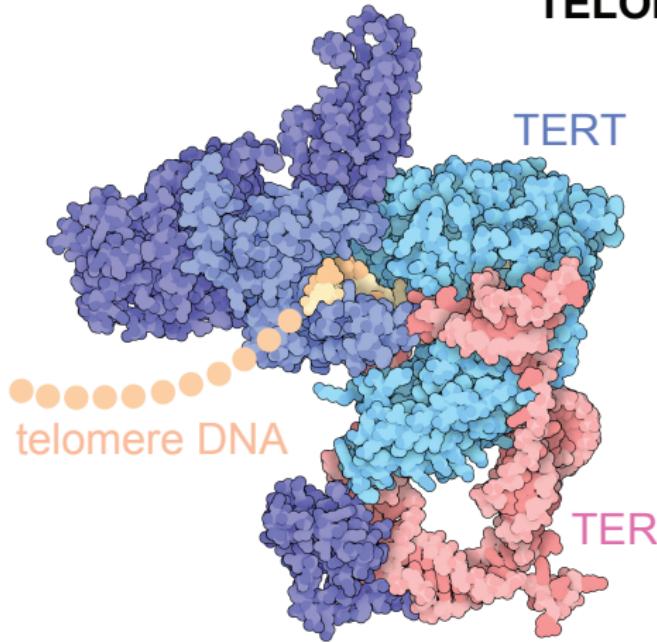


# **RNA structure prediction using positive and negative evolutionary information**

E Rivas, Harvard University  
[rivaslab.org](http://rivaslab.org)

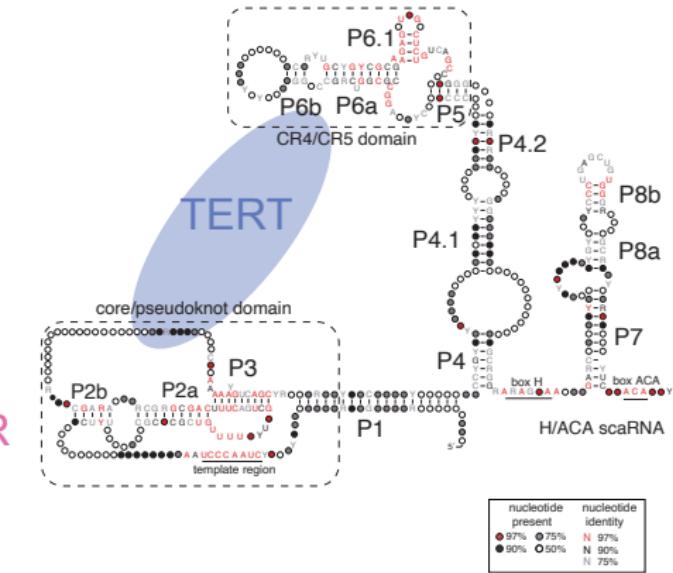
# Many functional RNAs have conserved structures

## TELOMERASE



3D structure (4.8 Å)  
Tetrahymena Telomerase

Jiang et al., Cell, 2018



2D structure  
vertebrate telomerase RNA (TER)

# vTER P4.1

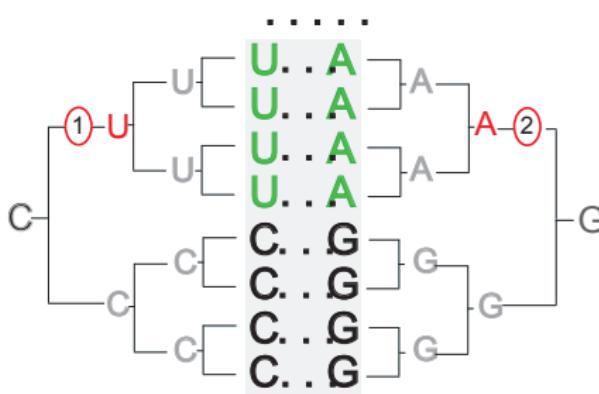
compensatory pair  
half-compensatory pair  
broken pair

## Pattern of sequence changes in a conserved RNA structure

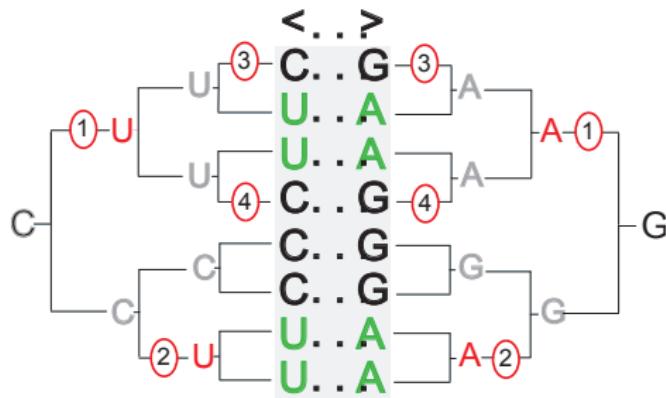
consensus	CCCC..GGGG
Human	GCCU..AGGU
Shark	GUCG..CGGC
Mustelus	GUCG..CGGC
Quoll	GUCU..AGGC
Stingray	CUCG..CGGG
Rhinoptera	CUCG..CGGG
Xenopus	CGGG..UCCG
Toad	-CUC..GAG-
Frog	-CUG..CAG-
Pyxicephalus	CGGG..CCGC
Dermophis	GCCC..GGGC
Herpele	GCCC..GGGC
Caecilian	GCCC..GGGC
Elephant	CC-C..GAGG
Manatee	CC-C..GAGG
Rabbit	CC-C..GAGG
Guinea_pig	UC-C..GAGU
Chinchilla	UC-C..GAGU
Gopher	CC-C..GC GG
Vole	GGCC..GGCC
Hamster	GGCC..GCC
Mus_musculus	GGCC..GGCC
Mus_spretus	GGCC..GGCC
Rat	GGCC..GGCC
Shrew_northern	CC-C..GAGG
Cat	CCUC..GAGG
Ferret	CC-C..GAGG
Raccoon	CC-C..GAGG
Bos	CC-C..-UGG
Pig	CC-C..GAGG
Shrew_house	CCG-..-CGG
Horse	CC-C..GAGG
Armadillo	UC-C..GAGG
Turtle	GGCC..GGUC
Macaw	GGCC..-GUC
	<<<..>>>

# Spurious pairwise covariations can appear from uncorrelated substitutions on a phylogenetic tree

two independent positions



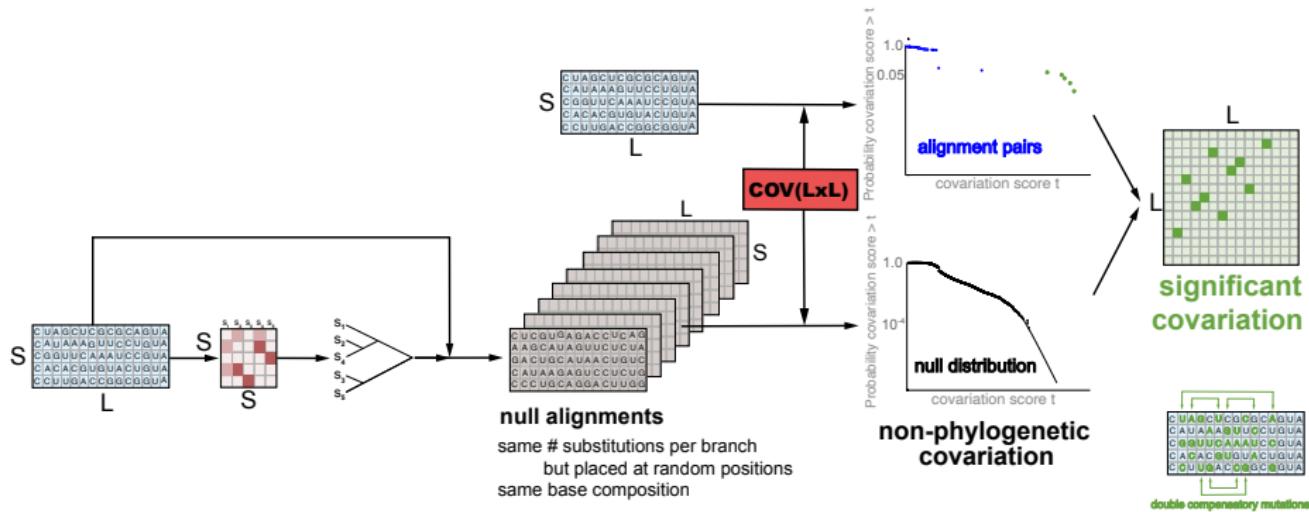
two base-paired positions



# A statistical test to simulate phylogenetic covariations



## RNA structural covariation above phylogenetic expectation

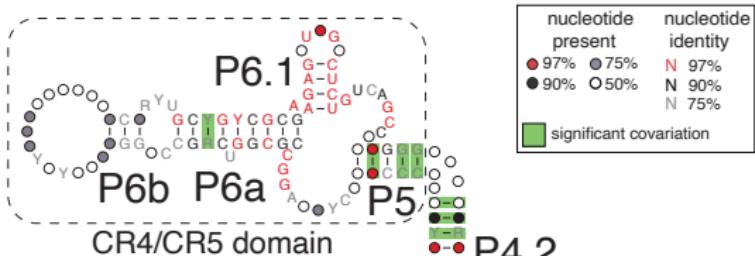


[eddylab.org/rscape](http://eddylab.org/rscape)

Lack of evidence for conserved secondary structure in long noncoding RNAs  
Rivas, Clements & Eddy, *Nat Methods*, 2017.

# R-scape finds RNA structural covariation above phylogenetic expectation

## telomerase RNA



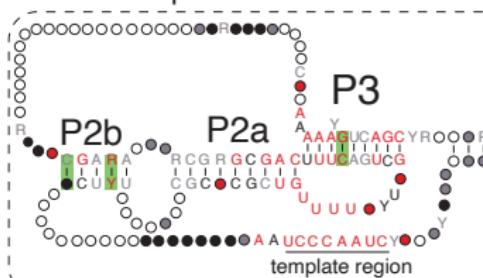
37 sequences

445 average length

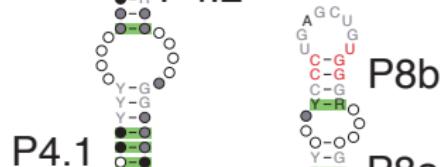
58% pairwise identity

27/107 base pair significantly covary

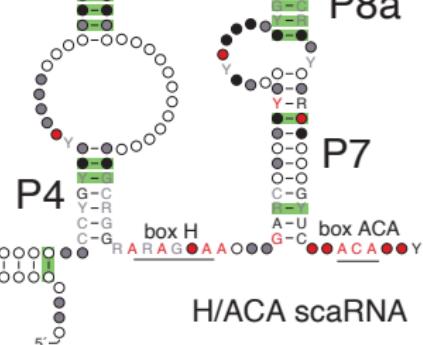
## core/pseudoknot domain



## P4.1



## P4

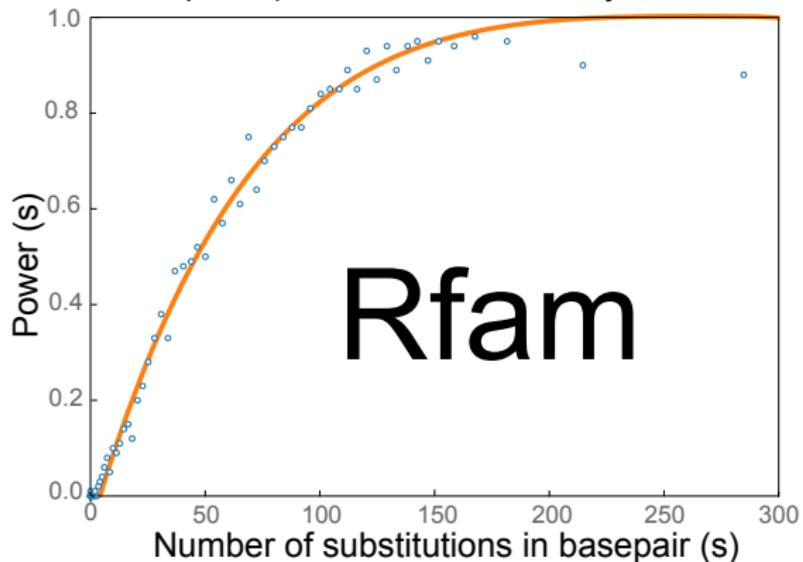


## H/ACA scaRNA

**Statistical Power = expected covariation given the variation observed**

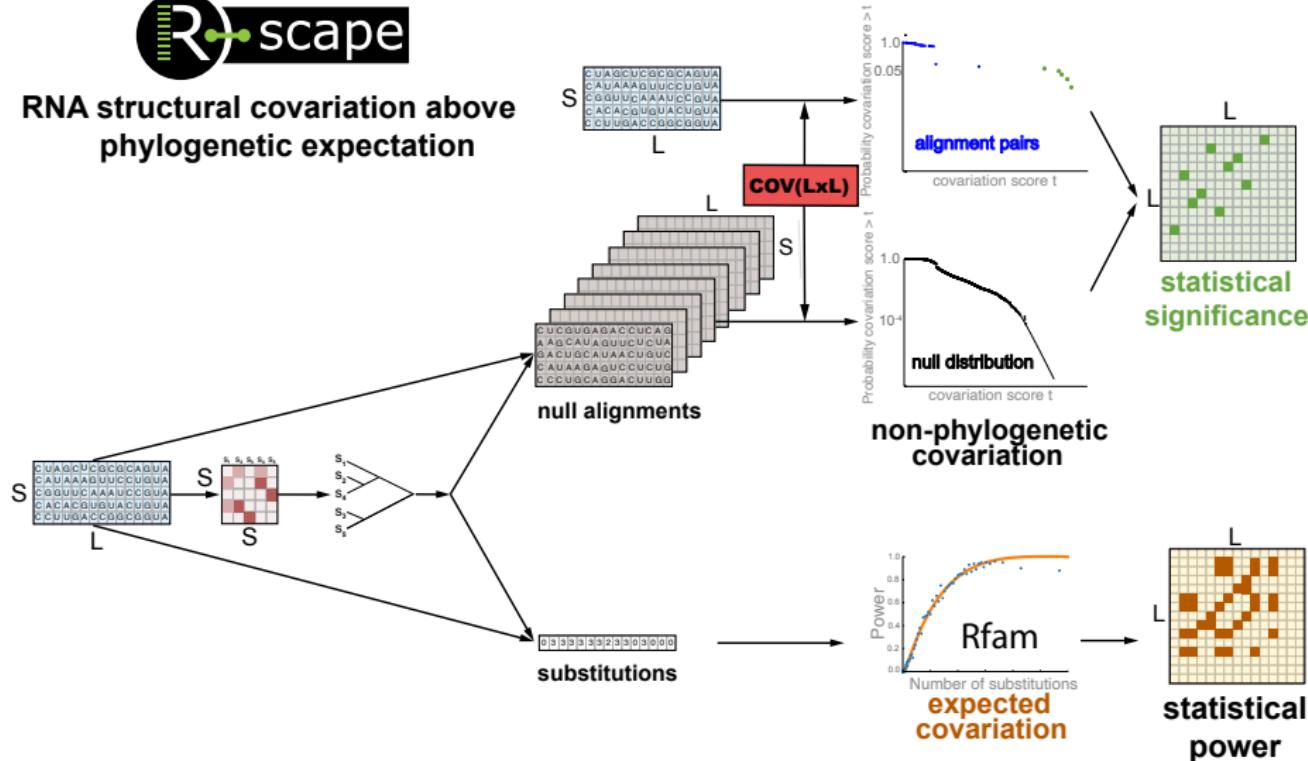
$$\text{power}(s) = P(\text{basepair with } s \text{ substitutions has an E-value} < 0.05)$$

7,012 basepairs (from 83 RNAs with crystal structures)

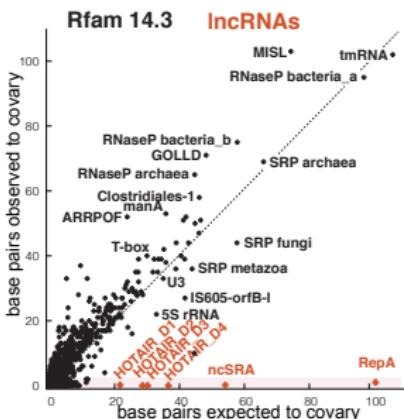
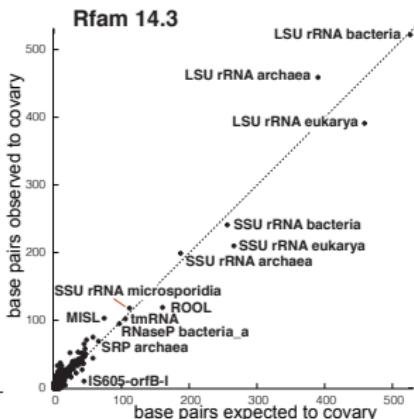
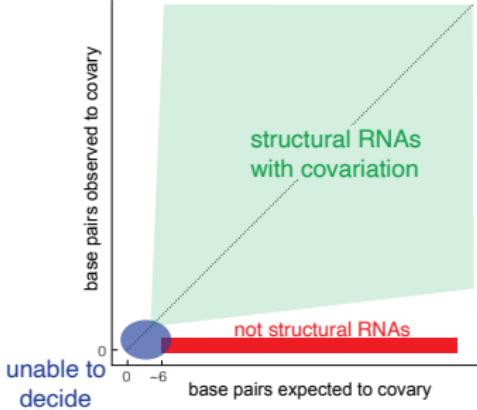




## RNA structural covariation above phylogenetic expectation

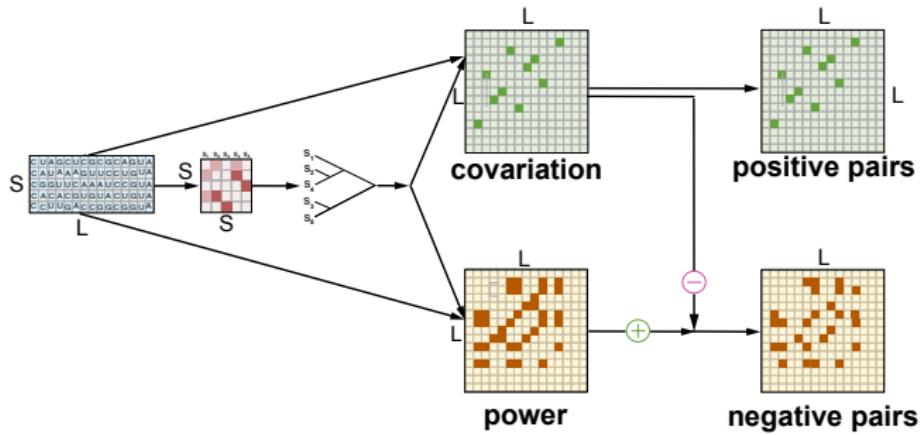


# How to tell when an RNA has an evolutionarily conserved structure ... or not?



Estimating the power of sequence covariation for detecting conserved RNA structure

Rivas et al., Bioinformatics, 2020.





## Cascade covariation/variation Constrained RNA Folding

- ▶ Produces an RNA structure that incorporates **all positive basepairs**
- ▶ Produces an RNA structure that forbids **negative basepairs**
- ▶ Uses a battery of probabilistic folding algorithms  
**(computationally efficient)**
- ▶ Visualization to critically analyze the structure
- ▶ In vivo structures

CaCoFold can incorporate pseudoknots, base triplets, other 3D interactions, even non RNA-structure related covariation

# RNA structure prediction using positive and negative evolutionary information

## CaCoFold

### a Input Alignment

5 sequences

50 consensus sequence length

76% average pairwise identity

```
CUGAAGUGACA-UCCUGCUGUUAUCUUAUCGAGCGGUUCCGAUAGCAUA  
CAGAAGUGACUUCUCCUAAGGUACUGUAUUGAUUGGUUCCAAUACCUGUA  
CGGAGGUAGCG-UCCUUUUCGUUACUUAUACGAAAGGUUCCGAUAAUCGUA  
CAG-UGUGACCUUCUACGGUUAUCUUAUCGAGUGGUUCCGAUAAUCGUA  
CCGAGGUACUU-CCUUGAGUUAUCUCAUUGACGGGUUCCGAUAGCGGUU
```

### c Cascade maxCov Algorithm

C0: 3/5 positive basepairs explained



C+: 2/5 positive basepairs explained



### e Alternative Helix Filtering

F0: The nested structure: keep unchanged



F+: One alternative positive helix: add to structure



### b Covariation Analysis

5 positive basepairs

E-value = 1e-4	E-value = 2e-6
CUGAAGUGACA-UCCUGCUGUUAUCUUAUCGAGCGGUUCCGAUAGCAUA	CAGAAGUGACUUCUCCUAAGGUACUGUAUUGAUUGGUUCCAAUACCUGUA
CGGAGGUAGCG-UCCUUUUCGUUACUUAUACGAAAGGUUCCGAUAAUCGUA	CAG-UGUGACCUUCUACGGUUAUCUUAUCGAGUGGUUCCGAUAAUCGUA
CCGAGGUACUU-CCUUGAGUUAUCUCAUUGACGGGUUCCGAUAGCGGUU	
	E-value = 6e-6
	E-value = 1e-5
	E-value = 3e-6

### d Cascade Constrained Folding

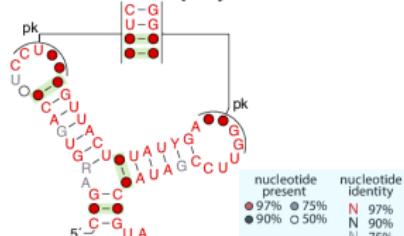
S0: Nested structure prediction: 3 forced/2 forbidden pairs

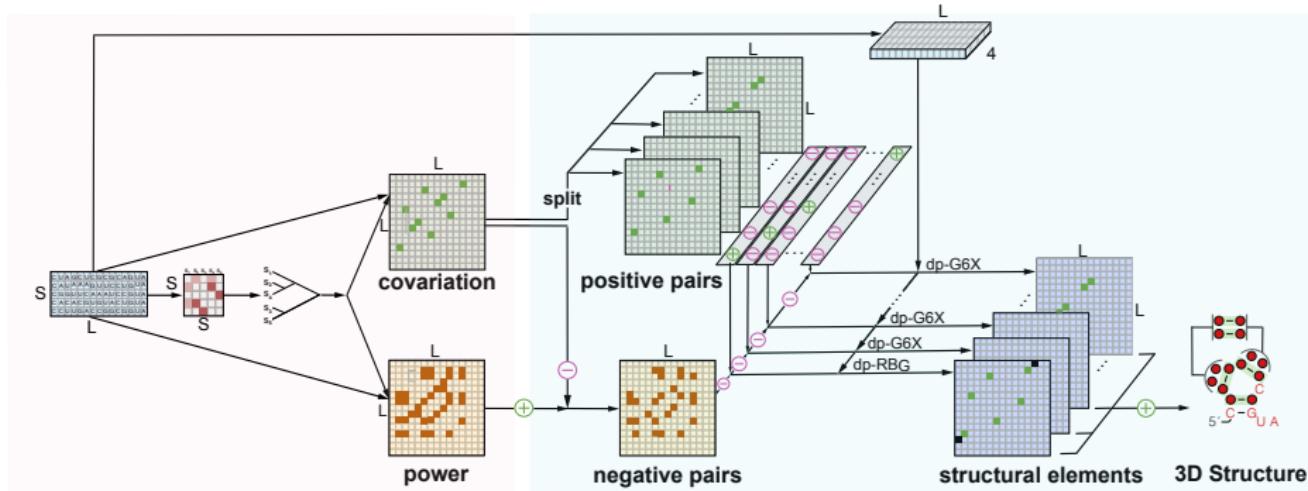


S+: Alternative helix prediction: 2 forced/3 forbidden pairs



### f Complete Structure Display





Rivas, *PLOS Comp Biol*, 2020.

a Model used by the maxCov algorithm

Nussinov Grammar

$S \rightarrow \circ S$  any non-covarying residue  
 $S \rightarrow \bullet S \bullet S$  a covarying basepair  
 $S \rightarrow S S$   
 $S \rightarrow \text{end}$

b Model used by the folding algorithm (first layer)

RNA Basic Grammar (RBG)

$S \rightarrow \circ S$  a free unpaired residue  
 $S \rightarrow L S$   
 $S \rightarrow \text{end}$

$L \rightarrow \bullet F \bullet$  a helix starts

$L \rightarrow \bullet P \bullet$  a one-basepair helix ends

$F \rightarrow \bullet F \bullet$  a helix adds one more basepair

$F \rightarrow \bullet P \bullet$  a helix ends

what can happen at the end of a helix...

$P \rightarrow \bullet \dots \circ$  a hairpin loop

$P \rightarrow \bullet \dots \bullet L$  a left bulge loop

$P \rightarrow \quad L \circ \dots \circ$  a right bulge loop

$P \rightarrow \bullet \dots \bullet L \circ \dots \circ$  an internal loop

$P \rightarrow M1 M$  a multiloop starts

$M \rightarrow M1 M$  multiloop adds one more branch

$M \rightarrow R$  multiloop about to add right residues

$R \rightarrow R \circ$  a right-unpaired residue in multiloop

$R \rightarrow M1$  multiloop about to add left residues

$M1 \rightarrow \circ M1$  a left-unpaired residue in multiloop

$M1 \rightarrow L$  multiloop starts another helix

c Model used by the folding algorithm (additional layers)

G6X Grammar

$S \rightarrow L$   
 $S \rightarrow L S$   
 $S \rightarrow \text{end}$

$L \rightarrow \bullet F \bullet$  a helix starts  
 $L \rightarrow \bullet \bullet$  a basepair of contiguous residues  
 $L \rightarrow \circ$  an unpaired residue

$F \rightarrow \bullet F \bullet$  a helix adds one more basepair  
 $F \rightarrow \bullet \bullet$  a helix ends without a hairpin  
 $F \rightarrow L S$  a helix ends, more stuff to come

- $\circ$  a non-covarying RNA residue
- $\bullet$  a covarying RNA basepair
- $\circ$  an RNA residue, not forming any basepairing
- $\dots \circ$  a set of contiguous unpaired RNA residues

$\bullet \bullet$  an RNA basepair; bases could be at arbitrary distance in the RNA backbone

$S, L, F, P, M, M1, R$  non-terminals that have to be transformed following one of the allowed rules

# transfer-messenger RNA (tmRNA)

## a Input Alignment

Rfam RF00023 seed alignment

477 sequences

354 consensus sequence length

357 average sequence length

42% average pairwise identity

## c Cascade maxCov Algorithm

121 positive basepairs explained in 6 layers

layer 1: 69      layer 2: 41

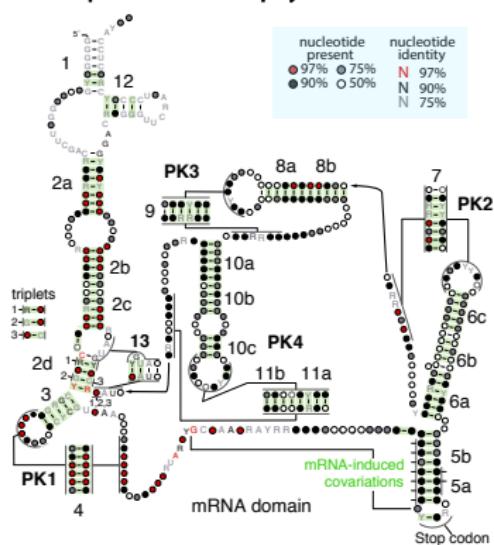
layer 3: 5      layer 4: 3

layer 5: 2      layer 6: 1

## e Alternative Helix Filtering

18 alternative helices    5 pseudoknots  
18 alternative helices    3 triplets  
10 mRNA-induced covariations

## f Complete structure display



## b Covariation Analysis

All possible pairs analyzed equally

119 annotated basepairs in alignment

(not used in analysis)

414 columns analyzed:

121 positive basepairs (significantly covary)

109 positive basepairs expected by power

31,027 negative basepairs

## d Cascade Constrained Folding

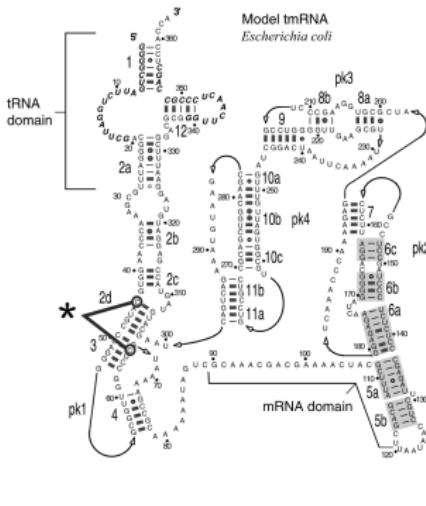
139 annotated pairwise interactions

121/139 positive basepairs

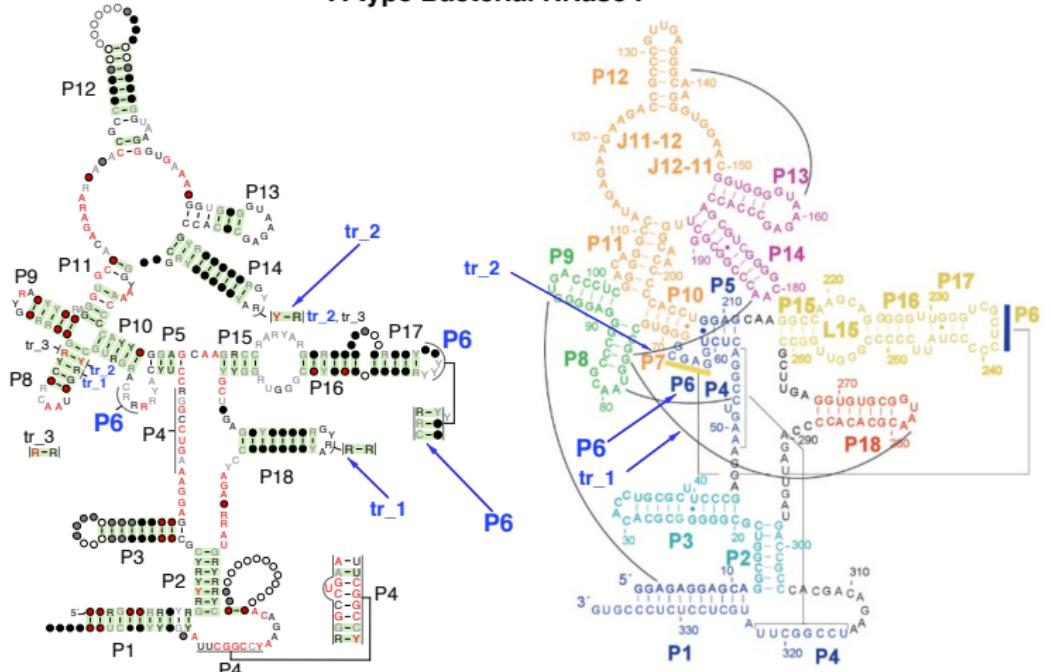
74 pairs not in final ss due to forbidden negative basepairs

## g Structure comparison

Kelley et al., RNA 2001, Fig 4



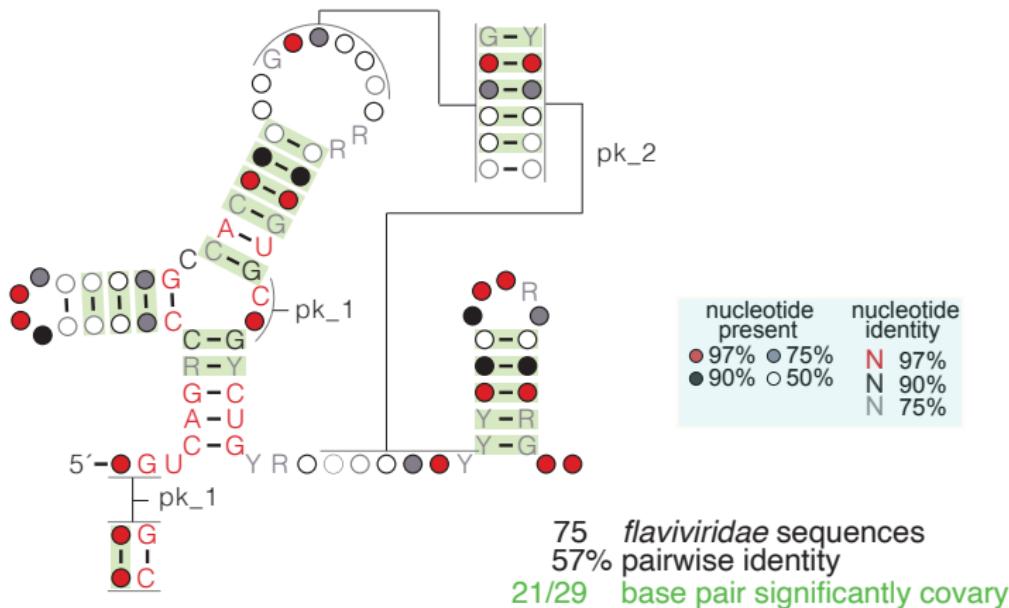
## A-type Bacterial RNase P



CaCoFold

Torres-Larios *et al.*, Nature 2005, Fig 2c

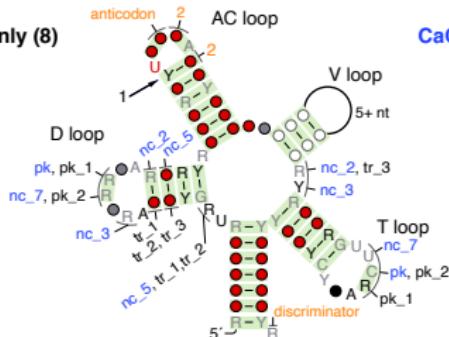
# exoribonuclease resistant RNA



Steckelberg, Vicens, Kieft, *mBio*, 2018  
Szucs, Nichols, Jones, Vicens, Kieft, *mBio*, 2020

**b****tRNA****CaCoFold only (8)**

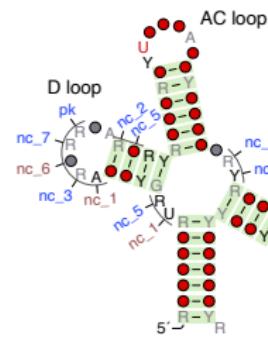
1 |Y=●  
 pk\_1 |R=R|  
 pk\_2 |R=Y|  
 tr\_1 |R=●|  
 tr\_2 |R=●|  
 tr\_3 |●=R|  
 2 |●=●|  
 anticodon/discriminator |●=R|



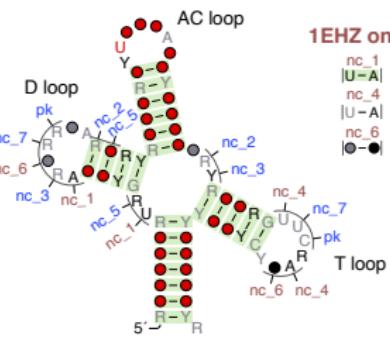
CaCoFold

**CaCoFold & 1EHZ (5)**

pk |R=Y| WWc  
 nc\_2 |R=R| HWt  
 nc\_3 |R=Y| WWt  
 nc\_5 |R=●| HHt  
 nc\_7 |R=U| WSt

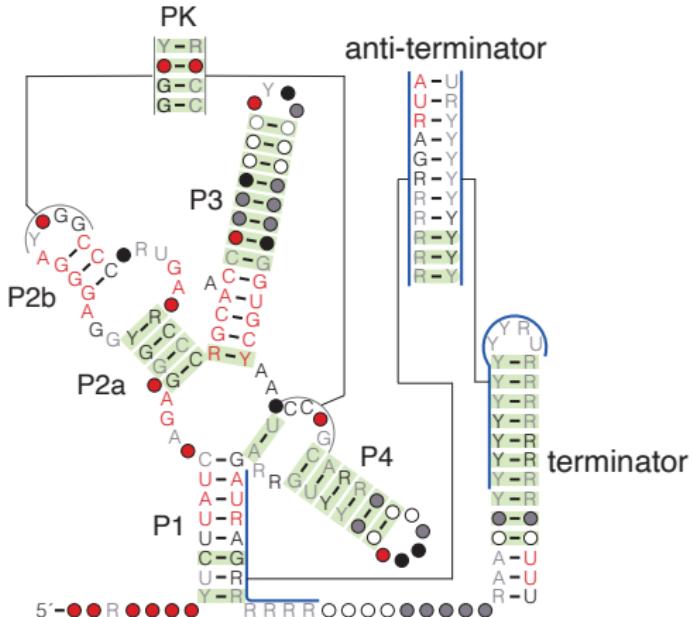
1EHZ  
Shi & Moore, Science 2020**1EHZ only (3)**

nc\_1 |U=A| WHt  
 nc\_4 |U=A| WHt  
 nc\_6 |●=●| SWt  
 nc\_7 |●=●| pk  
 nc\_4 |U=C| 1EHZ  
 nc\_7 |●=●| 1EHZ



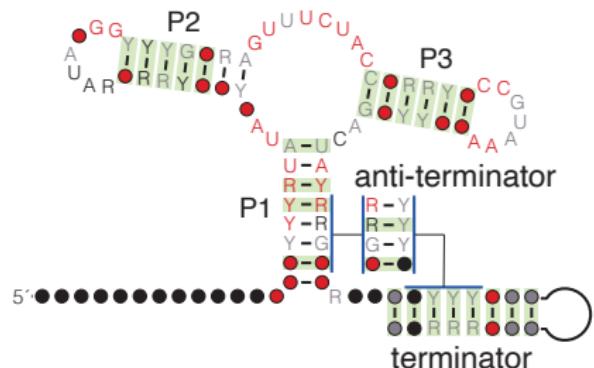
# Alternative/overlapping structures

## SAM-I riboswitch



Zhu & Meyer, *RNA Biology*, 2015

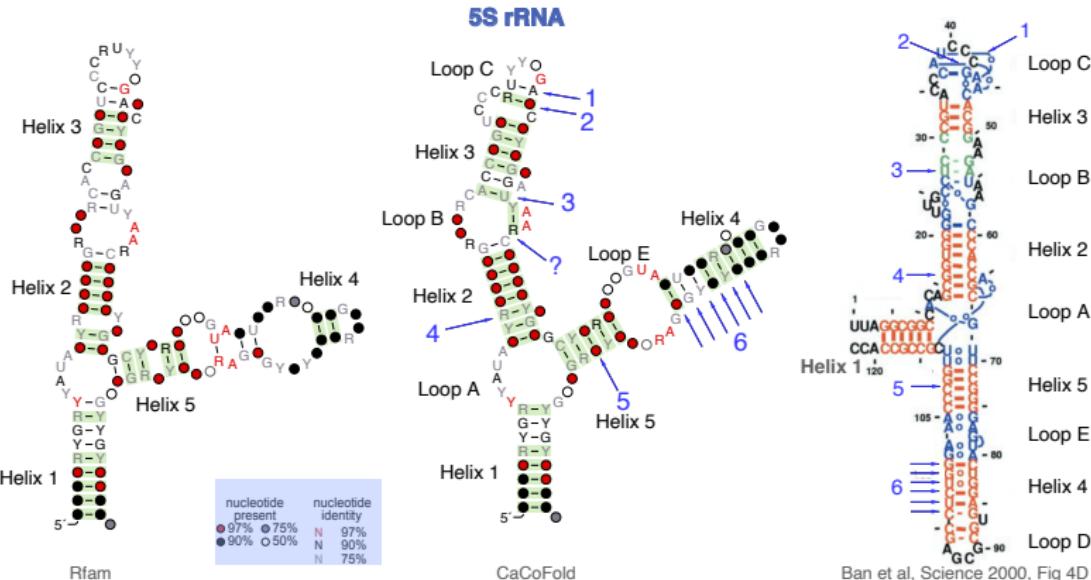
## Purine riboswitch



Ritz, Martin, Laederach, *PLCB*, 2013

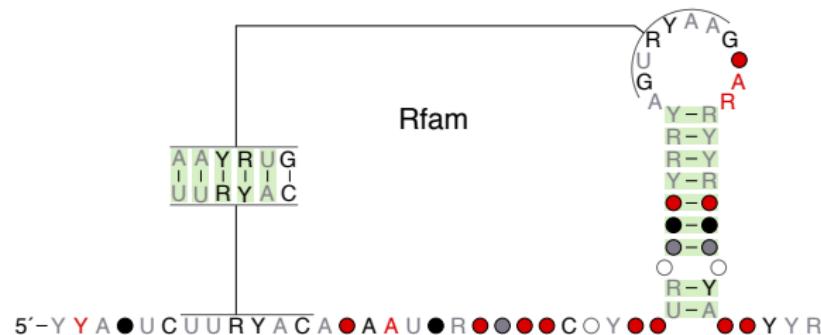
nucleotide present	nucleotide identity	
97% ● 75%	N 97%	■ significant covariation
● 90% ○ 50%	● 90%	○ variable length hairpin loop
● 75%	N 75%	

# CaCoFold helps improve Rfam structures

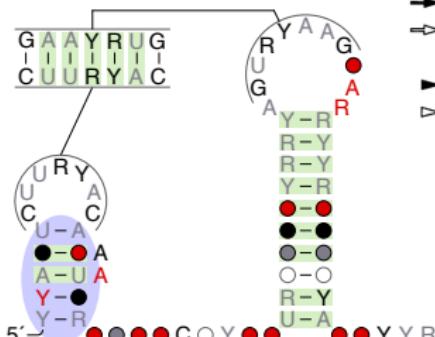


## Type 2

## Coronavirus 3'UTR pseudoknot

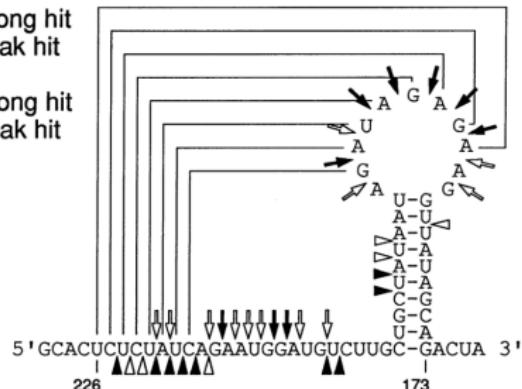


## CaCoFold



Williams *et al.*, J. Virol. 1999, Fig 4B

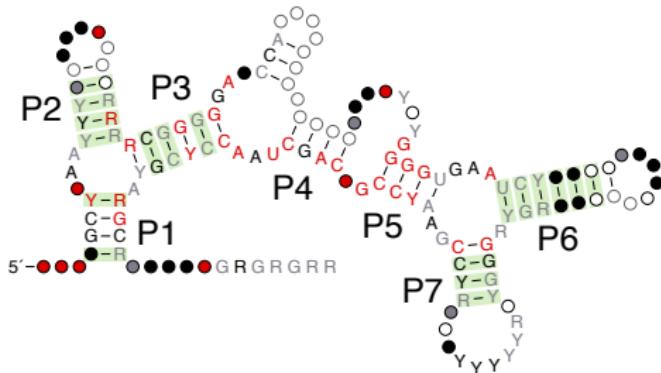
- ss strong hit
- ⇒ ss weak hit
- ds strong hit
- ▷ ds weak hit



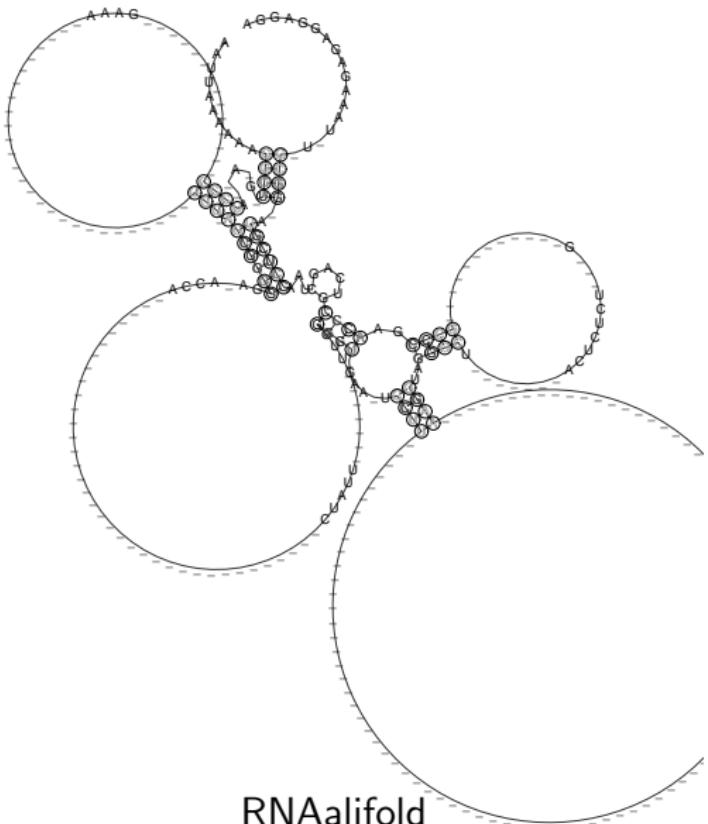
new helix with covariation support

### Cyclic di-AMP riboswitch

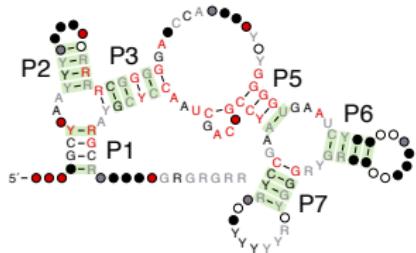
nucleotide identity  
N 97%  
N 90%  
N 75%  
nucleotide present  
● 97% ● 75%  
● 90% ○ 50%



CaCoFold

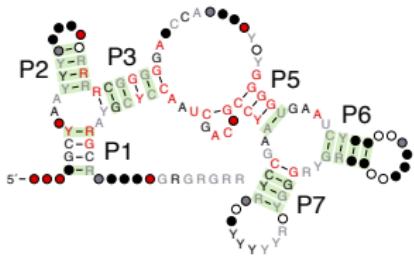


## Cyclic di-AMP riboswitch

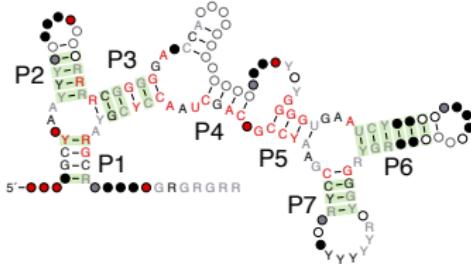


RNAalifold-R-scape

## Cyclic di-AMP riboswitch

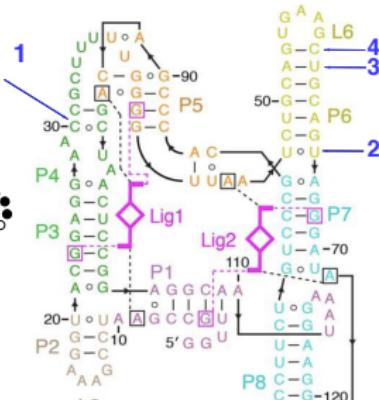
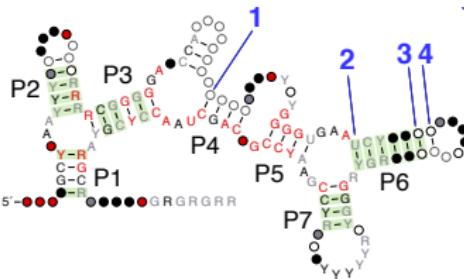
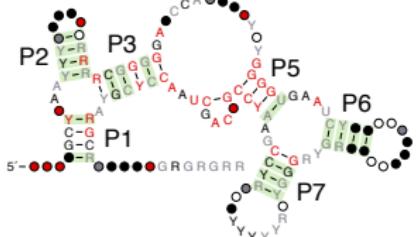


RNAalifold-R-scape



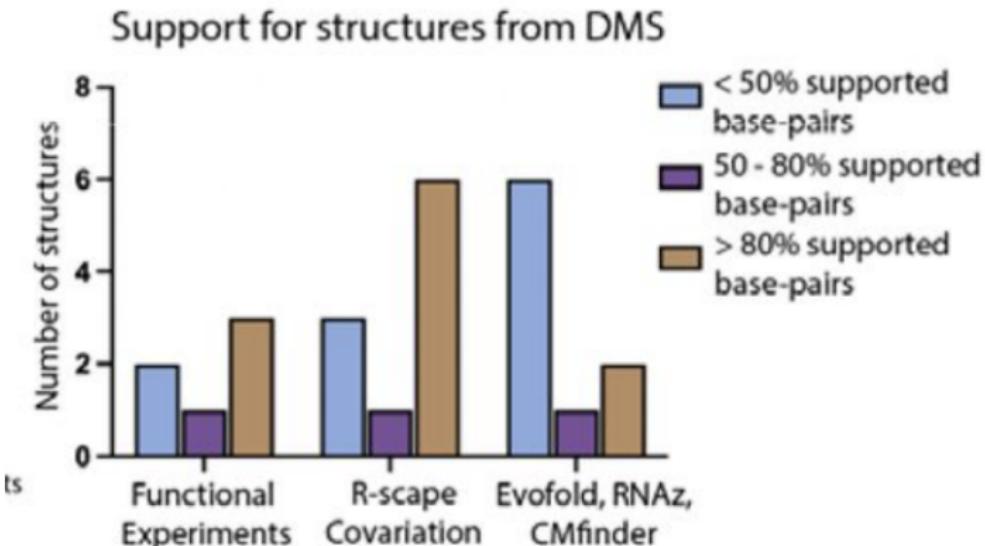
CaCoFold

### Cyclic di-AMP riboswitch



# CaCoFold find in vivo structures

## G structures with support from DMS reactivity for in vivo formation



"RNA structure landscape of *S. cerevisiae* intron"  
Rangan, Hunter, Pham, Ares Jr., Das

