


# Hunting RNA motifs

Paul Gardner

July 23, 2012

# Classification problems


- ▶ What is this?
  - ▶ The second most abundant *M. tuberculosis* transcript
  - ▶ No hits to Rfam, no homologs have been identified, no known function



# Classification problems


- ▶ What about this?

- ▶ The seventh most abundant *N. gonorrhoeae* transcript
- ▶ No hits to Rfam, no homologs have been identified, no known function



# Classification problems


- ▶ And this?
  - ▶ The ninth most abundant *C. difficile* transcript
  - ▶ No hits to Rfam, no homologs have been identified, no known function



# Classification problems

- ▶ And this?

- ▶ The eleventh most abundant *C. difficile* transcript
- ▶ No hits to Rfam, no homologs have been identified, no known function



Rank	Reads	Gene
1	12K	RNaseP RNA
2	12K	6S RNA
3	12K	SRP RNA
4	11K	ldhA
5	10K	tmRNA
6	6K	spoVG
7	6K	Gly reductase
8	5K	slpA
9	4K	sRNA/RUF?
10	4K	hypothetical prot.
11	4K	sRNA/RUF?

Deakin, Lawley  
et al. (2012)

*Unpublished.*

# What is our next big challenge?

- ▶ Past:
  - ▶ How many non-coding RNA genes are there?
- ▶ Future:
  - ▶ Can we determine the functions, if any, for large sets of given RNAs?




# Will a “periodic table of RNA” be useful?

- ▶ My aim is to build an analog to the Periodic Table for classifying RNA families and motifs, enabling researchers to predict function.

base	basepair
ANYA	GNRA
sarric1	sarric2
TRIT	IRE
SECIS	mir-TAR
mir-30	mir-9
lin-4	mir-5
mir-8	mir-8
mir-1	mir-1
mir-2	mir-6
let-7	Y_RNA
UAA_GAN	Csr
Cloop	domV
ktrn1	ktrn2
term1	term2
tRNA	RNaseP
SAM_v	symR
CPeB3	FinP
sroB	msr
SAM_u	HH_3
Vmtnn3	Vmtnn3
l1vK	DsrA
CAESAR	CAESAR
isrK	sroD
isrB	isrB
6C	6C
rspL	subB
dimB	dimB
suCA	SraD
	*
sky	*
RNAl	RNAl
Purine	SAM-Chl
cdGMP_2	Anti-Q
GadY	RnkJdr
PrfA	PrfA
OmrA-B	RyeB
traJ2	traJ2
SraH	23Smeth
23Smeth	DS-pep
Ps-Rho	rnk-ps
	**
Mgsns	Mgsns
tRNAs	Qrr
isrC	HH_1
SNR_24	TrpJdr
greA	greA
preQ12	HAR1F
TermLeu	TermLeu
MicC	C4
RsmY	RsmY
Ribosome	Ribosome

# What is a motif?





- ▶ Wiktionary definitions for **motif**:
  - ▶ A recurring or dominant element; a theme.
- ▶ For the purposes of this work:
  - ▶ an RNA motif is a recurring RNA sequence and/or secondary structure found within larger structures that can be modelled by either a covariance model or a profile HMM.

Escher, MC (1938) Sky and Water 1.

# What is a motif?



- ▶ The kink-turn RNA motif
  - ▶ Found in rRNA, riboswitches, RNase P, snoRNAs, ...
  - ▶ An asymmetric internal loop
  - ▶ Causes a sharp kink between two flanking helical regions



Cruz & Westhof (2011) Sequence-based identification of 3D structural modules in RNA with RMDetect. *Nature Methods*.




# Can we build a seed alignment and CM resource of these motifs?

- ▶ For the hairpin motifs?
  - ▶ YES
- ▶ For the internal loop motifs?
  - ▶ MAYBE
- ▶ For the sequence motifs?
  - ▶ Not certain yet



Doshi et al. (2004) Evaluation of the suitability of free-energy minimization using nearest-neighbor energy parameters for RNA secondary structure prediction. *BMC Bioinformatics*.

# The RNA Motifs: RMfam



5' - ● ● ● ● ● ● ● ● ● ● R R A R G G R G R R ● ● ● ● ● Y A U G R ● ● R ●

5' - ● ● ● ● ● ● ● ● ● ● R R G G R R ● ● ● R Y ● A U G A R ● R ●

5' - ● ● ● ● ● ● ● ● ● ● A A G G R R ● ● ● ● A U G ● ● R ●

# Breaking Rules

- ▶ Prefer motifs found in multiple families &/or multiple locations in the same family
    - ▶ Aligning **analogous** structures
  - ▶ Motifs are allowed to overlap
  - ▶ Models are non-specific
    - ▶ High false-positive rates
  - ▶ Sequences are not edited, spliced or otherwise interfered with
- 
- ▶ Sequence sources are diverse (PDB, EMBL & publications)




# An example entry

```
# STOCKHOLM 1.0

#=GF ID      twist_up
#=GF AU      Gardner PP
#=GF SE      Gardner PP
#=GF SS      Published; PMID:21976732
#=GF GA      19.00
#=GF DR      URL;http://rnabrabase.cs.put.poznan.pl/?act=pdbdetails&id=1S72
#=GF RN      [1]
#=GF RM      21976732
#=GF RT      Clustering RNA structural motifs in ribosomal RNAs using
#=GF RT      secondary structural alignment.
#=GF RA      Zhong C, Zhang S
#=GF RL      Nucleic Acids Res. 2012;40:1307-17.
```

2zkr\_Y/27-46                   CCGAUCUCGUC-UGAUCUCGG  
#=GR\_2zkr\_Y/27-46 SS <<<<---<----->--->>>  
2gv3\_A/3-20                   CCAGACUC---CCGAAUCUGG  
#=GR\_2gv3\_A/3-20 SS (((((. (.----....))))))  
3iz9\_C/29-48                CGGAUCCCGUC-AGAACUCCG  
#=GR\_3iz9\_C/29-48 SS (((((...(. ...-.)...))))))  
3bbo\_B/31-50                CCAA-UCAUCCCGAACUUGG  
#=GR\_3bbo\_B/31-50 SS .(.....<.....>.....).  
1nkw\_9/33-53                CCCACCCCAUGCCGAACUGGG  
#=GR\_1nkw\_9/33-53 SS (((.....(.....))))  
1c2x\_C/31-51                CUGACCCCAUGCCGAACUCAG  
#=GR\_1c2x\_C/31-51 SS ((((...<.....>....))))  
1giy\_B/33-53                CCGUUCCCAUCCGAACACGG  
1S72\_9/29-50                CCGUACCACAUCCGAACACGG  
#=GR\_1S72\_9/29-50 SS (((((...(. ....)....))))  
#=GR\_1S72\_9/29-50 MT ...XXXXXX.....XXXXXX...  
#=GC SS\_cons  
//




# An example entry

```
# STOCKHOLM 1.0

#=GF ID      twist_up
#=GF AU      Gardner PP
#=GF SE      Gardner PP
#=GF SS      Published; PMID:21976732
#=GF GA      19.00
#=GF DR      URL;http://rnabrabase.cs.put.poznan.pl/?act=pdbdetails&id=1S72
#=GF RN      [1]
#=GF RM      21976732
#=GF RT      Clustering RNA structural motifs in ribosomal RNAs using
#=GF RT      secondary structural alignment.
#=GF RA      Zhong C, Zhang S
#=GF RL      Nucleic Acids Res. 2012;40:1307-17.
```


2zkr\_Y/27-46                   CCGAUCUCGUC-UGAUCUCGG  
#=GR\_2zkr\_Y/27-46 SS        <<<<---<          >--->>>  
2gv3\_A/3-20                   CCAGACUC---CCGAAUCUGG  
#=GR\_2gv3\_A/3-20 SS        (((((.(-.-.----))))))  
3iz9\_C/29-48                   CGGAUCCCGUC-AAACUCCG  
#=GR\_3iz9\_C/29-48 SS        ((((...(.--.))....)))  
3bbo\_B/31-50                   CCAA-UCCAUCCCGAACUUGG  
#=GR\_3bbo\_B/31-50 SS        . .... <.....>..... .  
1nkw\_9/33-53                   CCCACCCCCAUUGCCTAACUGGG  
#=GR\_1nkw\_9/33-53 SS        ((.....-.....-.....))  
1c2x\_C/31-51                   CUGACCCCCAUGCCAACUCAG  
#=GR\_1c2x\_C/31-51 SS        ((((...<.....>....))))  
1giy\_B/33-53                   CCGUUCCCCAUUCCGAACACCGG  
#=GR\_1giy\_B/33-53 SS        CCGUUACCCAUUCCGAACACCGG  
1S72\_9/29-50                   ((((...(.....)....))))  
#=GR\_1S72\_9/29-50 SS        ...XXXXXX.....XXXXXX...  
#=GR\_1S72\_9/29-50 MT        ((((...(.....)....))))  
#=GC SS\_cons  
//



# An example annotated Rfam alignment

```
# STOCKHOLM 1.0
#=GF ID      5S_rRNA
#=GF AC      RF00001
#...
#=GF WK      5S_ribosomal_RNA
#=GF SQ      712
#=GF MT.G    GNRAA
#=GF MT.s    sarcin-ricin-2
#=GF MT.t    twist_up
```


X62858.1/62-115	GCUGCG-U-UCUCCGUGUGUACUGC GG UUUU-UUG-CUGUGGGAA--GCCCA CUUCA-CUG
X01588.1/64-115	UCACGU--UAGUGGG---GCCGUGGAUACCGUGAGGAUCC-GCAG-CCCCACU-AAG-CUG
#=GR X01588.1/64-115 MT.0	.....GGGGGGGGGGGGGGGGGGG.....
X05870.1/363-414	UCACGU--UGGUGGG---GCCGUGGAUACCGUGAGGAUCC-GCAG-CCCACU-AAG-CUG
#=GR X05870.1/363-414 MT.0	.....GGGGGGGGGGGGGGGGG.....
L27170.1/63-116	CCAGCG--UCCGGCA--AGUACUGGAGUGC CGGAGCCUCUGGGAA-AUCCGGU-UCG-CCG
#=GR L27170.1/63-116 MT.0	..ssss..ssssss..sssssssssssssssssssssssssssssssssssssss.sssssss.sss.ss..
L27163.1/61-115	CCAGCGGUUCGGGCA--GUACUGGAGUGC CGGAGCCUCUGGGAA-ACCGGGUUCG-CCG
#=GR L27163.1/61-115 SS	...((((((.(((.))))))))(((((.((.(((.)))))))))))....
#=GR L27163.1/61-115 MT.0	.....GGGGGGGGGGGGG..GGGGGGG.....
#=GR L27163.1/61-115 MT.1	..ssssssssssss..ssssssssssssssssssssssssssssss..ssssssssssss.ss..
L27343.1/60-114	CCAGCGAAC CAGCUA--GUACUGAGUGGGAGACCCUCUGGGAG-CGCUGGU-UCG-CCG
#=GR L27343.1/60-114 MT.0	.....GGGGGGGGGGGGGGGGG.....
#=GR L27343.1/60-114 MT.1	...ssssssssss..ssssssssssssssssssssssssssss..ssssssss.sss.s..
M33891.1/66-117	CCAGCG--CCGGAGA--GUACUGCGC-GGGCAACCGCGUGGGAG--GCGAGG-CCGCUCG
#=GR M33891.1/66-117 MT.0	.....GGGG.GGGGGGGGGGGG.....
X01484.1/62-114	GUCAGG--CC CAGUU--AGUACUGAGGUGGGCGACCACUUGGGAA-CACUGGG-U-G-CUG
#=GR X01484.1/62-114 MT.0	.....GGGGGGGGGGGGGGGGG.....
#=GR X01484.1/62-114 MT.1	...sss..ssssss..ssssssssssssssssssssssssss..ssssssss.s.s.ss..
#=GC SS_cons	:::::::<-<-----<<-----<<----->>----->>----->>----->>----->
#=GC RF	cCaGcG..cccgga...GUAcUggGGUggGcgAcCcucUGGgAA.CaCcgg.uCG.CUG
//	




# An example annotated Rfam alignment

```
# STOCKHOLM 1.0
#=GF ID    5S_rRNA
#=GF AC    RF00001
#...
#=GF WK    5S_ribosomal_RNA
#=GF SQ    712
#=GF MT.G  GNRA
#=GF MT.s  sarcin-ricin-2
#=GF MT.t  twist_up
```

X62858.1/62-115	GCUGCG-U-UCUGGGUGU	GUACUC CGG UUUU-UUG-CUGUGGGAA-	GCCCACUUCA-CUG
X01588.1/64-115	UCACGU--UAGUGGG--	CCCGUGGAUACC GUGAGGAUCC-GCAG-CCCCACU-AAG-CUG	
#=GR X01588.1/64-115 MT.0	.....	.....GGGGGGGGGGGGGGGGG.....	.....
X05870.1/363-414	UCACGU--UGGUGGG--	CCCGUGGAUACC GUGAGGAUCC-GCAG-CCCCACU-AAG-CUG	
#=GR X05870.1/363-414 MT.0	.....	.....GGGGGGGGGGGGGGGGG.....	.....
L27170.1/63-116	CCAGCG--UC CCG CA--A	GUACUGGAGUGC GCGAGCC UCGGGAA-AUCCGGU-UCG-CCG	
#=GR L27170.1/63-116 MT.0	....SSSS...SSSSSS...SSSSSSSSSSSSSSSSSSSSSSSSSSSS...SSSSSS.SSS.SS.		
L27163.1/61-115	CCAGCGGUUC CGGG GA--	GUACUGGAGUGC GCGACC CUCUGGGAA-ACC GG UCG-CCG	
#=GR L27163.1/61-115 SS	....(((((.((.)))))))	(((((((.((.))))))))..))..)	....
#=GR L27163.1/61-115 MT.0	.....	.....GGGGGGGGGGGGG..GGGGGGG.....	.....
#=GR L27163.1/61-115 MT.1	....SSSSSSSSSSSS...SSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...SSSSSSSSSS.SS.		
L27343.1/60-114	CCAGCGAAC CAG CUA--	GUACUAGAGUGGAG ACC CUCUGGGAG-CGCUGGU-UCG-CCG	
#=GR L27343.1/60-114 MT.0	.....	.....GGGGGGGGGGGGGGGG.....	.....
#=GR L27343.1/60-114 MT.1	....SSSSSSSSSS...SSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...SSSSSS.SSS.S.		
M33891.1/66-117	CCAGCG--CCGGAGA--	GUACUGCGC-GGGCAACC GCGUGGGAG-GCGAGG-CCGCU CG	
#=GR M33891.1/66-117 MT.0	.....	.....GGGG.GGGGGGGGGGG.....	.....
X01484.1/62-114	GUCAGC--CCCAGUU--A	GUACUGAGGUGG GCGACCACU UGGGAA-CAC UGGG-U-G-CCG	
#=GR X01484.1/62-114 MT.0	.....	.....GGGGGGGGGGGGGGGG.....	.....
#=GR X01484.1/62-114 MT.1	....SSS...SSSSSS...SSSSSSSSSSSSSSSSSSSSSSSSSSSS...SSSSSS.S.S.SS.		
#=GC SS_cons	::::::::::<-<-<-<-<-<-<->>->>->>->:::::		
#=GC RF	cCaGcG..cccgga...GUAcUggGGuggGcgAcCcucUGGgAA.CaCcgg.uCG.CUG		
//			




# Uses for a motif DB (I): Improving a RNA alignments and structure prediction constraints




## Uses for a motif DB (II): RNA structural motifs

### ► The Lysine Riboswitch




# Uses for a motif DB (II): RNA structural motifs

## ► The Lysine Riboswitch




# Uses for a motif DB (III): Functionally characterising novel RNAs



- ▶ TwoAYGGAY sRNA: 216 homologues found in Human gut metagenome sequences,  $\gamma$ -Proteobacteria and Clostridiales species
- ▶ Weinberg Z *et al.* (2010)  
"Comparative genomics reveals 104 candidate structured RNAs from bacteria, archaea and their metagenomes". *Genome Biol.*


# Csr-Rsm background




Toledo-Arana *et al.* (2007) Small noncoding RNAs controlling pathogenesis. *Current Opinion in Microbiology*.

## Uses for a motif DB (IV): Identifying and ranking ncRNAs of interest in transcriptome results

Distribution of motif scores on *S. typhi* ncRNAs




CDF of motif scores on *S. typhi* ncRNAs




# Uses for a motif DB (IV): Identifying and ranking ncRNAs of interest in transcriptome results


Distribution of motif scores on *S. typhi* ncRNAs alignments




CDF of motif scores on *S. typhi* ncRNAs alignments



# 80 highest scoring matches between RMfam & Rfam



# What about the highly expressed ncRNAs?




RNA	RNA code	RNAz	Motifs	complexity	G + C sRNA (gen.)
M.tub. RUF3	No	RNA (prob=1.00)	Terminator,k-turn	No	59% (66%)
N.gon. RUF7	No	RNA (prob=0.94)	Terminator	No	43% (52%)
C.dif. RUF9	? (P=0.044)	RNA (prob=1.00)	Terminator	No	30% (29%)
C.dif. RUF11	? (P=0.079)	RNA (prob=1.00)	Terminator	Yes (23/ 170nts)	26% (29%)




# What about the highly expressed ncRNAs?

- Are they antisense to any 5' UTRs

**RNAup RNA–RNA energies between sRNAs and 5' UTRs**




**Difference between native and shuffled CDFs**




# What about the highly expressed ncRNAs?

- Are they antisense to any 3' UTRs

**RNAup RNA–RNA energies between sRNAs and 3' UTRs**



**Difference between native and shuffled CDFs**





## Discussion points

- ▶ A useful curation tool for building alignments and structures
- ▶ Can provide testable function hypotheses, sometimes
- ▶ Limited use outside of alignments based on the *Salmonella Typhi* results
- ▶ Even in alignments, false-positives & negatives are common
- ▶ More motifs (Terminators, Shine-Dalgarno, ...)
- ▶ Limit trivial motifs and prevent rediscovery of old motifs
- ▶ Snapshots of the data are available from  
<https://github.com/ppgardne/RMfam>
- ▶ Can anyone here do better on MTB3, NGON7, CDIF9 & CDIF11?

# Thanks!

- ▶ Lars Barquist, Alex Bateman & Rob Knight


# RNAbiology



PPG is supported by a Rutherford Discovery Fellowship from Government funding, administered by the Royal Society of New Zealand.



# RNA classification



Backofen et al. (2007) RNAs everywhere: genome-wide annotation of structured RNAs. *Journal of Experimental Zoology*.