

# SNOBOARD

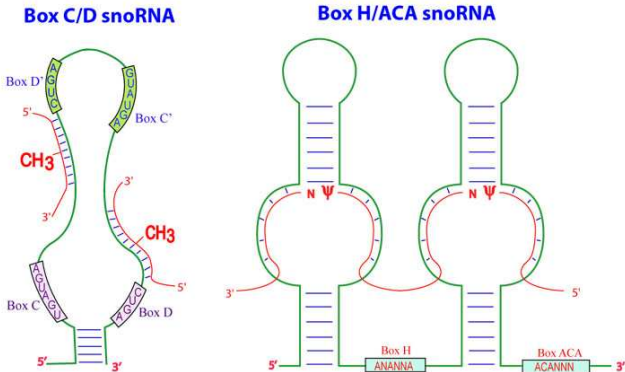
Stephanie Kehr, Sebastian Bartschat

Bioinformatics  
University of Leipzig

Bled, Slovenia, 2010

# SNORNAS

- CD-box: guide methylation of target RNA
- H/ACA-box: guide pseudouridylation of target RNA
- scaRNAs: hybrid of H/ACA and C/D domains



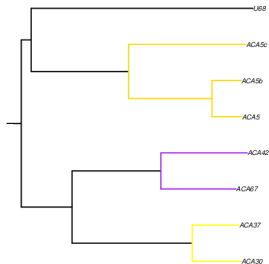
## MOTIVATION

- snoRNA with microRNA most often newly detected ncRNA
- pseudouridylation and methylation are highly abundant post-transcriptional modifications
- guide modifications in rRNA, snRNA, tRNA and some mRNAs
- precursors of microRNAs
- orphan snoRNAs with unknown function



















⇒ *snoRNA set which is as complete as possible to get further insights into snoRNAs and their evolution*

## STARTINGPOINT

- human, chicken, platypus, rhesus, C.elegans
- try to find homologs and paralogs among metazoan
- name conflicts
- paralogs and families??











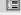
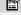

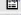






# PROPERTIES

Table ▲		
CD-boxes		
CD-targets		
genome-source		
HACA-boxes		
HACA-targets		
homology		
snoRNA-source		
snoRNAs		
target-source		
9 table(s)		

# PROPERTIES

Table

<b>CD-boxes</b>		
CD-targets		
genome-source		
HACA-boxes		
HACA-targets		
homology		
snoRNA-source		
snoRNAs		
target-source		
9 table(s)		

name	C-box	D-box	position_C	position_D	distance(C/D)
hsa_CD_1	ATGATGA	CTGA	8	66	50
hsa_CD_10-24	TTGATGG	CTGA	17	61	36
hsa_CD_2-1	GTGATGA	CTGA	11	65	46
hsa_CD_37	GTGATGA	CTGA	11	106	87
hsa_CD_47	ATGATGA	CTGA	10	77	59

# PROPERTIES

The image shows a database interface with a table list on the left and a table view on the right. Red circles highlight 'CD-boxes' and 'HACA-targets' in the table list, with red arrows pointing to the 'hsa\_CD\_37' row in the table view.

name	C-box	D-box	position_C	position_D	distance(C/D)
hsa_CD_1	ATGATGA	CTGA	8	66	50
hsa_CD_10-24	TTGATGG	CTGA	17	61	36
hsa_CD_2-1	GTGATGA	CTGA	11	65	46
hsa_CD_37	GTGATGA	CTGA	11	106	87
hsa_CD_47	ATGATGA	CTGA	10	77	59

Field	Type
<b>name</b>	varchar(50)
<b>hairpin_number</b>	int(11)
<b>target</b>	varchar(50)
<b>target-source</b>	int(11)
<b>modification_site</b>	int(11)
<b>target_sequence</b>	varchar(20)
<b>left_binding_site</b>	varchar(20)
<b>dist(LeftEnd/RightStart)</b>	int(11)
<b>right_binding_site</b>	varchar(20)
<b>dist(RightEnd/Box)</b>	int(11)
<b>binding_structure</b>	varchar(150)
<b>binding_energy</b>	double
<b>hairpin_sequence</b>	varchar(150)
<b>status</b>	varchar(100)

# PROPERTIES

name	C-box	D-box	position_C	position_D	distance(C/D)
hsa_CD_1	ATGATGA	CTGA	8	66	50
hsa_CD_10-24	TTGATGG	CTGA	17	61	36
hsa_CD_2-1	GTGATGA	CTGA	11	65	46
hsa_CD_37	GTGATGA	CTGA	11	106	87
hsa_CD_47	ATGATGA	CTGA	10	77	59

Table ▾

CD-boxes		
CD-targets		
genome-source		
HACA-boxes		
HACA-targets		
homology		
snoRNA source		
snoRNAs		
target-source		

9 table(s)

Field	Type
name	varchar(50)
hairpin_number	int(11)
target	varchar(50)
target-source	int(11)
modification_site	int(11)
target_sequence	varchar(20)
left_binding_site	varchar(20)
dist(LeftEnd/RightStart)	int(11)
right_binding_site	varchar(20)
dist(RightEnd/Box)	int(11)
binding_structure	varchar(150)
binding_energy	double
hairpin_sequence	varchar(150)
status	varchar(100)

name	sequence	chromosome	position	strand	length
ptr_HACA_43	GCTGTCCTGGACCTGTTGGACACACAGACAGTTGCTGCTGCTGCCTGTG...	chr9	136970559,136970694	2	136

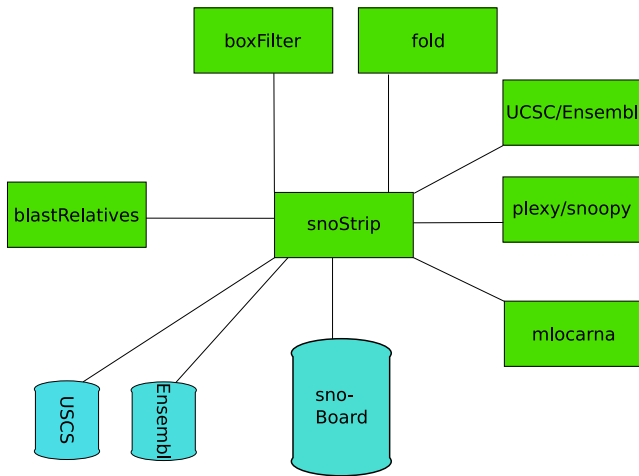
  

hostgene	localisation	intronnumber	transcript ID	structure	minimum_free_energy	found_with
SNHG7 ()	intronic	2,1,1	NR_024543.NR_024542.NR_003672	.(((((((.(.((((.( (((.....))))))))).....	-51.80	hsa_HACA_43



# HOMOLOGY

name	H.sapiens	P.troglodytes	P.pygmaeus	M.mulatta
HACA_1	hsa_HACA_1	ptr_HACA_1-1,ptr_HACA_1-2	ppy_HACA_1-1,ppy_HACA_1-2,ppy_HACA_1-3	mml_HACA_1-1,mml_HACA_1-2
HACA_10	hsa_HACA_10	ptr_HACA_10	ppy_HACA_10	mml_HACA_10
HACA_2	hsa_HACA_2-1,hsa_HACA_2-2	ptr_HACA_2	ppy_HACA_2	mml_HACA_2-1,mml_HACA_2-2
HACA_23	hsa_HACA_23	ptr_HACA_23	ppy_HACA_23	mml_HACA_23
HACA_38	hsa_HACA_38-1,hsa_HACA_38-2	ptr_HACA_38	ppy_HACA_38	mml_HACA_38
HACA_39	hsa_HACA_39	ptr_HACA_39	ppy_HACA_39	mml_HACA_39
HACA_64	hsa_HACA_64	ptr_HACA_64-1,ptr_HACA_64-2,ptr_HACA_64-3,ptr_HACA...	ppy_HACA_64	mml_HACA_64
HACA_68	hsa_HACA_68-1,hsa_HACA_68-2	ptr_HACA_68-1,ptr_HACA_68-2	ppy_HACA_68	mml_HACA_68-1,mml_HACA_68-2



## BLAST-SEARCH

- stepwise approach using all paralogs
- best blasthit as anchor for cutoff values
- decrease penalties for mismatches, gap opening and gap extensions

Alternatives:

- christian's chainer
- infernal search runs

# BOX-FILTER

FASTA



**clustalw**

```
TTGGCTAGGTTTCATGATGACACAGGACCTTGCTGATCATAATGATTTCAAAAATTGAGCTTAAAAATGACACTCTGAAATCCAGTCAG...
TTGGCTAGGTTTCATGATGACACAGGACCTTGCTGAACATAATGATTTCAAAAATTGAGCTTAAAAATGACACTCTGGAATCCAGTCAG...
TTGGCTAGGTTTCATGATGACACAGGACCTTGCTGAACATAATGATTTCAAAAATTGAGCTTAAAAATGACACTCTGAAATCCAGTCAG...
TTGGCTAGGTTTCATGATGACATAGGACCTTGCTGAACATAATGATTTCAAAAATTGAGCTTAAAAATGACGCTCTGAAATCCAGTCAG...
TTGGCTAGGTTTCATGATGACATAGGACCTTGCTGAACATAATGATTTCAAAAATTGAGCTTAAAAATGAGCTTAAAAATGACACTT.GGAAGCCAGTCAA...
TTGACTGGATTTATGATGATATAGGACCTTGCTGAAATATAATGATTTCAAAAATTGAGCTTACAG.TGAC...TGAAATCTAGTCAATGT
```

# BOX-FILTER

FASTA

**clustalw**

```
TTGGCTAGGTT CATGATGAC ACAGGACCTTGCTGATCATAATGATTTCAAAAATTGAGCTTAAAAATGACACTCTGAAATCCAGTCAG...
TTGGCTAGGTT CATGATGAC ACAGGACCTTGCTGAAACATAATGATTTCAAAAATTGAGCTTAAAAATGACACTCTGAAATCCAGTCAG...
TTGGCTAGGTT CATGATGAC ACAGGACCTTGCTGAAACATAATGATTTCAAAAATTGAGCTTAAAAATGACACTCTGAAATCCAGTCAG...
TTGGCTAGGTT CATGATGAC ATAGGACCTTGCTGAAACATAATGATTTCAAAAATTGAGCTTAAAAATGACGCTCTGAAATCCAGTCAG...
TTGGCTAGGTT CATGATGAC ATAGGACCTTGCTGAAACATAATGATTTCAAAAATTGAGCTTAAAAATGACGCTCTGAAATCCAGTCAG...
TTGACTGGATT CATGATGAC ATAGGACCTTGCTGAAATATAATGATTTCAAAAATTGAGCTTACAG.TGAC...TGAAATCTAGTCAATGT
```

**isCorrectBox?**

# BOX-FILTER

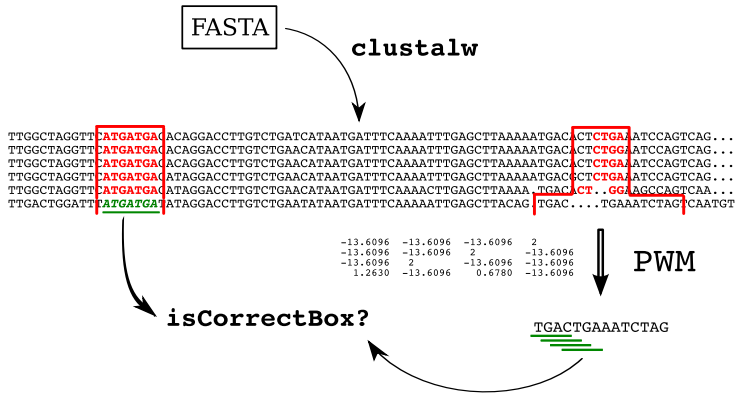
FASTA

clustalw

```
TTGGCTAGGTT CATGATGA ACAGGACCTTGCTGATCATAATGATTTCAAAATTTGAGCTTAAAAATGAC ACTCTGA ATCCAGTCAG...
TTGGCTAGGTT CATGATGA ACAGGACCTTGCTGAACATAATGATTTCAAAATTTGAGCTTAAAAATGAC ACTCTGG ATCCAGTCAG...
TTGGCTAGGTT CATGATGA ACAGGACCTTGCTGAACATAATGATTTCAAAATTTGAGCTTAAAAATGAC ACTCTGA ATCCAGTCAG...
TTGGCTAGGTT CATGATGA ATAGGACCTTGCTGAACATAATGATTTCAAAATTTGAGCTTAAAAATGAC ACTCTGA ATCCAGTCAG...
TTGGCTAGGTT CATGATGA ATAGGACCTTGCTGAACATAATGATTTCAAAACTTGAGCTTAAAA TGACACT . GGAGCCAGTCAA...
TTGACTGGATT CATGATGA ATAGGACCTTGCTGAATATAATGATTTCAAAATTTGAGCTTACAG TGAC . . . . TGAAATCTAG CAATGT
```

isCorrectBox?

## BOX-FILTER



# STRUCTURE

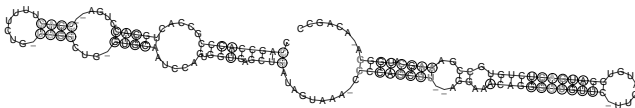
- fold single sequences with RNAsubopt

```

GUGCCUUUUAAAGGUUGACCCAGUGCUUUAAAGAGGCUAACACAGAGGGUAAAGUAAGUCUCCAUAAAACCCAGAGAAGAGACUGGAAAGCUCUUCUUGGAUCCUGUCUGGAGUCACAACU -3524
..(((((((.....)).....)).....)).....).....(.....((((.....((((.....(((.....(((.....))).....))).....))).....)).....)).....)..... -35.24
..(((((((.....)).....)).....)).....).....(.....((((.....((((.....(((.....(((.....))).....))).....))).....)).....)).....)..... -35.24
..(((((((.....)).....)).....)).....).....(.....((((.....((((.....(((.....(((.....))).....))).....))).....)).....)).....)..... -35.24
..(((((((.....)).....)).....)).....).....(.....((((.....((((.....(((.....(((.....))).....))).....))).....)).....)).....)..... -35.20

```

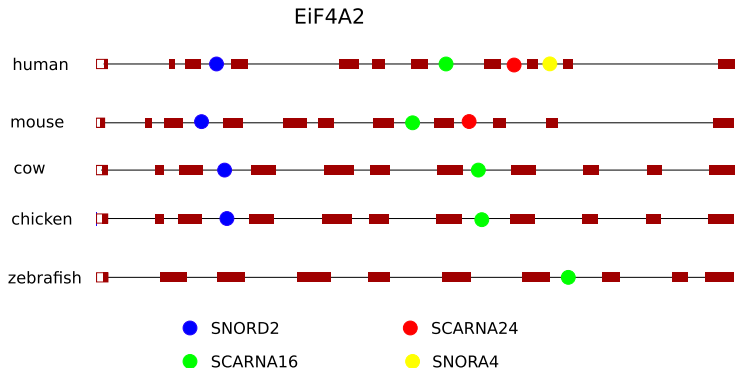
- fold alignments with RNAalifold





## GENOMIC ORGANISATION & HOSTGENE

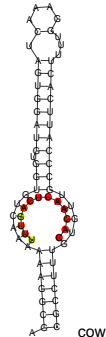
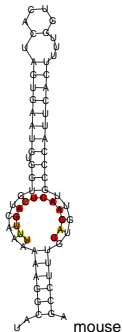
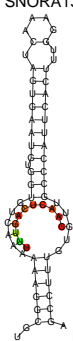
- known gene annotation from UCSC and ENSEMBL
- in vertebrata mostly located in introns
- recognize evolutionary events like intron or hostgene switches
- problem alternative splicing



# TARGET PREDICTION

- HACA: RNAsnoop for each hairpin
- apply svm to evaluate results
- CD: compute RNAduplex for 15nts upstream of D and D'-box
- use accessibility-profiles of target RNAs by RNAup

SNORA13-hairpin1





# ALIGNMENT CONSTRUCTION

- use mlocarna to compute final alignments
- structural constraints
- anchors constraints to align boxes

```

>H.sapiens
AGCCTTTGTGTTGCCATTCACTTTGGAACTAGTGAATGTGGTGTCAAAAAGGCGTAAATTAACGCTTTGCAGCCTTTCTGCGCCTTAAATTTGATACCTTTGGTGTAGGAGCTGCATAAGTAACAGTT
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....#S
.....AAAAAAAAA.....BBBBB.....#1
.....123456.....#2
>M.musculus
AGCCTTTGTGTTGCCATTCACTTTGGTCACTAGTGAATGTGGTGTCAAAAAGGCATAAATAATGCTTTGCGGCCTTCTGCGCCTGGAGTTTGGTATCTGGGTGTACGAGCTGCATAAGTAACAGTA
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....#S
.....AAAAAAAAA.....BBBBB.....#1
.....123456.....#2

>H.sapiens
AAGTGAATGATGGCAATCATCTTTGGGACTGACCTGAAATGAAGAGAATACTATTGCTGATCACTT
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....#S
.....AAAAAAA.....BBBBB.....#1
.....1234567.....#2
>M.musculus
AAGTGAATGATGGCAATCATCTTTGGGACTGACCTGAAATGAAGAGAATACTATTGCTGATCACTTA
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....#S
.....AAAAAAA.....BBBBB.....#1
.....1234567.....#2

```



# THANK YOU

Hakim, Jana, Peter, Christian