

PRESENTATION

## Forbidden motifs and the cardinality of secondary structure space

YAO, Hua-Ting

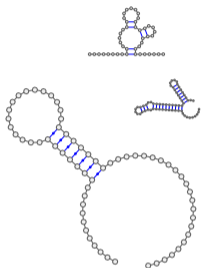
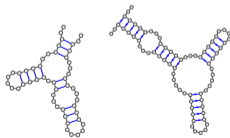
Ecole Polytechnique, France  
McGill University, Canada



TBI, University of Vienna, Austria

Benasque — Aug 11, 2022

## Phenotype Space

Secondary structure  $S^*$ 

## Genotype Space

CACGGCUAUUCAACCUUGC ...

UAGCGAGCUGAAUUCGACUCGAA ...

UAAUUUAAGAUGGCGGUGAA ...

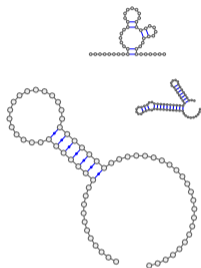
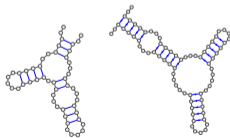
UUUAAGAUAAACUGGGCGAA ...

RNA sequence  $w$ 

AAAAAA ...

GAACUAGCUAAAGCUUGGCGU ...

## Phenotype Space

Secondary structure  $S^*$ 

## Genotype Space

CACGGCUAUUCAACCUUGC ...  
 UAGCGAGCGUGAAUUCGACUCGAA ...

UAAUUUAAGAUGGCGGUGAA ...  
 UUUAAGAUAAACUGGGCGAA ...

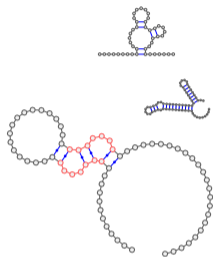
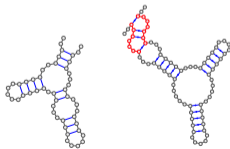
RNA sequence  $w$ 

AAAAAA ...  
 GAACUAGCUAAAGCUUGGCGU ...

$$\text{MFE}(w) = S^*$$



## Phenotype Space

Secondary structure  $S^*$ 

## Genotype Space

CACGGCUAUUCAACCUCUGC ...

UAGCGAGCGUGAAUUCGACUCGAA ...

X

RNA sequence  $w$ 

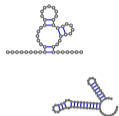
AAAAAAA ...

GAACUAGCUAAAGCUUGGCGU ...

$$\text{MFE}(w) = S^*$$

(Aguirre-Hernández *et al.*, 2007)

## Phenotype Space



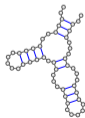
$$\text{MFE}(w) = S^*$$



## Genotype Space

CACGGCUAUUCAACCUCUGC ...  
UAGCGAGCGUGAAUUCGACUCGAA ...

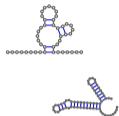
Secondary structure  $S^*$



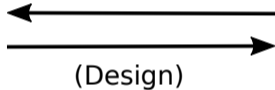
RNA sequence  $w$

AAAAAAA ...  
GAACUAGCUAAAGCUUGGCGU ...

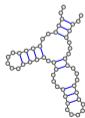
## Phenotype Space



$$\text{MFE}(w) = S^*$$



Secondary structure  $S^*$



## Genotype Space

CACGGCUAUUCAACCUUGC ...

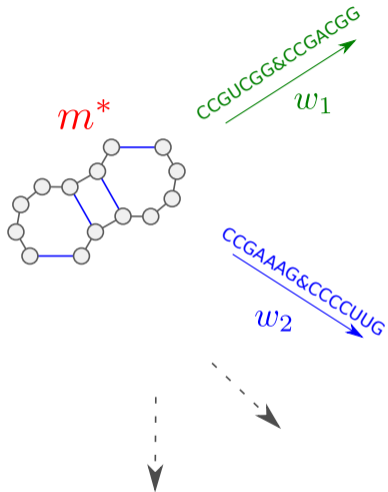
UAGCGAGCUGAAUUCGACUCGAA ...

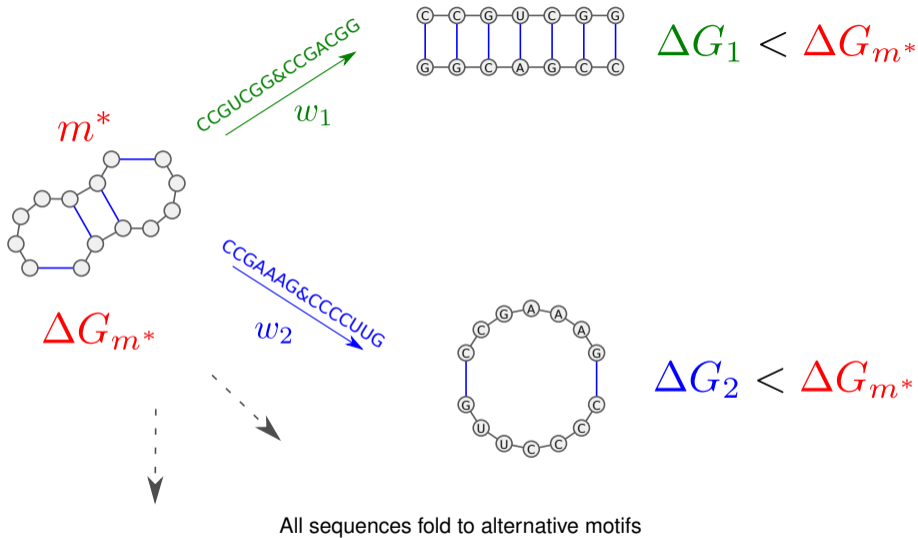
RNA sequence  $w$

AAAAAA ...

GAACUAGCUAAAGCUUGGCGU ...

# Definition of Forbidden motif

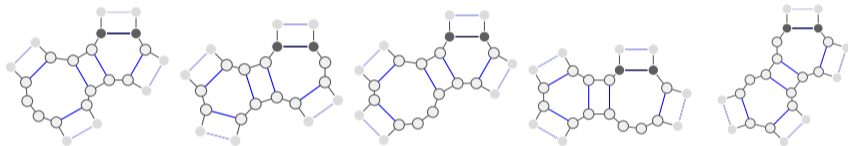




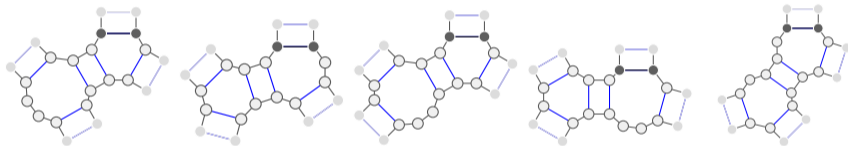


# Forbidden Motifs of size up to 14

- Almost half of motifs are forbidden of size up to 14 (4 561 out of 10 886)
- 2 323 forbidden motifs are minimal
- 63 do not contain isolated base pair



- Almost half of motifs are forbidden of size up to 14 (4 561 out of 10 886)
- 2 323 forbidden motifs are minimal
- 63 do not contain isolated base pair



- Impact on the combinatorics of secondary structures

→ (Asymptotic) **Proportion**  $P_n$  of structures avoiding forbidden motifs decreases exponentially with the size  $n$

Grammar generates secondary structures

$$S \rightarrow \varepsilon + \bullet S + (S) S$$

Grammar generates secondary structures

$$S \rightarrow \varepsilon + \bullet S + (S) S$$



Grammar generates secondary structures  $\mathcal{D}$  avoiding forbidden motif set  $\mathcal{F}$

$$S \rightarrow \varepsilon + \bullet S + (S)S$$

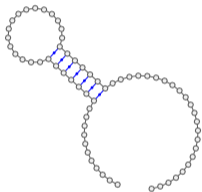


Grammar generates secondary structures  $\mathcal{D}$  avoiding forbidden motif set  $\mathcal{F}$

$$S \rightarrow \varepsilon + \bullet S + (S)S$$



		Proportion			
$ \mathcal{F} $	$ \mathcal{D}_n $	$P_n$	$P_{100}$	$P_{500}$	$P_{1000}$
387	$0.67 \frac{2.242^n}{n\sqrt{n}}$	$0.980^n$	$1.19 \cdot 10^{-1}$	$2.98 \cdot 10^{-5}$	$9.40 \cdot 10^{-10}$



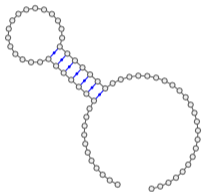
Secondary structure  $S^*$

$$\text{MFE}(w) = S^*$$



UAAUUUAAGAUGGCGGUGAA ...  
UUUAAGAUAAACUGGGCGAA ...

RNA sequence  $w$



Secondary structure  $S^*$

$$\mathbb{E}_w(\delta(S^*, \hat{S})) \leq \varepsilon$$

$$\mathbb{P}_w(S^*) \geq \varepsilon$$

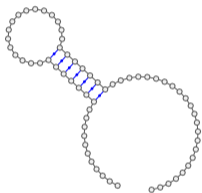
$$\text{MFE}(w) = S^*$$



UAAUUUAAGAUGGCGGUGAA ...  
UUUAAGAUAAACUGGGCGAA ...

RNA sequence  $w$



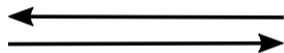


Secondary structure  $S^*$

$$\mathbb{E}_w(\delta(S^*, \hat{S})) \leq \varepsilon$$

$$\mathbb{P}_w(S^*) \geq \varepsilon$$

$$\text{MFE}(w) = S^*$$

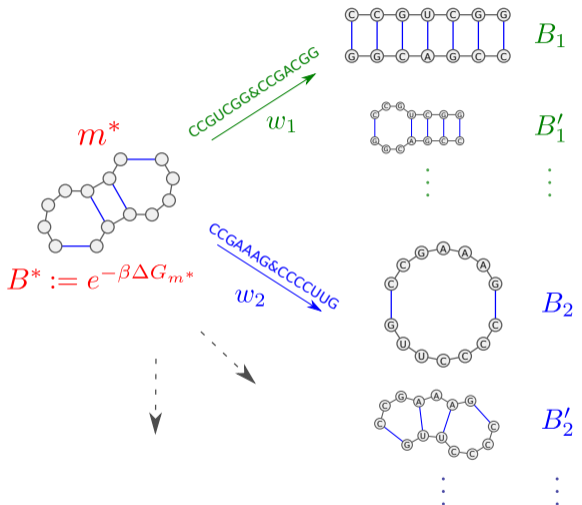


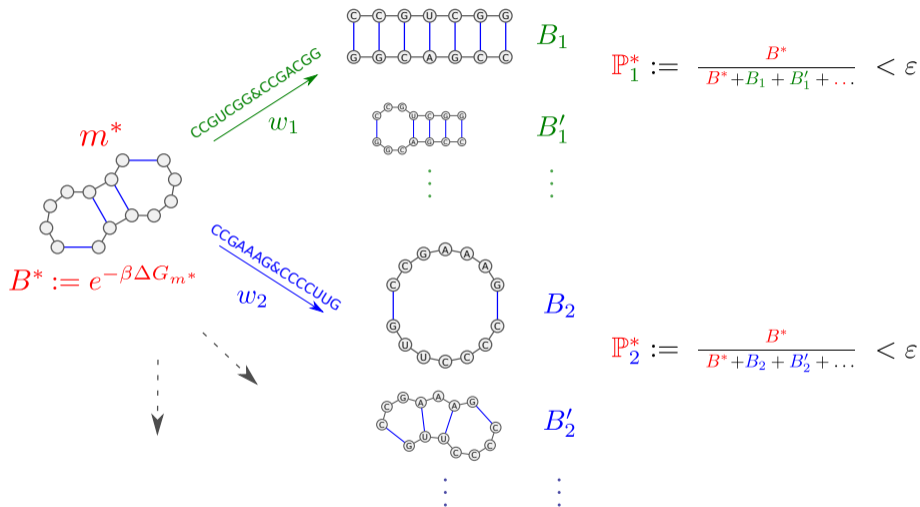
Probability defect

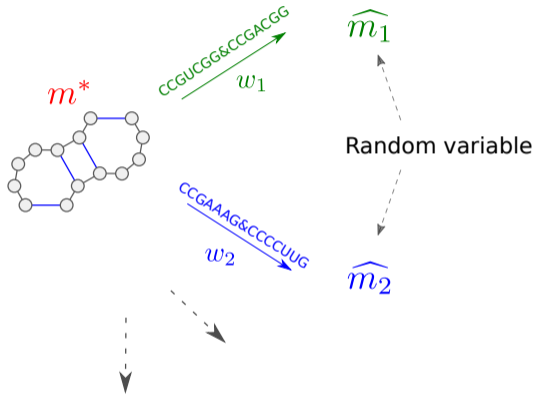
Ensemble defect

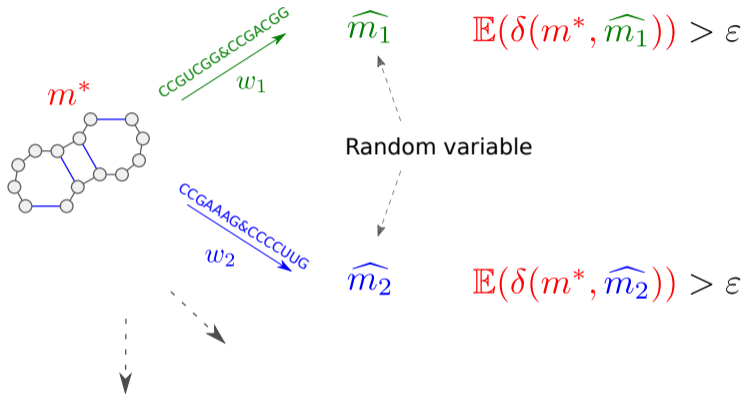
UAAUUUAAGAUGGCGGUGAA ...  
 UUUAAGAUAAACUGGGCGAA ...

RNA sequence  $w$



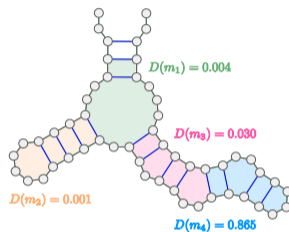






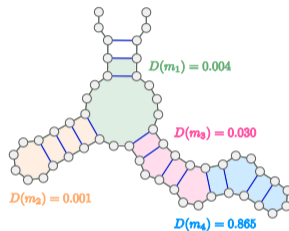
	$ \mathcal{F} $	$ \mathcal{D}_n $	$P_n$	Proportion		
				$P_{100}$	$P_{500}$	$P_{1000}$
not MFE	387	$0.67 \frac{2.242^n}{n\sqrt{n}}$	$0.980^n$	$1.19 \cdot 10^{-1}$	$2.98 \cdot 10^{-5}$	$9.40 \cdot 10^{-10}$
$\mathbb{P} < 50\%$	401	$0.66 \frac{2.239^n}{n\sqrt{n}}$	$0.978^n$	$1.03 \cdot 10^{-1}$	$1.53 \cdot 10^{-5}$	$2.49 \cdot 10^{-10}$
$\mathbb{E} > 1$	411	$0.65 \frac{2.236^n}{n\sqrt{n}}$	$0.977^n$	$9.08 \cdot 10^{-2}$	$8.52 \cdot 10^{-6}$	$7.86 \cdot 10^{-11}$

- Lower bound of structural ensemble defect



$$\mathcal{D}(S) := \min_w(w, S) \geq \mathcal{D}(m_1) + \mathcal{D}(m_2) + \mathcal{D}(m_3) + \mathcal{D}(m_4) = 0.9$$

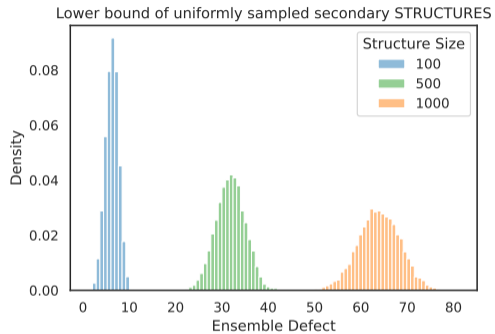
- Lower bound of structural ensemble defect



$$\mathcal{D}(S) := \min_w(w, S) \geq \mathcal{D}(m_1) + \mathcal{D}(m_2) + \mathcal{D}(m_3) + \mathcal{D}(m_4) = 0.9$$

- Variable tolerance grows with size  $n$ , e.g.  $0.01n$  recommended in NUPACK (Zadeh *et al.*, 2011)

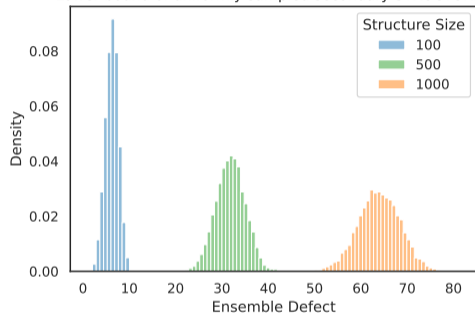




Random Structures

→ With  $\varepsilon = 0.01n$ ,  $P_{1000} \approx 10^{-33}$

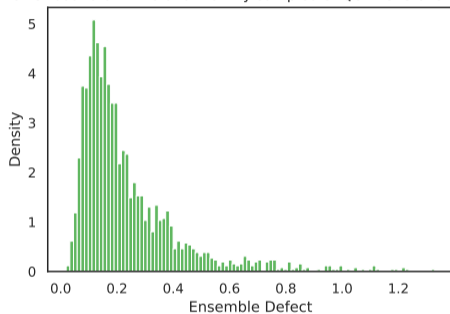
Lower bound of uniformly sampled secondary STRUCTURES



Random Structures

→ With  $\varepsilon = 0.01n$ ,  $P_{1000} \approx 10^{-33}$

Lower bound of MFEs of uniformly sampled SEQUENCES of size 500



MFE Structures

- Almost half of motifs are forbidden
- The cardinality of phenotype (structure) space is much smaller, but still exponentially grows
- Occurrences in PDBs shows a selection pressure on forbidden motifs

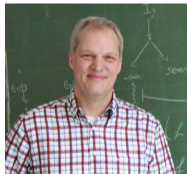
- Almost half of motifs are forbidden
- The cardinality of phenotype (structure) space is much smaller, but still exponentially grows
- Occurrences in PDBs shows a selection pressure on forbidden motifs
- Estimate how hard a phenotype can be realized
- Estimate the neuTral network size (number of designs)



AMIBio  
Ecole Polytechnique



TBI  
University of Vienna



Cedric Chauve  
Simon Fraser University

