

---

# Prediction of functional siRNA structure

---

Zuzana Hořejší  
1.4.2005



Ústav Molekulární  
Genetiky  
odd. Molekulární  
Virologie

Czech **FOBIA**

---

# Previous and current work

Institute of Molecular Genetics, dep. Of Molecular Virology  
- projects focused on transcription factor c-Myb

2002-2004 Copenhagen Institute of cancer Biology, dep. of  
Cell Cycle and Cancer  
- cell cycle checkpoint pathways (DNA damage checkpoint  
pathways in G1/S phase)

---

---

# RNA interference

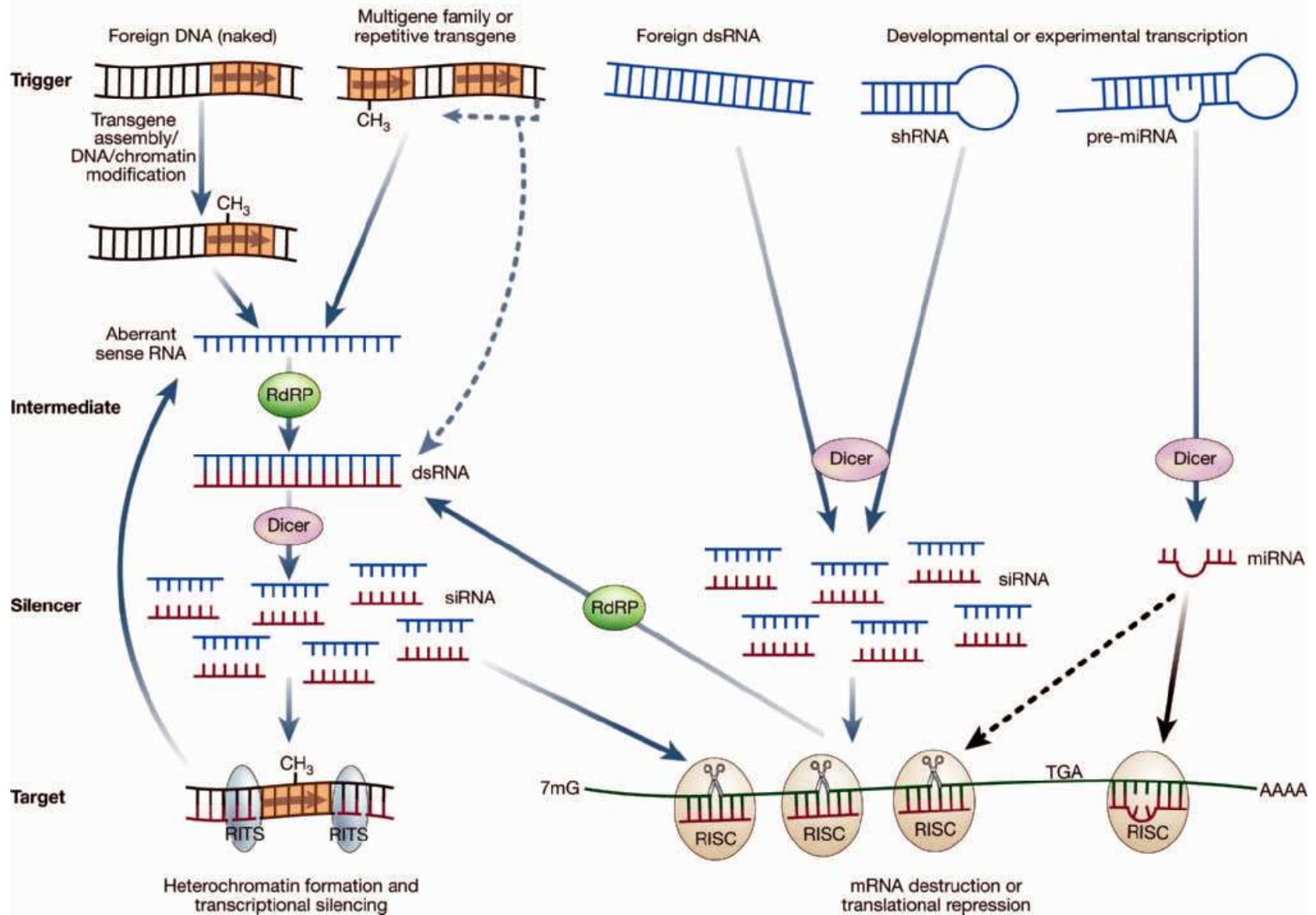
Cleavage of specific mRNA triggered by presence of complementary dsRNA

- first recognized in *Caenorhabditis Elegans*, evolutionary conserved
- antiviral mechanisms - protecting from RNA viruses
- prevention of random integration of transposable elements

Final product involved in recognition and cleavage is short (19-23 bp) dsRNA – short interfering RNA (**siRNA**)

Other related short RNAs:

- repeat associated short interfering RNAs (**ra-siRNAs**) – chromatin changes
  - microRNAs (**miRNAs**) - natural hairpin dsRNAs, regulation of protein translation
-



---

# Mechanism of recognition and cleavage

**RISC** = RNA-induced silencing complex  
Argonaute protein family (Ago proteins)

- components of RISCs
- tightly bind ssRNAs
- contain PAZ and PIWI domains

PIWI domain – binds Dicer

RNA - similar to RNase H family domain (cleaves the strand of RNA/DNA duplexes)

PAZ domain – recognizes the 2-nucleotide 3' overhang of siRNA

RNA- - prevents processing of unrelated RNA or turnover products

---

---

# Features important for siRNA function

dsRNA should form  $\alpha$ -helix

duplex-end stability - whichever strand is most easily unwound in a 5'-3' direction will be preferentially assembled within RISC

GC contents between 30-60%

Avoid GGGG, CCCC, UUUU, AAAA

particular nucleotides at specific positions

mRNA conformation

---

# Differences between the siRNA design algorithms

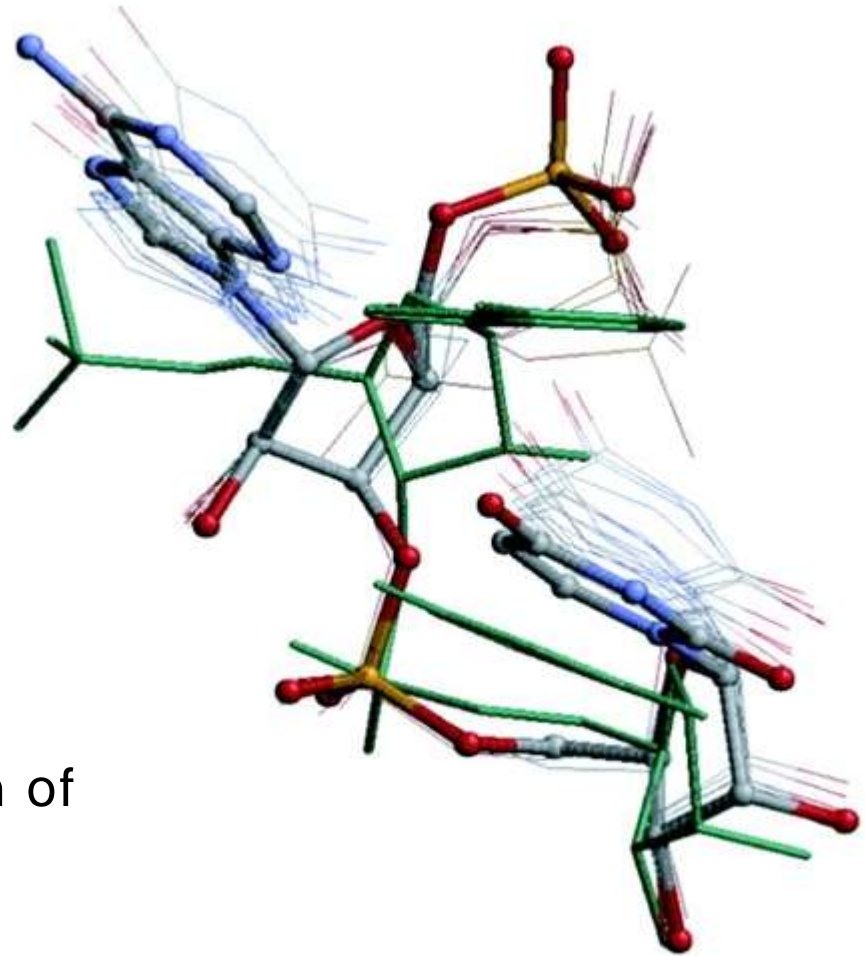
Algorithm	Description
GPboost motifs/patterns	Weighted sum of sequence
Ui-Tei	Sequence features
Amarzguioui	Sequence features
Hsieh	Sequence features
Takasaki	Sequence features
Reynolds 1	Hairpin potential, sequence features
Reynolds 2	Sequence features
Schwarz	Difference between 3' and 5' stability
Khvorova	Duplex stability profile
Stockholm 1	Energy features
Stockholm 2	Energy features
Tree	Sequence features in decision tree
Luo	mRNA secondary structure features

# Sequence characteristics used by different algorithms

Algorithm	siRNA sense strand position																																						
	1	2			3	6			7	8		9	10		11	13		15	16		17		18		19														
	A	C	G	U	A	U	A	U	A	C	G	U	A	G	G	U	U	C	G	A	G	A	U	A	G	U	A	U	A	U	A	C	G	U					
Reynolds 1 and 2					1										1					-1	1	1	1	1	1	1	1	1	1	1	2	-1	-1	1					
Ui-Tei	-1	1	1	-1																																			
Amarzguioui	-1	1	1	-2	-1	-1	-1	-1	1																														
Hsieh										-1																													
Takasaki	-3.97		7.4	-3.75																																			



# Our approach



Find any rules for prediction of functional siRNAs based on their 2/3 D structure

# Available data sets

Author	Amarzguioui	Harborth	Hsieh	Khvorova	Reynolds	Ui-
Tei Vickers						
Count	46	44	108	14	240	53

---

76

## Data from Reynolds

90 siRNAs each targeting every other position of the regions

Human cyclophilin B (M60857) 193-390 [90]

efficacy  $\geq 80$  [9]

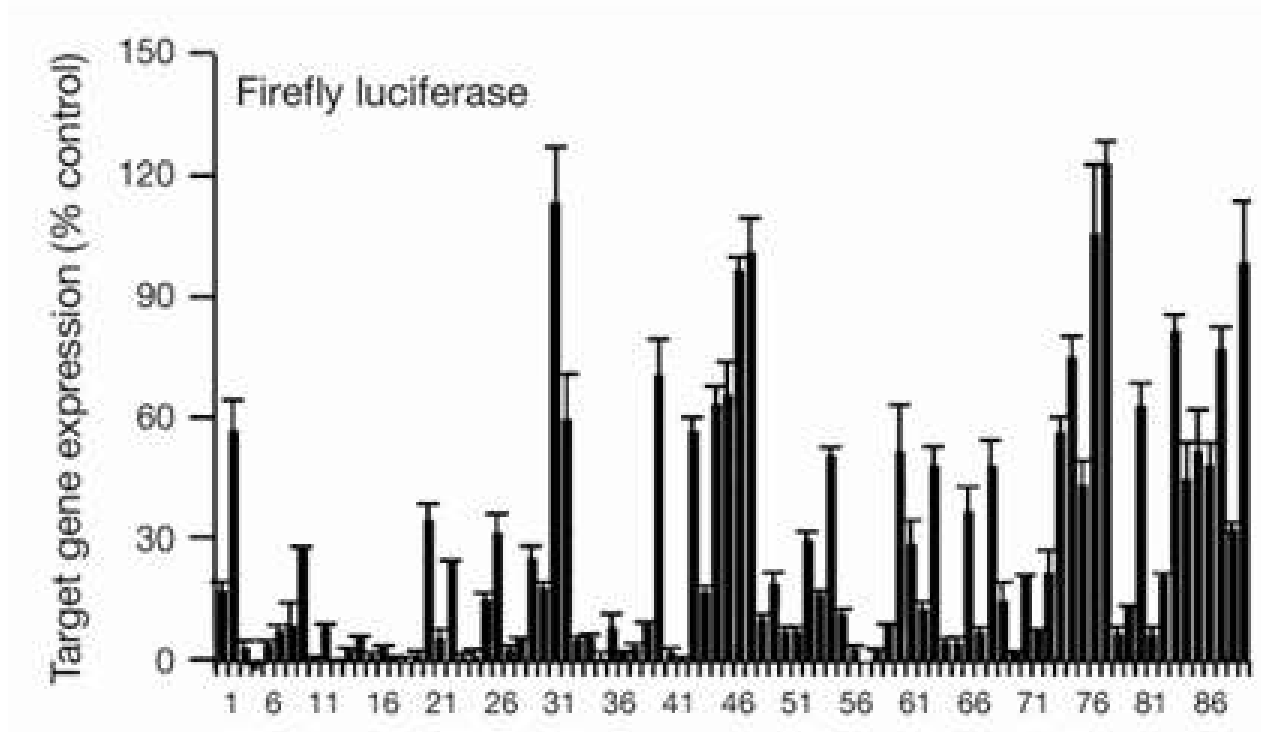
efficacy  $\leq 20$  [41]

Firefly luciferase (U47298) 1434-1631 [90]

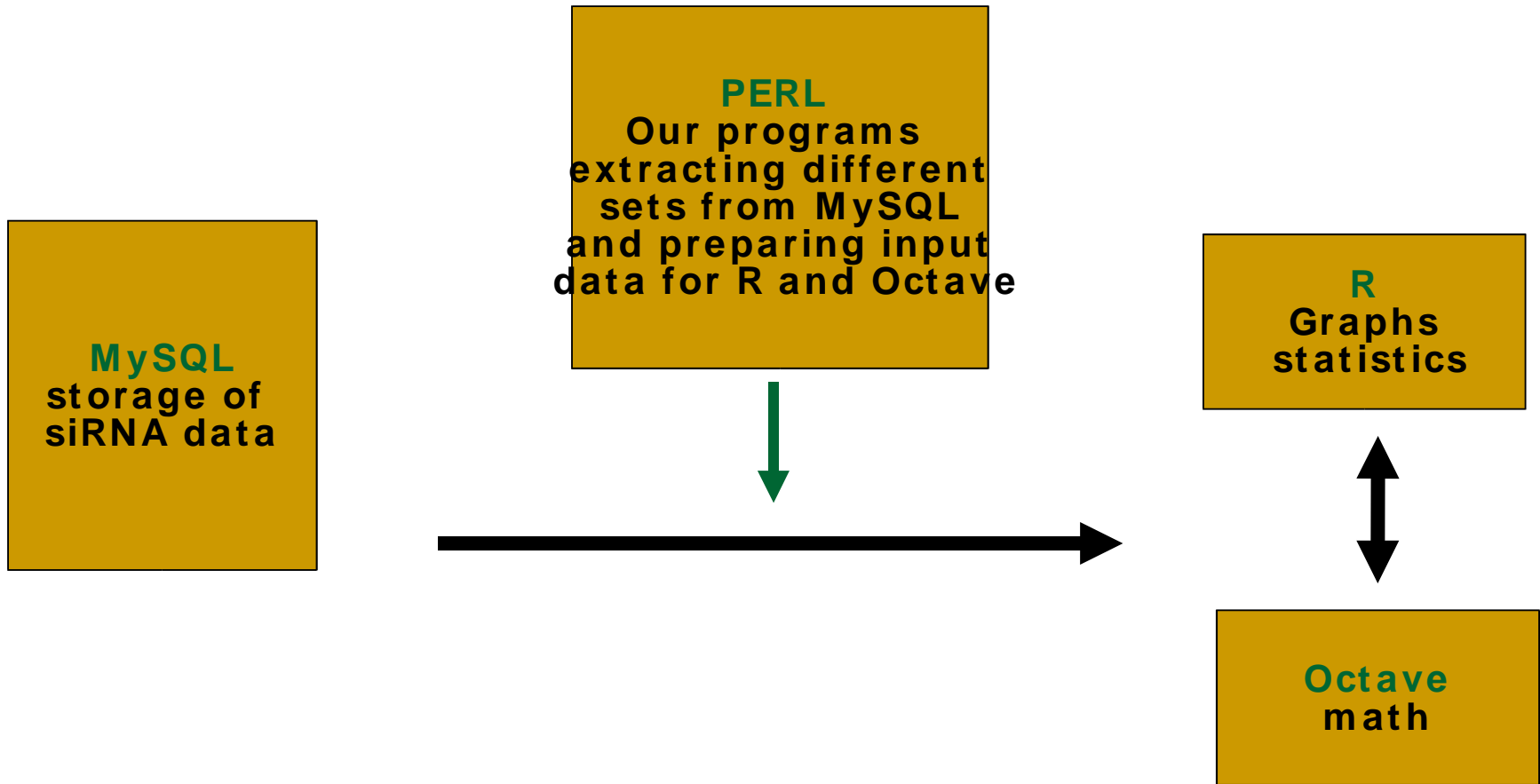
efficacy  $\geq 80$  [8]

efficacy  $\leq 20$  [49]

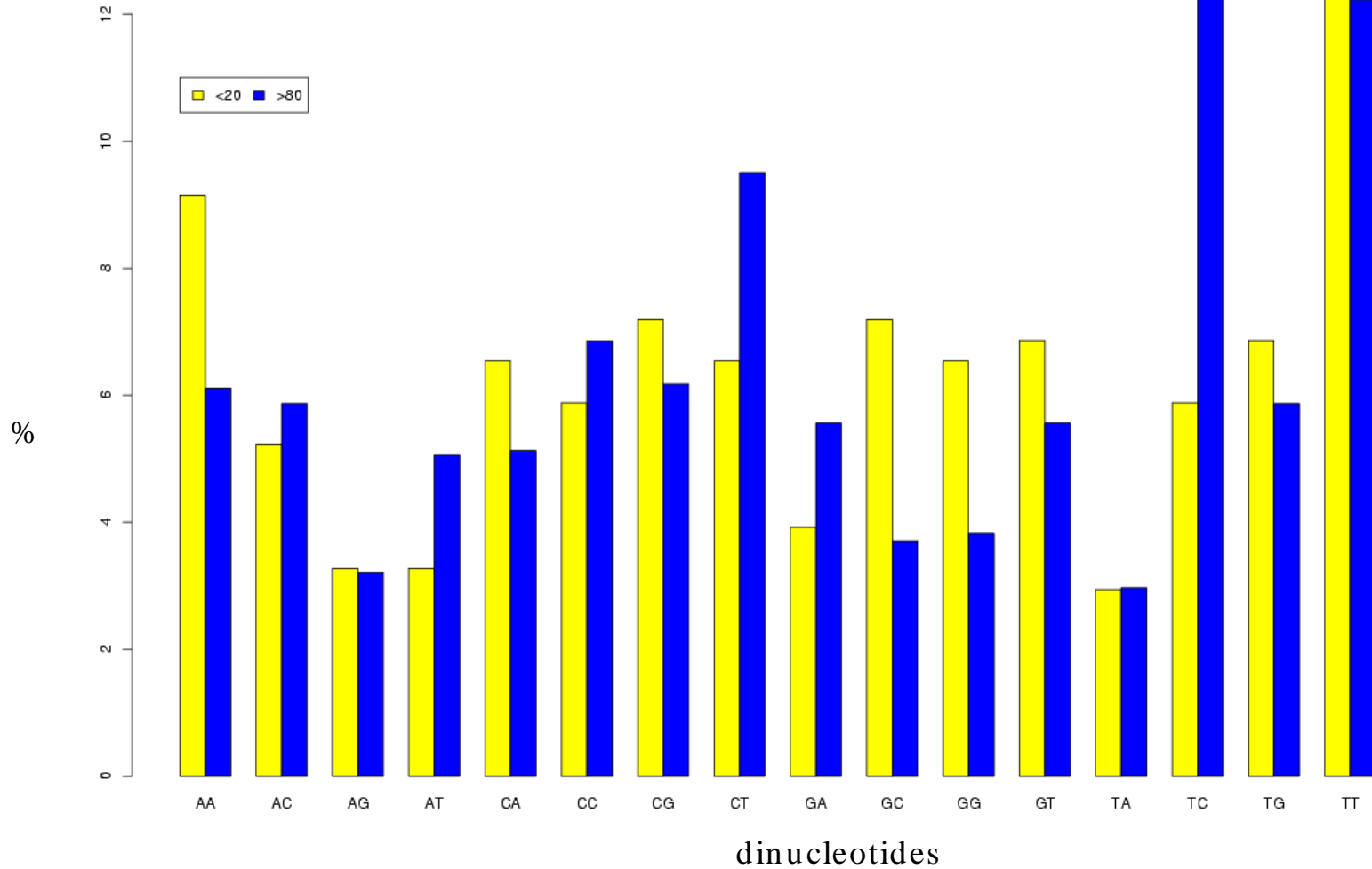
# Reynolds cyclophylin B



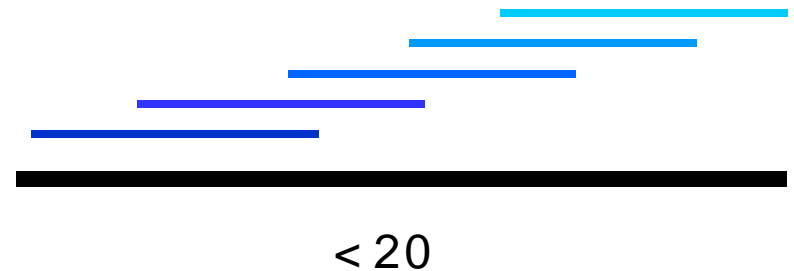
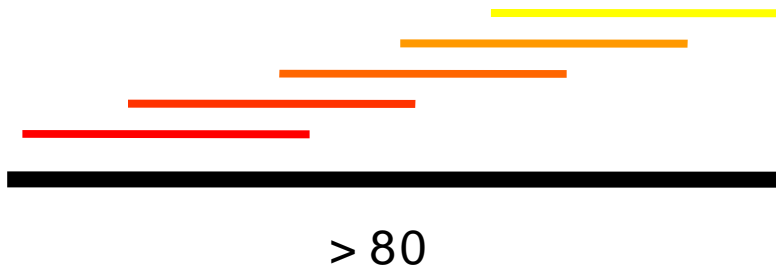
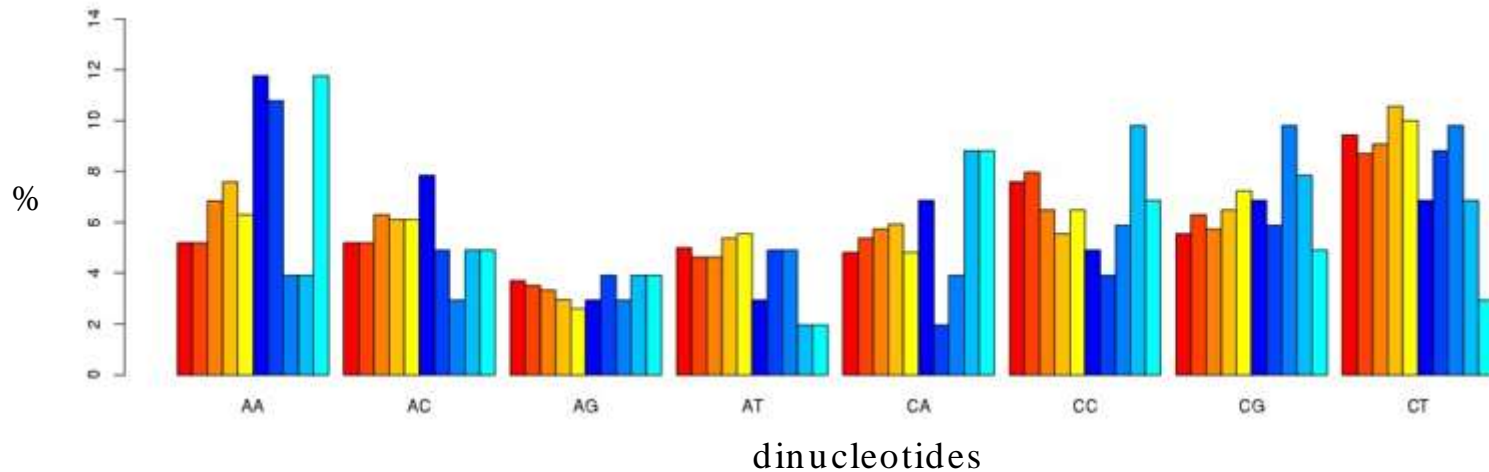
# Bioinformatic workflow



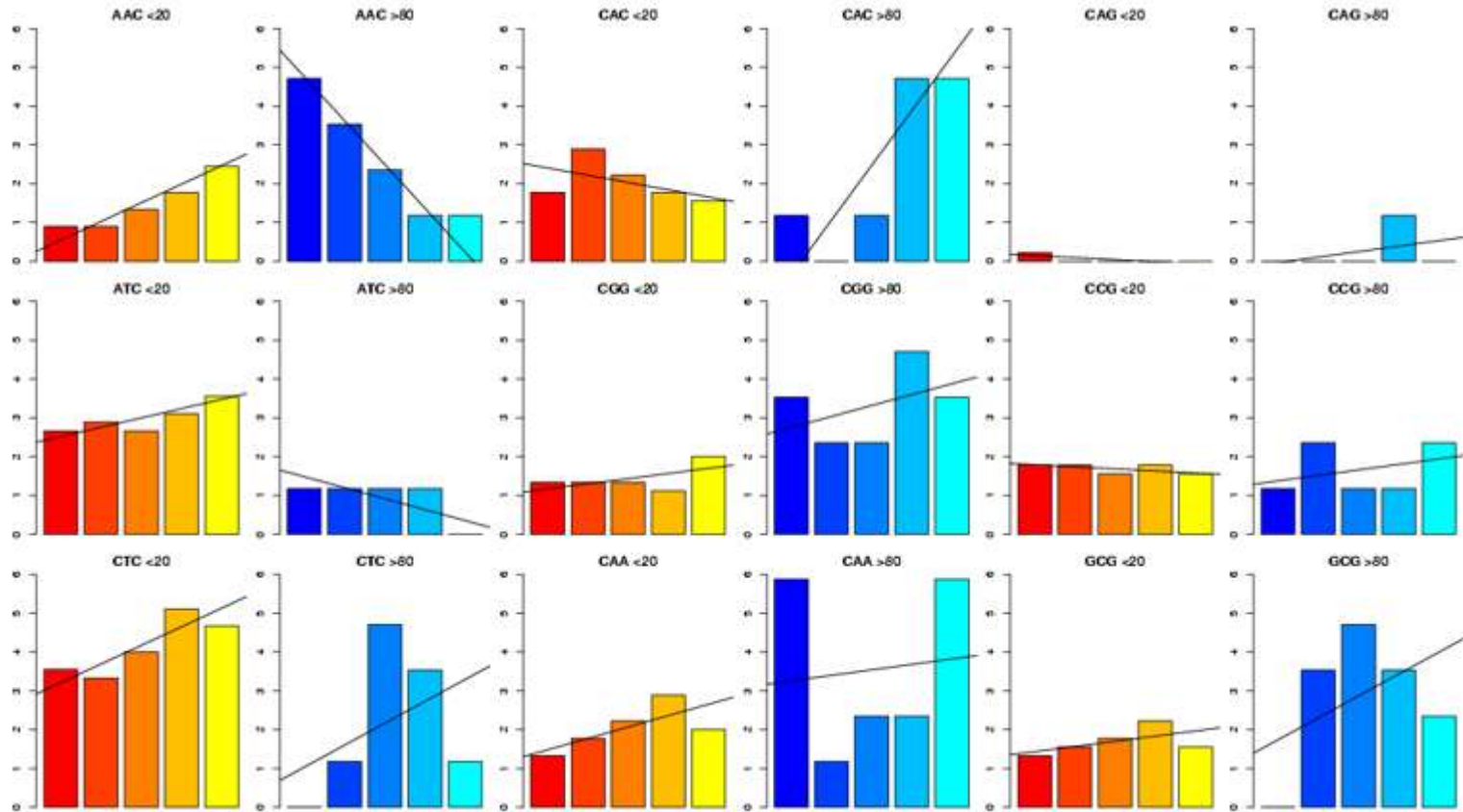
# Number of dinucleotides in both Reynolds sets



# Number of dinucleotides in both Reynolds sets with 3 nucleotide shift



# Number of trinucleotides in both Reynolds sets with 3 nucleotide shift



# Principle of prediction

data sets :

Several different sets of chosen di or trinucleotides

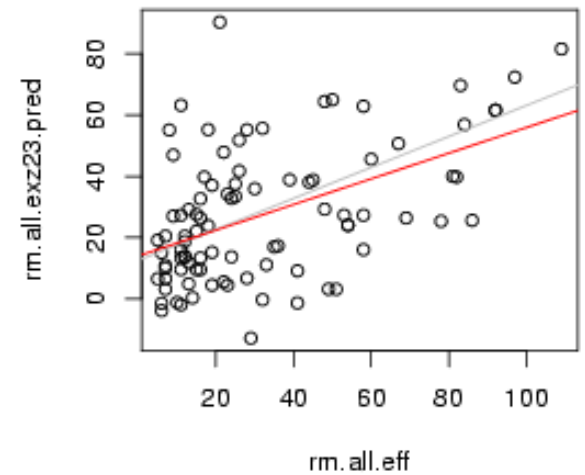
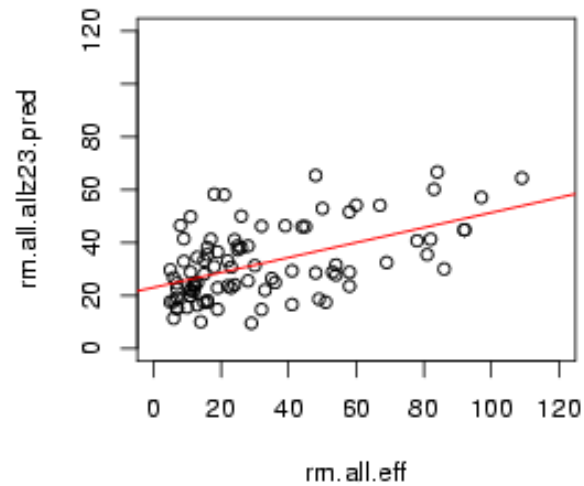
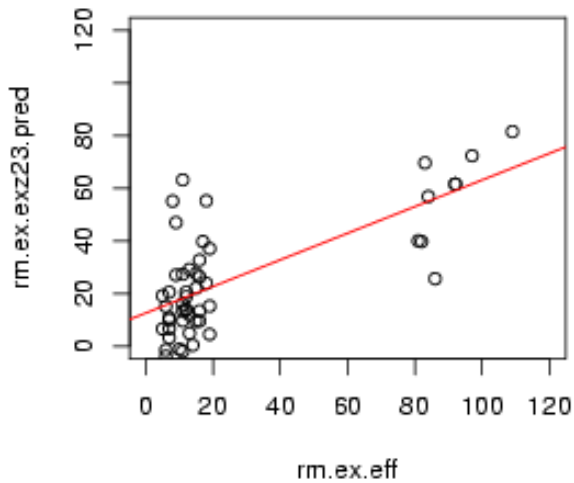
Set of functions that allow obtaining of parameters which show importance of each di/trinucleotide for resulting efficacy of concrete siRNAs

- eg. for sequence ATGTGGCATC with efficacy 82 for ATG [1..10], GGG [1..10], ATC [1..5]:
- $1 \cdot \text{ATG} + 0 \cdot \text{GGG} + 0 \cdot \text{ATC} = 82$

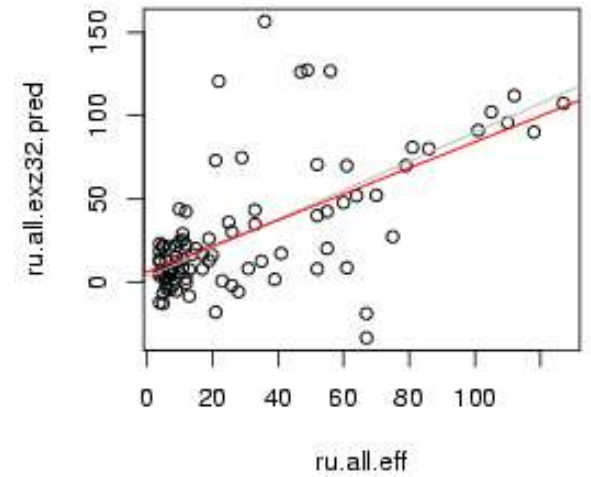
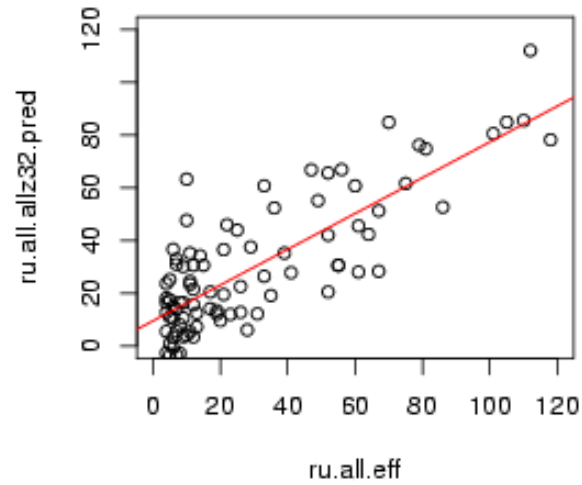
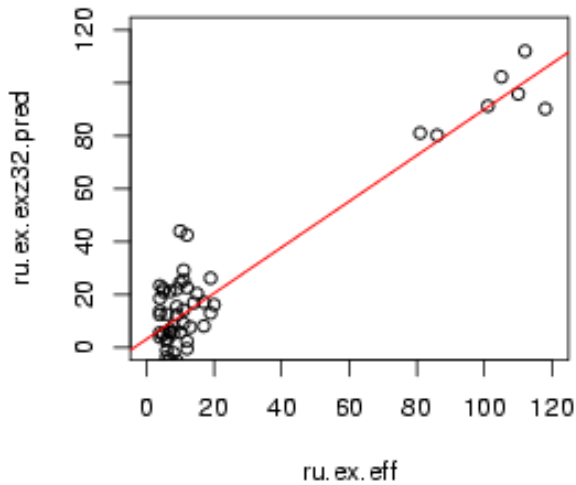
...system of linear equations of all siRNAs with efficacy  $< 20$  and  $> 80$  from one half of Reynolds data



# Prediction for a set of dinucleotides



# Prediction for a set of trinucleotides



---

# mechanism

dsRNA

- produced by RNA-templated RNA polymerization (viruses)  
siRNA

- hybridization of overlapping transcripts (repetitive sequences - transgene arrays or transposons) rasiRNA

-endogenous transcripts with complementary 20-50 bp inverted repeats miRNA

dsRNA processed by dsRNA specific RNase-III-type endonucleases - Drosha and Dicer

---