# The Brave New World of Non-Coding RNAs
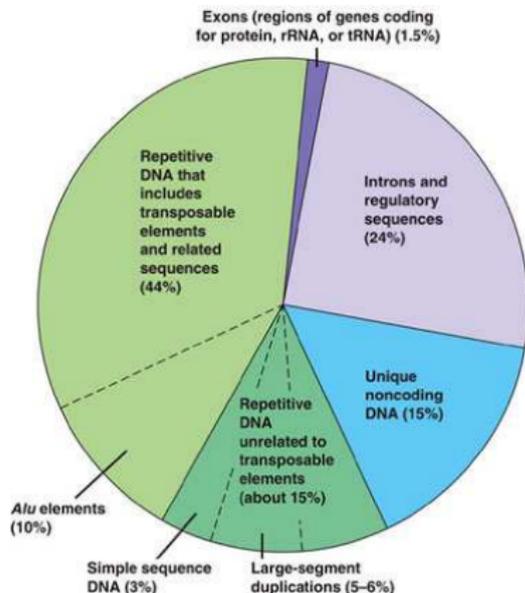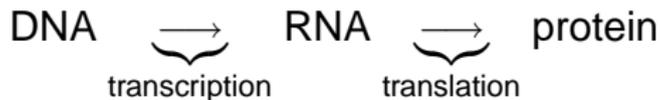
## Peter F. Stadler

Bioinformatics Group, Dept. of Computer Science &
Interdisciplinary Center for Bioinformatics,
**University of Leipzig**
Max-Planck-Institute for Mathematics in the Sciences
RNomics Group, Fraunhofer Institute for Cell Therapy and Immunology
Institute for Theoretical Chemistry, Univ. of Vienna (external faculty)
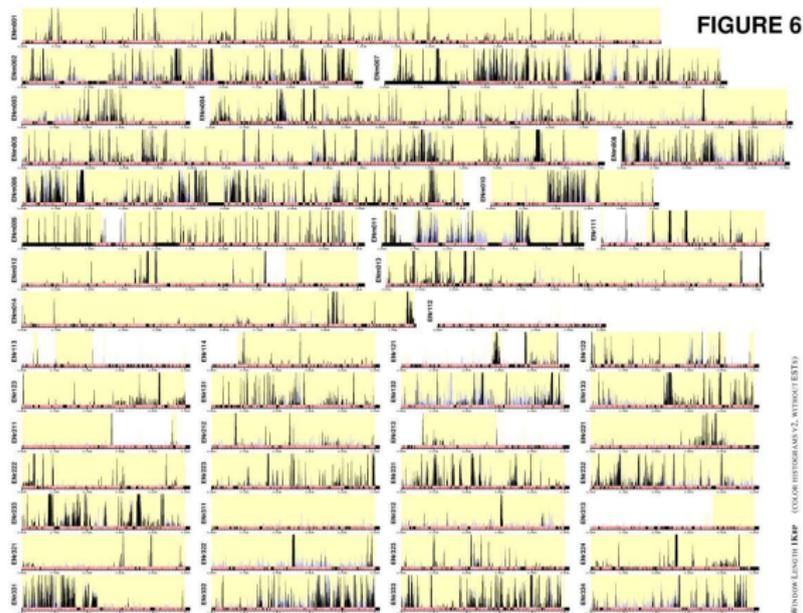The Santa Fe Institute (external faculty)

Jena, Aug 2010

DNA $\longrightarrow$ RNA $\longrightarrow$ protein

transcription     translation



Exons (regions of genes coding for protein, rRNA, or tRNA) (1.5%)

Repetitive DNA that includes transposable elements and related sequences (44%)

Introns and regulatory sequences (24%)

Unique noncoding DNA (15%)

Repetitive DNA unrelated to transposable elements (about 15%)

Alu elements (10%)

Simple sequence DNA (3%)

Large-segment duplications (5–6%)

- only 3% of the non-repetitive part of genome codes for proteins
- Is all the rest **junk DNA**?
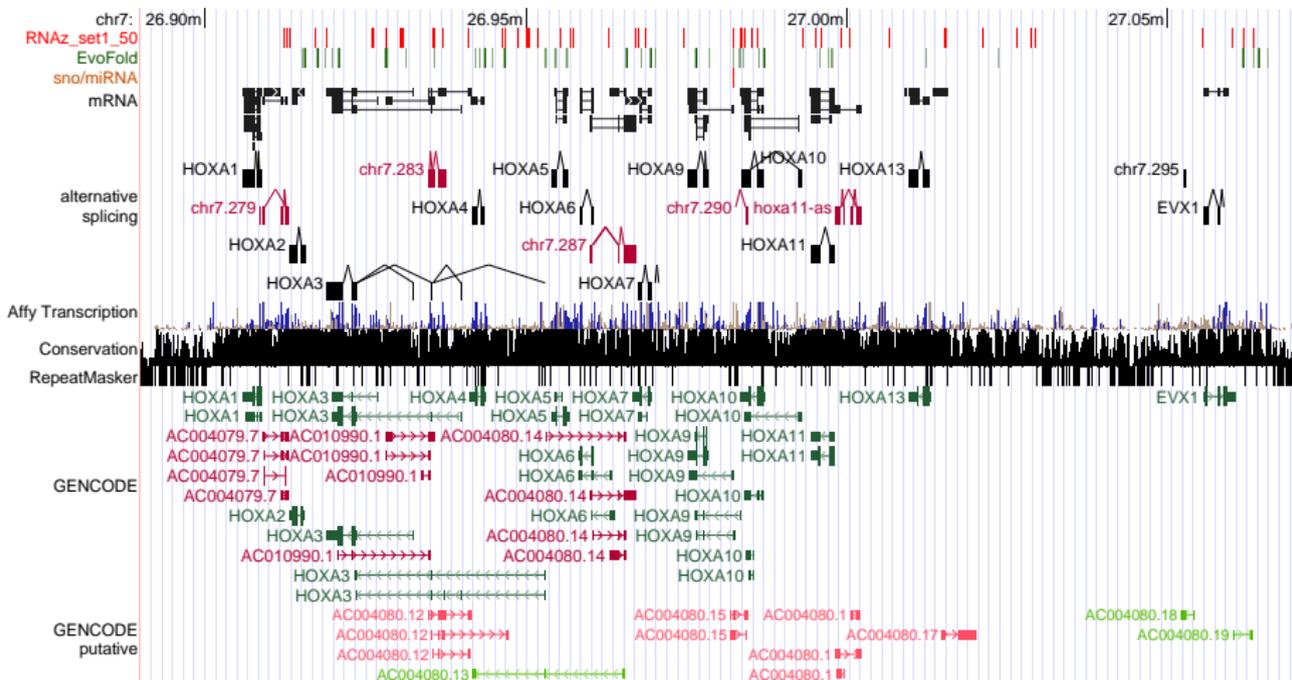- Are all the repeats just genomic parasites?

# Pervasive Transcription



FIGURE 6

More than 90% of the non-repetitive genome shows evidence for transcription in at least one direction
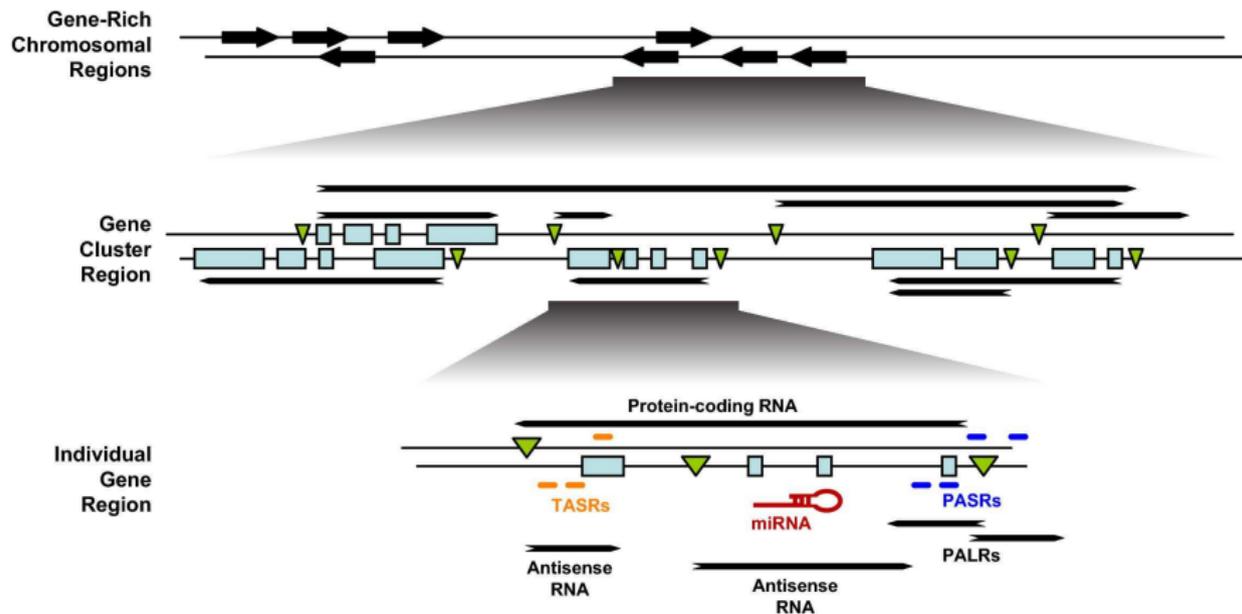
The ENCODE Consortium, *Nature*, 447: 779-816 (2007).

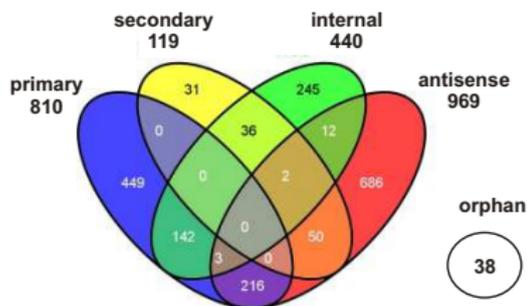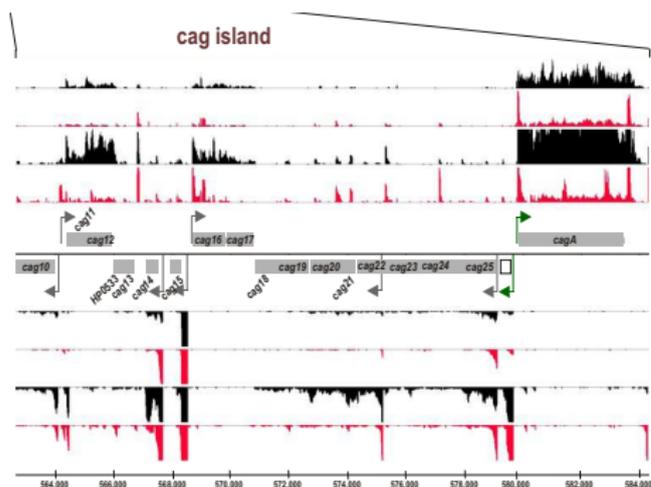# Transcriptome Complexity



*Hox A* cluster.

# Transcriptome Complexity



Science 316: 1484-1488 (2007)

# *H. phylori* doesn't read textbooks

mapping of transcription start sites in *Helicobacter pylori*
secondary start-sites and pervasive antisense transcription



Nature **464**: 250-255 (2010)

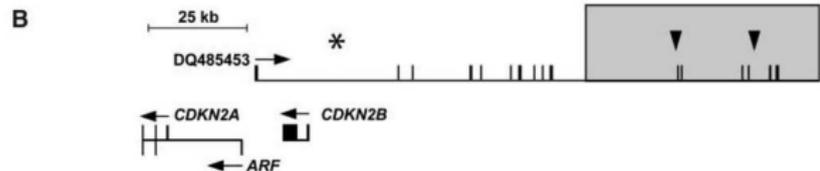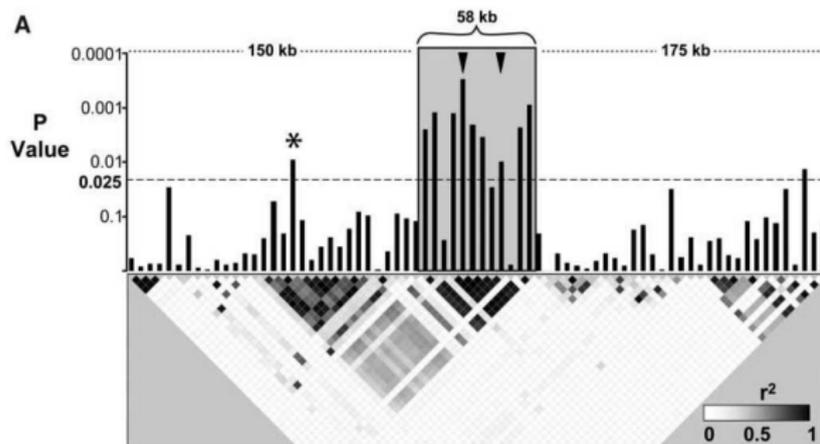# A New Paradigm of Molecular Biology!

- There is no junk!
  Most of the human genome is transcribed, and there are good reasons to believe to most of the transcripts have function
- Most "genes" do not code for proteins
  We have to re-think — and maybe even abandon — the very notion of a gene
- Are these ncRNAs really functional????

## Evidence for ncRNA function

- A small number of well-studied transcripts have functions identifyable by genetic methods (e.g. deletion/complementation)
- Statistical arguments:
    - differential regulation
    - Conservation at sequence level
    - Conservation of RNA structure
    - Conservation of splicing patterns
    - Association with (disease) phenotypes
    - Specific processing

# CHD QTL Locus

The majority of QTLs for complex multi-genic diseases fall into non-coding regions



Association of coronary heart disease (CHD) with a 58kb region on chr. 9p21

non-coding locus, produces the ANRIL transcript(s) ANRIL expression is associated with the atherosclerosis risk
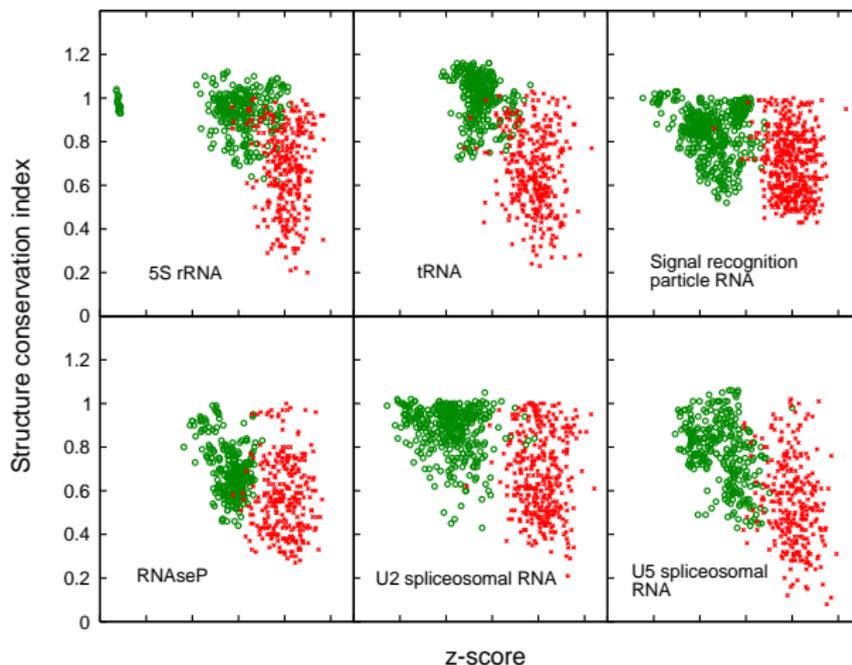
Holdt *et al. Arterioscler Thromb Vasc Biol.*, **30**, 620-627 (2010)

McPherson *et al.*, Science (2007)

# Computational RNA Gene Finding

- Many (but by no means all known functional RNAs are structured, i.e. certain base pairing patterns must be conserved
- This implies that substitutions are not random, but must be consistent with (GC→GU) or even compensate for base pairs (GC→AU)
- Empirical Observation: Known ncRNAs are (a little bit) more stable than genomic background with the same base composition.

IDEA: use this to build a gene finder
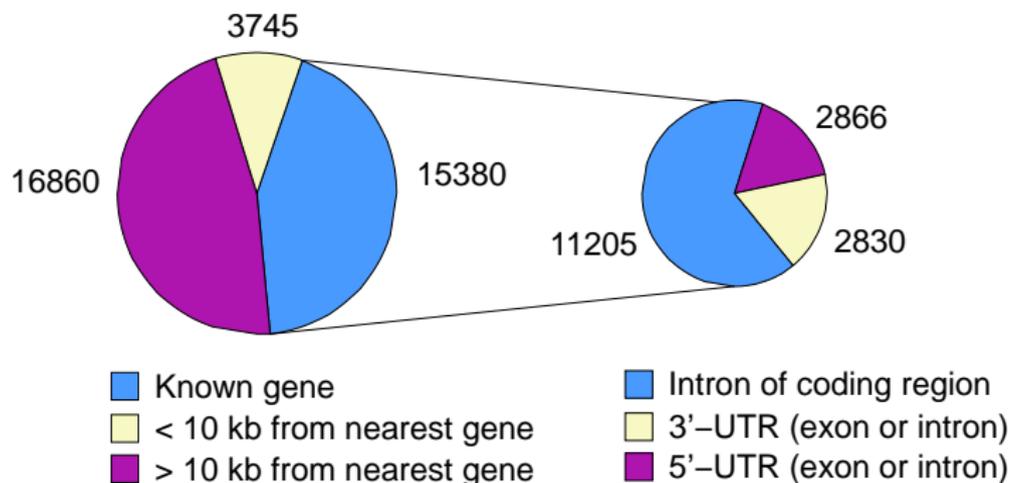
# RNAz: a gene finder for structured RNA



Separation of native ncRNAs from random controls in two dimensions

# Structured RNAs in the Human Genome

# Structured RNAs in the Human Genome

Mammalian genomes contain $\sim 10^5$ structured RNA motifs
Statistics of the highest-confidence fraction ($\sim 36000$):
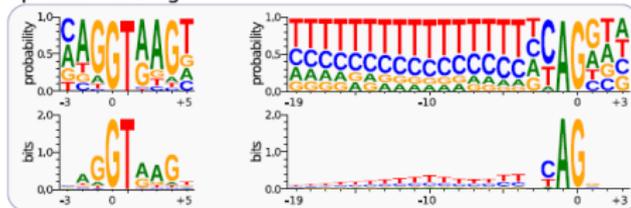


3745

16860

15380

11205

2866

2830

- Known gene
- < 10 kb from nearest gene
- > 10 kb from nearest gene

- Intron of coding region
- 3'–UTR (exon or intron)
- 5'–UTR (exon or intron)

# Finding mRNA-like ncRNAs

- long = contains at least one intron
- predict non-coding transcripts by predicting **conserved** **short** introns
- Why introns?
  - intron evolution is slow and essentially independent of the evolution of the mature sequence
  - splice sites are often conserved
  - disruption of correct splicing usually destroys function
  - ! non-coding transcrips do not have randomly placed large in/dels.
- Why short introns?
  - Most *Drosophila* introns are short.
  - Can be accurately predicted (94% with both splice sites correct)
- Intron prediction (Lim & Burge 1999): machine learning using patterns of donor, acceptor, intron length, branch point, intron composition

# mlncRNAs – splice sites

# Intron-prediction pipeline
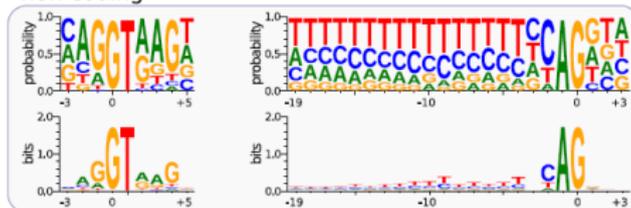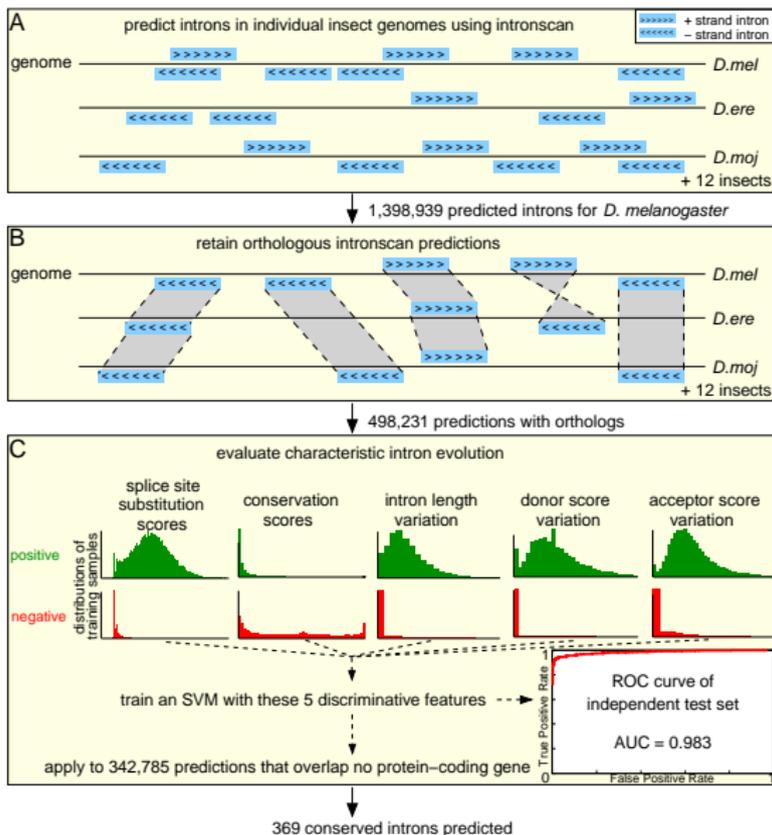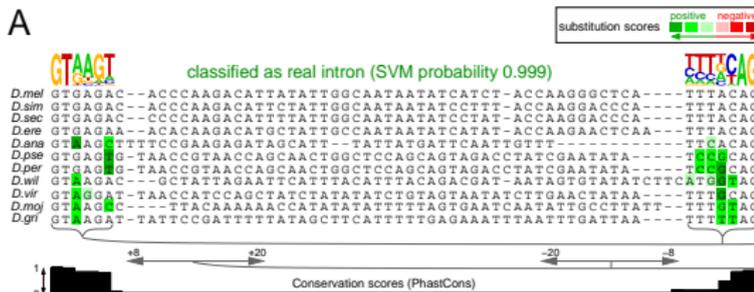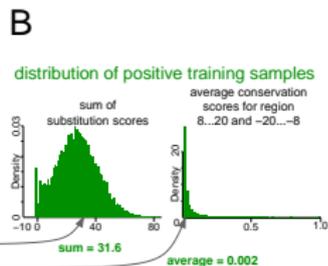
11 out of 17 predictions verified by PCR and sequencing

Expression of transcripts

and existance of introns

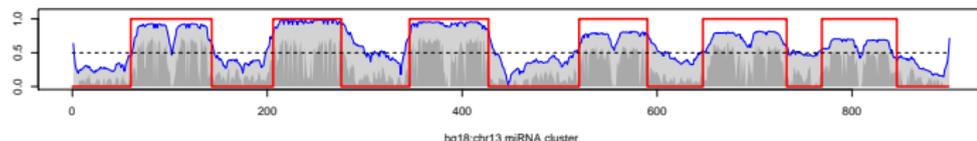also verified in 3 other fly species

Embryo

Larva

Pupa

male
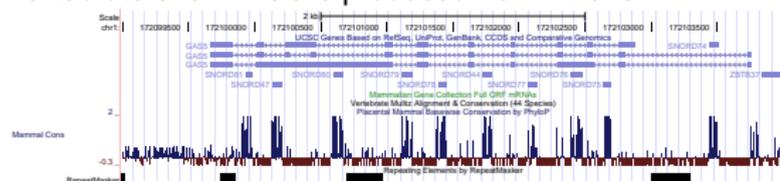
female

# Generation of small RNAs

Many types of small RNAs are produced from longer precursors. In many cases, these precursors are mRNA-like pol-II transcripts.

**1** primary microRNA precursors.
miRNAs a processed out of either exons or introns



hg18:chr13 miRNA cluster

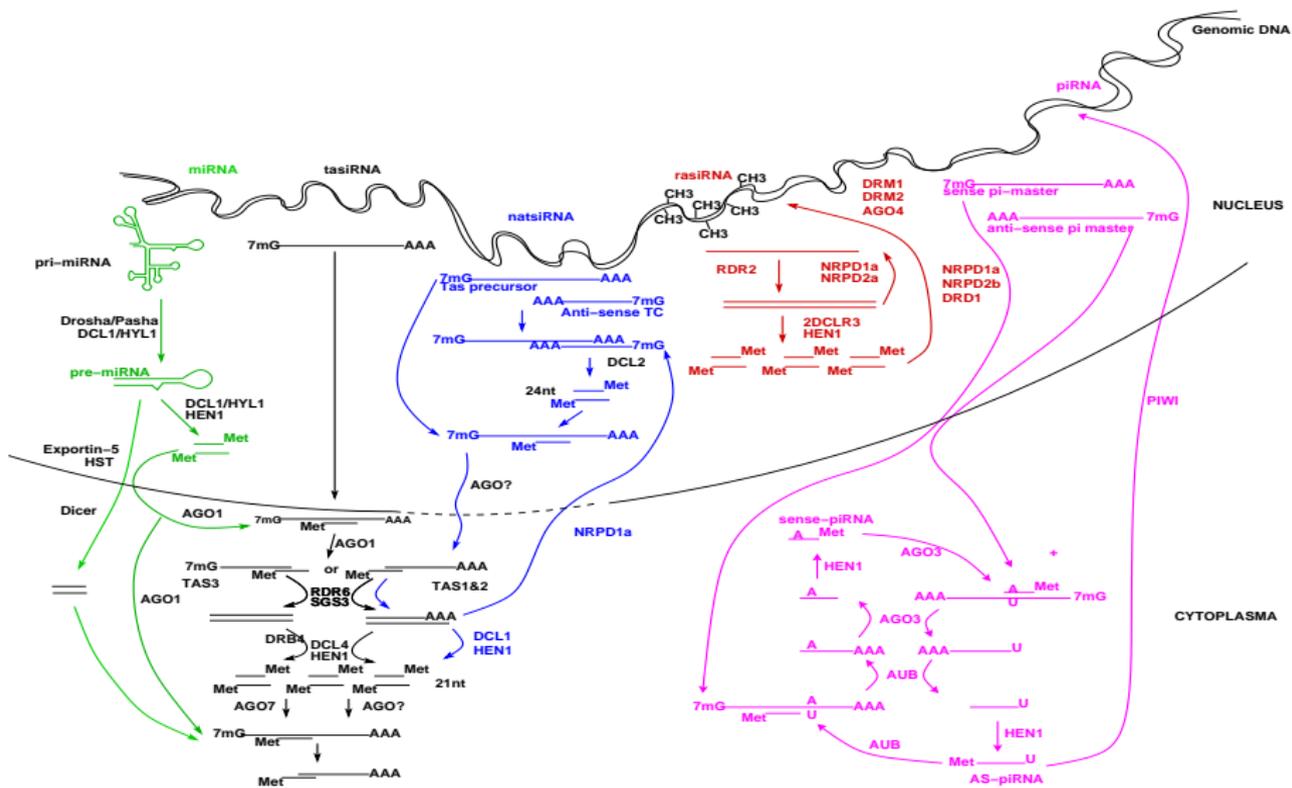primary precursor of the *mir-17* cluster

**2** snoRNA host genes
vertebrate snoRNAs are produced from introns



**3** some piRNA precursors

**4** Affymetrix high-density arrays showed that at least 1% of the human genome produces small RNAs
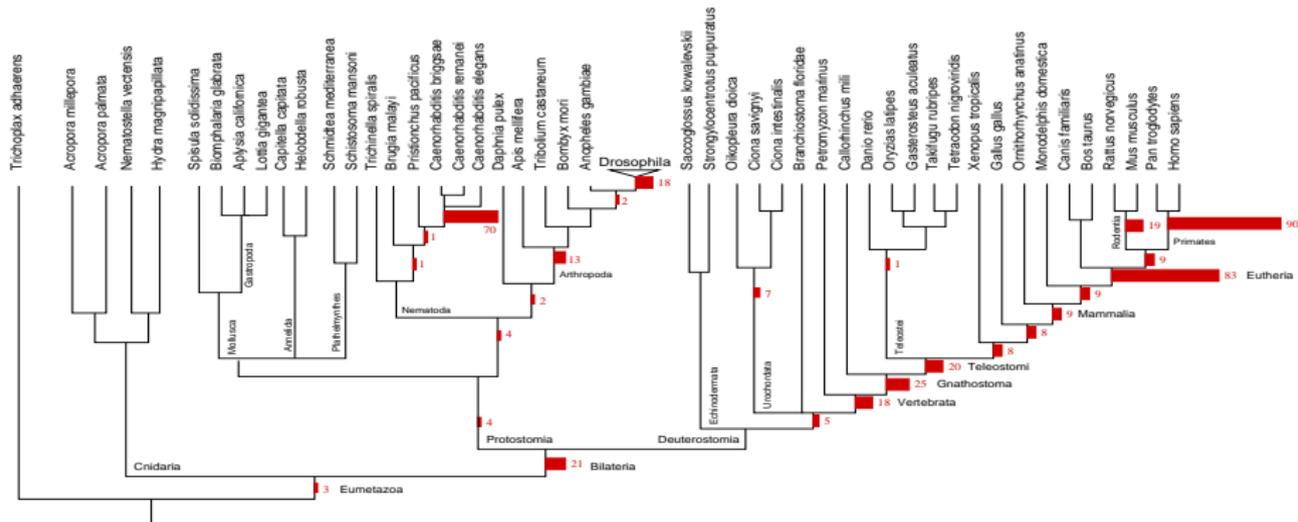Science **316**: 1484-1488 (2007) (joint work with Affymetrix)

## MicroRNAs: Innovation

- Most protein-coding genes are evolutionarily old. For instance, there are no or very few new transcription factor that were invented throughout vertebrate evolution
- Small RNAs, in particular microRNAs, however, are readily created *de novo*.
- Is there a link between ncRNA innovation and novelty at the organismal level?

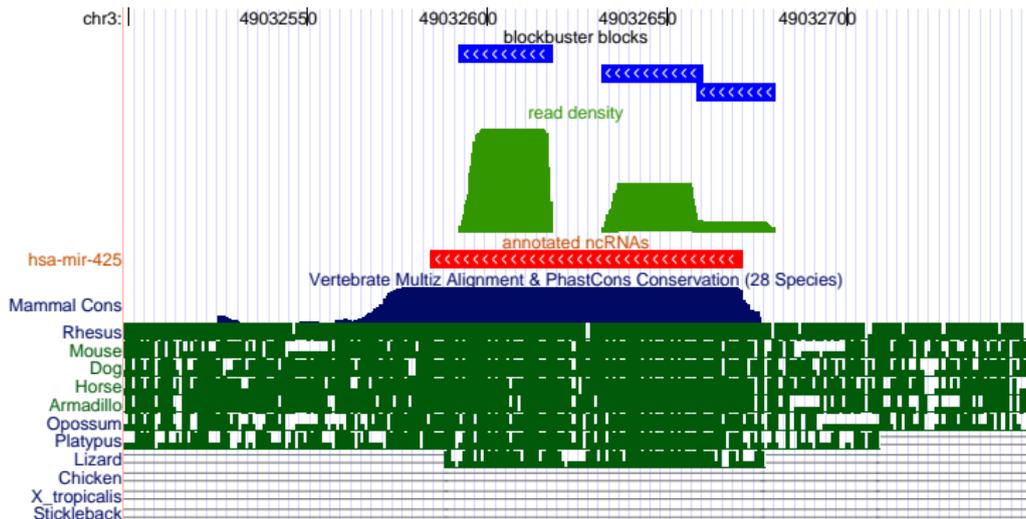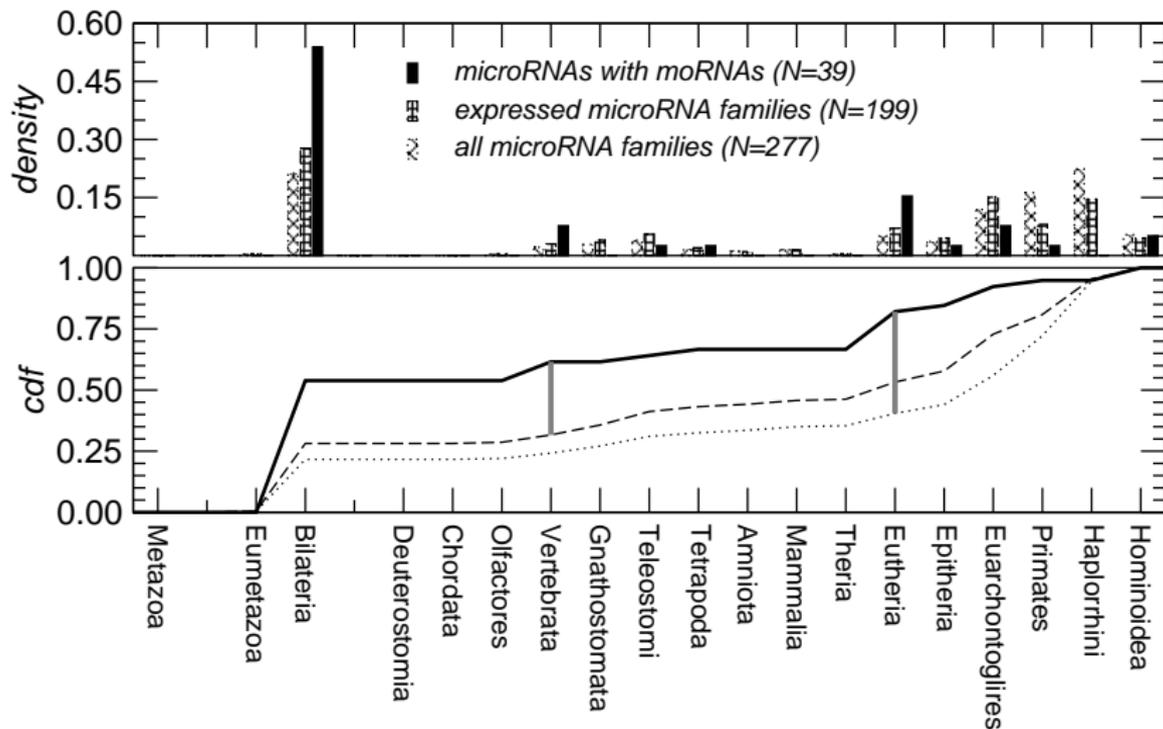BMC Genomics **7**: 15 (2006), updated

# MicroRNA Offset RNAs



Distribution of short reads at the *hsa-mir-425* locus. There are three clearly distinct blocks of reads: the two more abundant ones correspond to miR and miR*, the third one to the 5'moRNA.

density / cdf plot legend:
- microRNAs with moRNAs (N=39)
- expressed microRNA families (N=199)
- all microRNA families (N=277)

x-axis categories: Metazoa, Eumetazoa, Bilateria, Deuterostomia, Chordata, Olfactores, Vertebrata, Gnathostomata, Teleostomi, Tetrapoda, Amniota, Mammalia, Theria, Eutheria, Epitheria, Euarchontoglires, Primates, Haplorrhini, Hominoidea

High-density tiling array screen of human small RNAs



More than 1% of the human genome is transcribed into small RNAs.

in *Aspergillus fumigatus*

recently discovered by several groups also in mammals

# Block patterns are source specific

A first attempt:
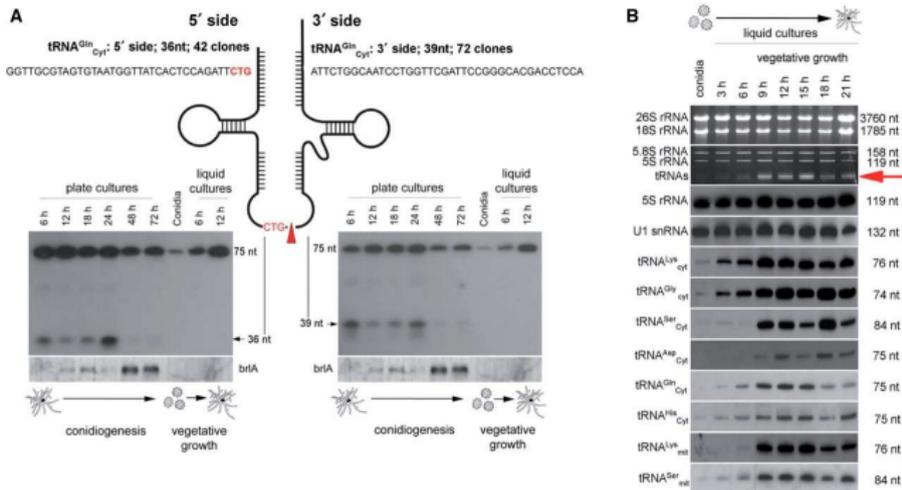random forrest classificator to distinguish the block patterns of
microRNAs, the two snoRNA classes, and tRNAs.
Confusion matrix: (10-fold crossvalidation)

| class | classified as | | | | |
|-------|-------|-------|------|------|-------|
|       | miRNA | H/ACA | C/D  | tRNA | other |
| miRNA | 249   | 2     | 6    | 8    | 21    |
| H/ACA | 6     | 8     | 5    | 2    | 4     |
| C/D   | 20    | 3     | 82   | 13   | 22    |
| tRNA  | 7     | 0     | 12   | 310  | 41    |
| other | 25    | 4     | 16   | 56   | 312   |

# New RNA Classes from Structural Clustering

Comparison of RNA secondary structures:

- Structure-enhanced alignments (e.g. `stral`)
- Tree-alignment or Tree-editing (e.g. `RNAforrester`, `MARNA`)
- `RNAshapes`
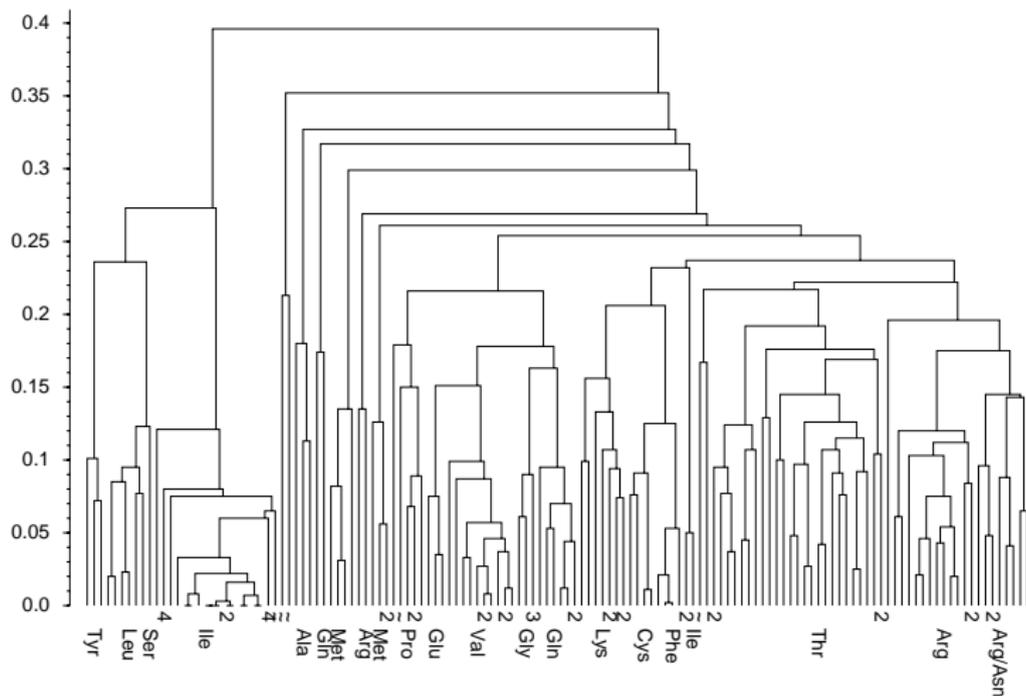- Variants of the *Sankoff* algorithm



`locarna`: a Sankoff-based local structure alignment tool

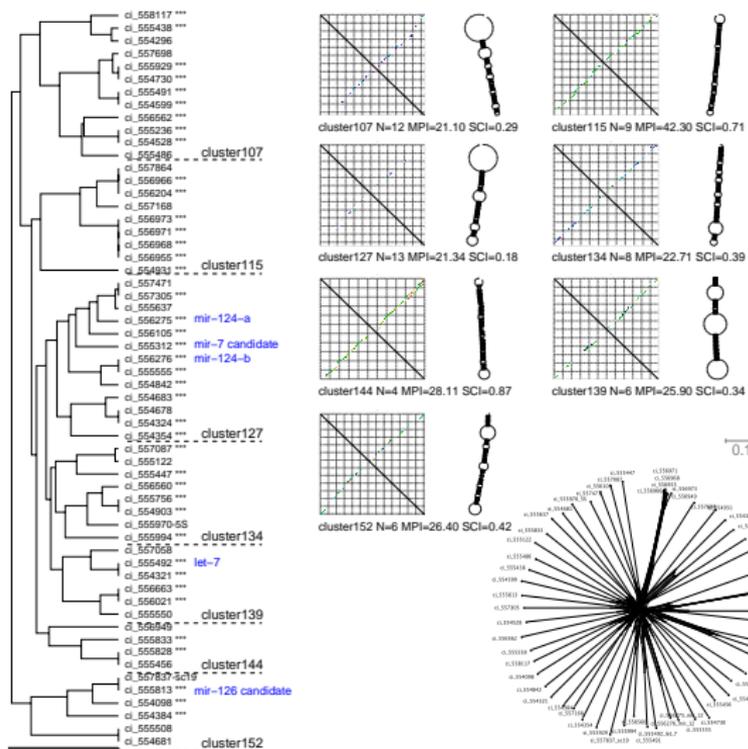Trick: use thermodynamically most plausible base-pairs only

Joint work with Rolf Backofen's group

# Clustering *Ciona intestinalis* `RNAz` Predictions


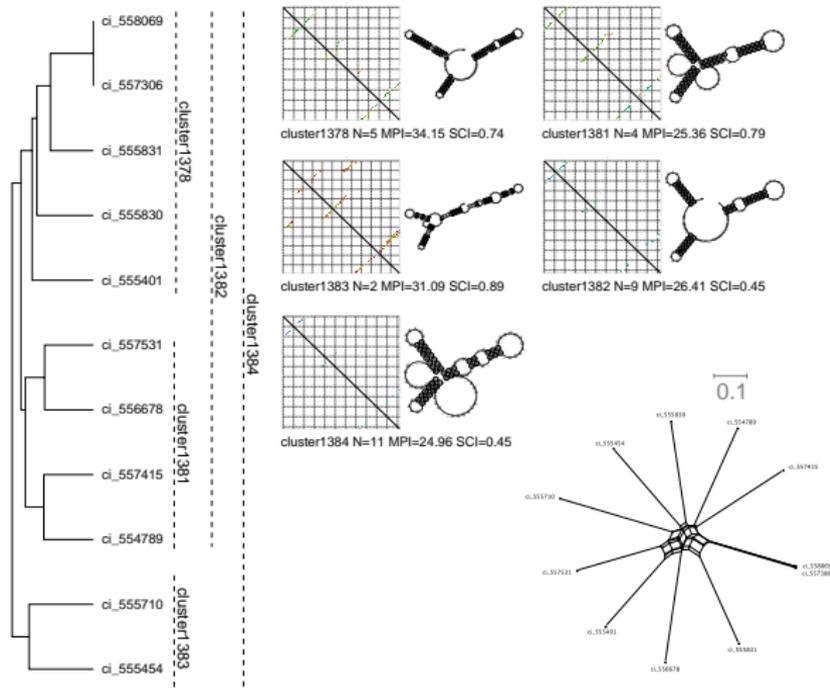
tRNAs subtree from a clustering 3332 ncRNA candidates

*Ciona intestinalis*: microRNA subtree

# Putative Novel RNA Classes



cluster1378 N=5 MPI=34.15 SCI=0.74

cluster1381 N=4 MPI=25.36 SCI=0.79

cluster1383 N=2 MPI=31.09 SCI=0.89

cluster1382 N=9 MPI=26.41 SCI=0.45

cluster1384 N=11 MPI=24.96 SCI=0.45

# Many, many thanks . . .

- Leipzig: Sonja J. Prohaska, Dominic Rose, Jana Hertel, Manja Marz, Claudia & Roman Stocsits, Sven Findeiß, . . .
  FH RNomics group: Jörg Hackermüller, Antje Kretzschmar, Kristin Reiche, Kathy Schutt, Kerstin Ullmann, Tine Schulz
- Vienna: Stefan Washietl, Ivo L. Hofacker, Christoph Flamm, Andrea Tanzer, Stefan Bernhart, Hakim Tafer, Susanne Rauscher, Caroline Thurner, Christina Witwer, . . .
- Halle: Günter Reuter's Lab
- Freiburg: Rolf Backofen, Sebastian Will
- Tübingen: Kay Nieselt's Group
- Freiburg: Rolf Backofen's Lab
- Würzburg: Jörg Vogel, Cynthia Sharma
- Copenhagen: Jan Gorodkin, Stefan Seemann, Peter Menzel
- Affymetrix: Tom Gingeras, Phil Kapranov, *et al.*
- CAS Beijing: Wei Deng and all the others in Runsheng Chen's Lab
- PICB Shanghai: Axel Mosig and Phil Khaitovich and their students (PICB/SIBS)
- ASU Tempe: Julian L. Chen and his lab
- Stanford: Michael Hiller
- ENCODE: Ewan Birney and $10^{2.5}$ coauthors