# Detecting non-Coding RNA in Genomic Sequences

I.    Overview of ncRNAs

II.    What's specific about RNA detection ?

III.   Looking for known RNAs

IV.   Looking for unknown RNAs

**Daniel Gautheret**

INSERM ERM 206 & Université de la Méditerranée

# Non-coding or non-messenger RNA

★ All organisms
  • rRNA 5S/5.8S 15S/18S 23S/28S (5-300 copies)
  • RNAse P/MRP (1 copy)
  • tRNA (20 diff., 50-1000 copies)

rRNA & RNAse P: catalytic!

★ Eukaryotes and Archaes
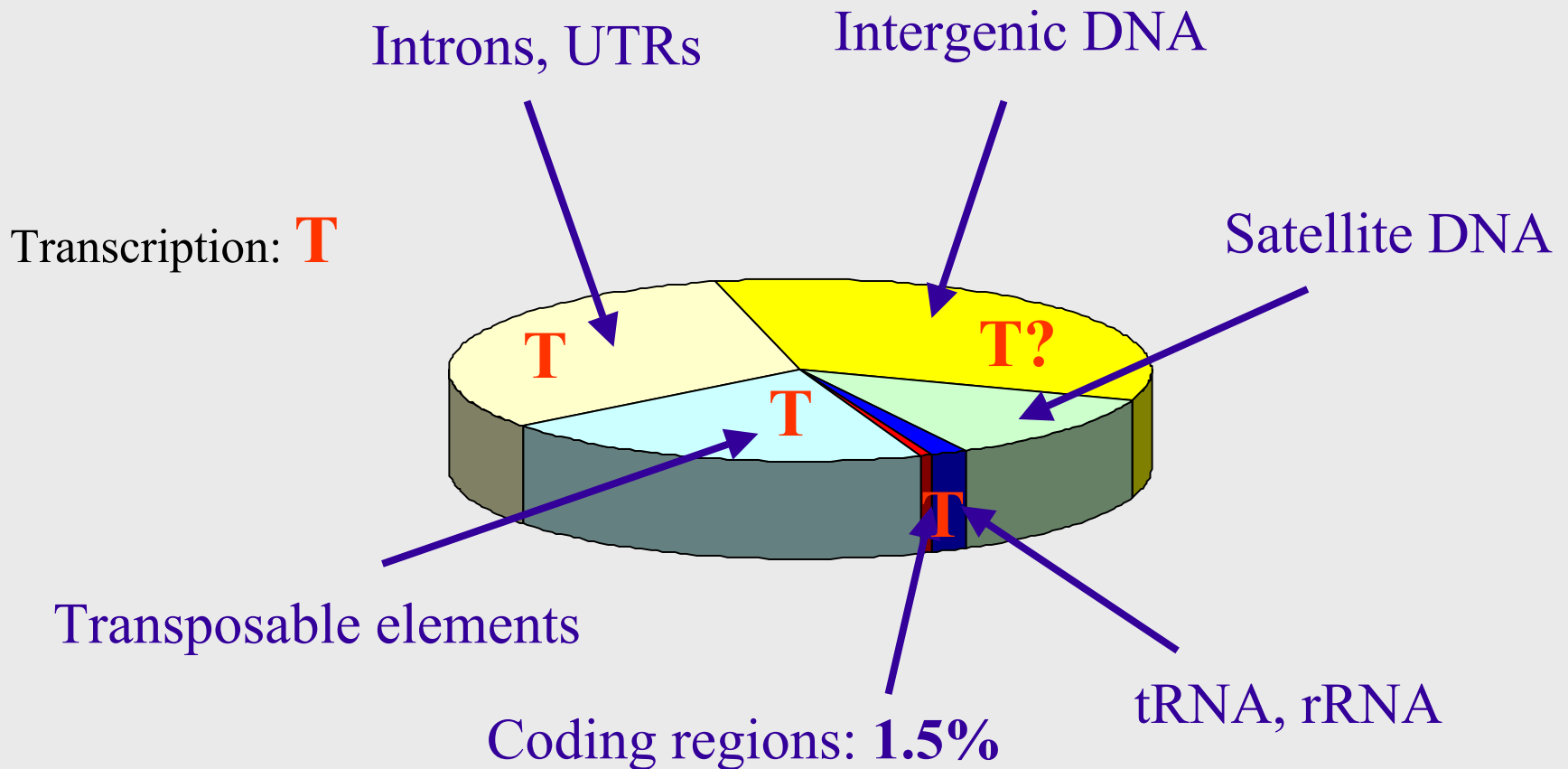  • snoRNAs (H/ACA and C/D: 10-100 different)

★ Eukaryotes
  • miRNAs (100-200 diff.)
  • XIST, H19, IPW (vertebrates)
  • snRNAs (U1, U2, U4, U5, U6)

★ Procaryotes
  • rprA, csrB, oxyS,…
  • tmRNA

# Vertebrate Genomes: > 50% transcribed!

Vertebrate gene: 30kb (coding: 1,5kb)

Introns, UTRs

Intergenic DNA

Transcription: **T**

Satellite DNA

**T**

**T?**

**T**

Transposable elements

**T**

Coding regions: **1.5%**

tRNA, rRNA

# How many other ncRNAs ?

★ Rnomics*
  - Extract total RNA
  - Isolate small RNAs (> mRNA size)
  - Tag & Reverse transcribe
  - Clone & Sequence

★ Success
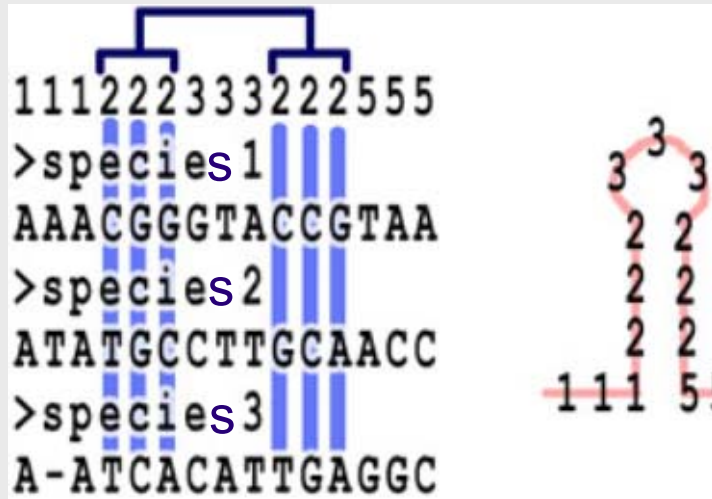  - 80 ncRNAs identified in mouse
  - Lots of new snoRNAs

★ Limitations
  - Access to nuclear RNAs
  - Tissue/time specific expression may be frequent

\* Huttenhofer et al., 2001

# Bioinformatics

★ Not sensitive to « rare » expression

★ Highly succesful in identifying protein-coding genes

★ Many complete genomes available

★ Large computational toolbox

- Statistics
- Thermodynamics
- Phylogeny

# What's Special About ncRNA Detection?



★ No ORF

★ No Markov model / sequence statistics

★ ncRNA is defined both by primary and secondary structure

★ « Substitution matrices » for nucleic acids are terrible compared to aminoacids counterparts

# III. Looking for known ncRNAs

★ How can you detect ncRNA genes from known families?

# Custom RNA search programs

Based on a variety of algorithms: seek regular expressions or base pairs, weight matrices, SCFGs, etc.

- tRNA
    - trnascan (Fichant & Burks 91)
    - trnascan-SE (Eddy 94)
- C/D Box snoRNAs (Lowe & Eddy 99)
    - One type of snoRNA
- miRNA
    - Not an automated procedure

# Descriptor-based programs

- RnaMot / Rnamotif (Gautheret 91, Macke '02)
- Palingol (Viari 96)
- Patscan (Overbeek '00)
- PatSearch (Pesole '01)

```
h1 s1 h1 s2 h2 s3 h2

h1 5:5 1
h2 5:5 NNNNR:YNNNN
s1 7:7 NUNNNNN
s2 4:40
s3 7:7 UUCNNNN
```

RnaMot descriptor for anticodon+TYC domain of tRNA

# Descriptor-based programs

| PROS | CONS |
|---|---|
| Draft descriptors can be **quickly sketched** and tested | Requires a **good prior knowledge** of secondary structure and sequence constraints |
| **No alignment is required**, although it is very helpful to have one | Requires **basic computer skills** to translate biological constraints into the descriptor language |
| **Biologists decide** what features are important or not (see also CONS!) | Biologists have the responsibility of correctly **weighting each important feature** |

# Probabilistic ncRNA search programs

Stochastic Context Free Grammars (first adaptation
of CFG to RNA: Searls 94; SCFG: Eddy & Durbin 94)

| Training set | → | Production rules |
|---|---|---|

describe how to generate
any structure in the
training set

★ Time cost = $O(N^4)$ for sequence of length N

★ Not « practical » for large alignments or genome-wide searches

★ Pseudoknots not allowed

# ERPIN: Profile-based search

Gautheret & Lambert, JMB,
2001, 313, p. 1005.

Training set



Helix profile
(16xN)

Single-strand
profile (5xN)

Search algorithm combines
**dynamic programming** for single
strands and **profile search** for
helices

# Profile-based search

| PROS | CONS |
|------|------|
| **All constraints** in the training set are **efficiently exploited**, resulting in highly specific detections | Alignment and secondary structure constraints must be accurate |
| After alignments and secondary structures are created, **no further programming is needed** | **Helices of variable length need to be reduced** to their shortest consensus |
| Scoring system is defined automatically | Program will not depart from initial alignment in terms of motif size |
| *E*-values are provided for each hit | **Users still have to decide on search order** and masked elements |

# Running a successful ncRNA search

Example: the Signal Recognition Particle (SRP) RNA

★ 172 sequences available

★ All 3 kingdoms

★ Signature: 50-nt domain IV

# Organize ncRNA information

★ Alignment is a must

★ Should be structure-based

★ ClustalW OK only as a first attempt

★ RNAalifold (Vienna package) can identify covarying basepairs



Secondary Structure annotation

Will help identify sequence/ structure constraints: helix sizes, conserved bases, etc.

# Want to publish your finds?
# Prepare Control Procedures

Sensitivity: SN = $\dfrac{TP}{TP+FN}$ ⟵ **Total « true » objects**

Specificity : SP = $\dfrac{TP}{TP+FP}$ ⟵ **Total predictions**

**TP** and **FN**: easy to obtain, using training set (*leave-one out*)

**FP**: harder! How do you know a hit is false?

*Hint*: express SP as: **FP / Mb in a random sequence**

    Make it large enough and of same composition (mono & di-nt) as search database (e.g. with the *shuffle* program)

# Using the ERPIN program

```
>structure
00000000000000000000000000000000000000000001111001000
22244333333366655588877777788899996661111443222
>AQU.AEO.
AGGGUGAACU-CCCCCAGGCCCGAA--AGGGAGCAAGGGUAAGC-CCG
>THE.THE.
GGCGUGAACC-GGGUCAGGUCCGGA--AGGAAGCAGCCCUAAGC-GCC
```

erpin srp.epn sequence.fasta -8,8 -nomask
erpin srp.epn sequence.fasta –2,2 -nomask
erpin srp.epn sequence.fasta -2,2 –umask 5 9 -nomask

**ERPIN results**

**Score:** based on profile values

**E-value:** How many hits expected at this score or higher?

**No need for random sequence tests!**

# The ERPIN Server

http://tagc.univ-mrs.fr/erpin/

**All searches parameterized to scan a bacterial genome in less than 5 minutes**

ncRNA detection

ncRNA detection

# IV. *De novo* ncRNA finding

★ How can we detect ncRNA genes when no prior sequence/structure data is available?

# Exciting times for comparative genomics

**Numerous potentially functional but non-genic conserved sequences on human chromosome 21**

Emmanouil T. Dermitzakis*, Alexandre Reymond*, Robert I
Nathalie Scamuffa*, Catherine Ucla*, Samuel Deutsch*,
Brian J. Stevenson†‡, Volker Flegel†‡, Philipp Bucher†§,
C. Victor Jongeneel†‡ & Stylianos E. Antonarakis*

**Research Update** — *TRENDS in Genetics* Vol.17 No.7 July 2001 — 373

**Selective constraint in intergenic regions of human and mouse genomes**

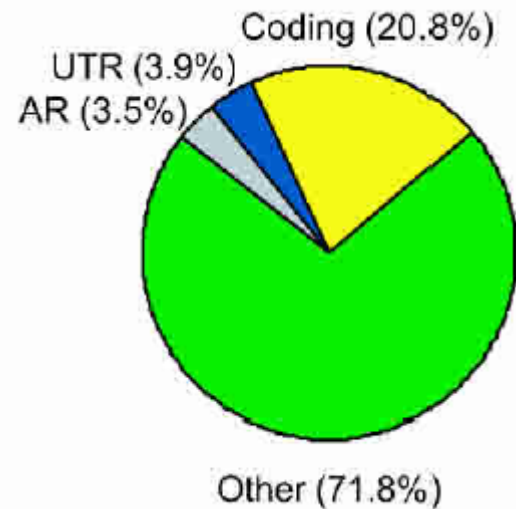Svetlana A. Shabalina, Aleksey Yu.Ogurtsov, Vasily A. Kondrashov and Alexey S. Kondrashov

★ 5-6% of mamalian genome under selection vs 1.5% coding

★ 3 times as much as in nematodes!

★ *« Intergenic regions might hold the key to the complexity of mammals »*

# Functional assignment of conserved regions

★ Coding exons

★ Regulatory non coding exons and introns

★ Promoters

★ ncRNA

★ Ancestral repeats

★ Others (matrix attachment, etc.)

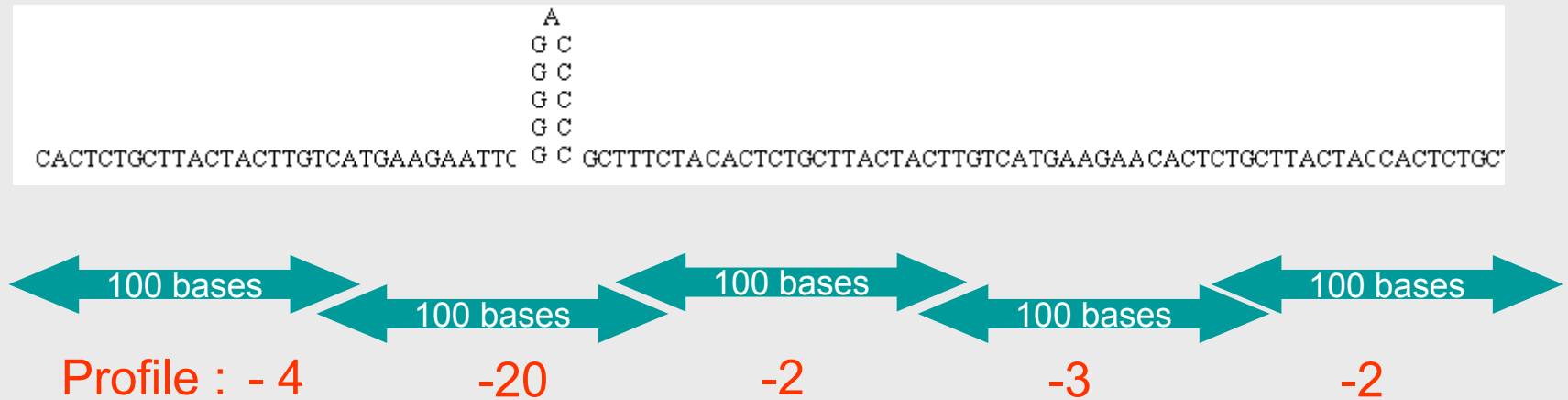**Detect this!**

**Fraction of conserved sequences in.. (AR=ancestral repeats)**



**Margulies et al, 2003**

# Thermodynamic Profiling (Le *et al*. 88)

```
              A
            G   C
            G   C
            G   C
            G   C
            G   C
CACTCTGCTTACTACTTGTCATGAAGAATTC G C GCTTTCTACACTCTGCTTACTACTTGTCATGAAGAACACTCTGCTTACTACCACTCTGC
```

100 bases   100 bases   100 bases   100 bases   100 bases

Profile :  - 4      -20        -2        -3        -2

$$Z\text{-score} = \frac{\text{window free energy} - \text{mean (energy of rnd seq.)}}{\sqrt{\text{Var(energy of rnd seq.)}}}$$

→ New software by Hofacker *et al*: (RNALfold)
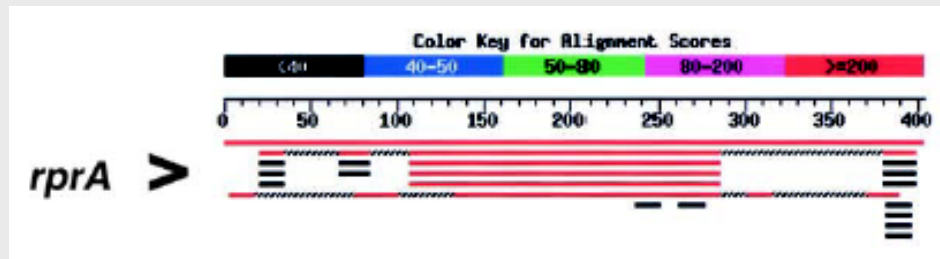
# The problem with thermodynamics

★ OK for strong local structures (some success in viral genomes)

★ However: <u>true ncRNA (tRNA, rRNA) do not display higher folding energy than random sequences of same composition</u> (di-nt: Rivas & Eddy 2000)

# G+C content

★ G+C content alone is a better ncRNA predictor than free energy

★ <u>In high A+T background</u> (thermophilic archaebacteria), <u>ncRNA stand out clearly</u>.

★ Combining (G+C)% and CpG% provides the best discriminant (Schattner '02).

★ Does not work in genomes with « normal » G+C contents, except as a complement to other methods (thermodynamics, etc.)

⟶ See software RNAGenie (Carter, Dubchak & Holbrook, 2001).
    - Combines energy and G+C contents

# Comparative Genomics + experiments

Bacteria: microarray + Northern in different growth conditions



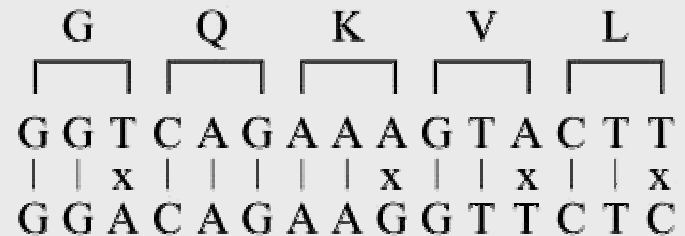From Wassarman et al. '01
Seq: Escherichia, Salmonella, Klebsiella

Wassarman *et. al.* '01: <u>60 ncRNA predicted</u>, 23 confirmed

★ Argaman et al. '01: 24 predicted, 14 confirmed

# Q-RNA (Rivas & Eddy 2001)

★ Analysis of Blast alignment (SCFG based)

• Model for protein coding gene



Synonymous mutations

• Model for ncRNA

(also include loop probabilities obtained from training set of real ncRNA)



Compensatory mutations

# Q-RNA results

★ Limited range for similarity (65%-85%)

- Too dissimilar: incorrect Blast alignments
- Too similar: no covariation

⟶ Problem: Human/mouse/rat ncRNAs not in this range!

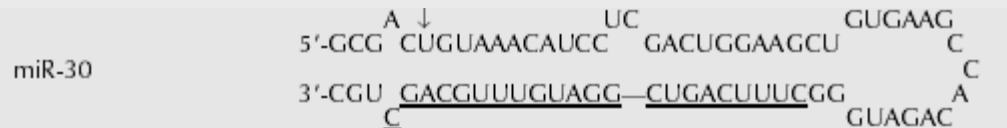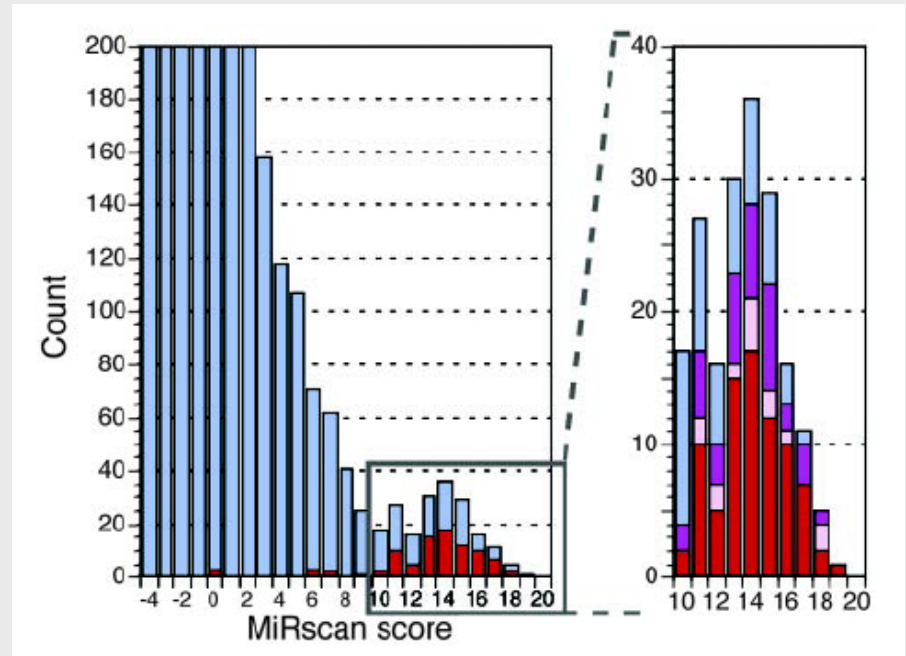★ *E.coli* vs *Salmonella typhy*: analysis of ~5000 Blast hits

- 115 true ncRNAs
- 33 with Blastn alignments in the 65%-85% range
- 33 detected as ncRNA
- 440 other candidates (half of them known elements: terminators, palindromic repetitive elements, etc.)

# Comparative Genomics & miRNA

*Lim et al. Science 2003*

★ Criteria:

- Loose conservation human/mouse/fugu
- Fall outside of protein coding gene
- Predicted to form stem-loop
  - ➤ 15000 hits

- Score based on resemblance to 21mer miRNA
- 107 potential new miRNAs





miR-30

# The right species for ncRNA detection?

➢Human/mouse ncRNA: ~98-100% id
➢18S fugu/xenopus/human: 95% id! Still too close
➢Obvious interest for older animals



**Human/mouse Asp tRNAs**

# Multiple species is the key

★ Multiple alignments will enable covariation detection

★ Covariation + GC-content + energy will provide enough evidence for ncRNA status