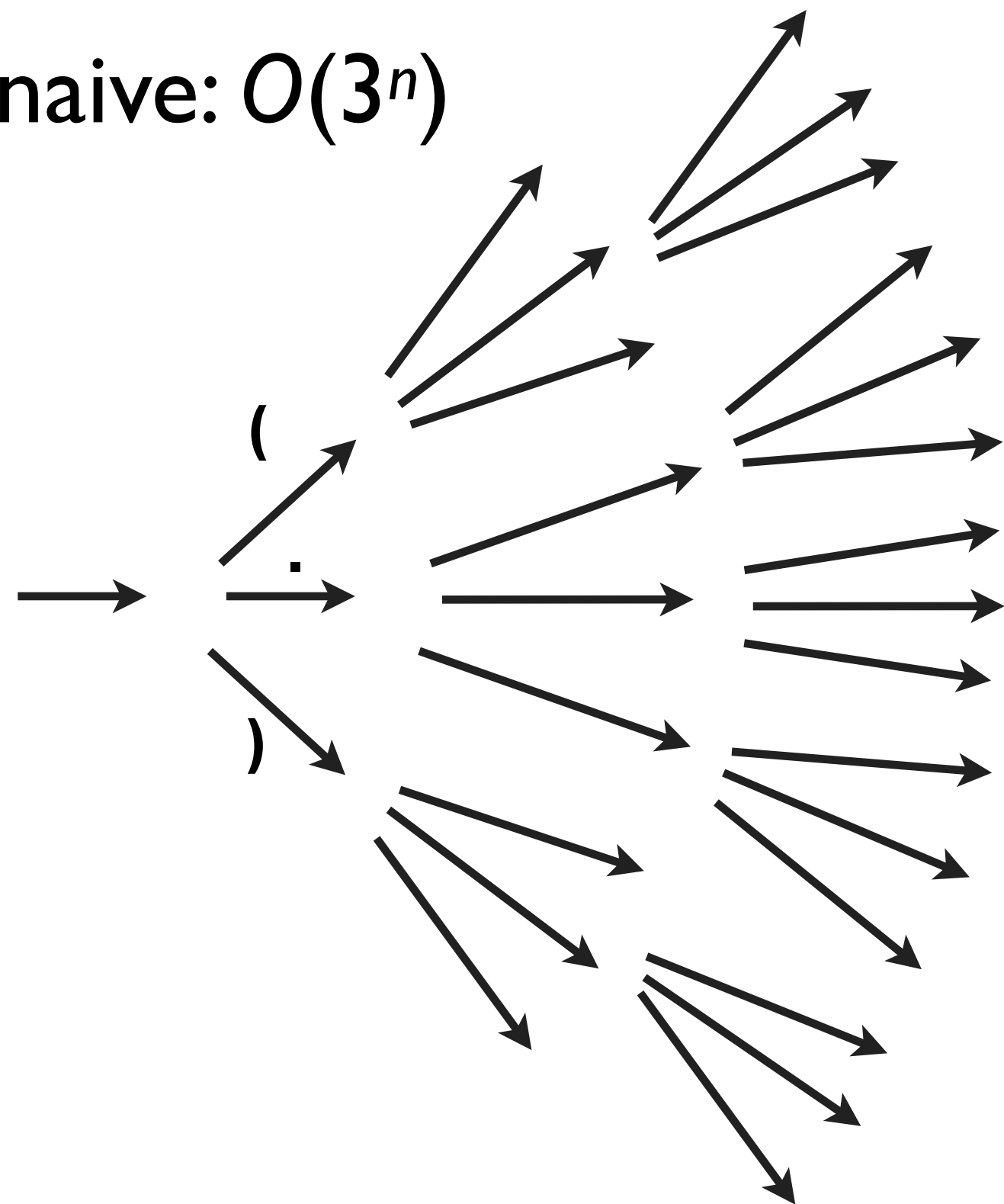
GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

(((((...(((...))))).(((...)))))...(((((...))))))...  


- idea 0: tag each nucleotide from left to right

- maintain a stack: push “(”, pop “)”, skip “.”

- naive:  $O(3^n)$



# How to Fold RNAs in Linear-Time?

GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

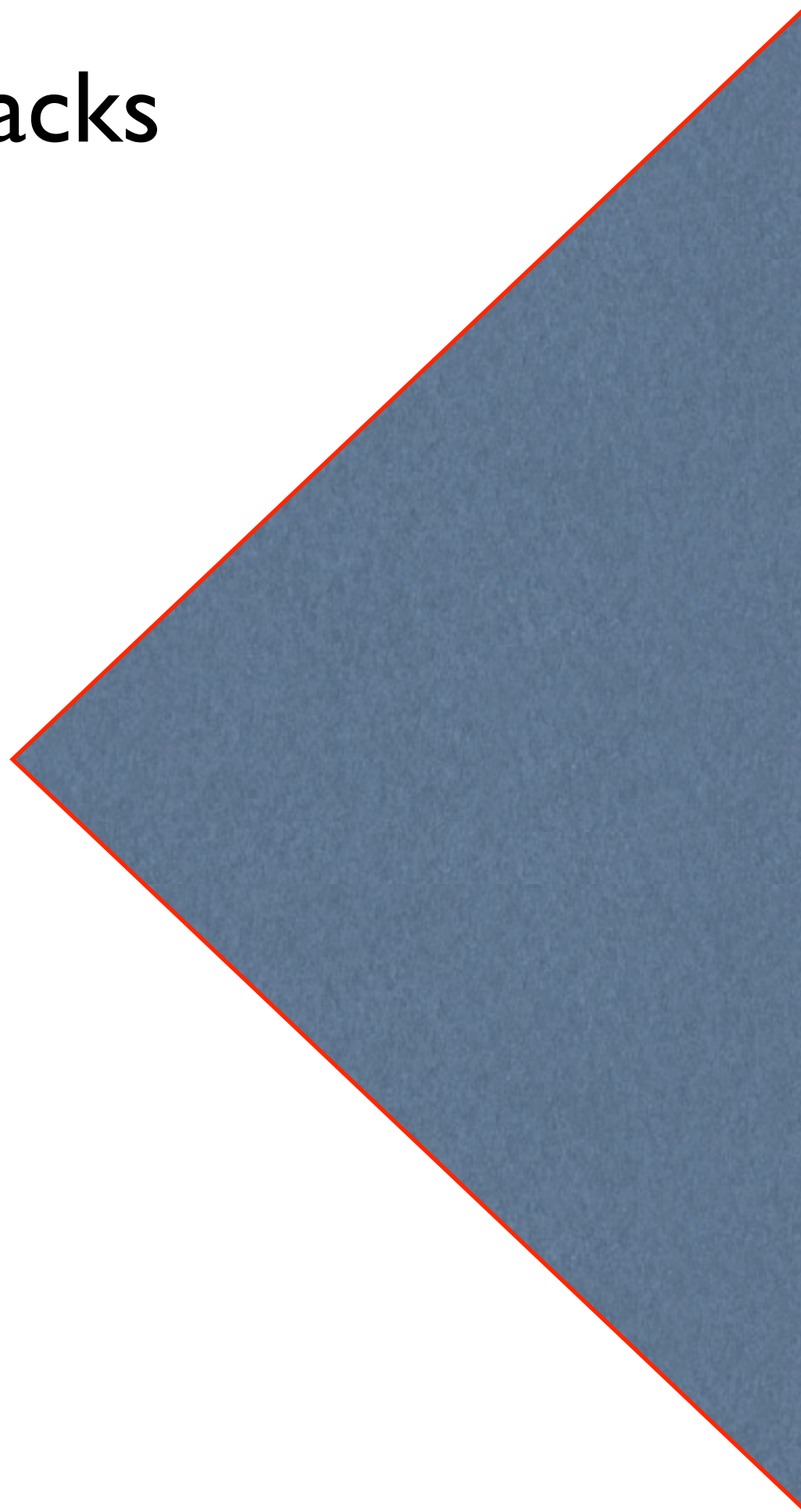
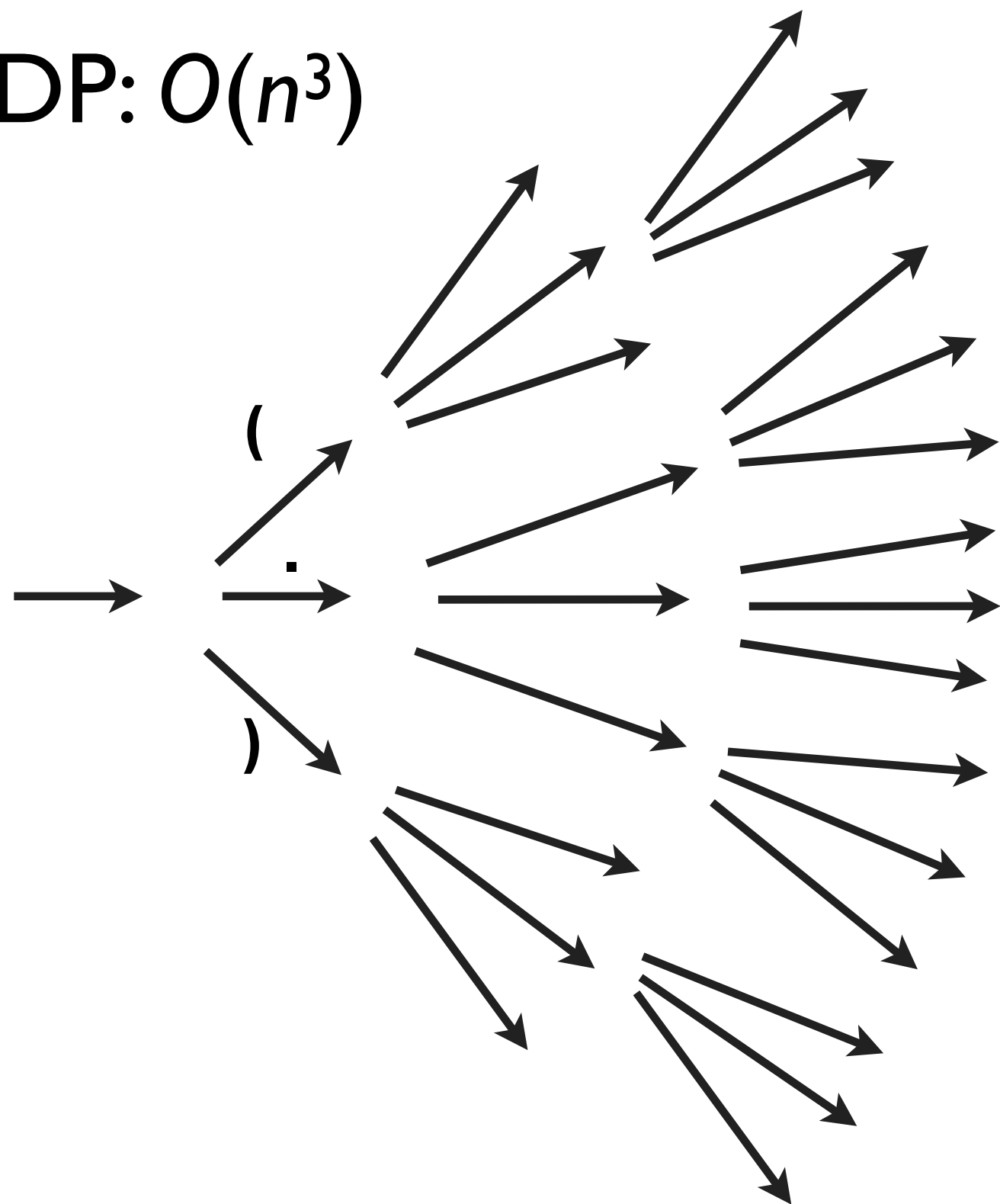
(((((...(((...)))))) . (((...)))) . . . . (((...))))))))) . . . .



- idea 1: DP by packing “equivalent states”


- maintain graph-structured stacks

- DP:  $O(n^3)$



# How to Fold RNAs in Linear-Time?

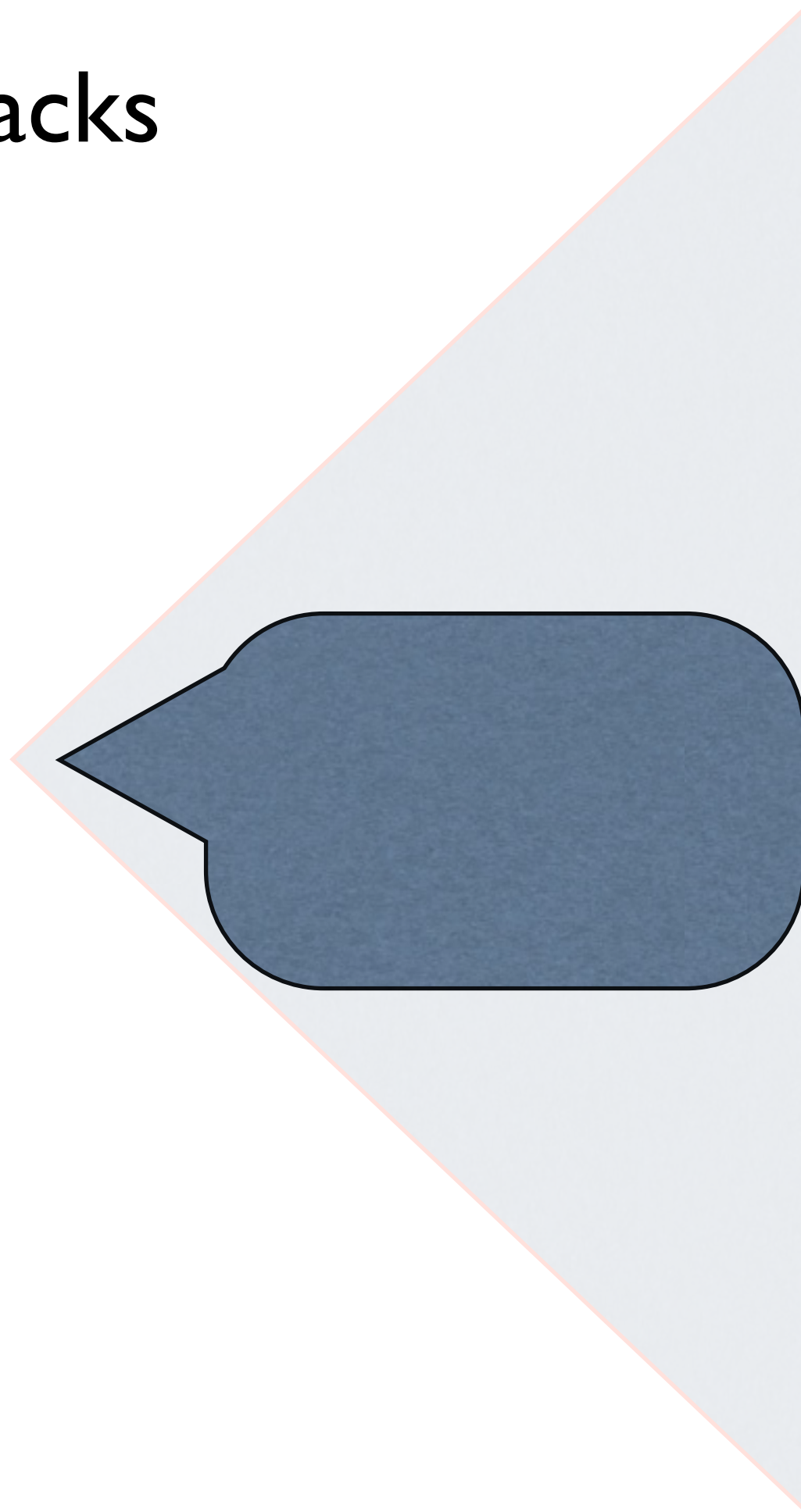
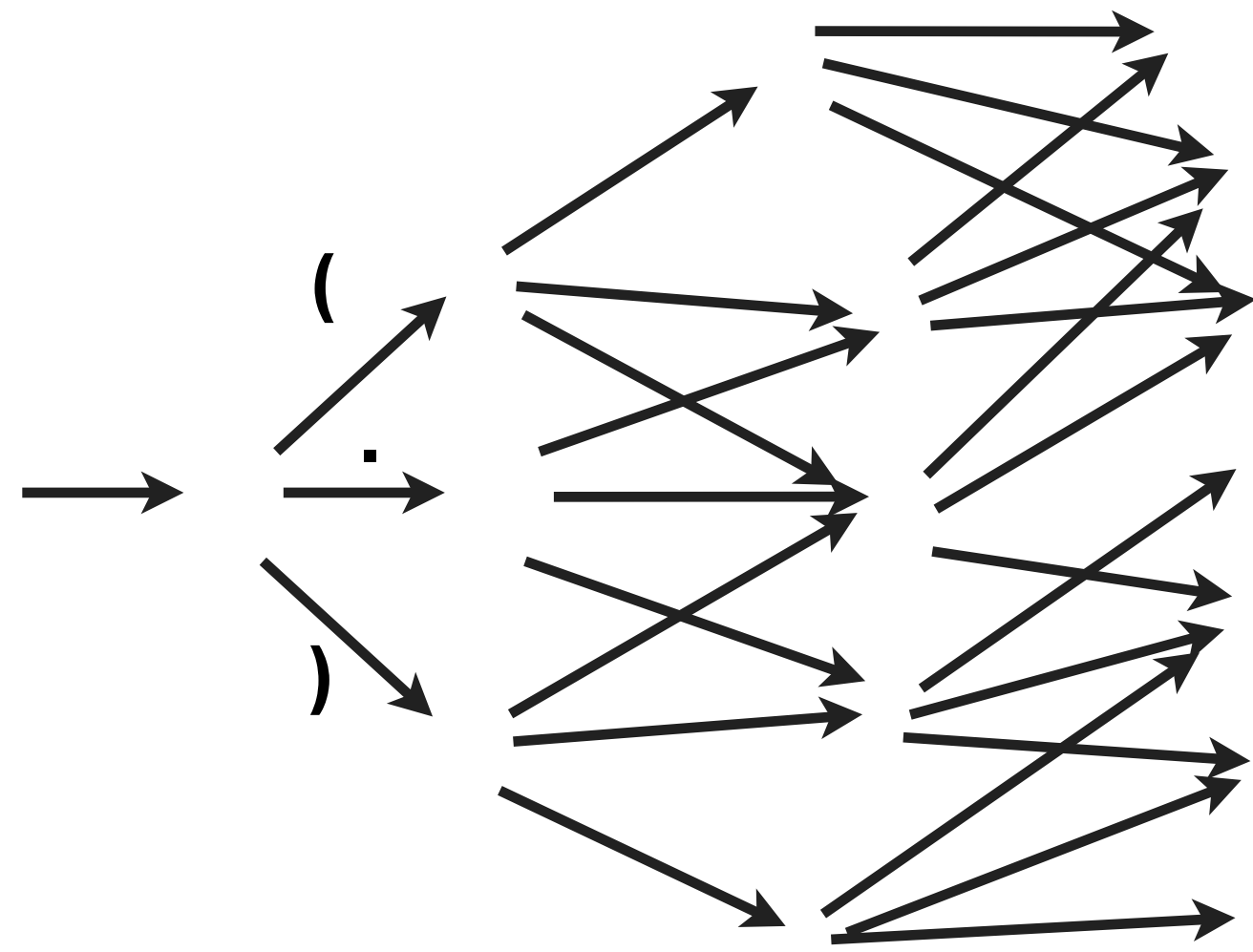
GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC~~CC~~GCUCCA

(((((...(((...))))).(((...)))))...(((...)))))))))...  


- idea 1: DP by packing “equivalent states”

- maintain graph-structured stacks

- DP:  $O(n^3)$



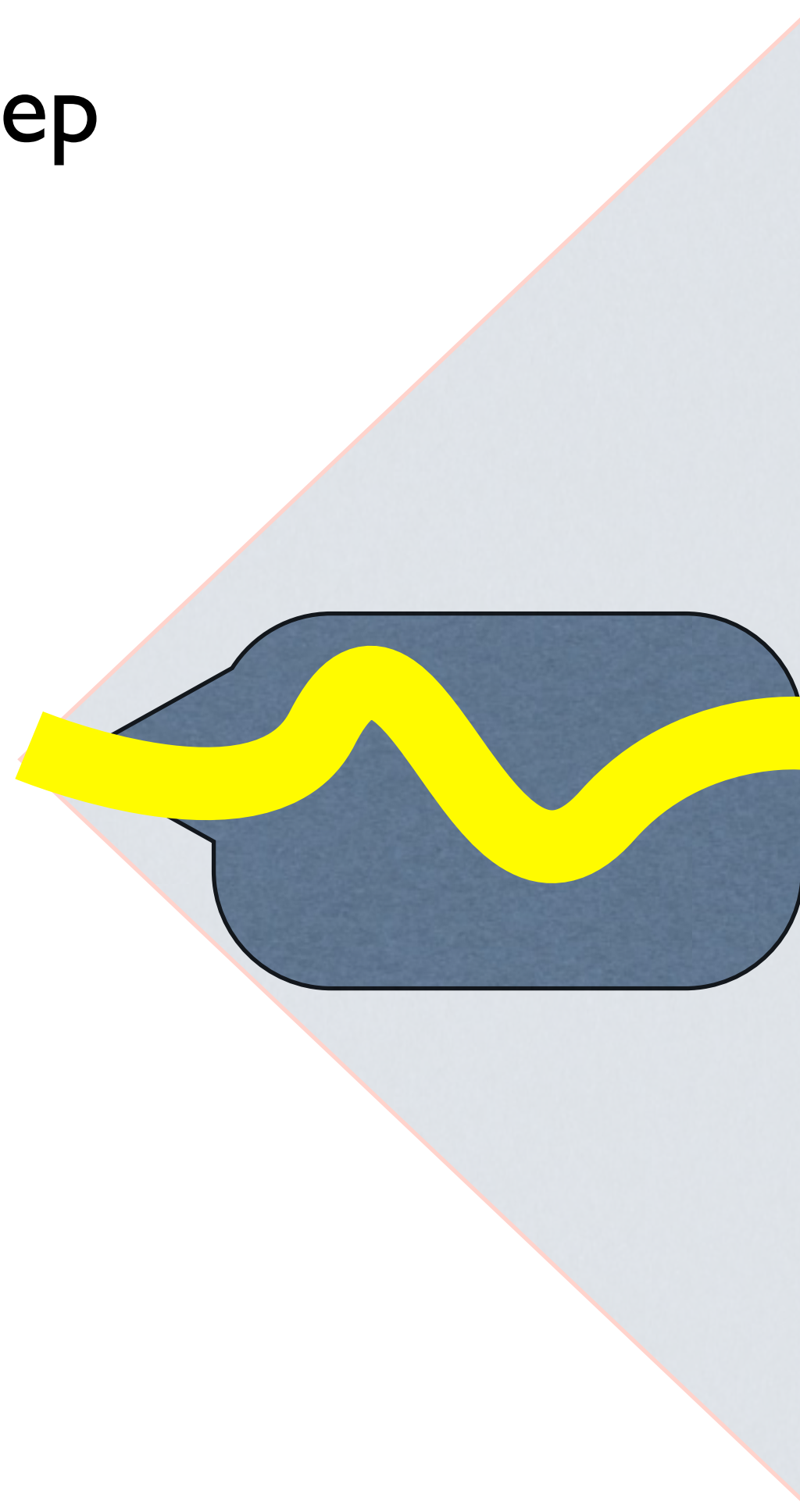
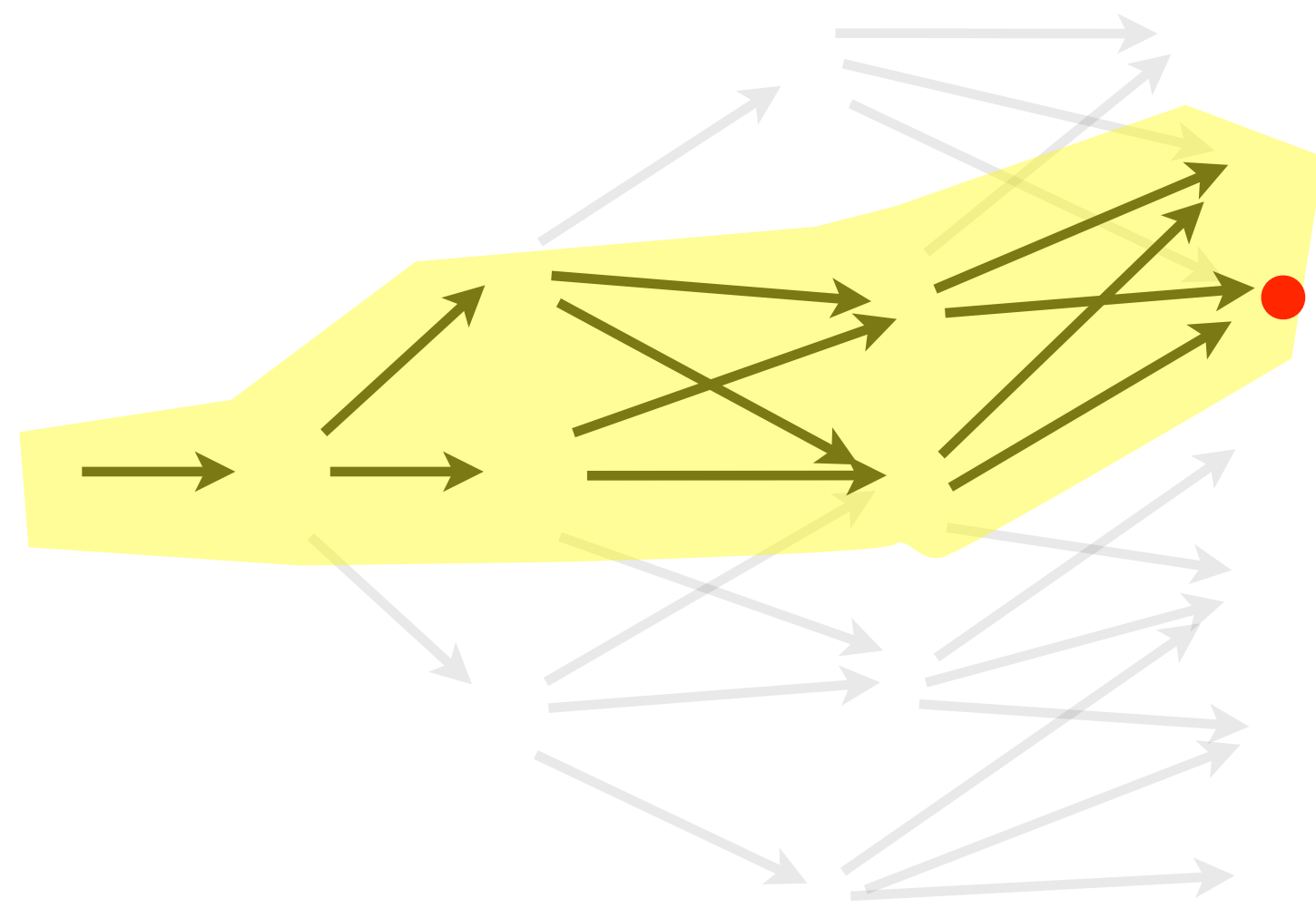
# How to Fold RNAs in Linear-Time?

GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

(((((...(((...)))))) . (((...)))) . . . . (((...))))))))) . . . .



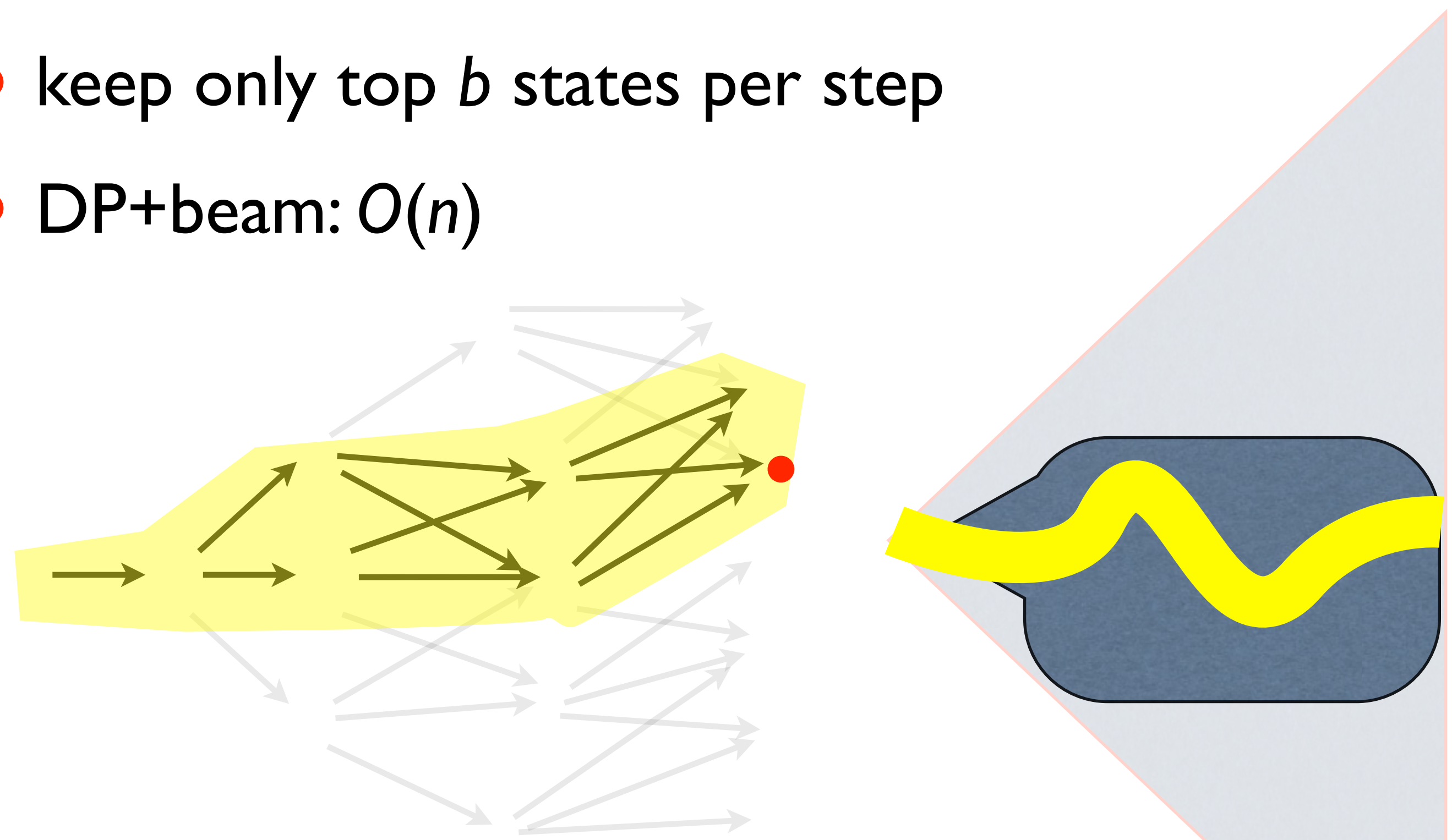
- idea 2: approximate search: beam pruning
  - keep only top  $b$  states per step
  - DP+beam:  $O(n)$



# How to Fold RNAs in Linear-Time?

GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC~~CCGCUCCA~~  
((((((...((((...)))))).((((...))))). . . . .((((((...))))))))) . . . .  
→

- idea 2: approximate search: beam pruning
  - keep only top  $b$  states per step
  - DP+beam:  $O(n)$



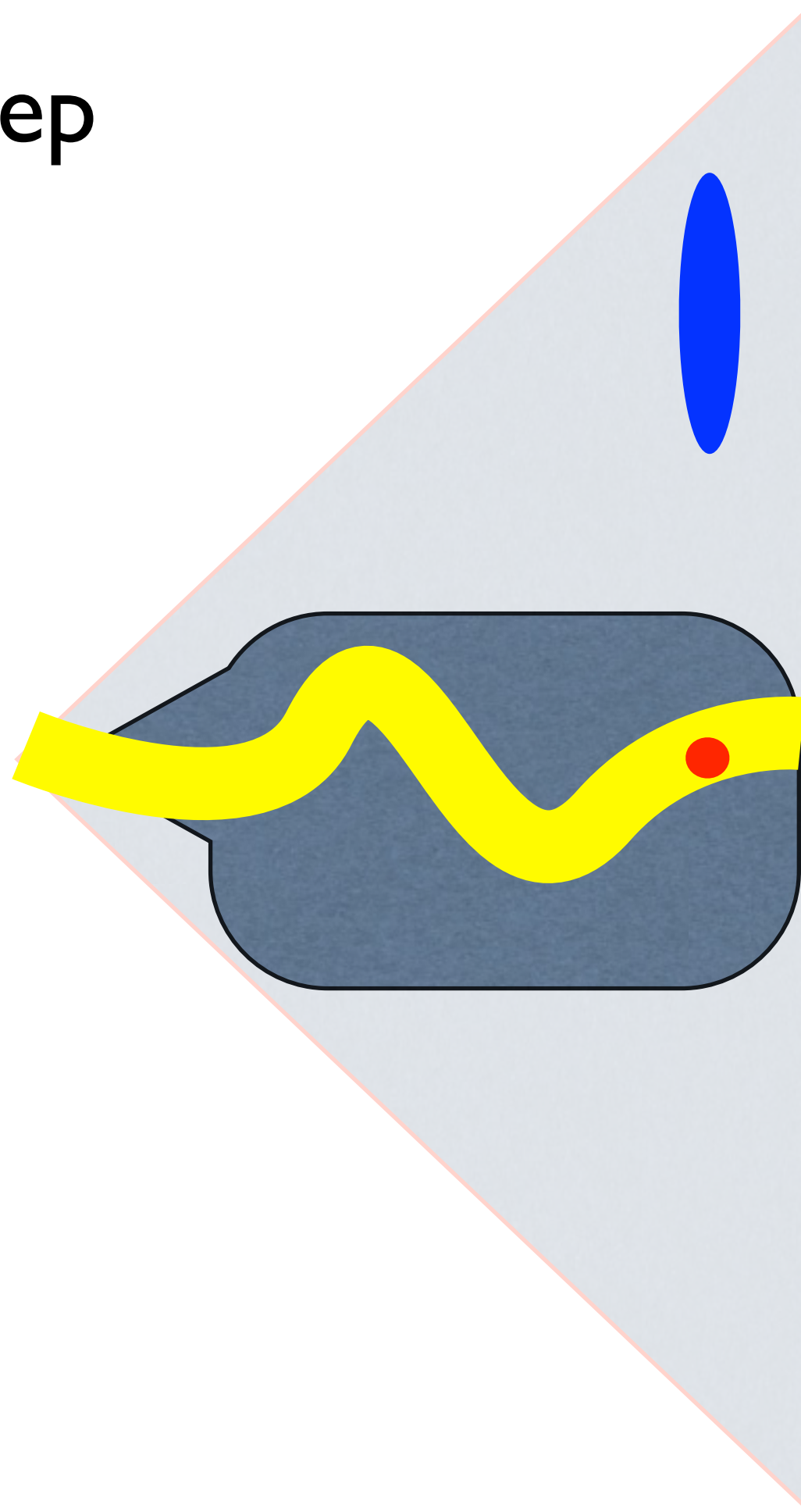
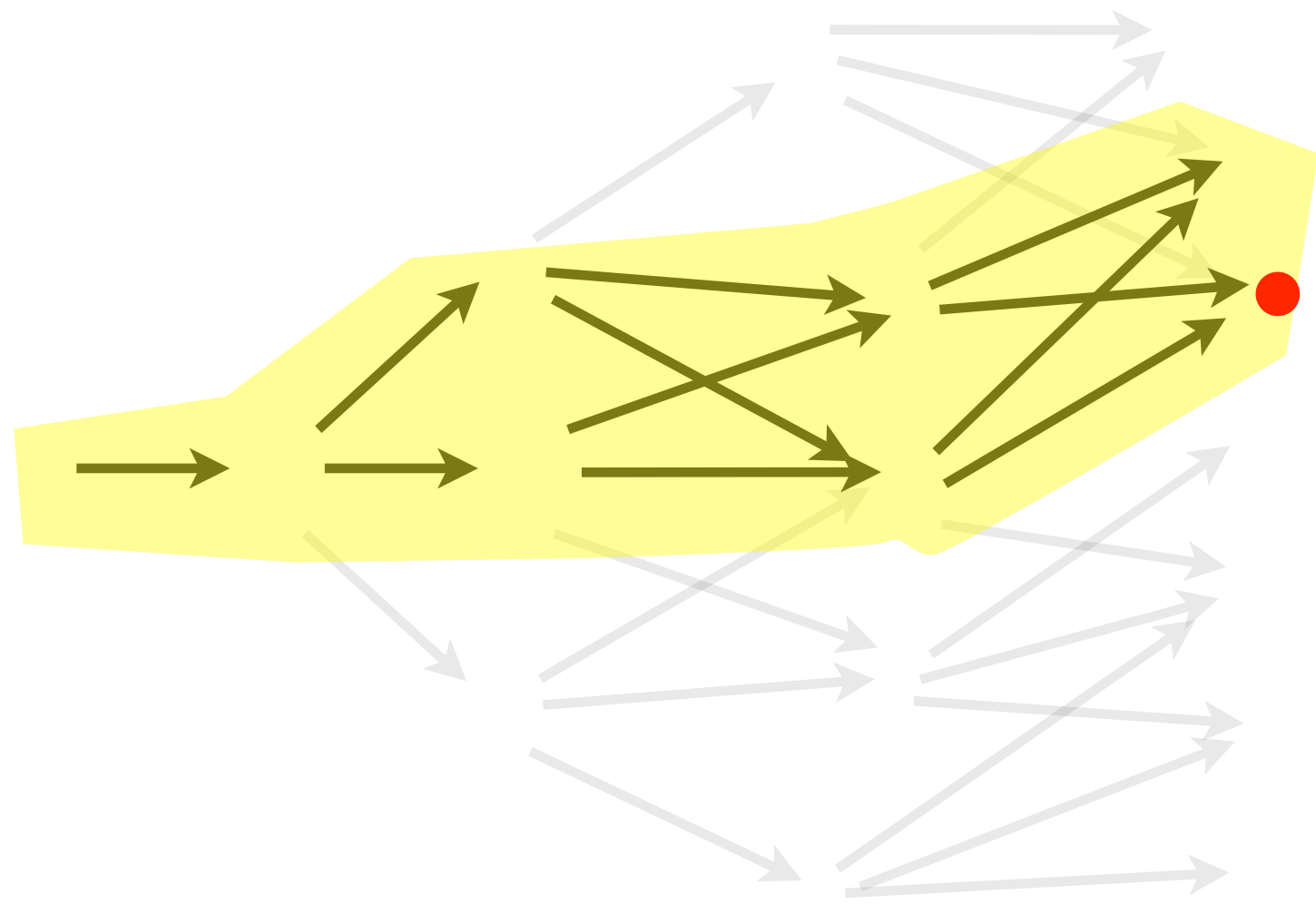
each DP state corresponds to exponentially many non-DP states

# How to Fold RNAs in Linear-Time?

GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC~~CCGCUCCA~~

(((((...(((...))))).(((...)))))...(((...)))))))))...  
\_\_\_\_\_→

- idea 2: approximate search: beam pruning
  - keep only top  $b$  states per step
  - DP+beam:  $O(n)$



each **DP state** corresponds to exponentially many **non-DP states**



# How to Fold RNAs in Linear-Time?

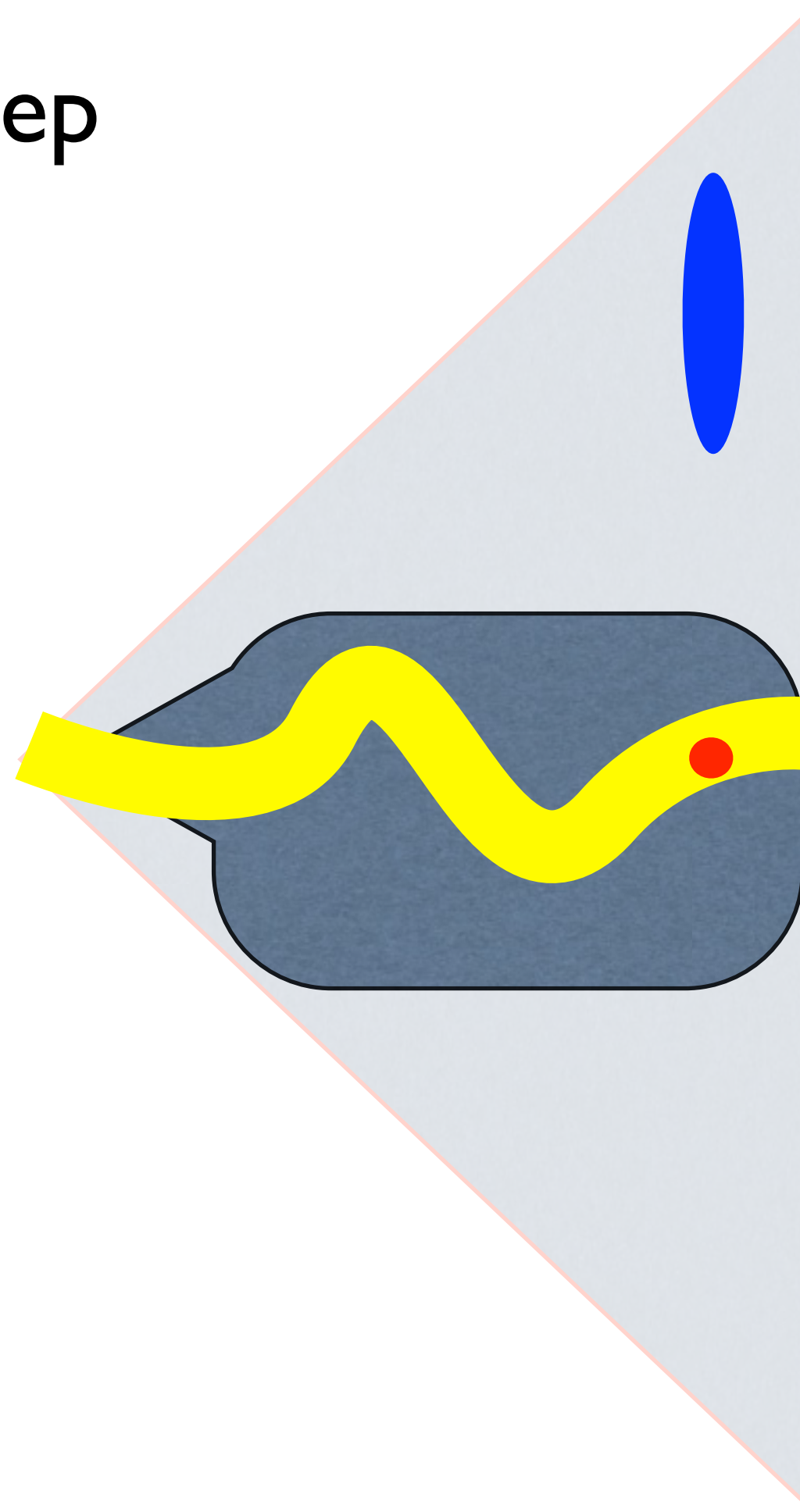
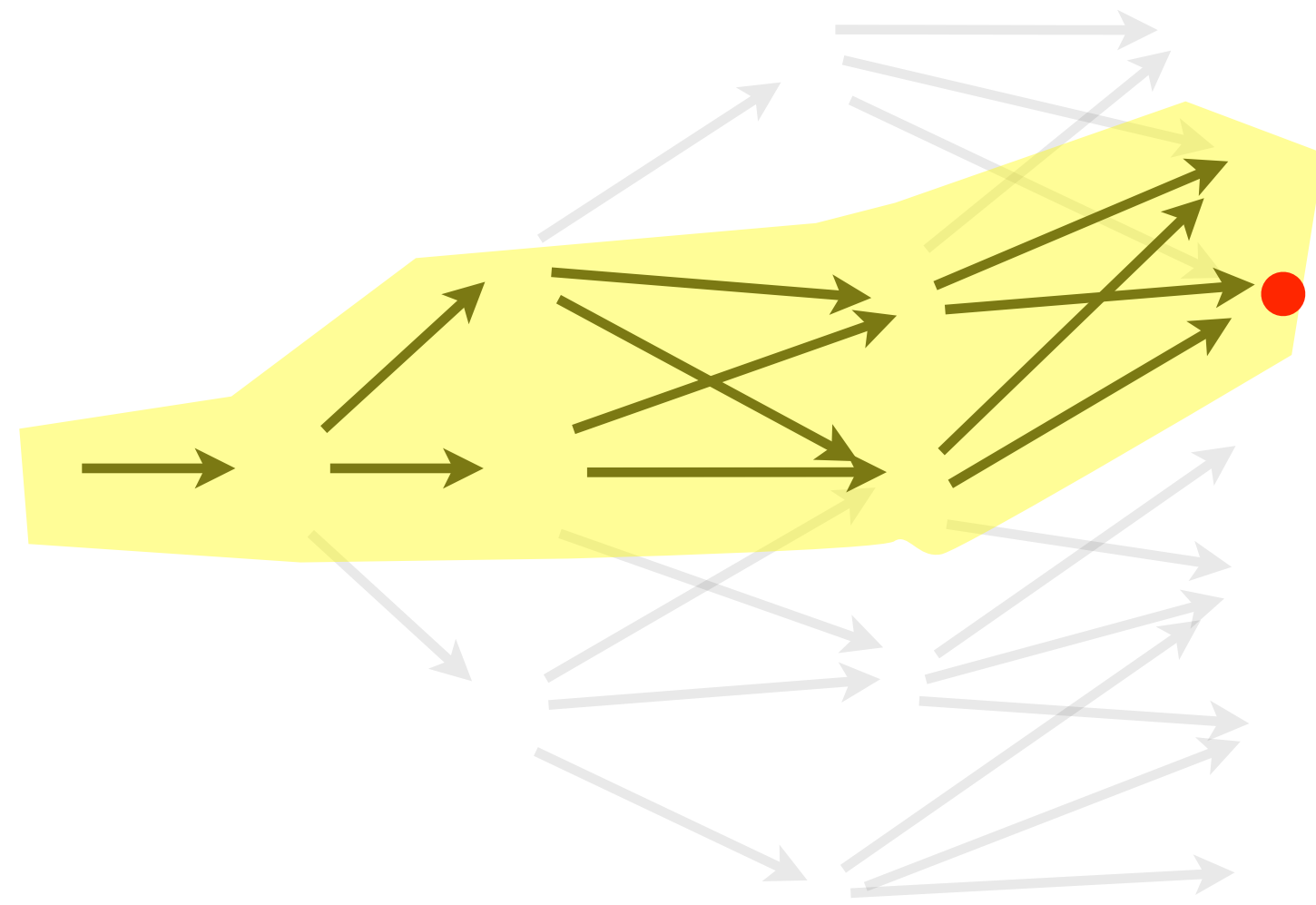
GCGGGAAUAGCUCAGUUGGUAGAGCACGACCUUGCCAAGGUCGGGGUCGCGAGUUCGAGUCUCGUUUC CCGCUCCA

(((((...(((...))))).(((...))))). . . . ((((((...))))))))) . . .

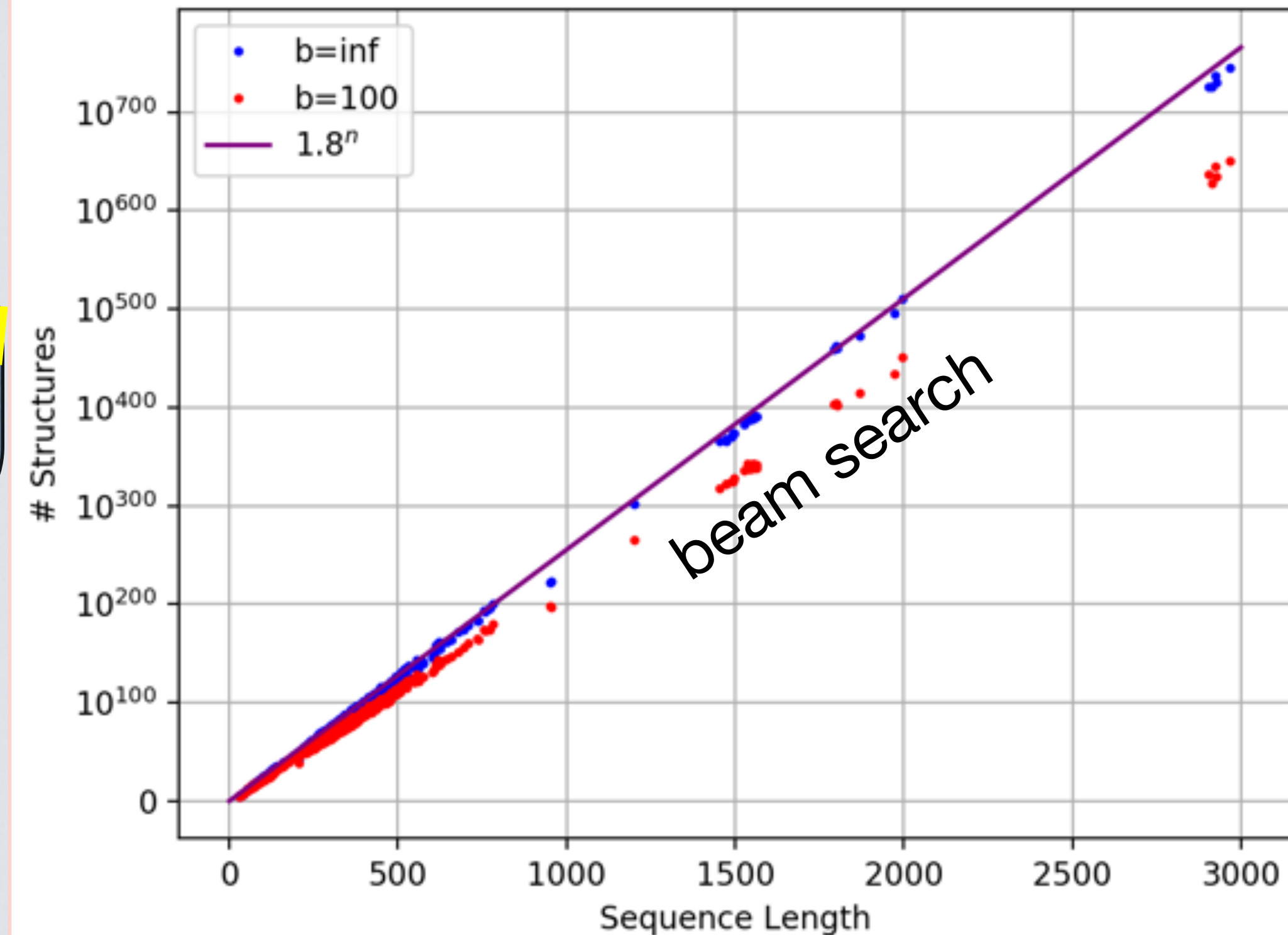


- idea 2: approximate search: beam pruning

- keep only top  $b$  states per step
- DP+beam:  $O(n)$

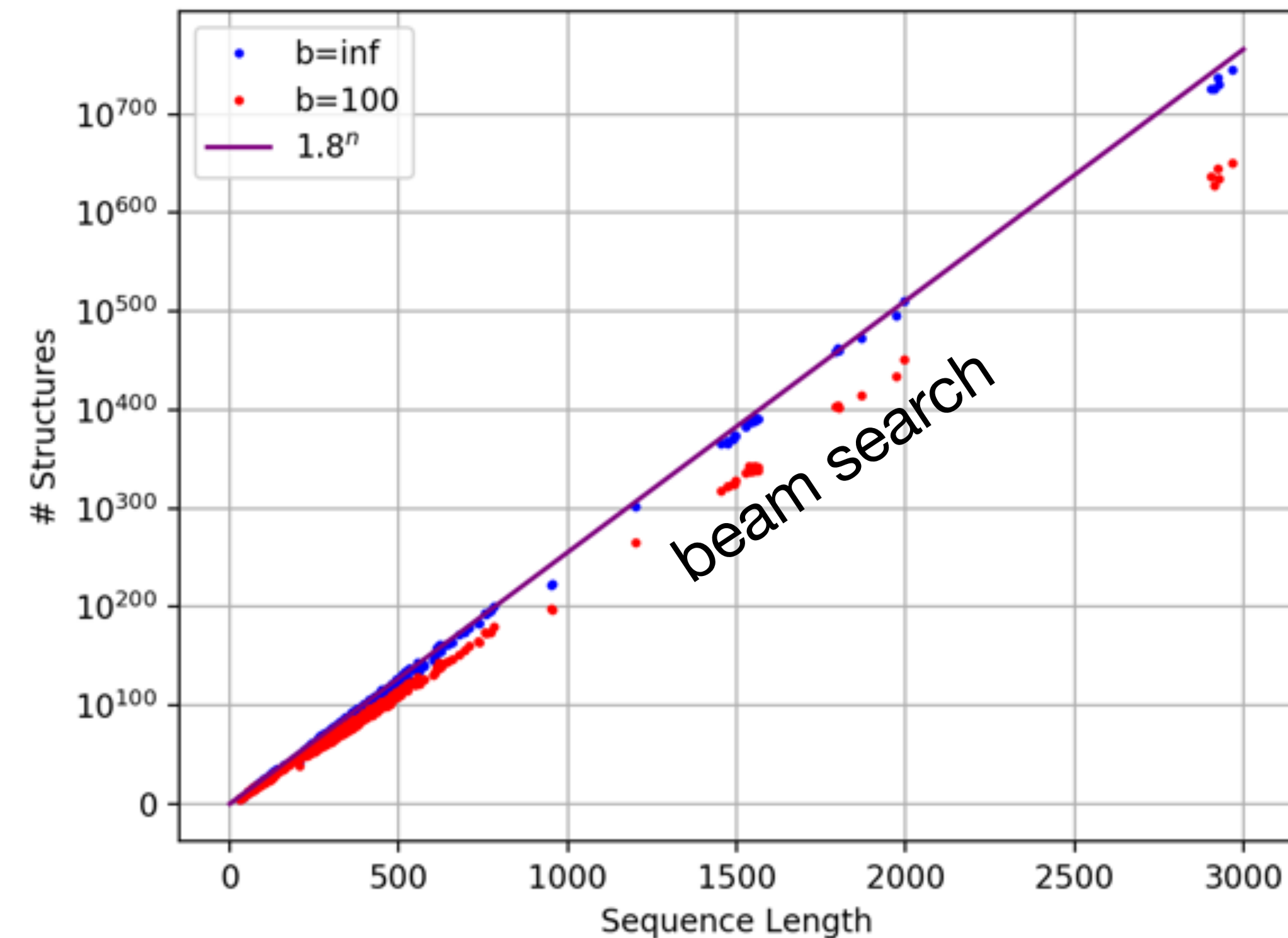
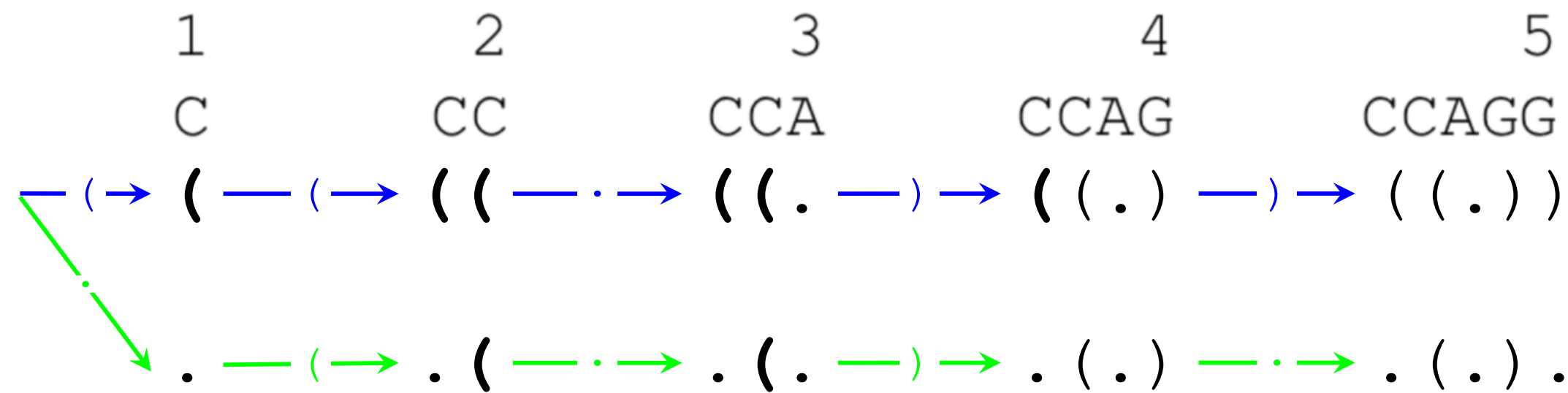


each **DP state** corresponds to exponentially many **non-DP states**



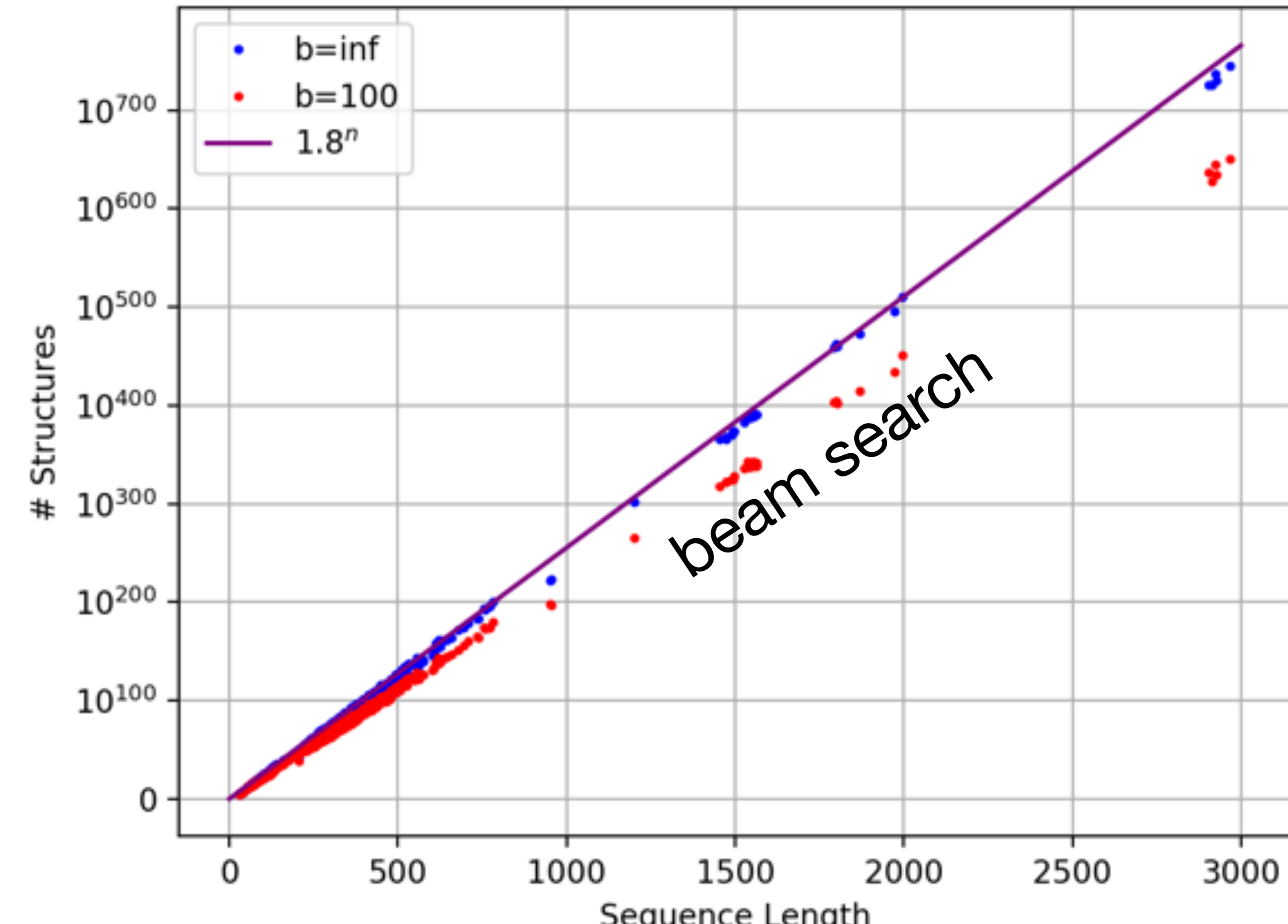
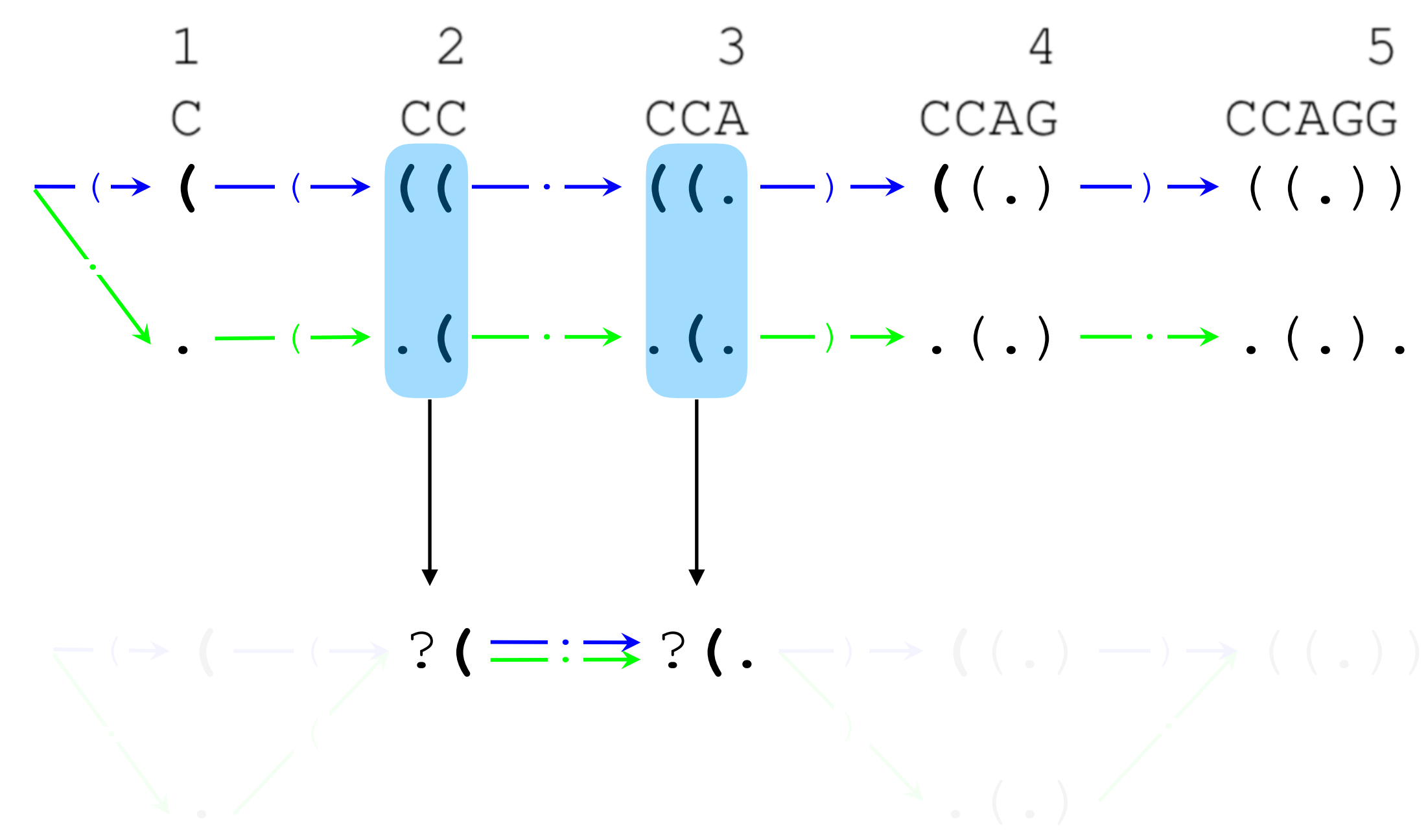
# Details: Packing “Temporarily Equivalent” States

- two states are “temporarily equivalent” if they share the same stack top index
- because they are looking for the same nucleotide(s) to match this stack-top nucleotide
- we pack them as a single state until a matching is found, when they unpack
- this is how we can explore exponentially many structures in linear time



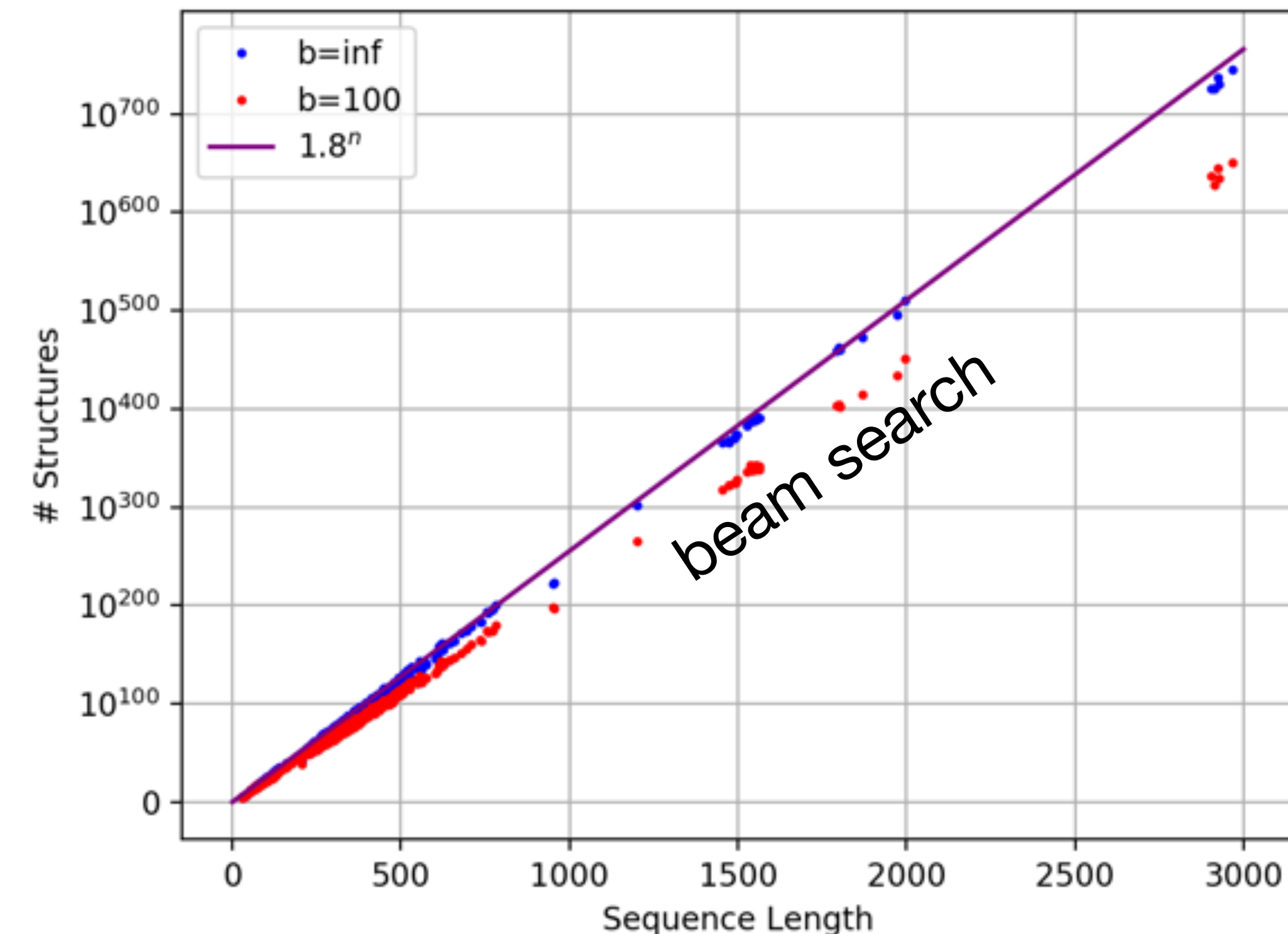
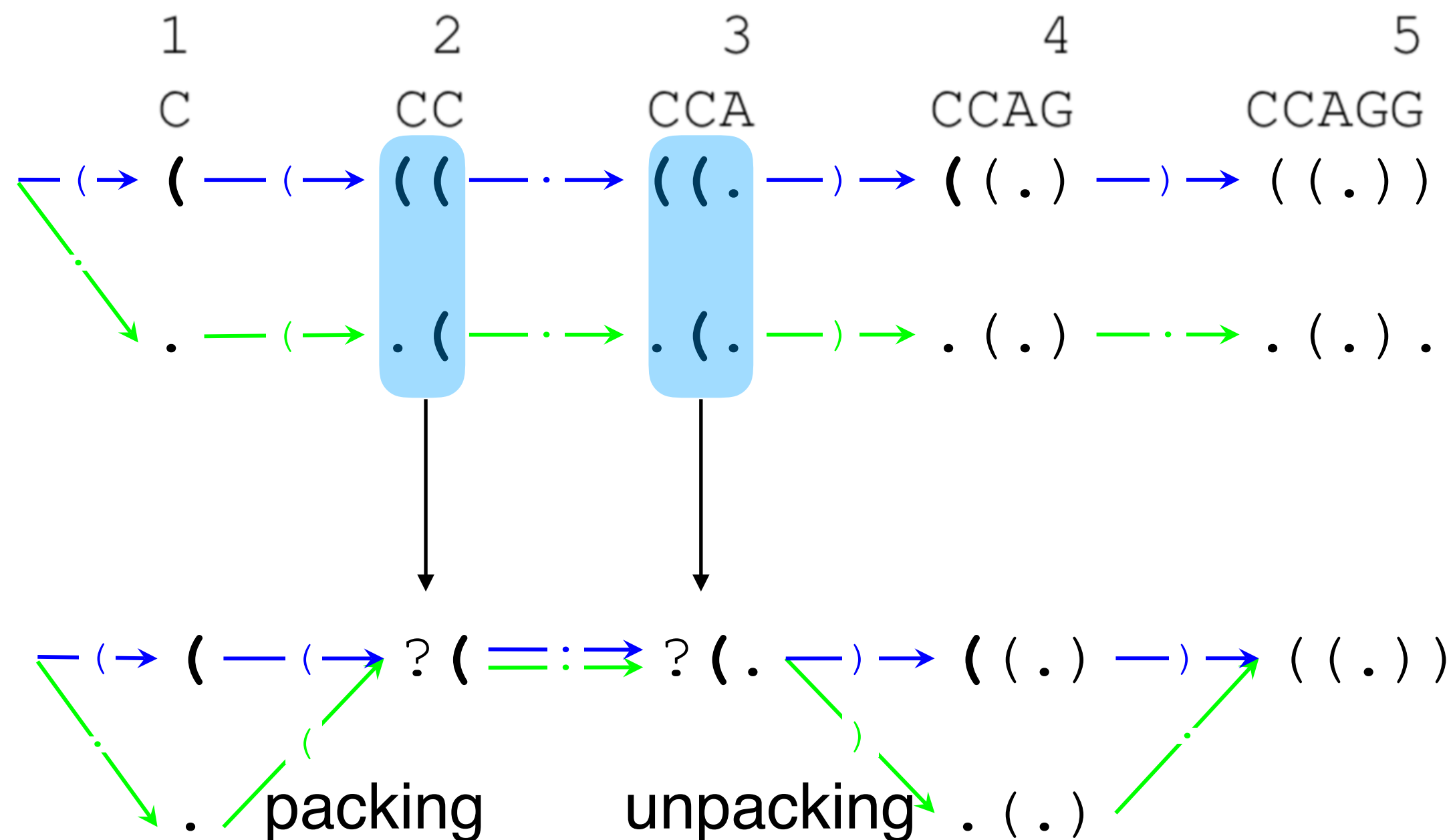
# Details: Packing “Temporarily Equivalent” States

- two states are “temporarily equivalent” if they share the same stack top index
- because they are looking for the same nucleotide(s) to match this stack-top nucleotide
- we pack them as a single state until a matching is found, when they unpack
- this is how we can explore exponentially many structures in linear time



# Details: Packing “Temporarily Equivalent” States

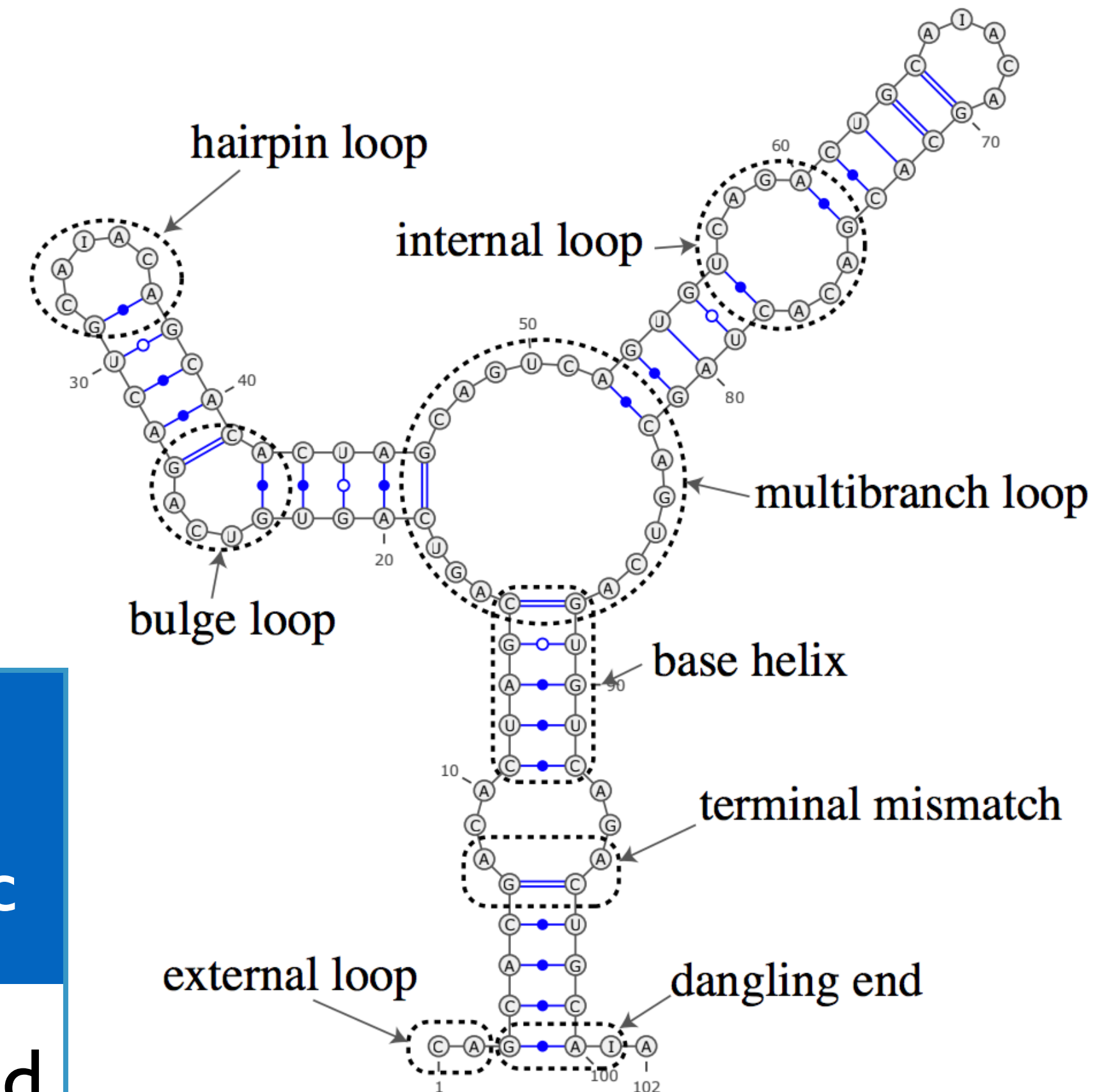
- two states are “temporarily equivalent” if they share the same stack top index
- because they are looking for the same nucleotide(s) to match this stack-top nucleotide
- we pack them as a single state until a matching is found, when they unpack
- this is how we can explore exponentially many structures in linear time



# Results

# LinearFold with SOTA Prediction Models

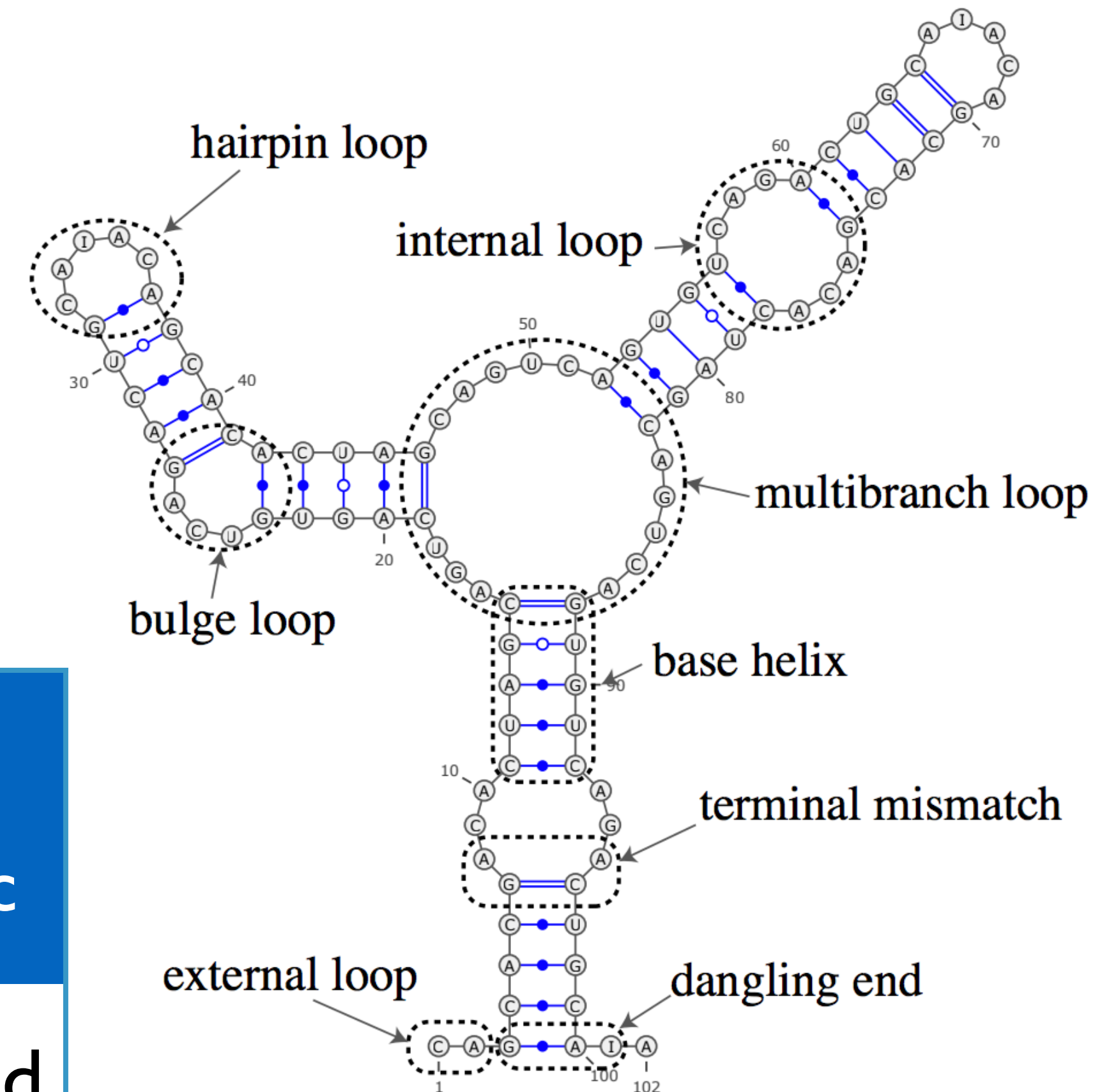
- models from two widely-used folding engines
  - CONTRAfold MFE (machine-learned)
  - Vienna RNAfold (thermodynamic)
- we linearized both systems from  $O(n^3)$  to  $O(n)$



	<i>efficiency</i>		<i>systems</i>	
	time	space	machine-learned	thermo-dynamic
baselines	$O(n^3)$	$O(n^2)$	<b>C</b> ONTRAfold	<b>V</b> ienna RNAfold
our work	$O(n)$	$O(n)$	LinearFold- <b>C</b>	LinearFold- <b>V</b>

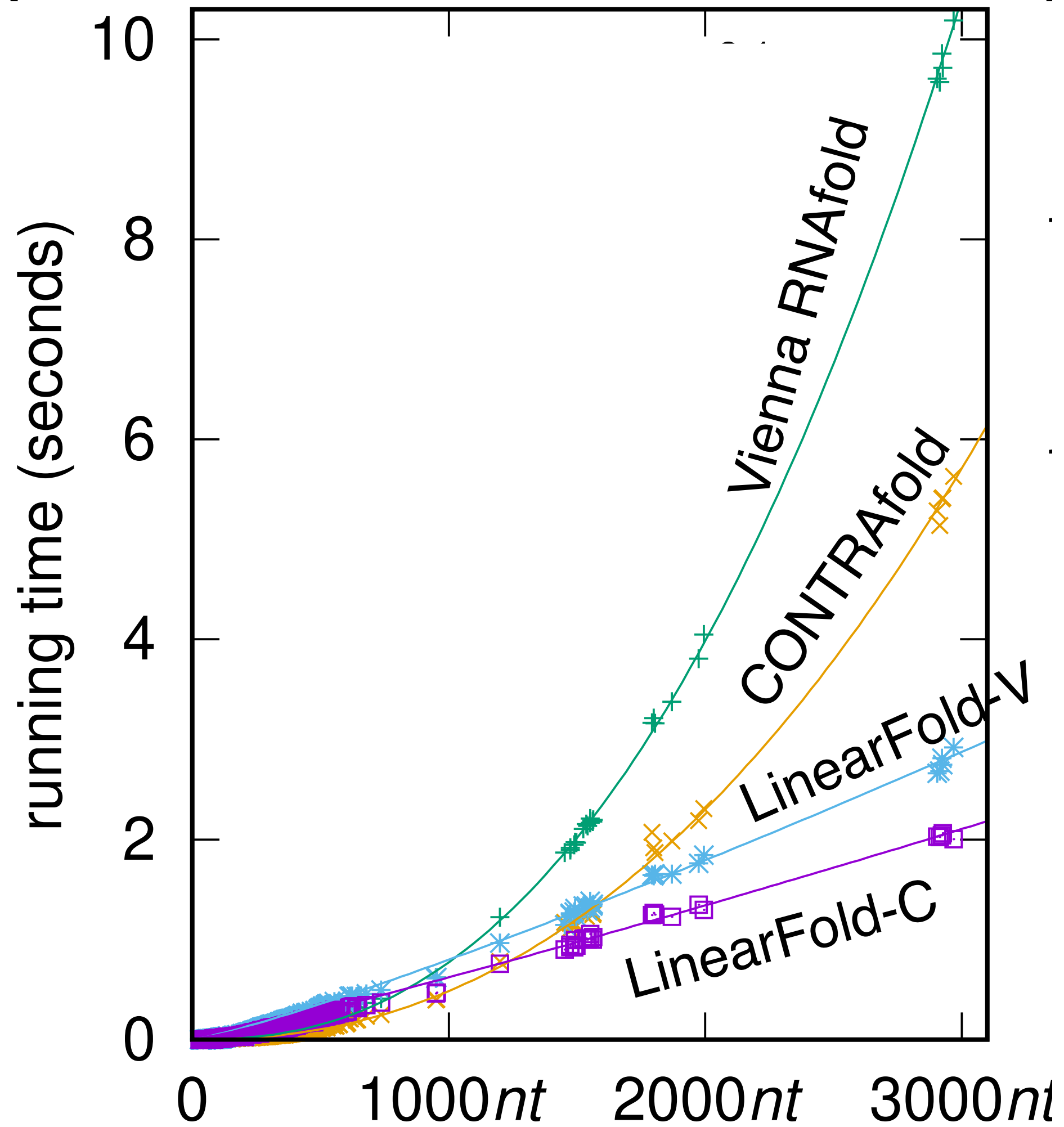
# LinearFold with SOTA Prediction Models

- models from two widely-used folding engines
  - CONTRAfold MFE (machine-learned)
  - Vienna RNAfold (thermodynamic)
- we linearized both systems from  $O(n^3)$  to  $O(n)$

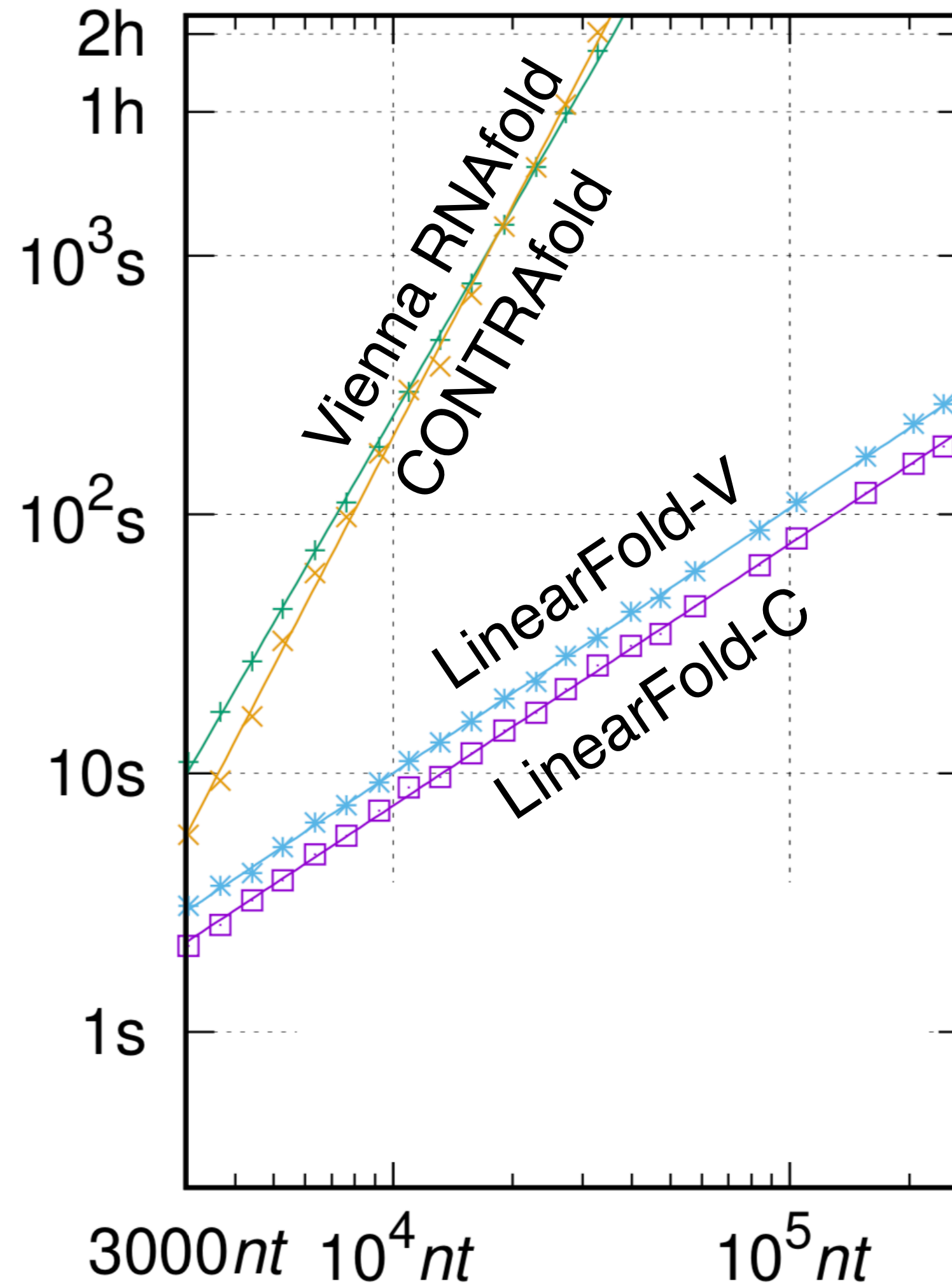


	<i>efficiency</i>		<i>systems</i>	
	time	space	machine-learned	thermo-dynamic
baselines	$O(n^3)$	$O(n^2)$	<b>CONTRAfold</b>	<b>Vienna RNAfold</b>
our work	$O(n)$	$O(n)$	<b>LinearFold-C</b>	<b>LinearFold-V</b>

# Efficiency & Scalability: $O(n)$ time, $O(n)$ memory



Archive II data set  
(~3,000 seqs, max len: ~3,000 nt)



RNACentral data set  
(sampled, max len: ~250,000 nt)

10,000nt (~HIV)  
4min  $\rightarrow$  7s

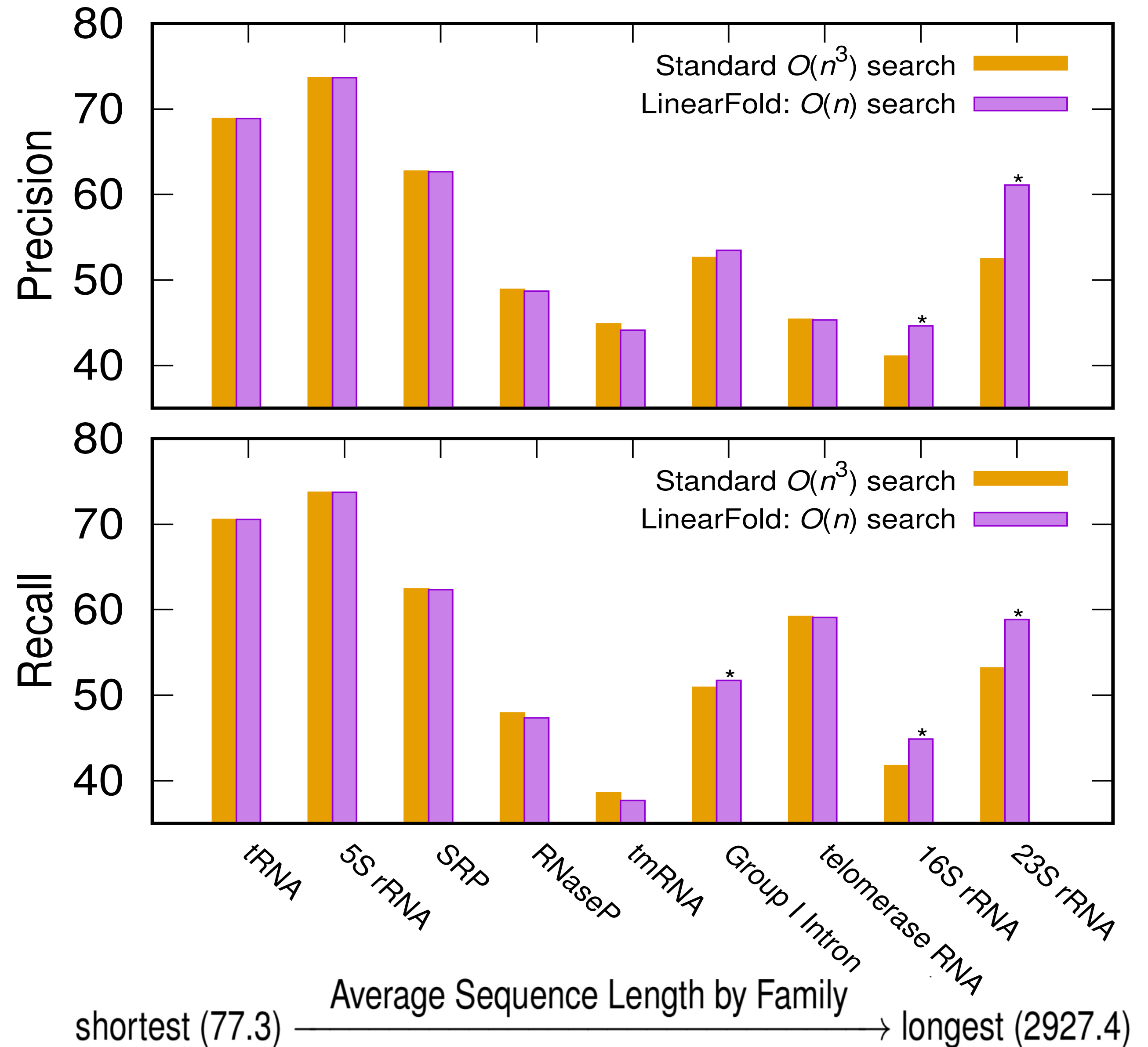
244,296nt  
(longest in RNACentral)  
 $\sim$ 200hrs  $\rightarrow$  120s



# Accuracy

- Tested on Archive II dataset (on a family-by-family basis)
- significantly better on 3 long families
- biggest boost on the longest families: 16S/23S rRNAs
- LinearFold-V vs. Vienna is similar

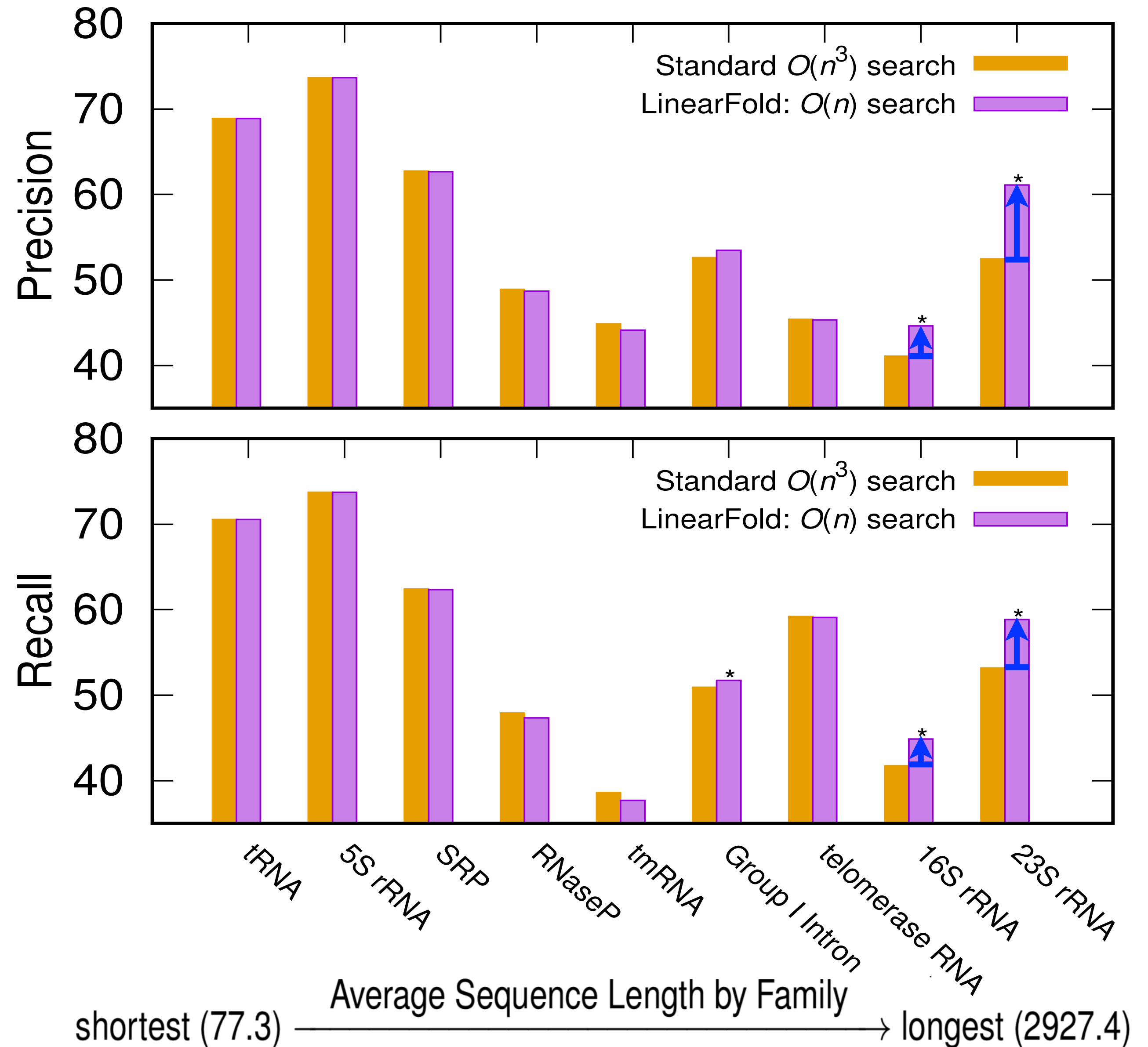
	(precision) PPV	(recall) Sensitivity
Overall		
CONTRAFold MFE	54.51	55.36
LinearFold-C	55.84 (+1.3)	56.24 (+0.9)
Vienna RNAfold	50.22	58.74
LinearFold-V	50.51 (+0.3)	58.97 (+0.2)



# Accuracy

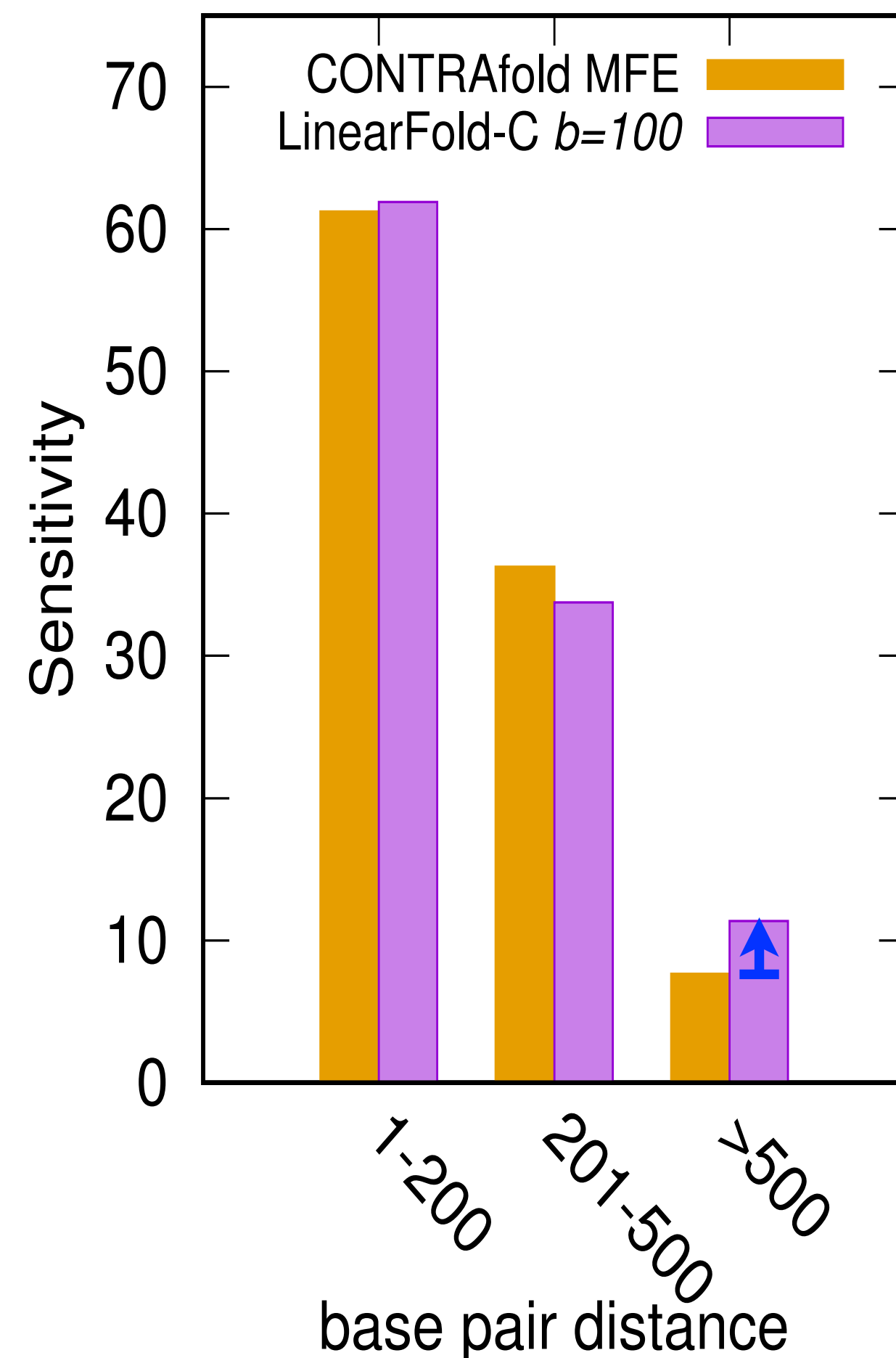
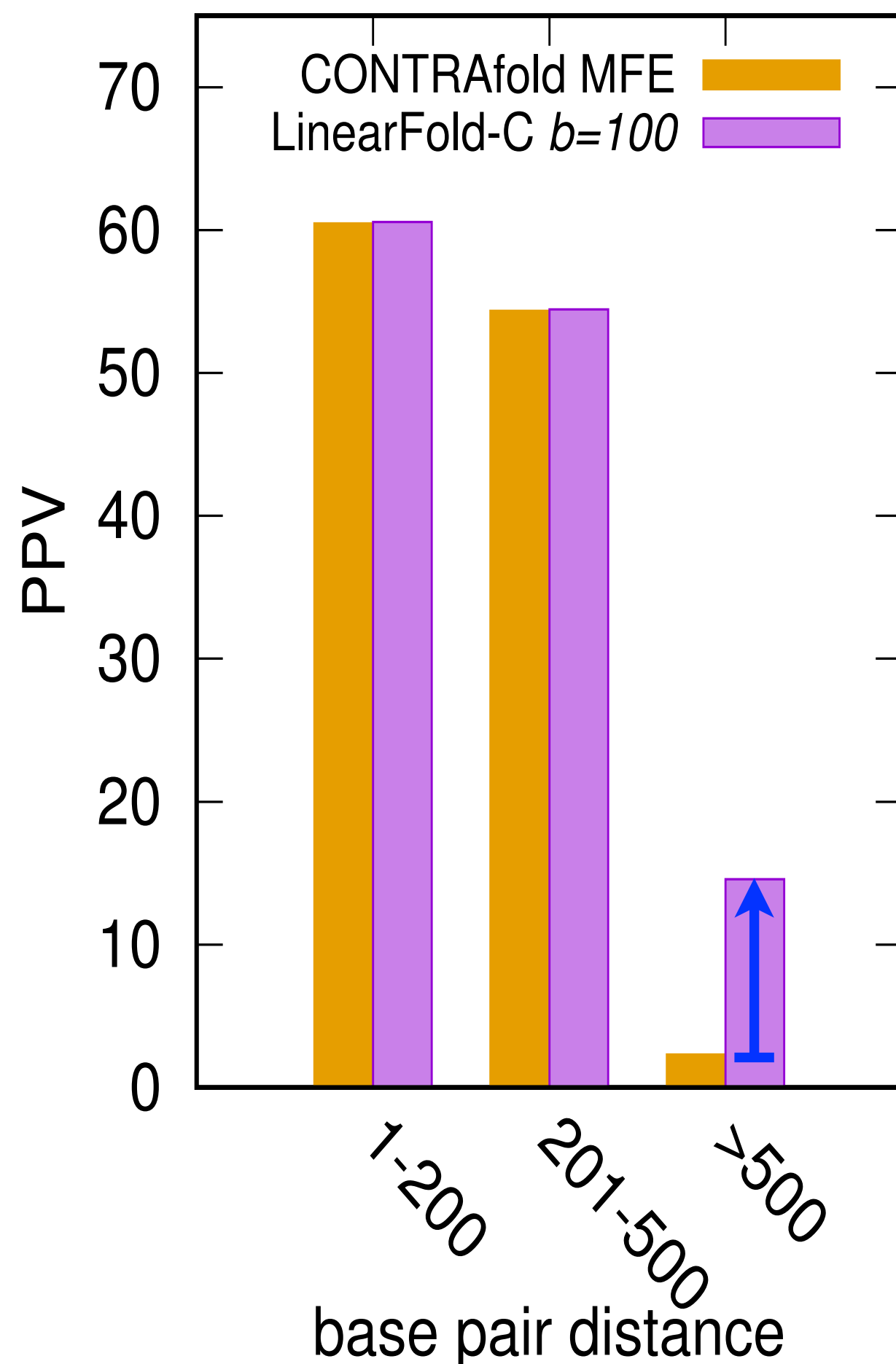
- Tested on Archive II dataset (on a family-by-family basis)
- significantly better on 3 long families
- biggest boost on the longest families: 16S/23S rRNAs
- LinearFold-V vs. Vienna is similar

	(precision) PPV	(recall) Sensitivity
Overall		
CONTRAFold MFE	54.51	55.36
LinearFold-C	55.84 (+1.3)	56.24 (+0.9)
Vienna RNAfold	50.22	58.74
LinearFold-V	50.51 (+0.3)	58.97 (+0.2)



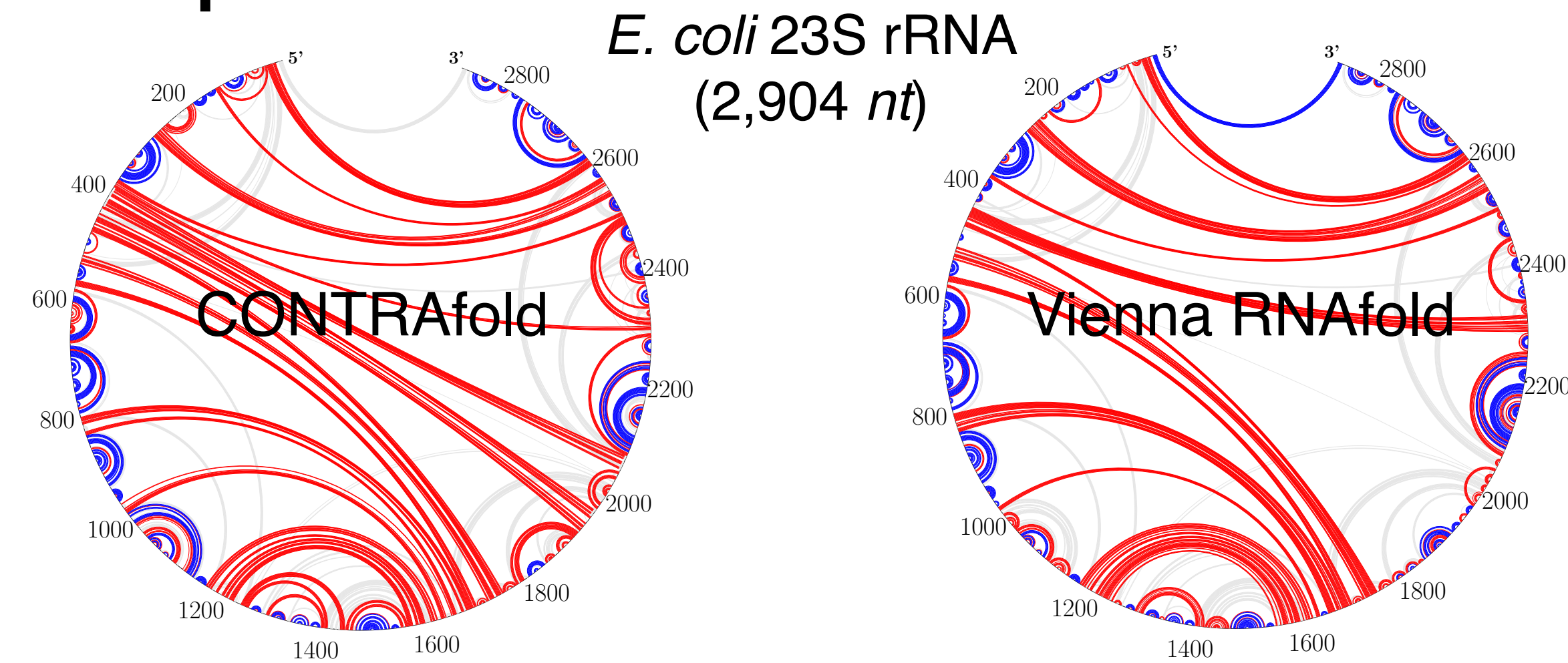
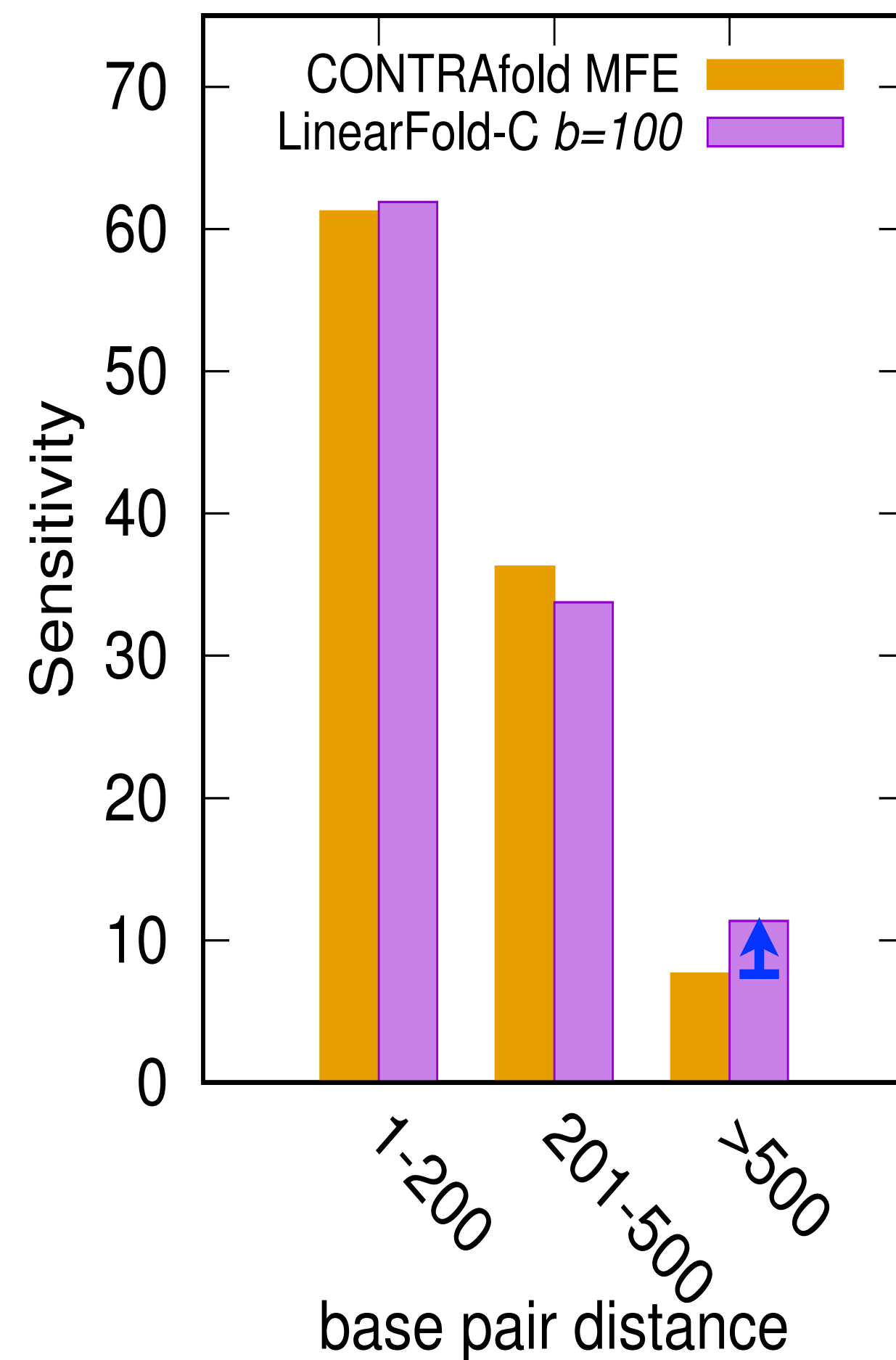
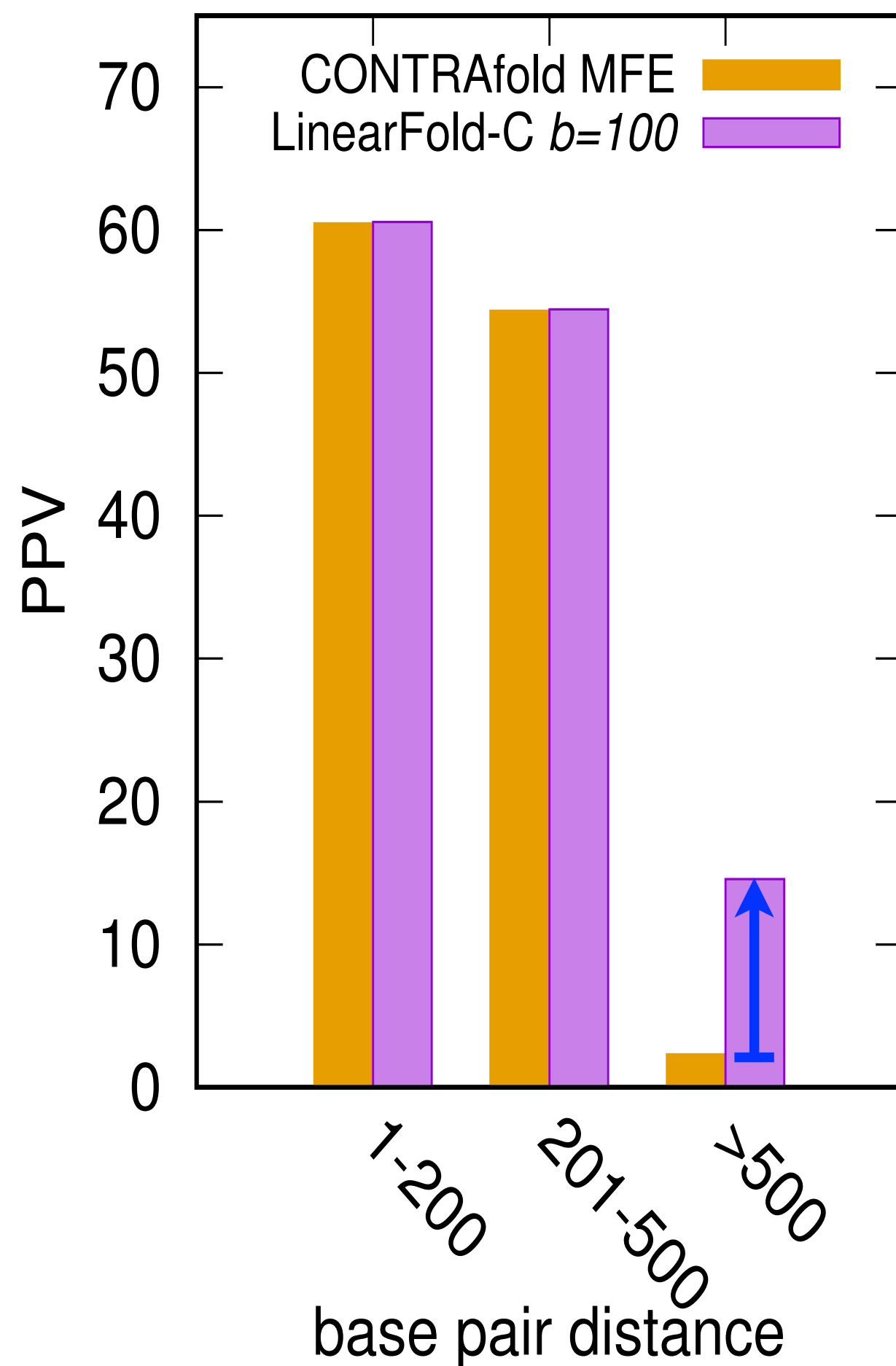
# Improvements on Long-Range Base Pairs

- long-distance pairs are well-known to be hard to predict
- LinearFold outputs **less** long-range base pairs, but **more correct** ones



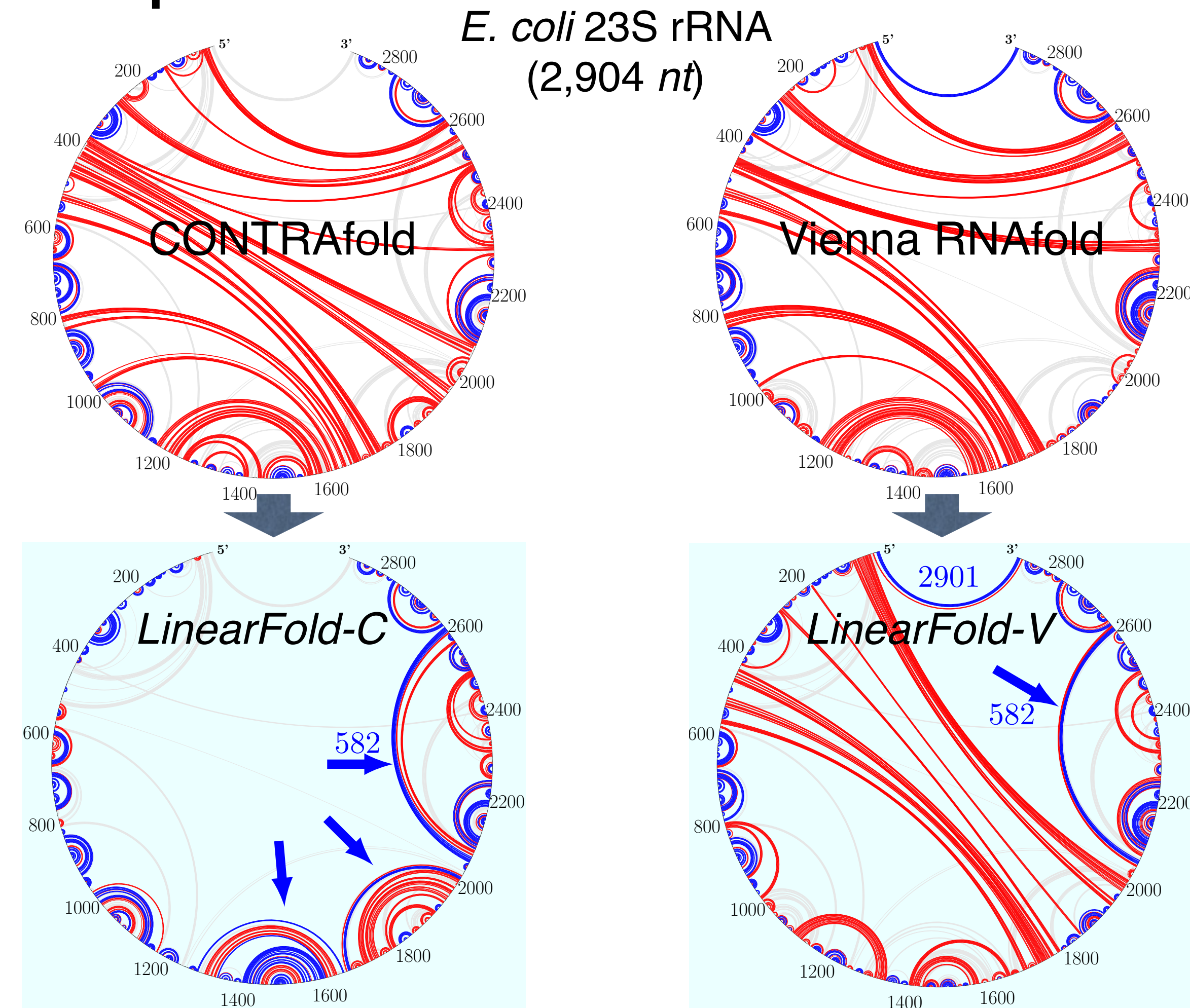
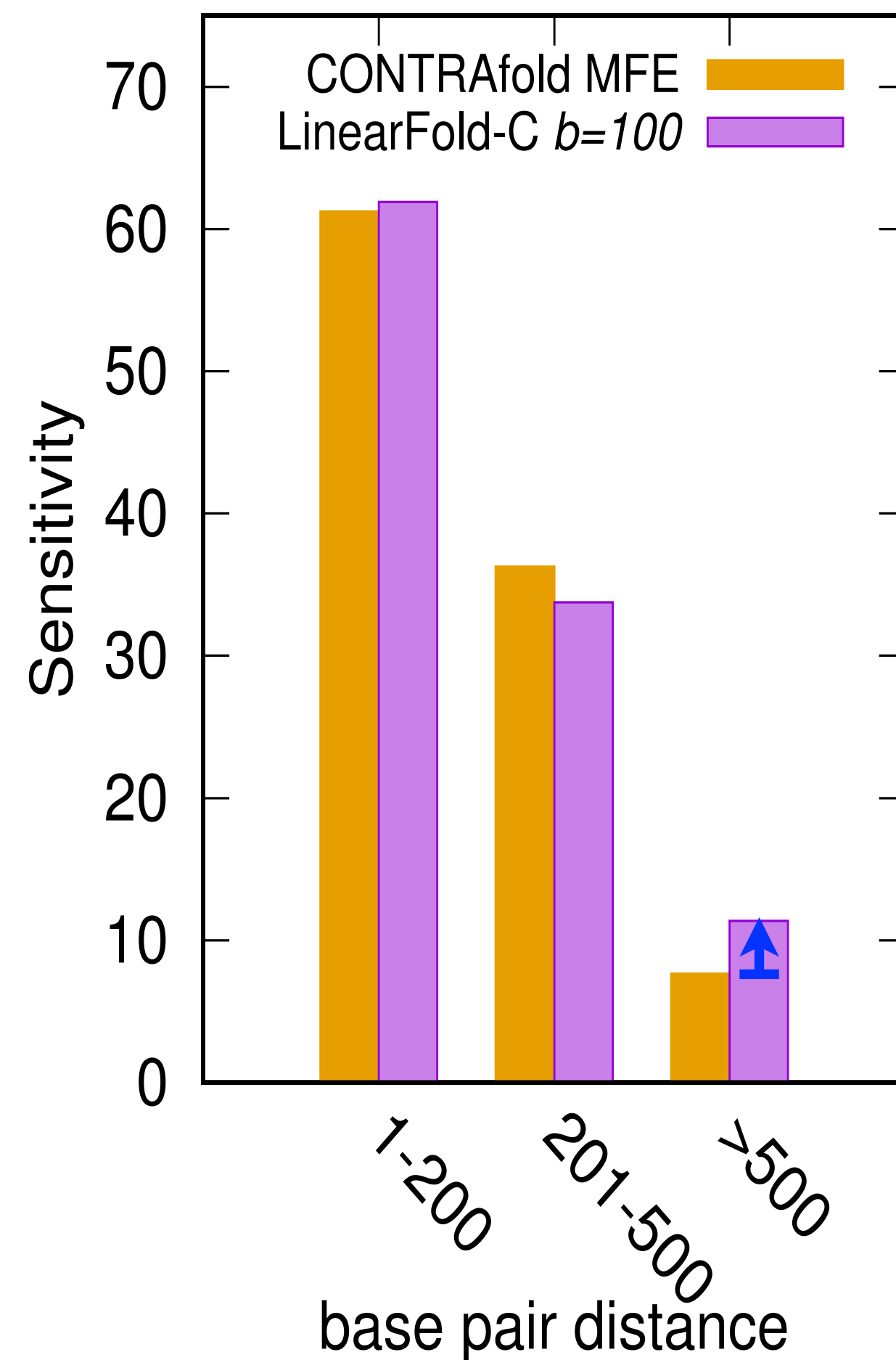
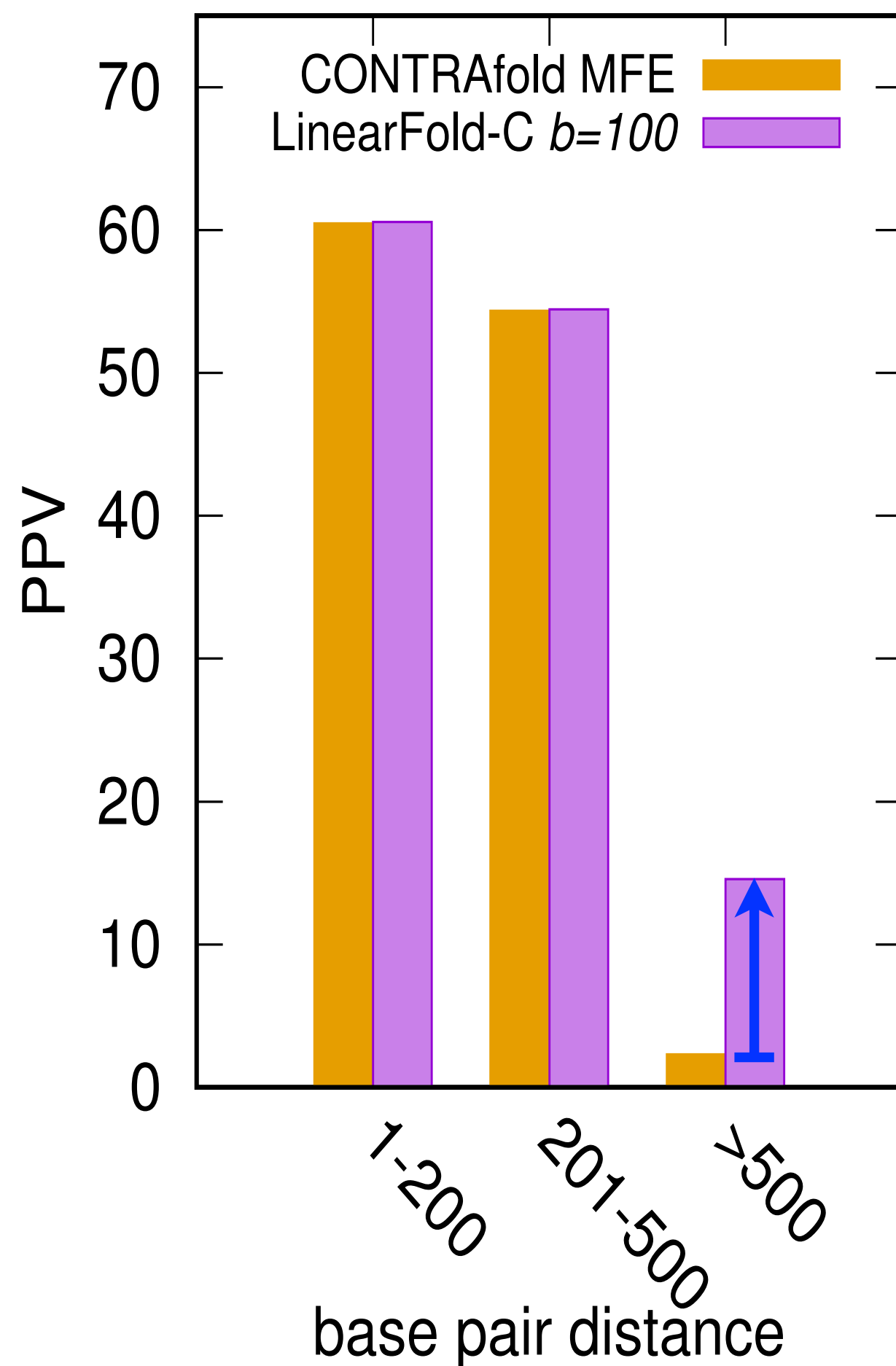
# Improvements on Long-Range Base Pairs

- long-distance pairs are well-known to be hard to predict
- LinearFold outputs **less** long-range base pairs, but **more correct** ones



# Improvements on Long-Range Base Pairs

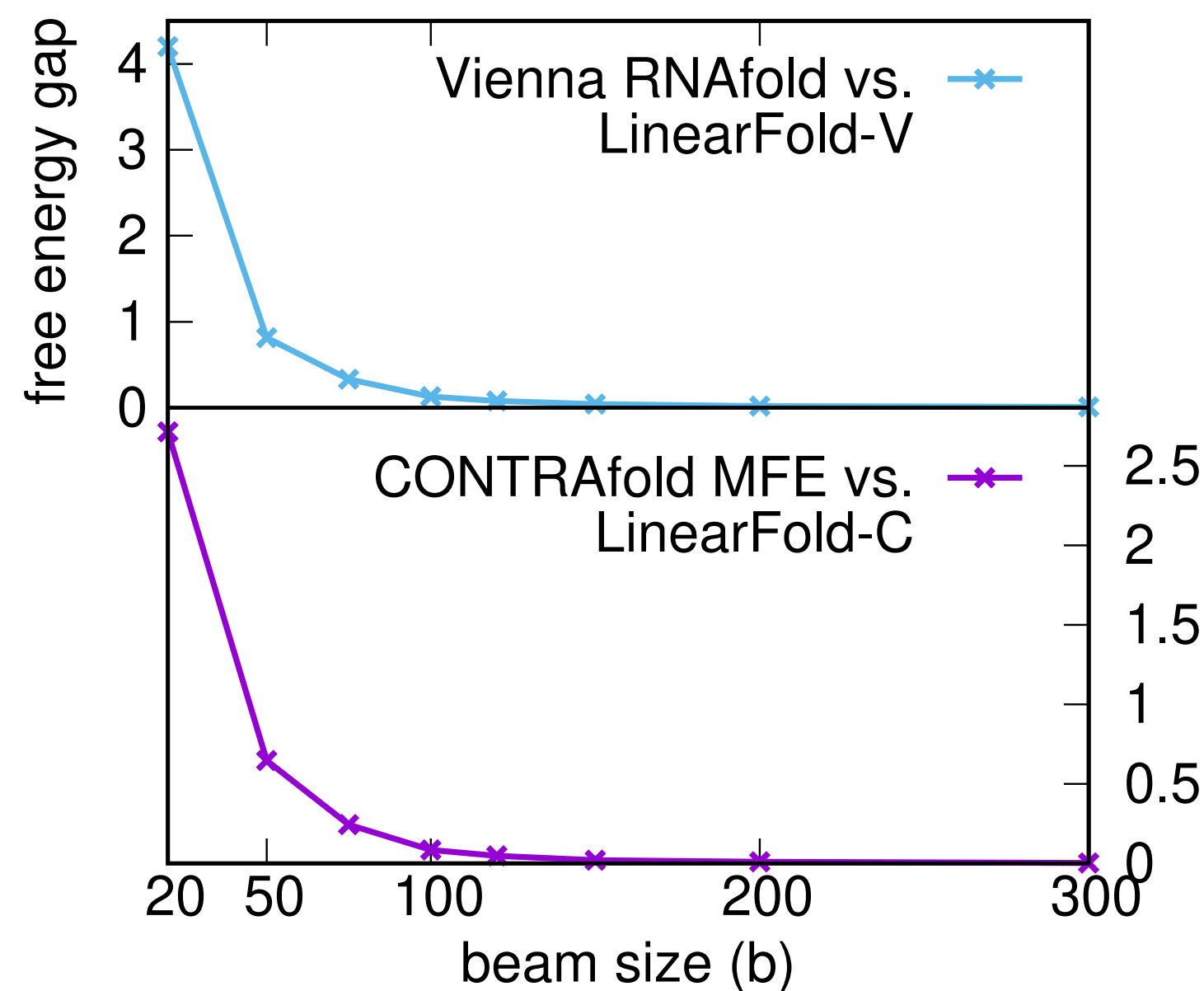
- long-distance pairs are well-known to be hard to predict
- LinearFold outputs **less** long-range base pairs, but **more correct** ones



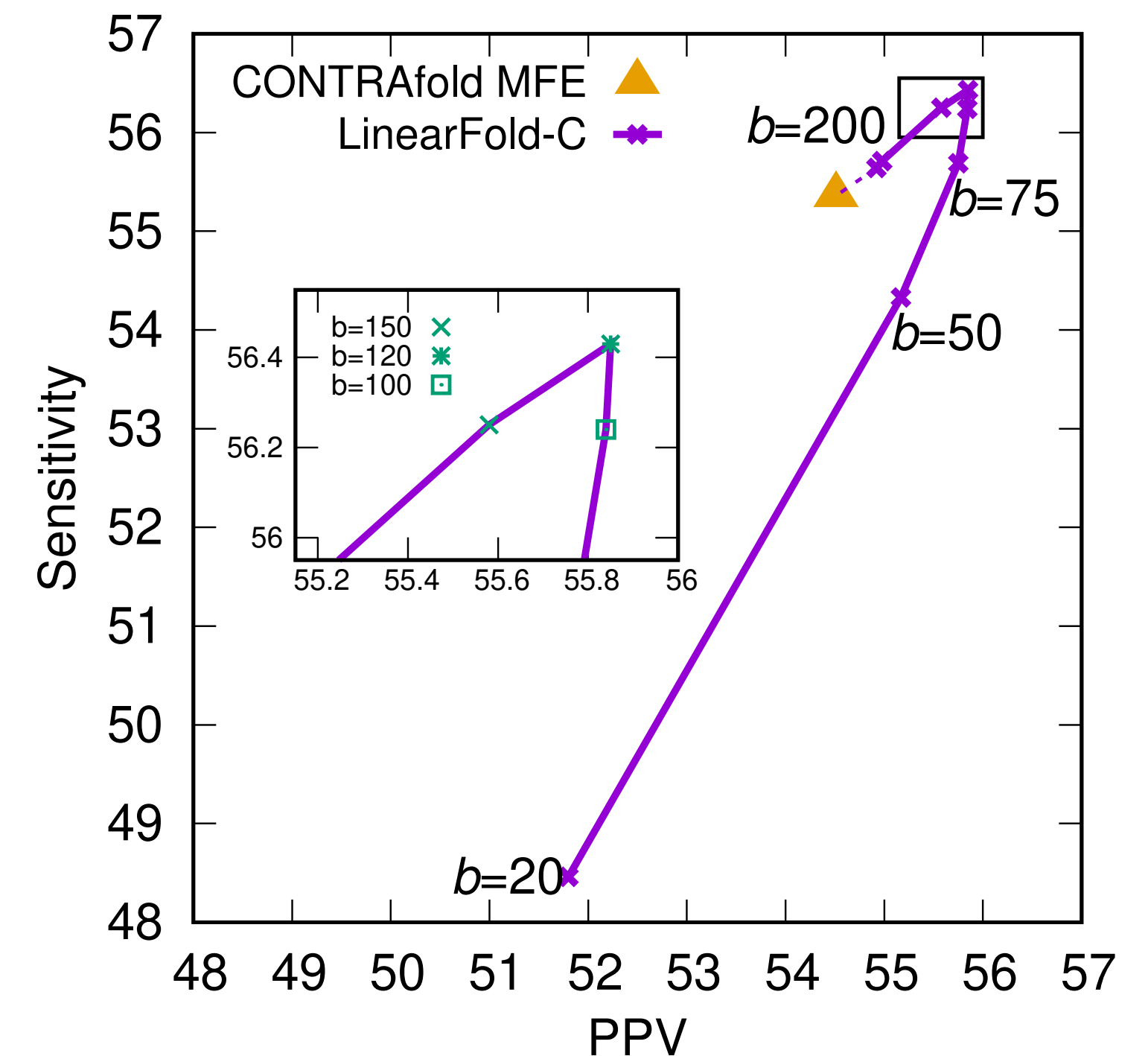
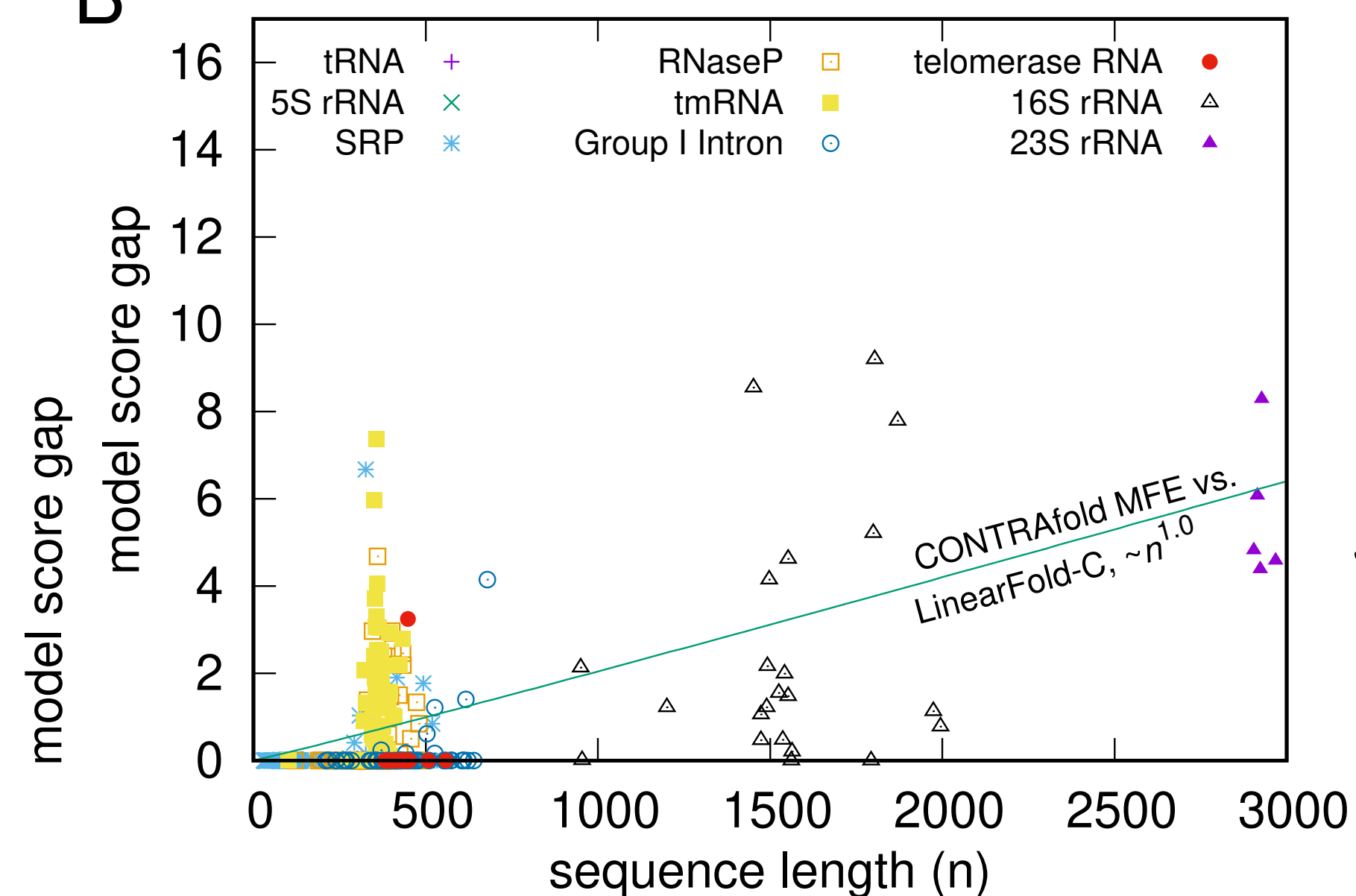
# Search (Approximation) Quality

- search error = energy (or model score) gap between exact search and our search
- search error quickly shrinks to 0 as beam size increases
- search error grows linearly with sequence length (constant search error per nucleotide)
- both PPV & Sensitivity peaks around beam=120, but we choose beam=100 for simplicity

A

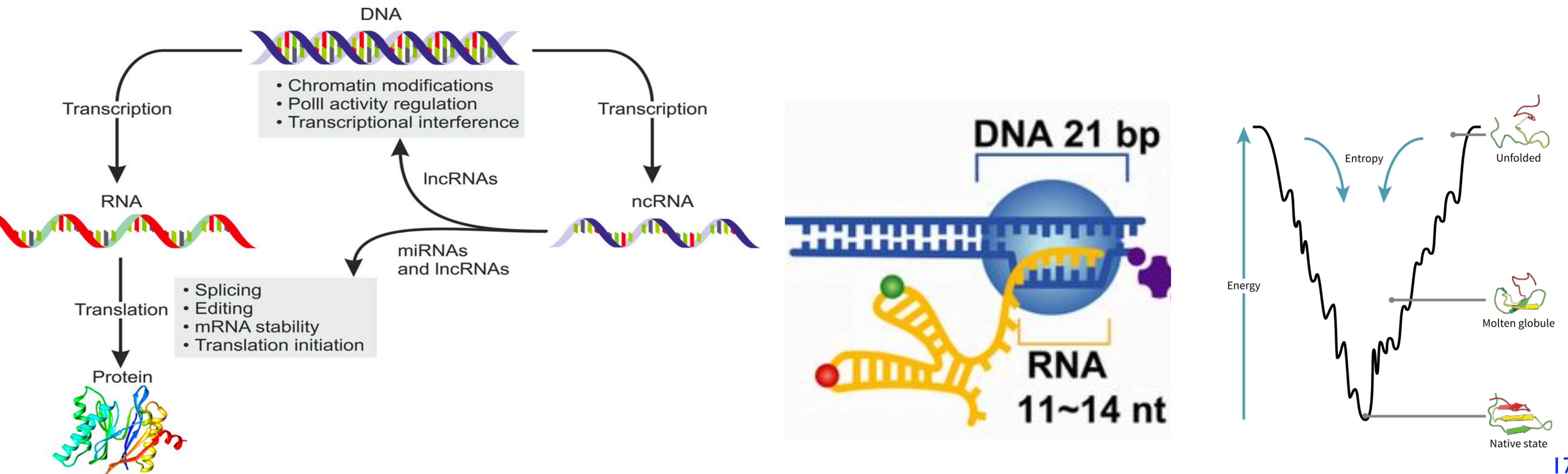


B



# Connections to Cotranscriptional Folding

- RNAs & proteins fold while being assembled
- RNA & protein sequences evolve to be incrementally foldable
- these might explain why LinearFold performs better than exact search



# LinearFold Server & GitHub

LinearFold Web Server (beta)

Interactive demo

Add a sequence

Paste or type your sequence here:

```
>tRNA_tdbR00000566-Homo_sapiens-9606-Tyr-9PA
CCUUCGAUAGCUCAGCUGGUAGAGCGGAGGACUGUAGAUCUAGGUCGCGUGGUUCGAUUCGCGCUCGAAGGACCA
```

Or upload a file in FASTA format:

Choose File No file chosen...

Set beam size

Beam size (1-200): 100

Choose model(s)

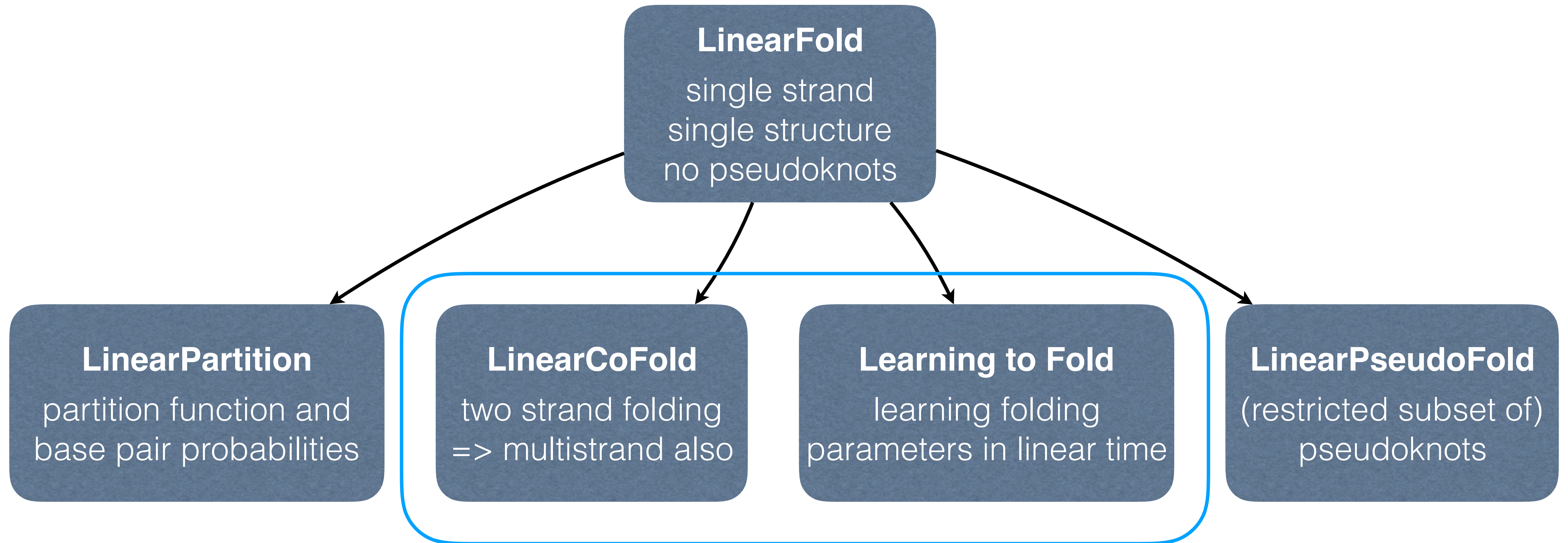
- LinearFold-C (using [CONTRAFold v2.0](#) machine-learned model, Do et al 2006)
- LinearFold-V (using [Vienna RNAfold](#) thermodynamic model, Lorenz et al 2011, with parameters from Mathews et al 2004)

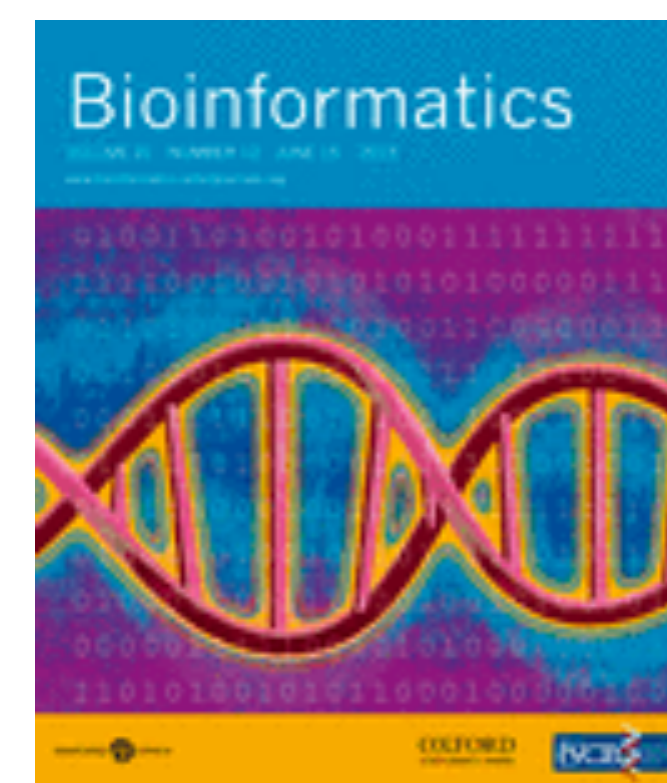
Run >> Reset

- fastest RNA folding server (by a very large margin)
- thanks to  $O(n)$  time
- longest sequence limit
  - 100,000 nt (10x Vienna)
  - thanks to  $O(n)$  space
- source code on GitHub
  - unified implementation of LinearFold-C & LinearFold-V



# Extensions of LinearFold

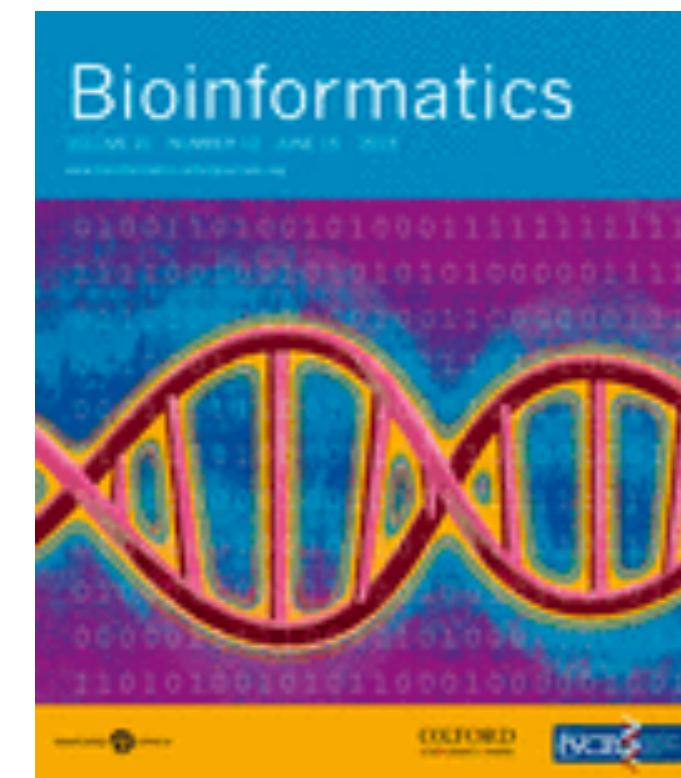




first linear-time (approximate) RNA folding algorithm  
better in accuracy (esp. long seqs and long-range pairs)  
<http://linearfold.org>



# Thank you very much !



first linear-time (approximate) RNA folding algorithm  
better in accuracy (esp. long seqs and long-range pairs)  
<http://linearfold.org>

# Backup Slide: Local Folding

<i>type</i>	<i>pair distance</i>	<i>time</i>	<i>space</i>	<i>systems/examples</i>
global	unbounded	$O(n^3)$	$O(n^2)$	RNAstructure, RNAfold, ...
local	$\leq L$	$O(nL^2)$	$O(nL)$	Rfold, RNApfold, LocalFold, ...
global	unbounded	$O(nb \log b)$	$O(nb)$	LinearFold ( <b>this work</b> )

23S rRNA *E. coli* (2904 nt, 830 pairs)

