

INTO THE DEEP

microRNA Detection using
Next-Generation Sequencing Data

David Langenberger
Bioinformatics Group
Department of Computer Science
University of Leipzig

UNIVERSITÄT LEIPZIG

Next-Generation Sequencing

Massively Parallel Sequencing technologies



**454 Life Sciences / Roche
(FLX Titanium Series)**

1 million reads
400 bp



**Solexa / Illumina
(Genome Analyzer)**

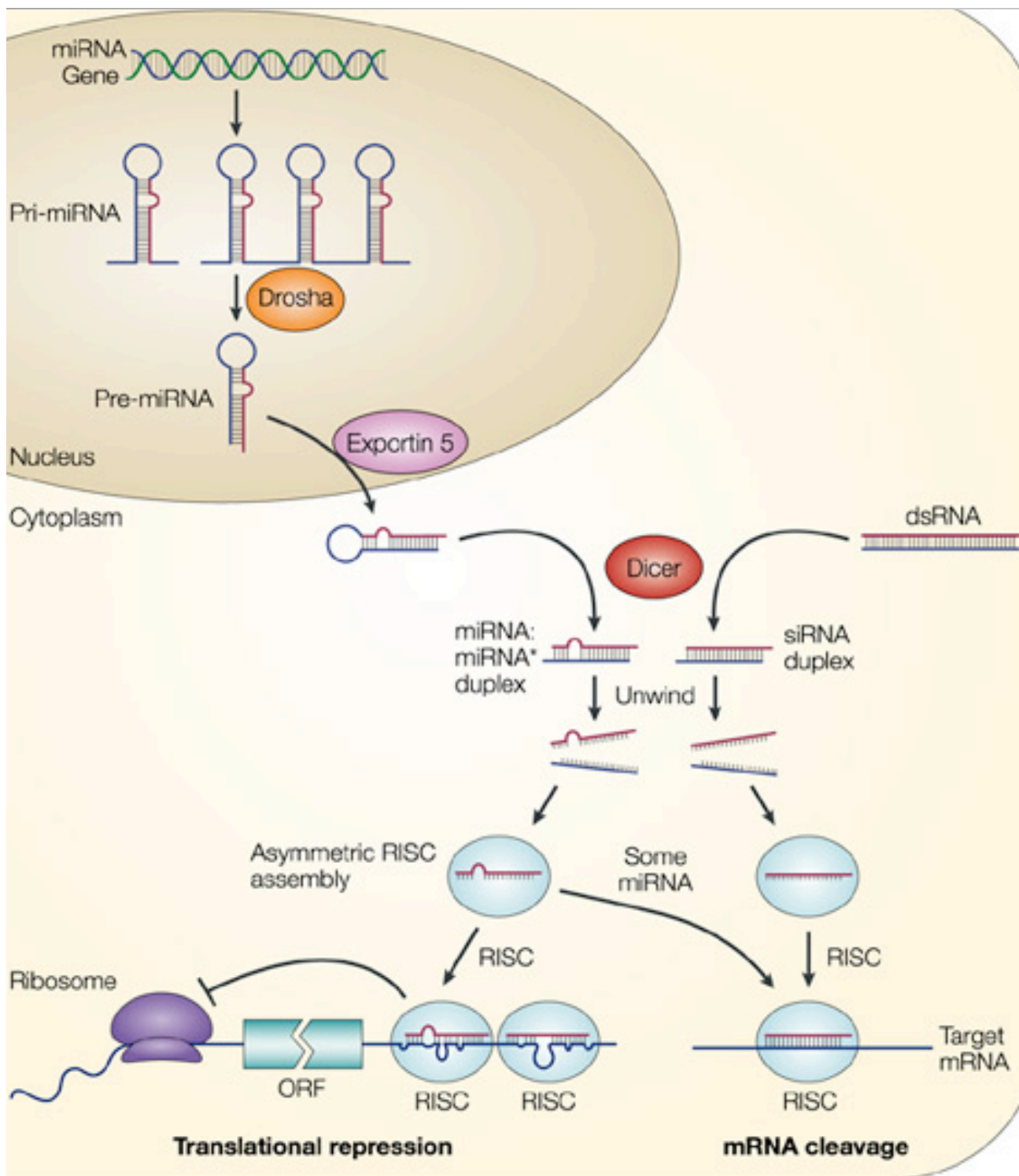
40 million reads
35 bp



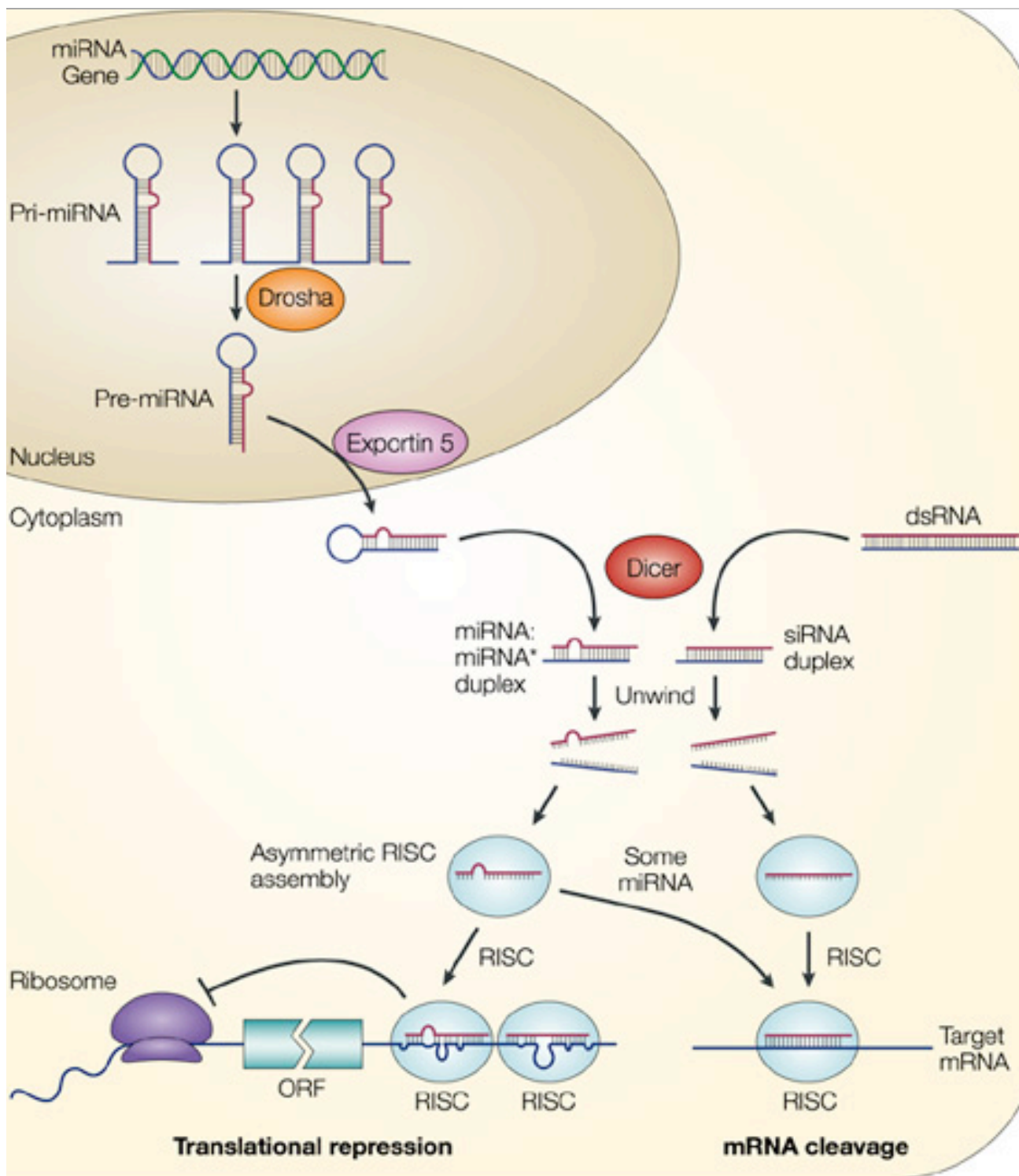
**Applied Biosystems
(SOLiD System)**

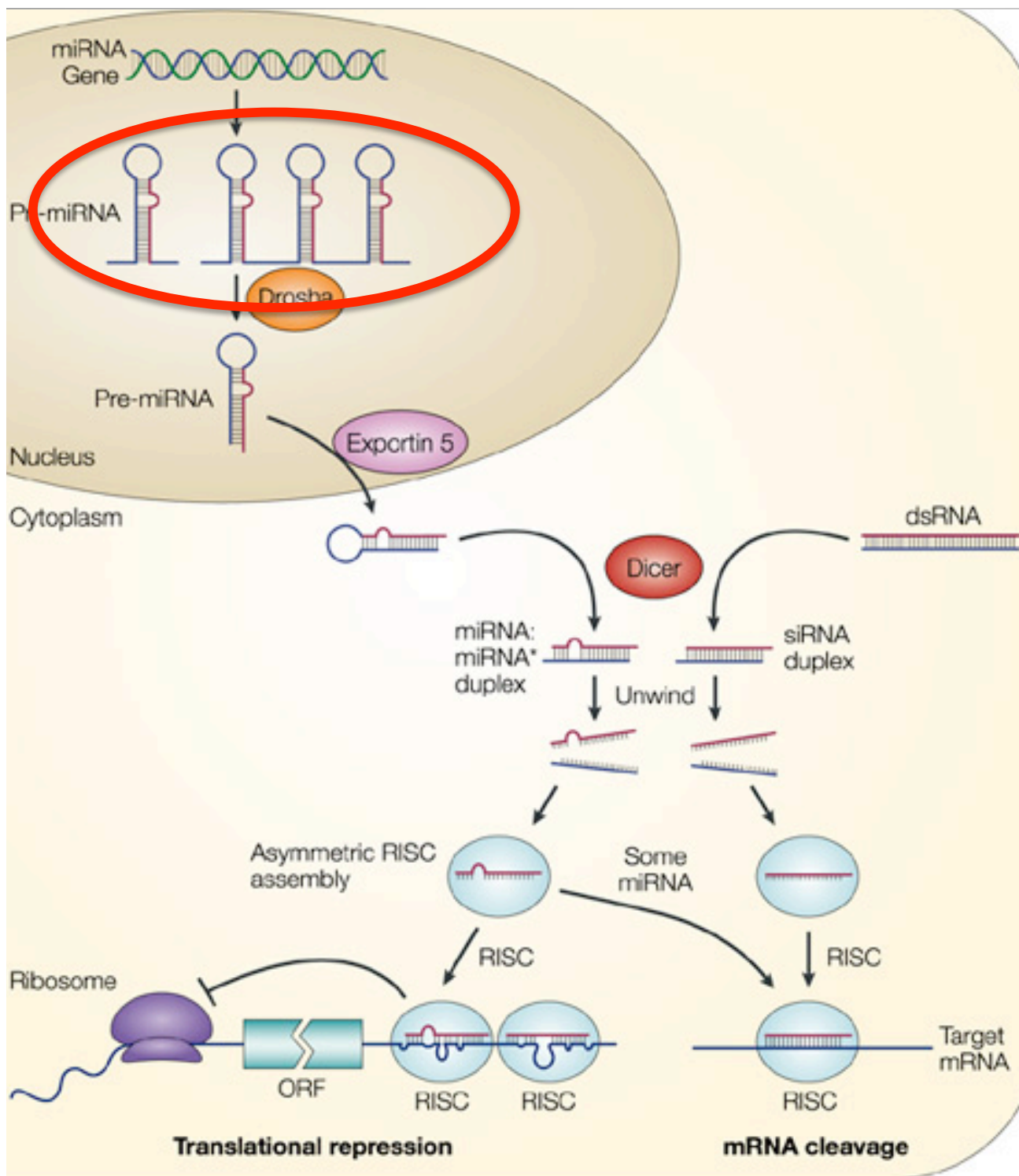
70 million reads
35 bp

microRNAs



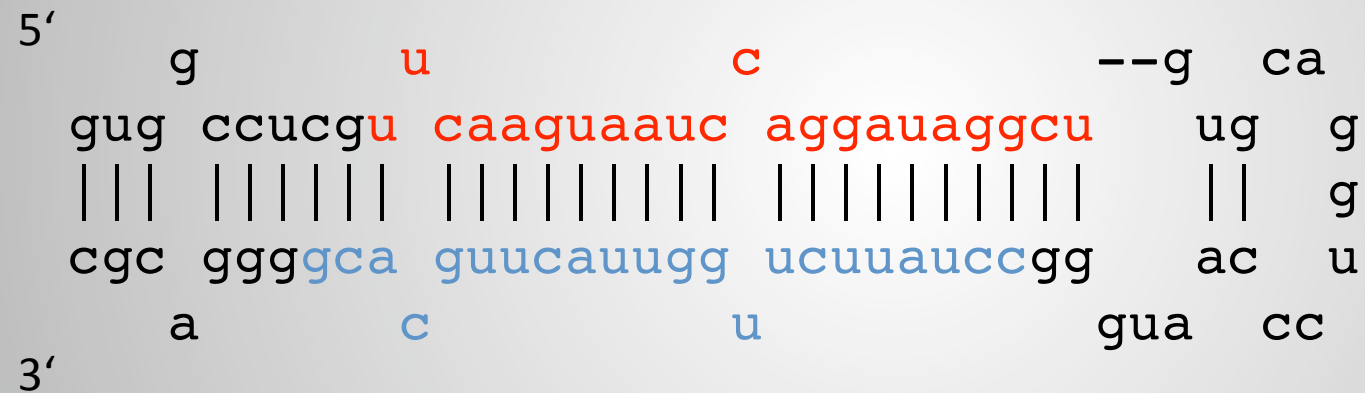
Information content



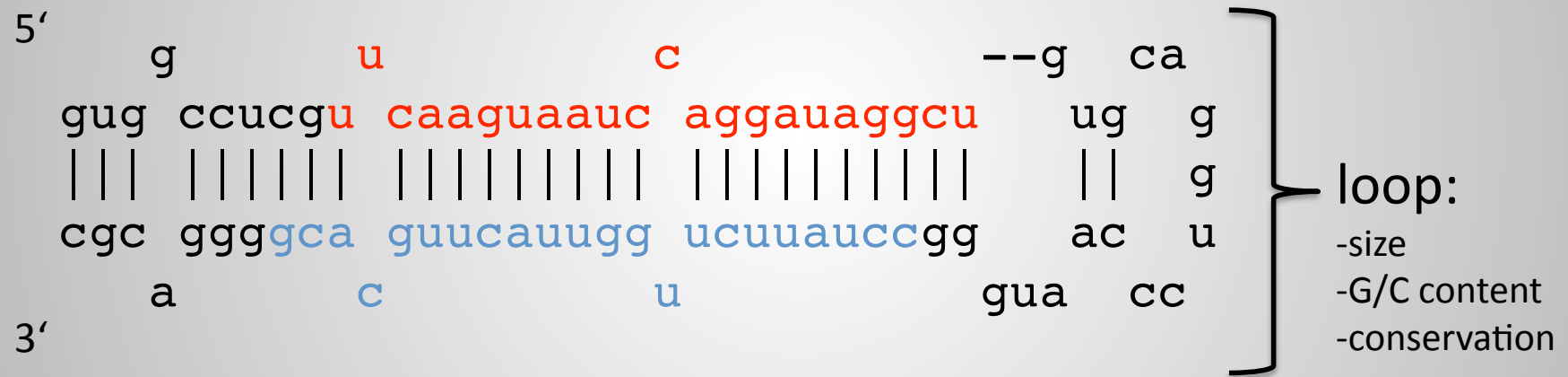


hairpin

miRNA

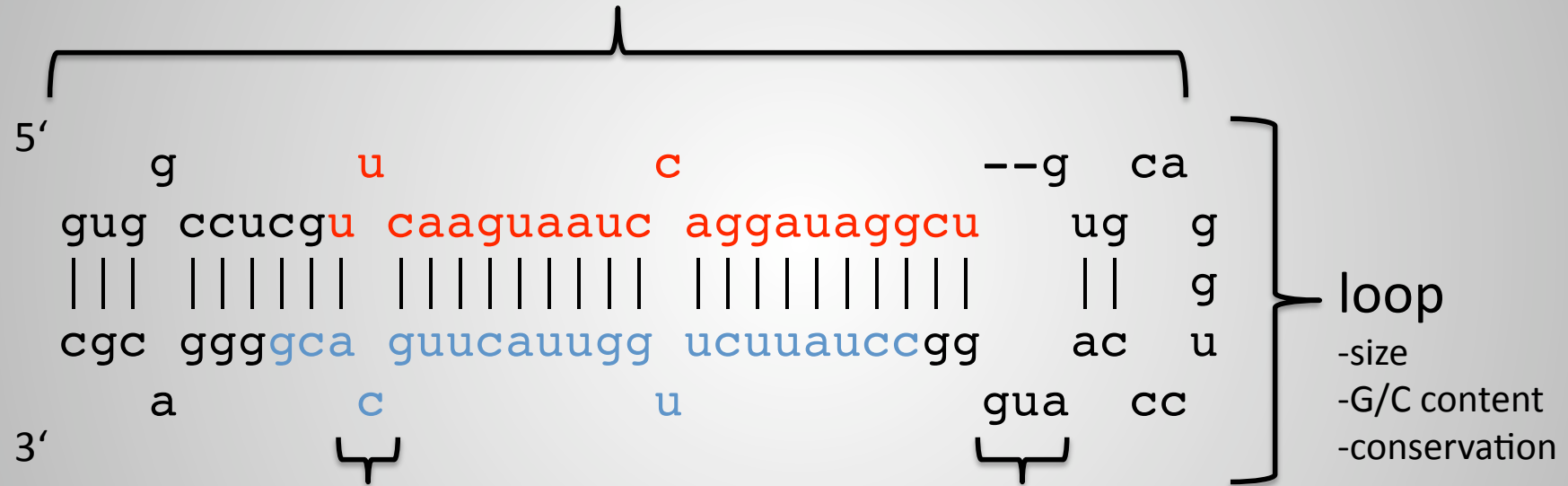


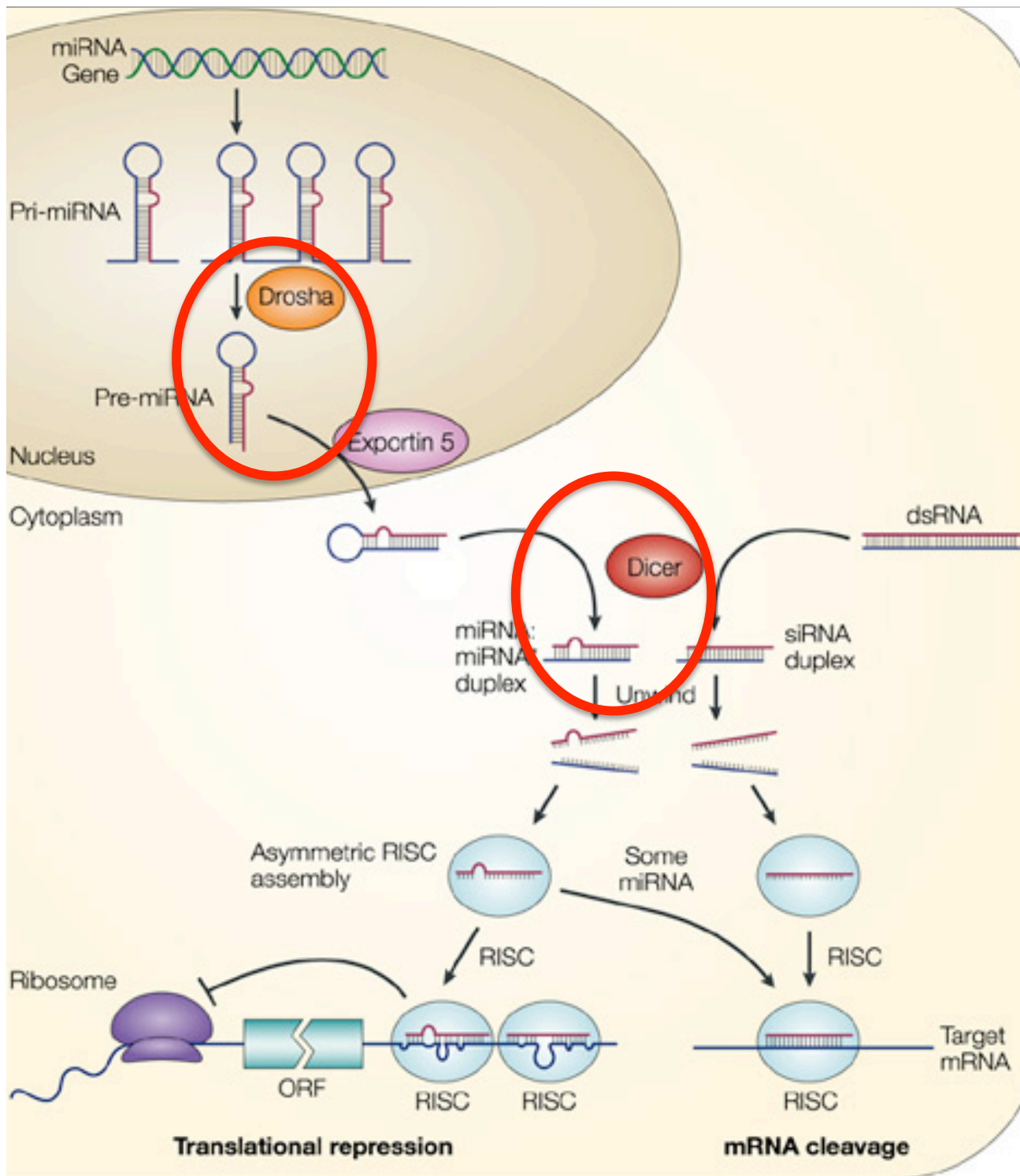
miRNA*



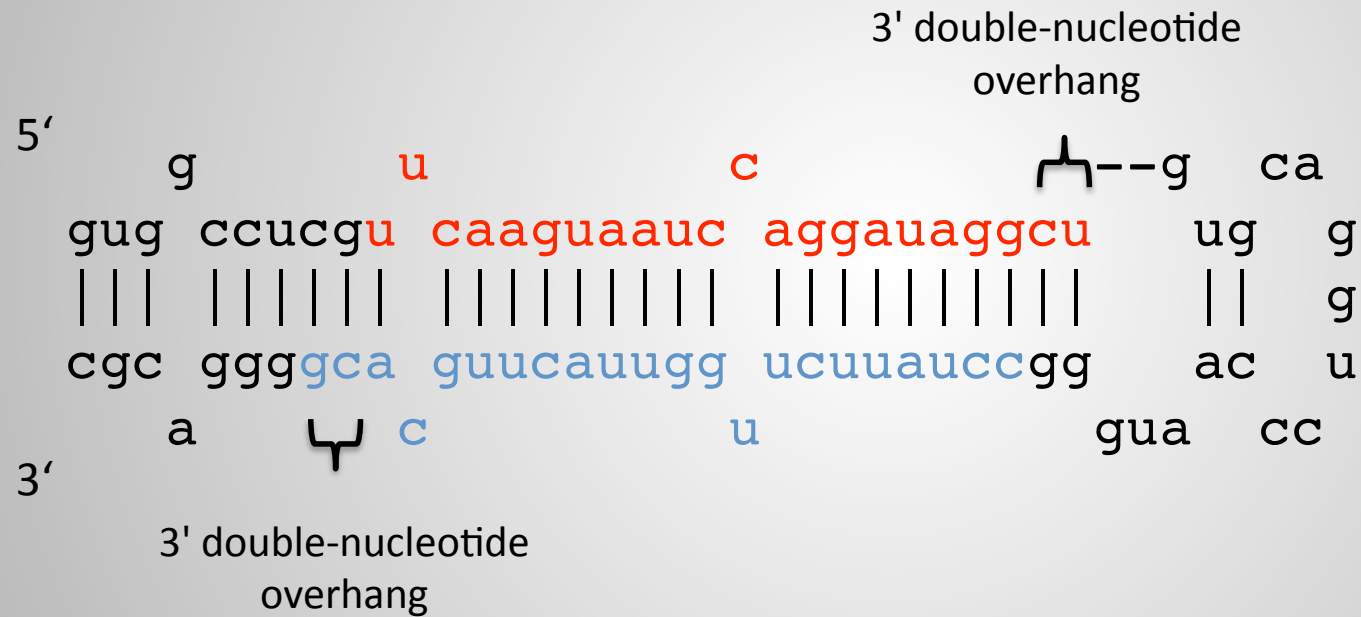
stem:

- size
- G/C content
- bound nucleotides
- conservation

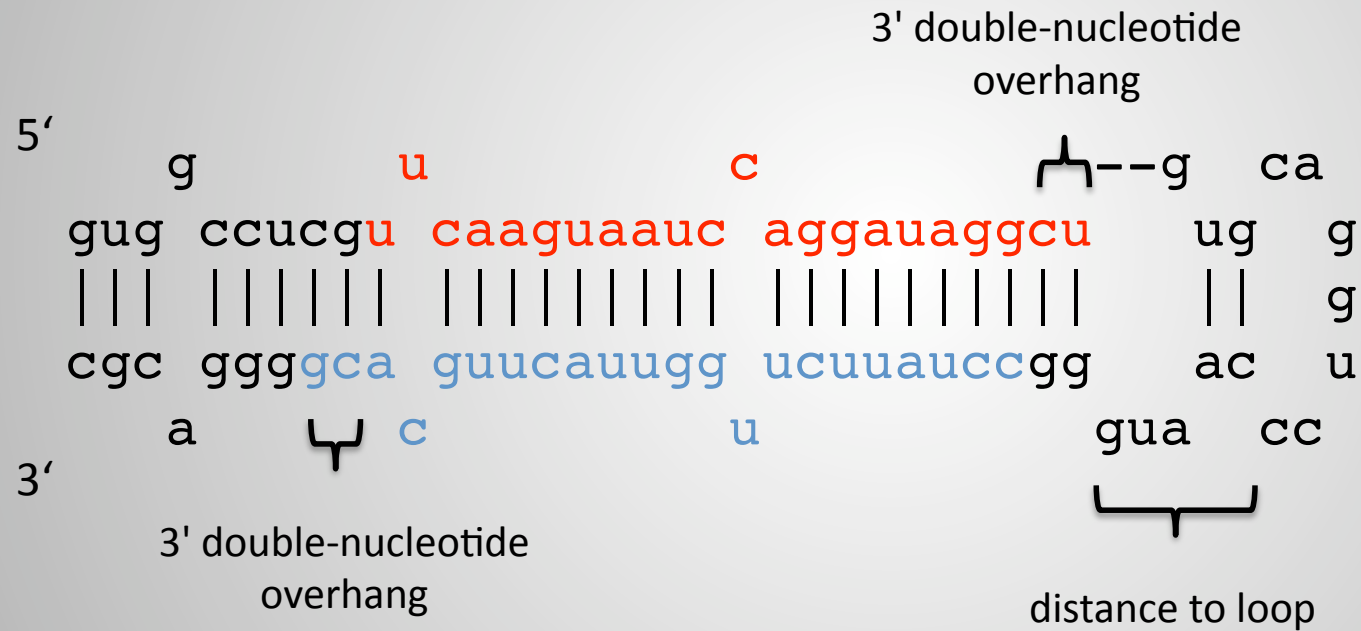


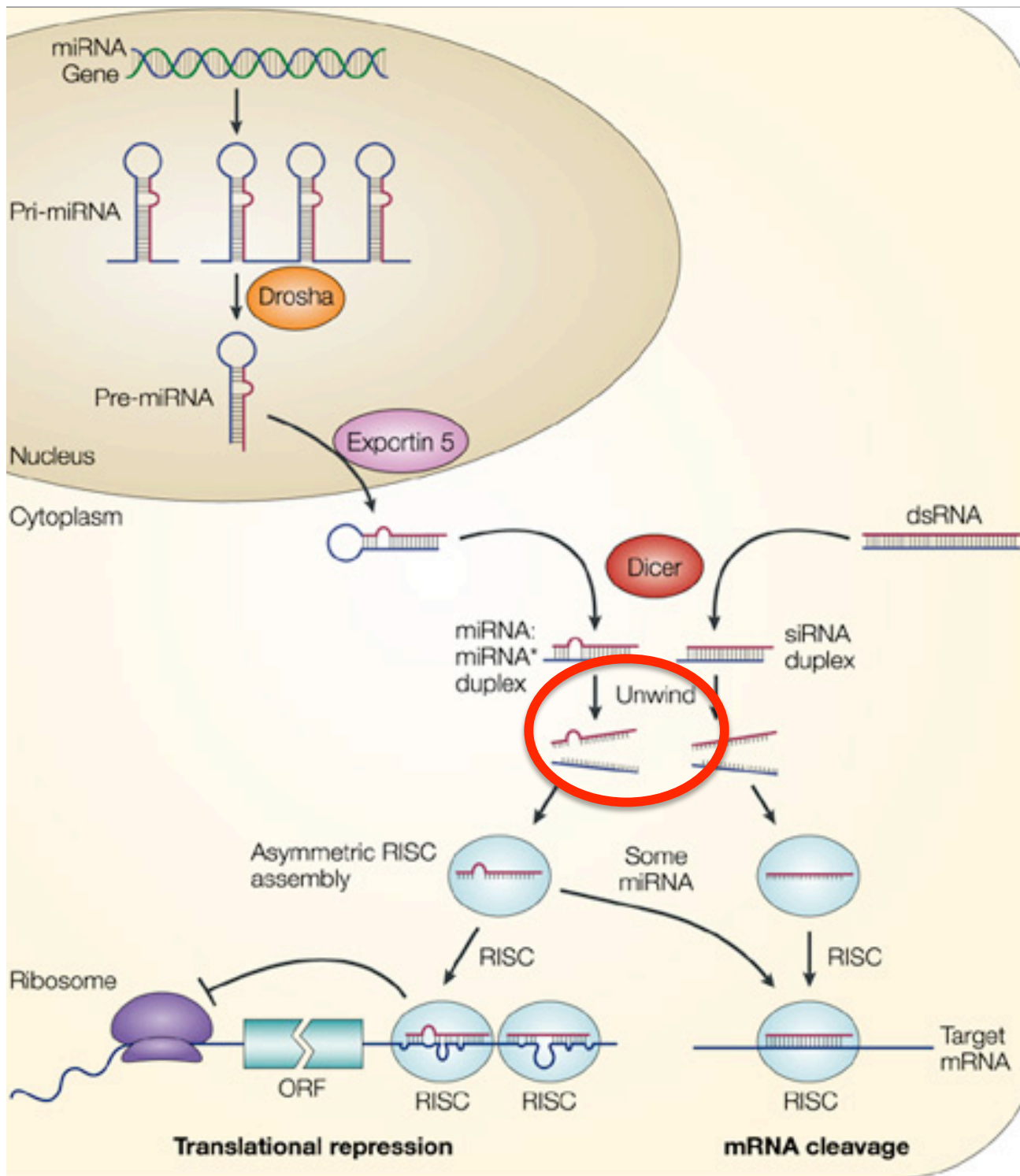


hairpin



hairpin





Data

355,453 reads



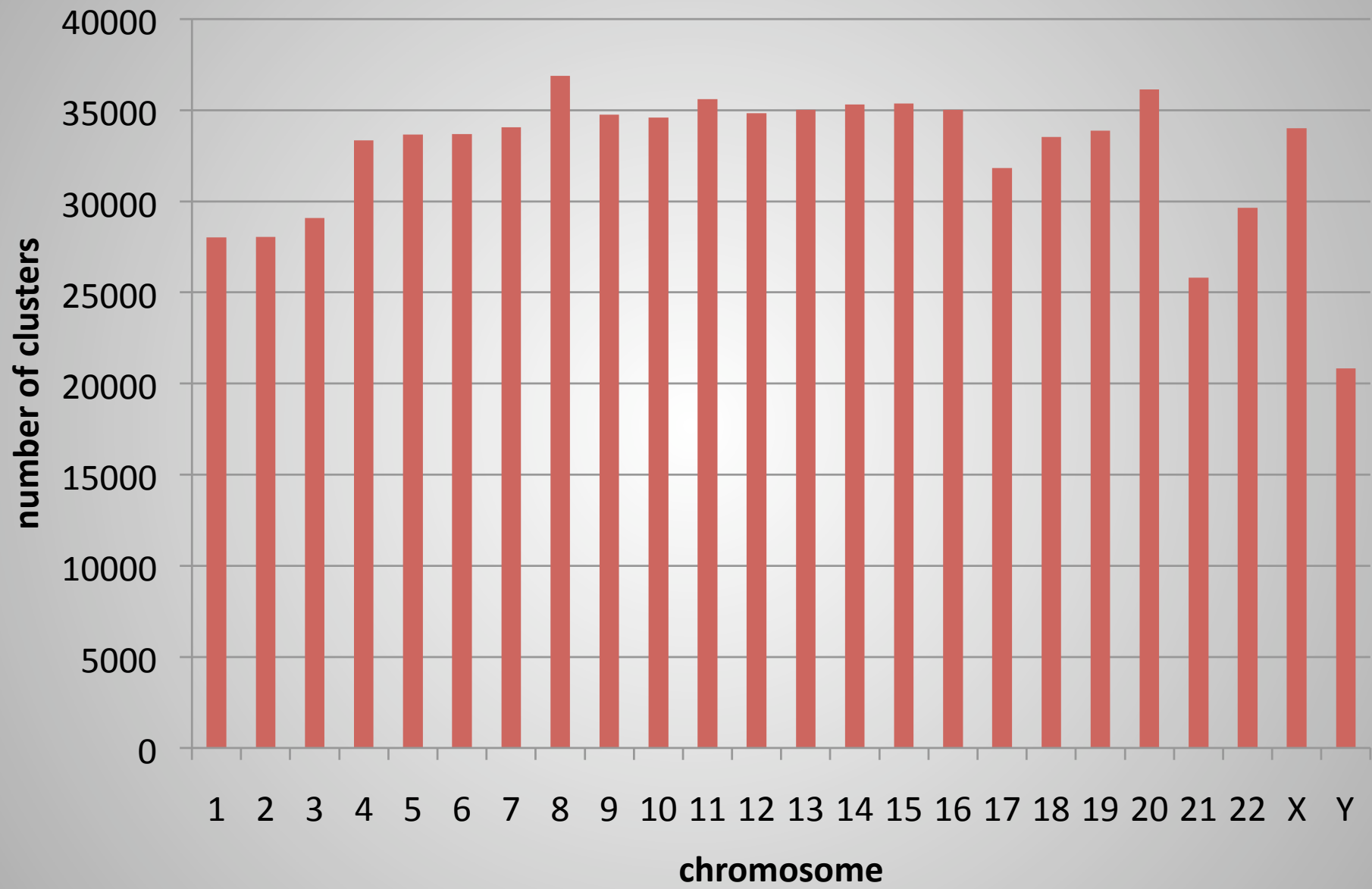
map reads to genome (segemehl)



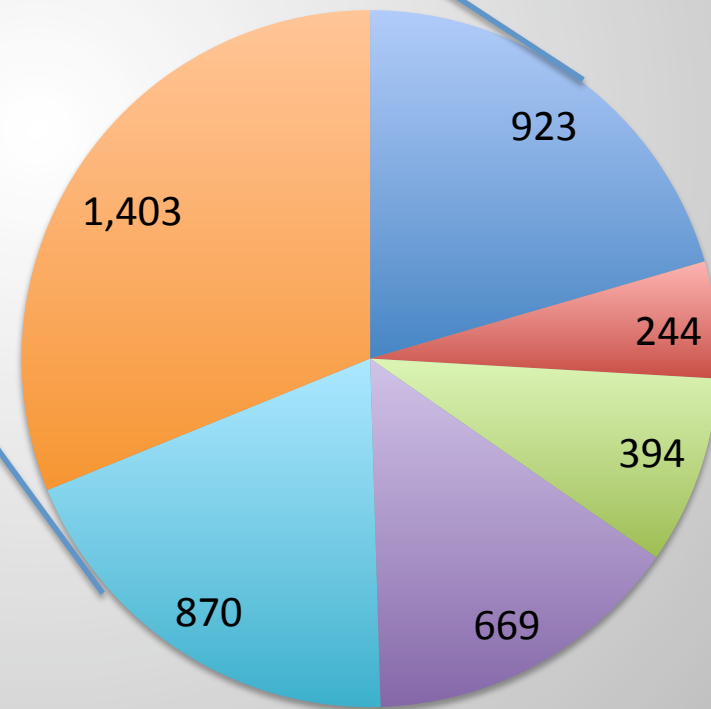
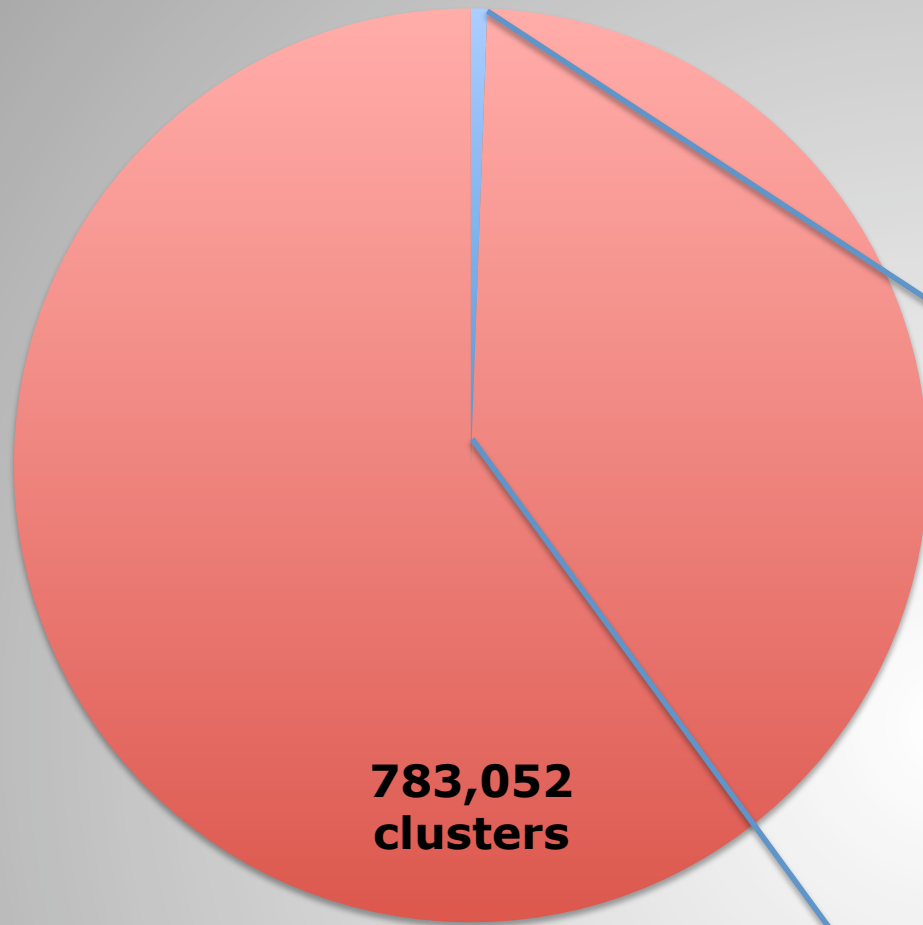
Cluster reads



783,052 clusters

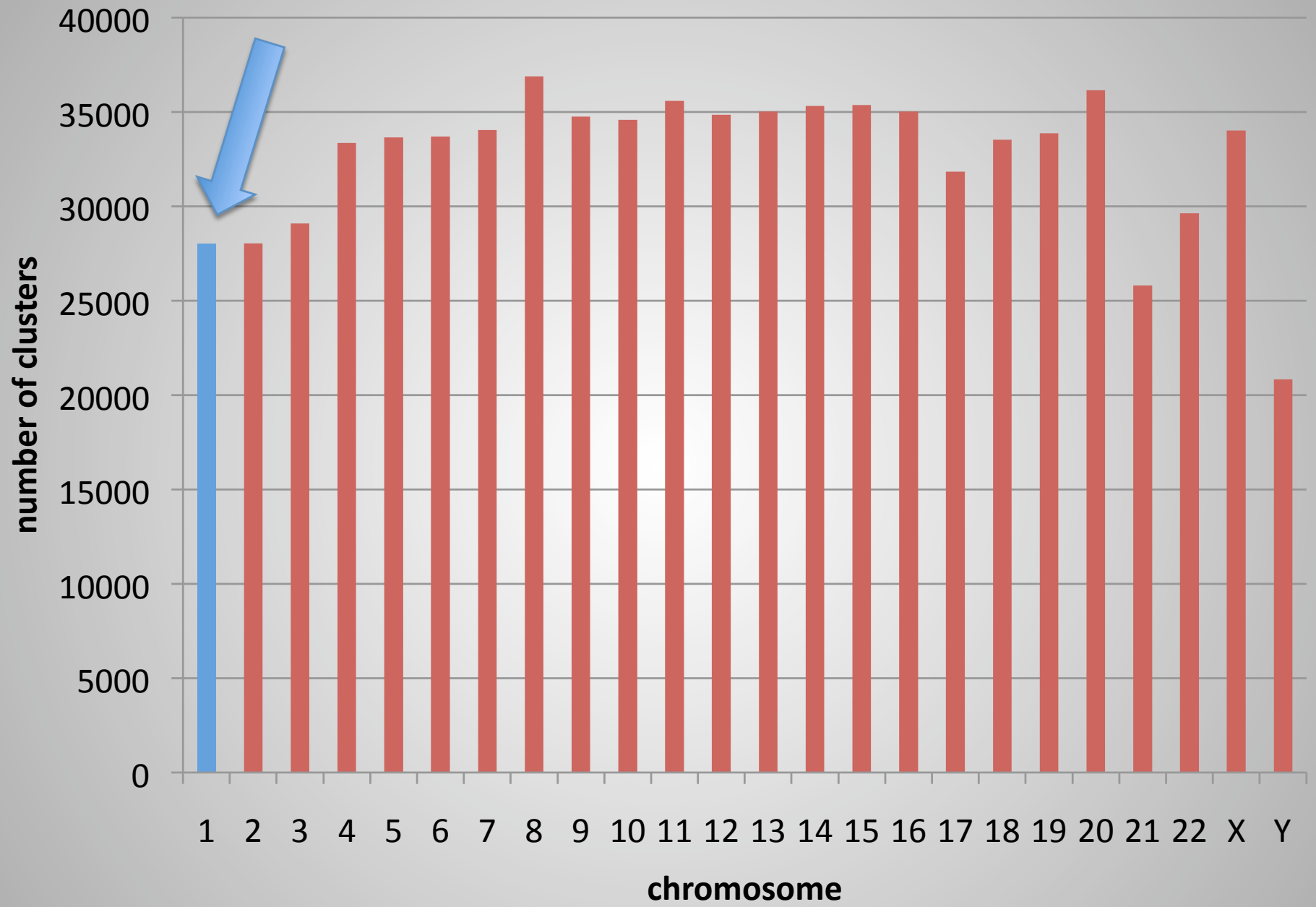


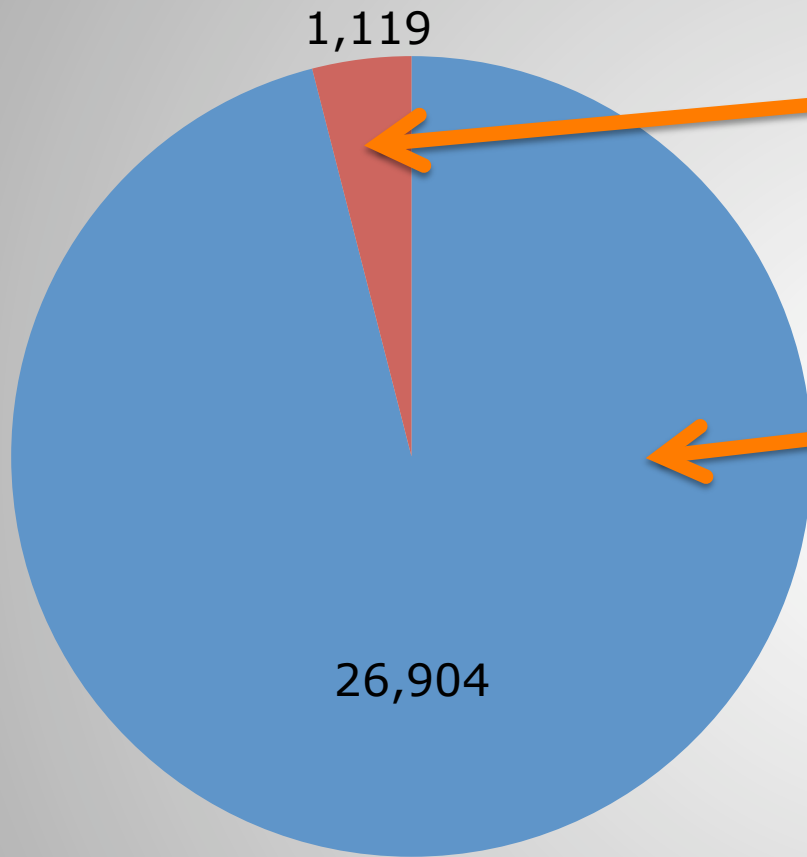
Annotated ncRNAs



4,503
known ncRNAs

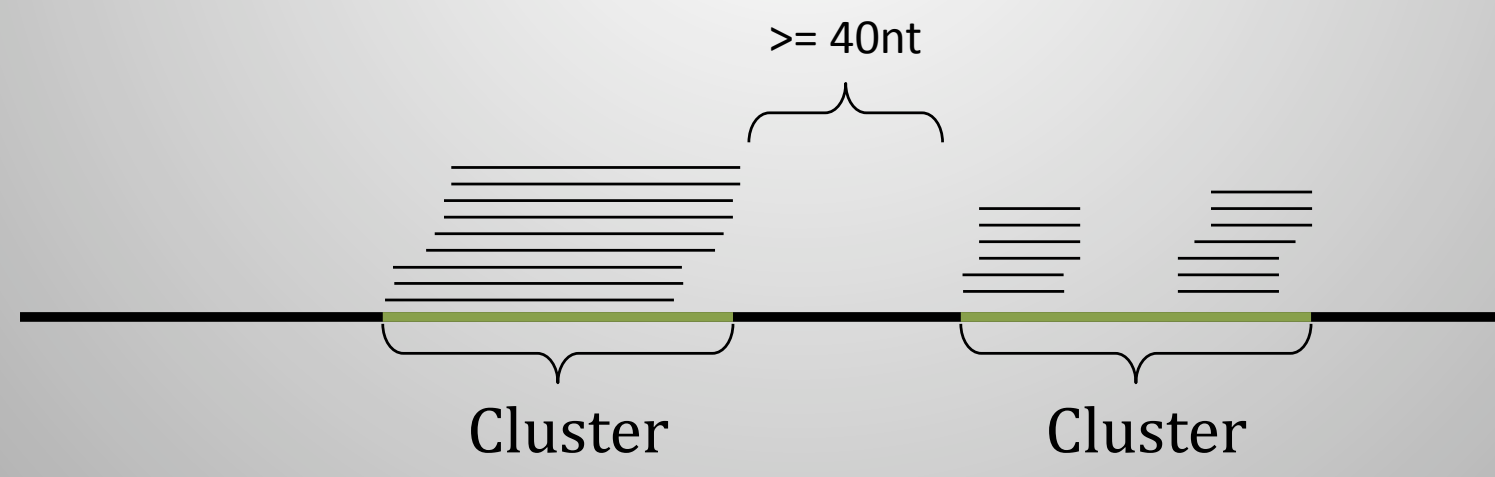
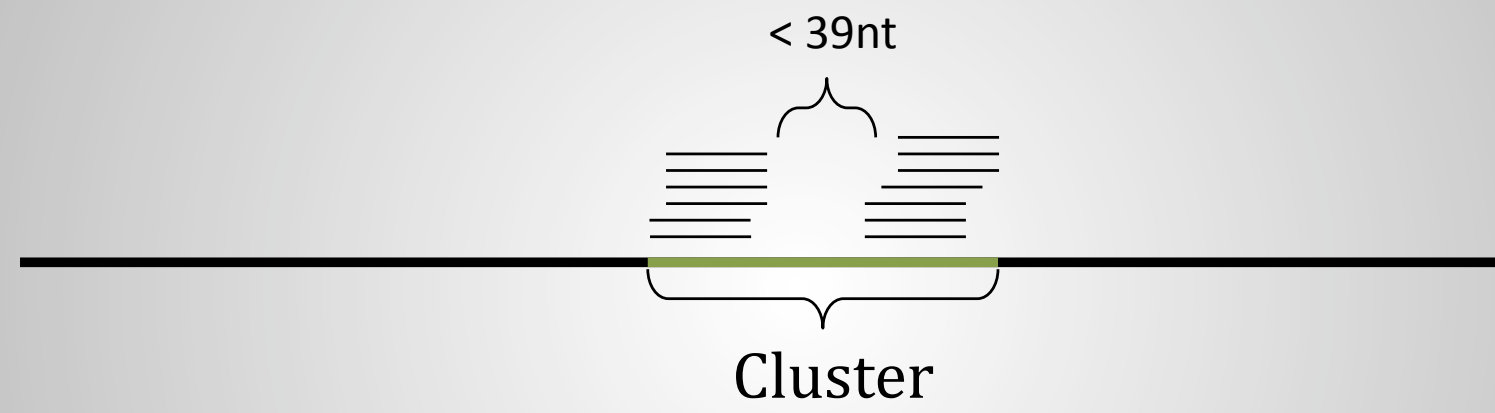
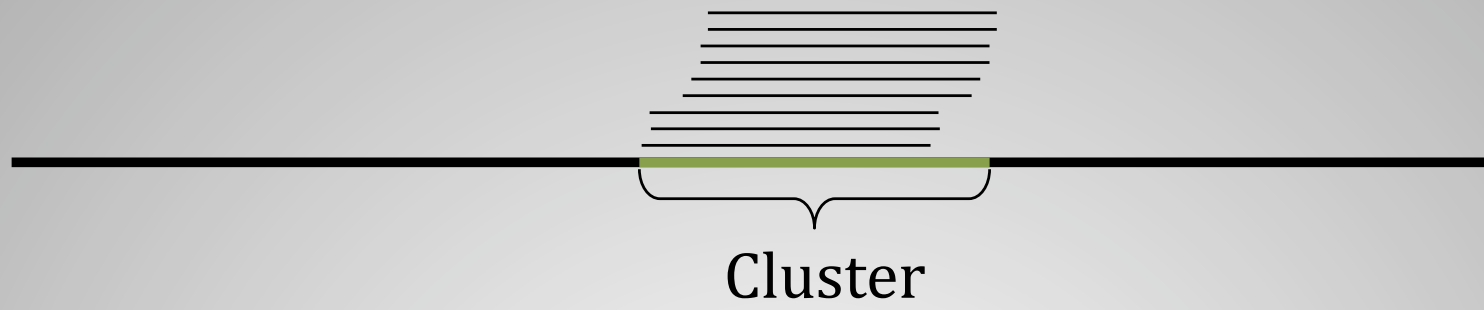
- miRNA
- rRNA
- snoRNA
- snRNA
- misc-RNA
- other





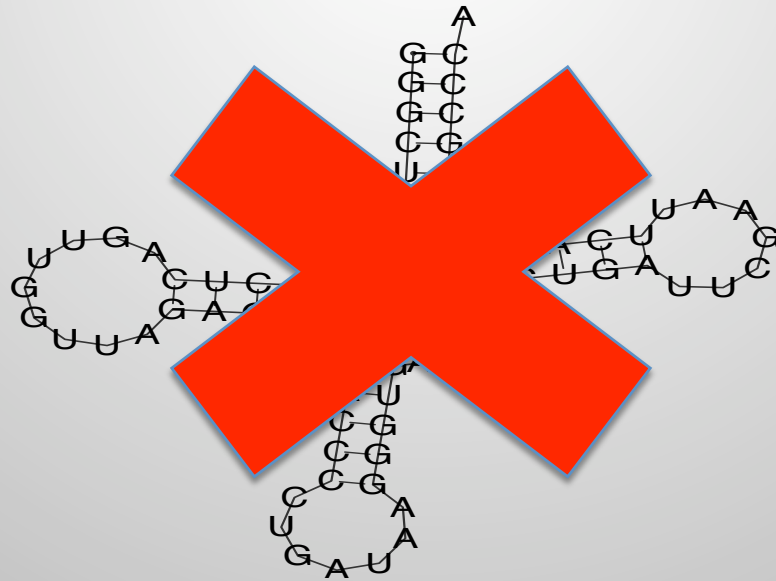
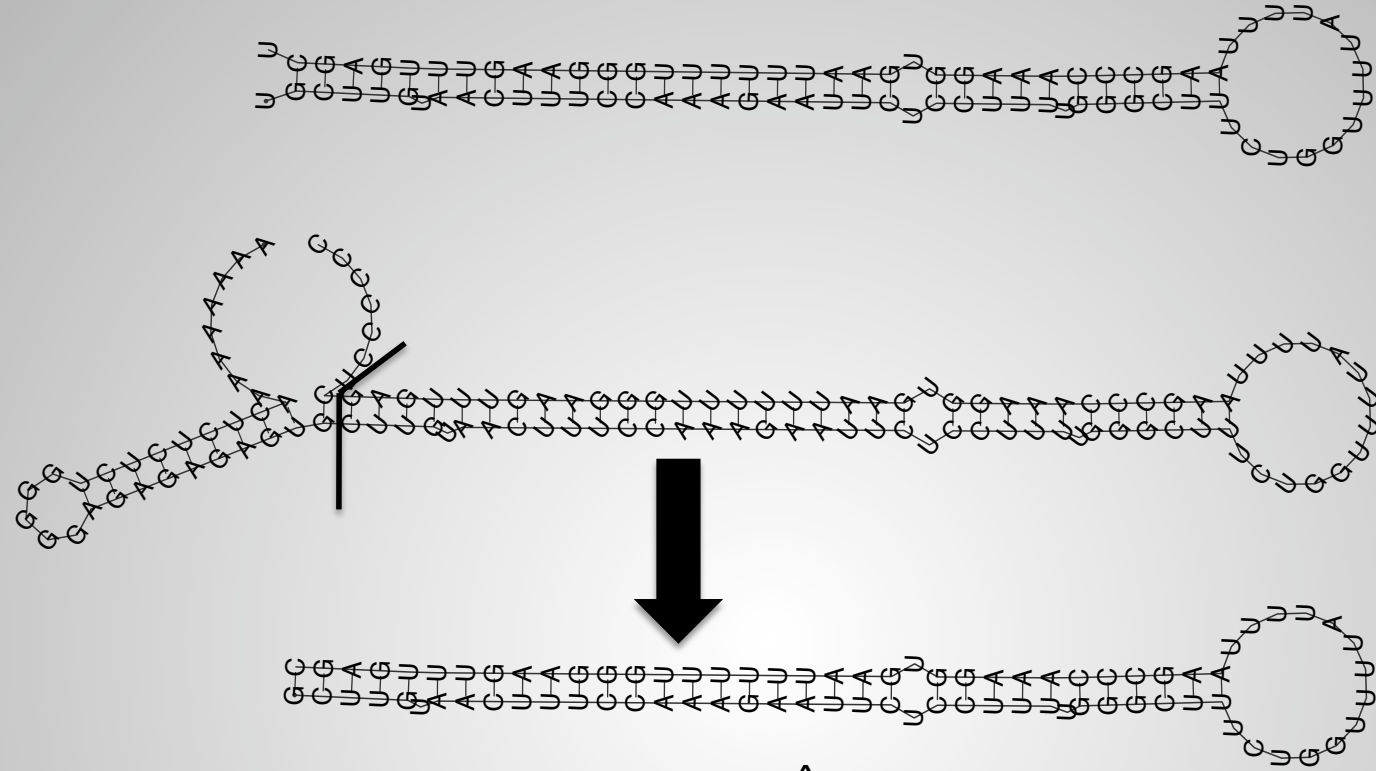
Interesting clusters

Clusters
in repeat-associated regions
or
overlapping with already annotated ncRNAs

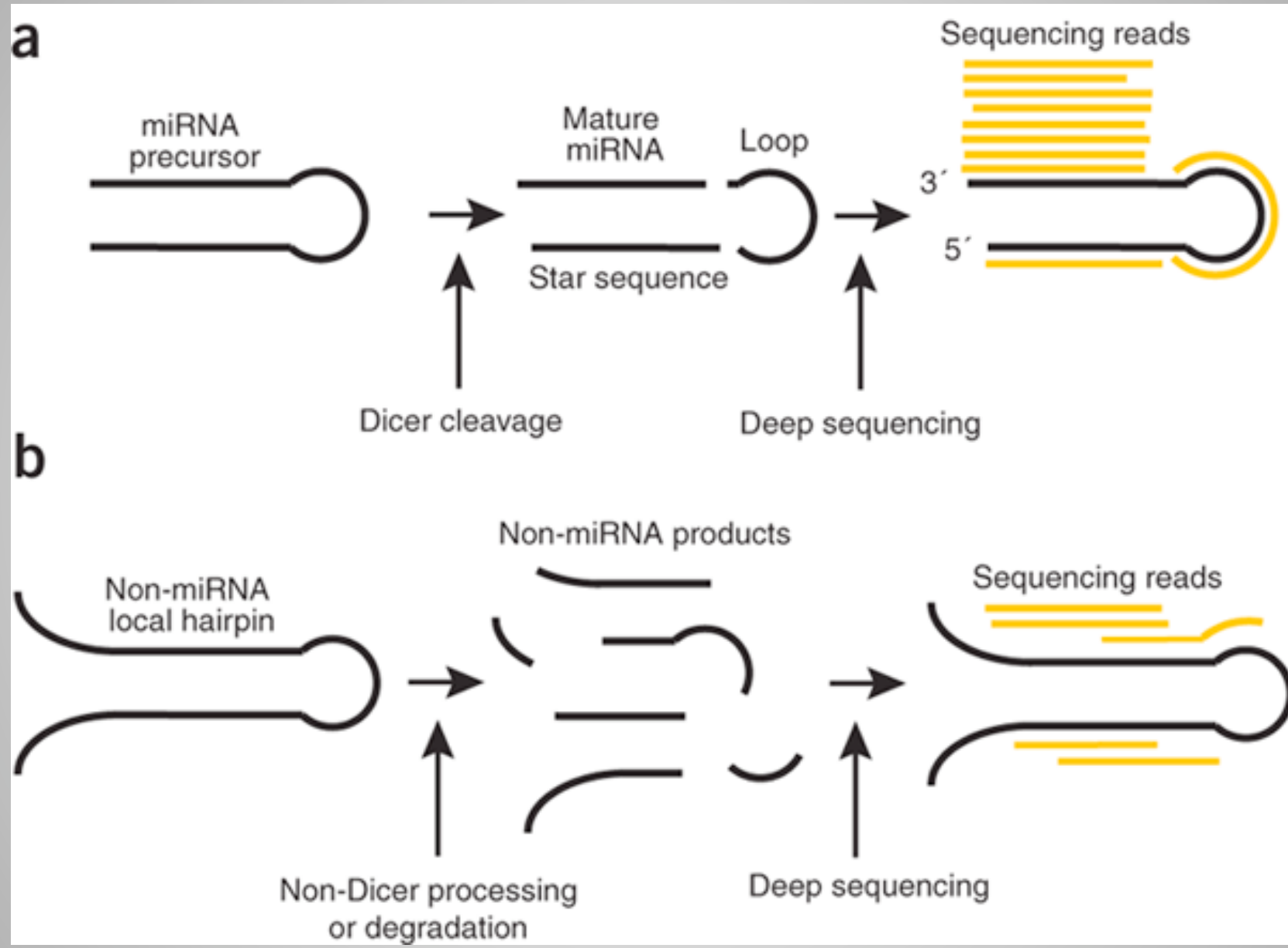


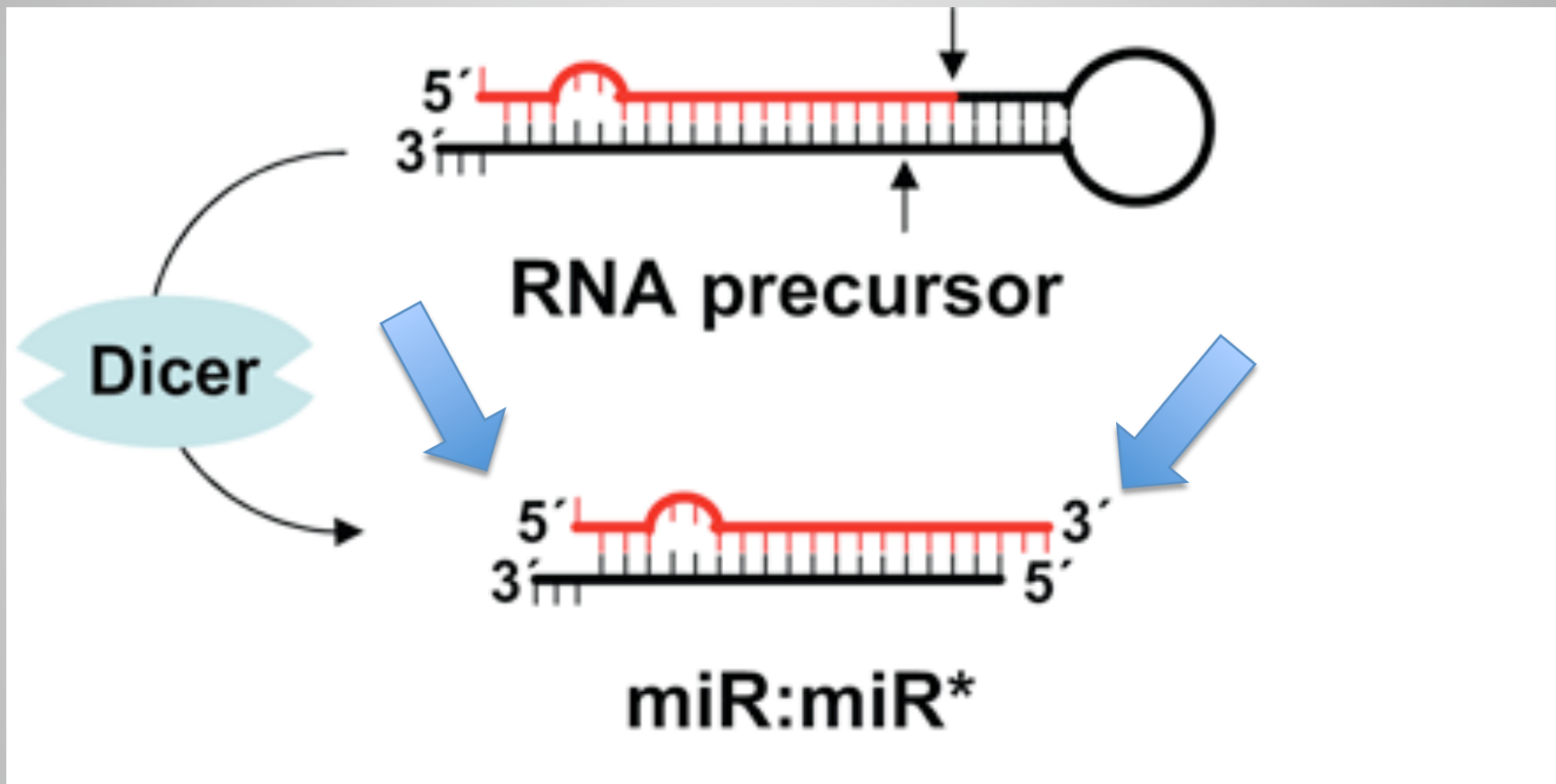
Bioinformatics pipeline

hairpin structure



dicer cleavage





Schwarz 2006

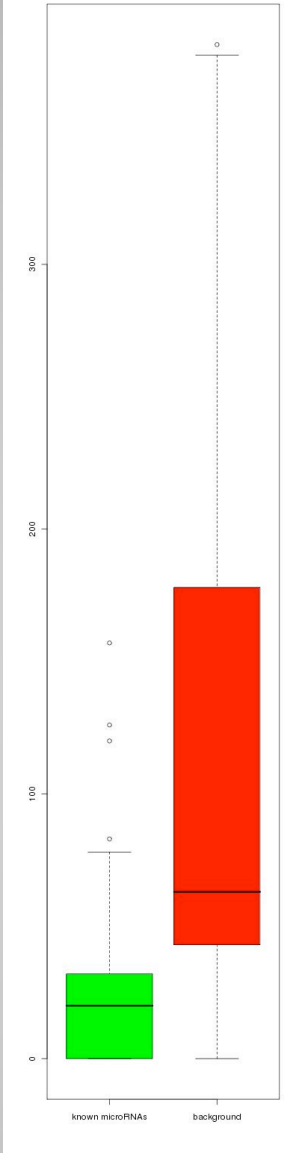
hairpin stability

```
(((((.(((((((((((((((((((..((((((((((((((....)))))))))))))))))))))))))))))))))))))))))))))
..(((((((((.(((((((((((((((((((..((((((((((((((....)))))))))))))))))))))))))))))))))))))))))))))
.(((...(((.(((.(((((((((((((((((((..((((((((((((((....))))))))))))))))))))))))))))))))))))))))))....))
....((((((...(((.(((.(((((((((((((((((((..((((((((((((((....))))))))))))))))))))))))))))))))))))))))))....)))))
(((.....(((...(((.(((.(((((((((((((((((((..((((((((((((((....))))))))))))))))))))))))))))))))))))))))))....)))))))))
```

Stable hairpin structure

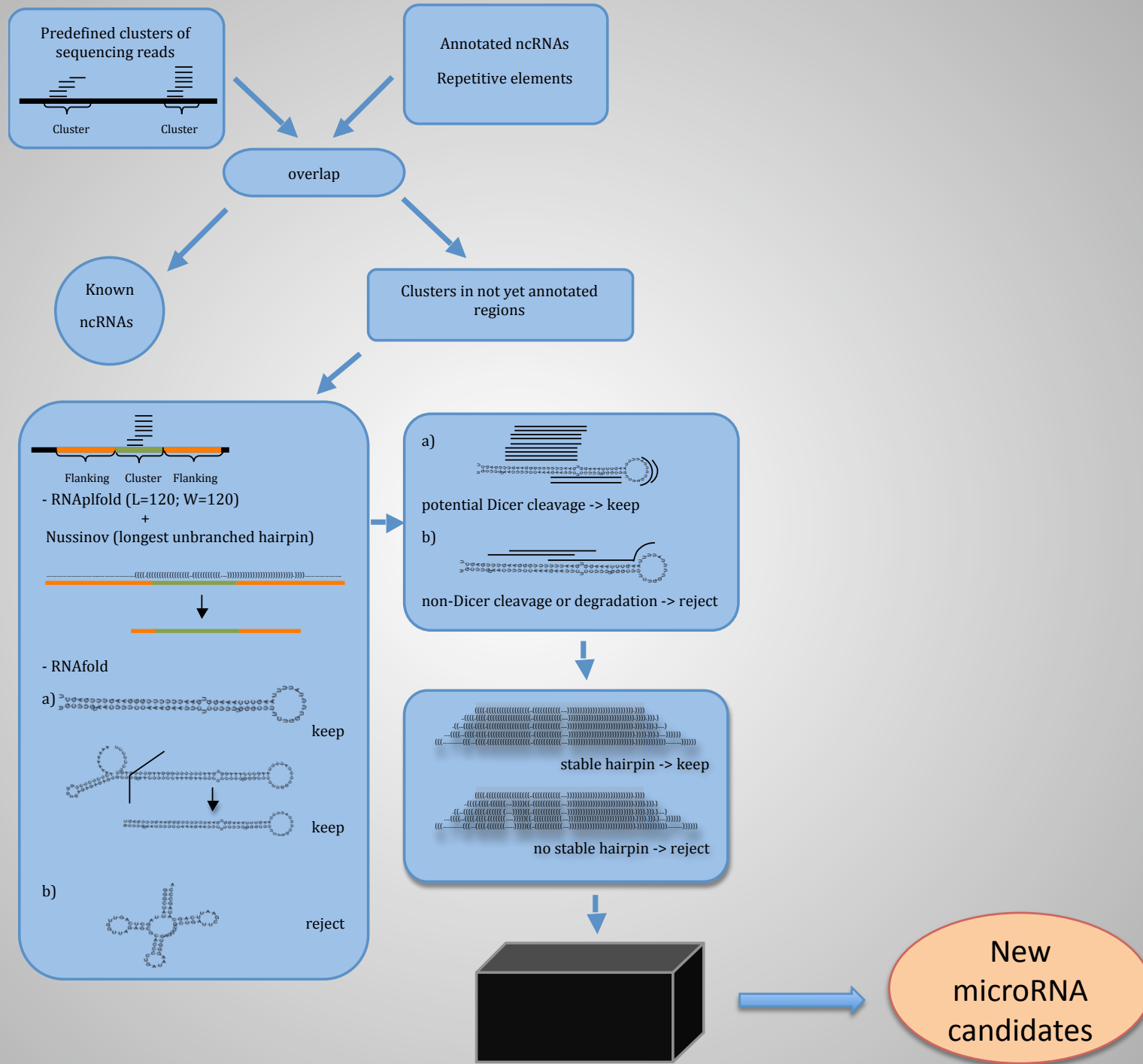
```
(((((.(((((((((((((((((((..((((((((((((((....)))))))))))))))))))))))))))))))))))))))))))))
..(((((((((.(((((((((((((((((((..((((((((((((((....)))))))))))))))))))))))))))))))))))))))))))))
.(((...(((.(((.(((((((((((((((((((..((((((((((((((....))))))))))))))))))))))))))))))))))))))))))....))
....((((((((((((((.....))))))))))))))..(((((((((((((((((((..((((((((((((((....))))))))))))))))))))))))))))))))))))))))))....)))))
.....((((((((((((((.....)))))))))))))..(((....(((.(((((((((((((((((((..((((((((((((((....))))))))))))))))))))))))))))))))))))))))))....))))))..(((..)))))
```

No stable hairpin



P-value: 4e-7

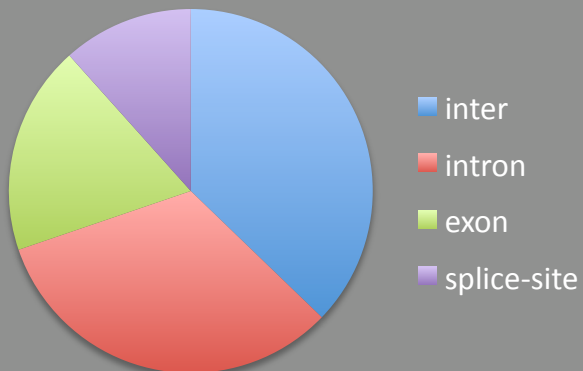
Workflow



RNAmicro

Recall of known
microRNAs (chrom1):
~50%

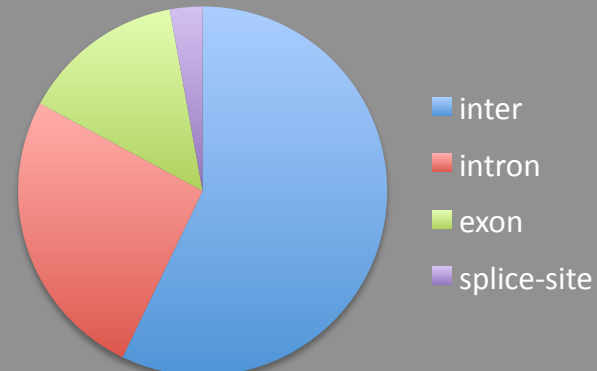
New candidates: **46**



tripletSVM

Recall of known
microRNAs (chrom1):
~66%

New candidates: **41**



Example of novel microRNA candidate

Outlook

microRNA detection:

- new machine learning approach
- new features

ncRNA Classification:

- using locaRNA to search for similar structures as tRNA, snRNA, etc.
- include preprocessed genome-wide predictions for snoRNAs and snRNAs, etc.

New machine learning approach for Deep-Sequencing data

Features:

- Structure
 - Energy
 - Loop
 - 5' and 3' Stems
 - Bulges
- Position of reads
 - Dicer cleavage (2nt overhang, valid products)
 - Distance to loop
 - Cluster size (shifted reads)
- Hairpin stability
-

Outlook

microRNA detection:

- new machine learning approach
- new features

ncRNA Classification:

- using locaRNA to search for similar structures as tRNA, snRNA, etc.
- include preprocessed genome-wide predictions for snoRNAs and snRNAs, etc.

Thanks To

Uni Leipzig

Peter F. Stadler

Clara I. Bermudez Santana (clustering, known ncRNAs)

Jana Hertel (microRNAs, RNAmicro)

Steve Hoffmann (segemehl)

TU Munich

Martin Sturm (machine learning)