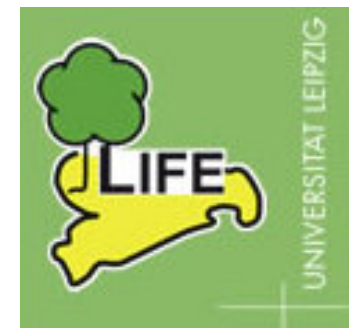# Detecting new miRNAs from deep sequencing data: A field study in worm

## David Langenberger

25th TBI Winterseminar

# Outline

- Data
- Pre-Analysis
- miRNA prediction
  - Classifier
  - miRDeep
  - miRanalyzer
- Results

# Data

# C. elegans small RNA sequencing

## Set 1 (MDC - set)

### Data

- SOLiD sequenzer
- 6,080,238 reads
- 18-27 nt in length

### Mapping

- 3,288,430 reads mapped with BWA (54%)
- remapping by segemehl (377,594 tags)

# C. elegans small RNA sequencing

**Set 2 (Bartel - set)**

**Data** (GSE5990)
- mixed-stage C. elegans
- 454 sequenzer
- 850,870 reads (181,668 tags)
- 18-36 nt in length

**Mapping**
- 138,868 tags mapped by segemehl (76%)

# C. elegans small RNA sequencing

## Set 3 (Berezikov - set)

**Data**(GSE15169)
- mixed-stage C. elegans
- 454 sequenzer
- 181,849 reads (23,327 tags)
- 18-38 nt in length

## Mapping
- 22,277 tags mapped by segemehl (95%)

# Pre-Analysis

# ncRNA mapping results (MDC – set)

| ncRNA type | # of annotated ncRNAs | # of annotated ncRNAs with read support | % of annotated ncRNAs with read support |
|---|---|---|---|
| all | 6606 | 4137 | 62.6 |
| microRNA | 174 | 165 | 94.8 |
| tRNA | 631 | 626 | 99.2 |
| snoRNA | 133 | 114 | 85.7 |
| snRNA | 90 | 84 | 93.3 |
| rRNA | 25 | 25 | 100 |
| 21U-RNA | 5356 | 2839 | 53.0 |
| SL2 splice leader | 8 | 8 | 100 |
| others | 189 | 149 | 78.8 |

# ncRNA mapping results (MDC – set)

| ncRNA type | # of tags | # of reads | % of all tags | % of all reads |
|:---:|:---:|:---:|:---:|:---:|
| **all** | 68,385 | 1,045,428 | 18,1 | 31,8 |
| **miRNAs** | 6,625 | 857,689 | 1,8 | 26,1 |
| **tRNAs** | 10,93 | 93,789 | 2,9 | 2,9 |
| **snoRNAs** | 1,163 | 3,774 | 0,3 | 0,1 |
| **snRNAs** | 2,186 | 5,652 | 0,6 | 0,2 |
| **rRNAs** | 33,519 | 80,152 | 8,9 | 2,4 |
| **21U-RNAs** | 3,524 | 6,181 | 0,9 | 0,2 |
| **others** | 6,946 | 1,626 | 1,8 | 0,0 |

# miRNA prediction

# Classifier

# Data preparation

1. Map small RNAs to the human genome, using `segemehl`



2. Cluster hits based on their genomic location (distance <100nt)



>100nt

3. Divide consecutive reads into blocks, using `blockbuster`



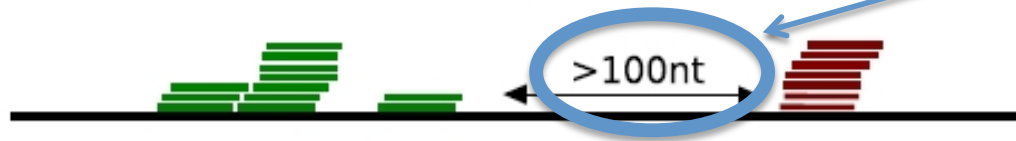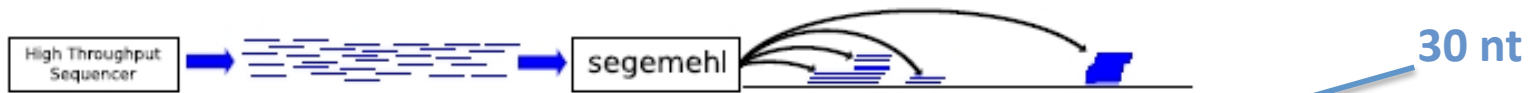4. Discard clusters with <2 blocks and/or <10 reads (small information content)
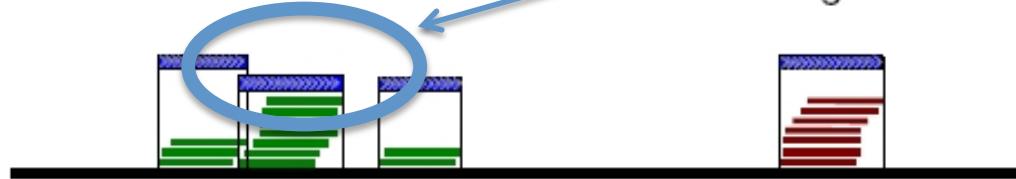
# Data preparation



1. Map small RNAs to the human genome, using `segemehl`

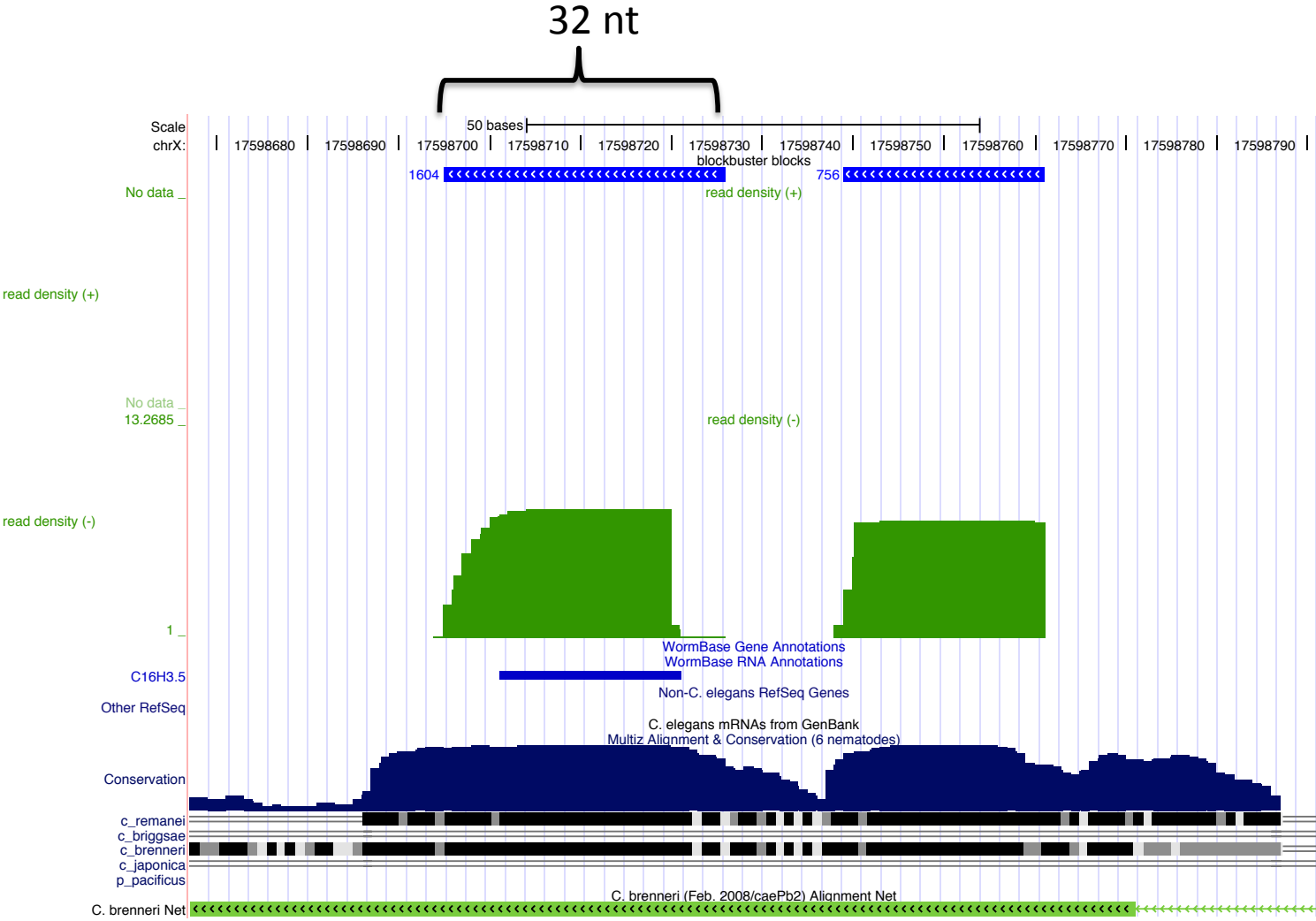2. Cluster hits based on their genomic location (distance <100nt)

**30 nt**

>100nt

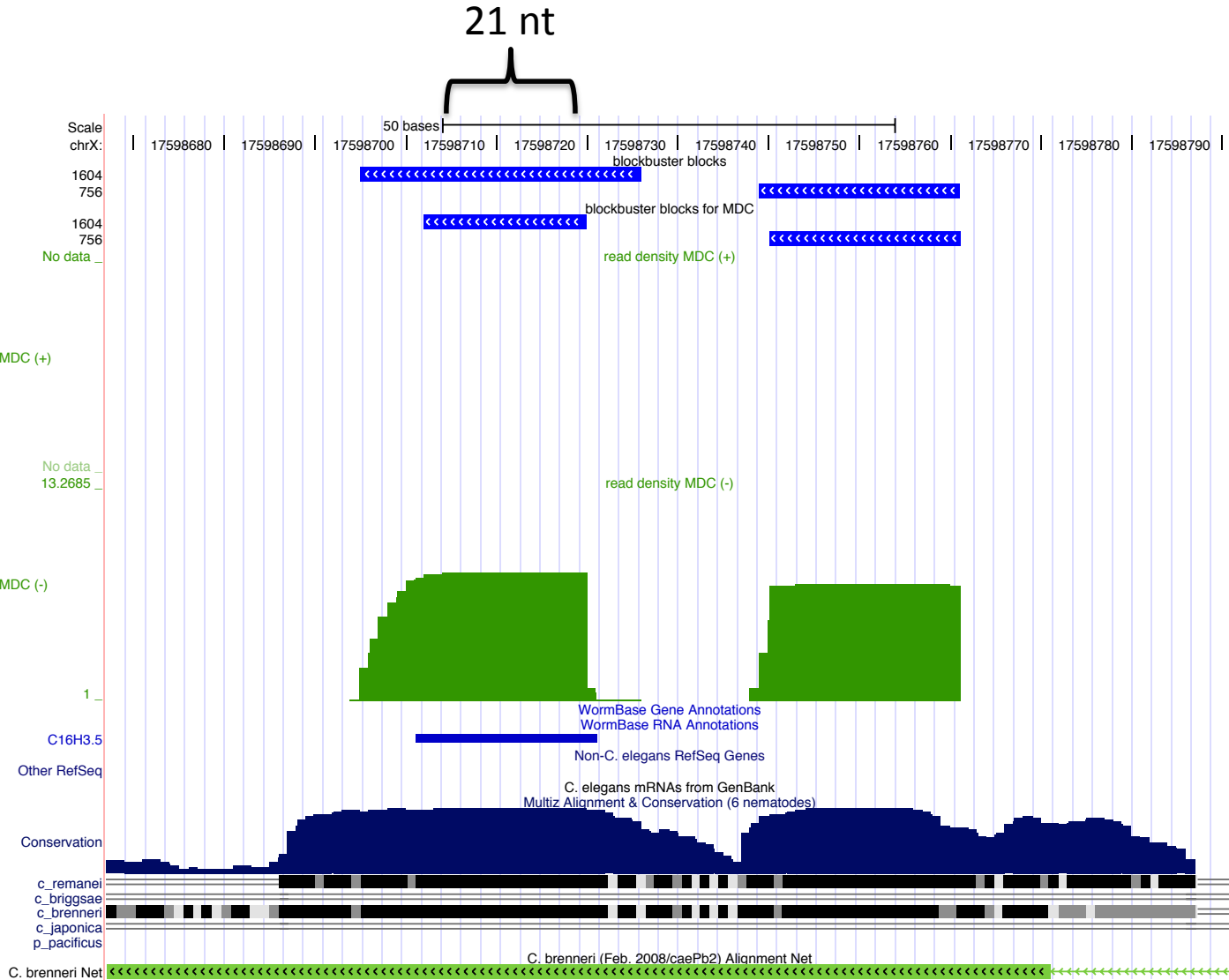3. Divide consecutive reads into blocks, using `blockbuster`

4. Discard clusters with <2 blocks and/or <10 reads (small information content)

# Data preparation

1. Map small RNAs to the human genome, using segemehl

2. Cluster hits based on their genomic location (distance <100nt)

>100nt

**30 nt**

**Blocks are cut at positions with less than 80% read support**

3. Divide consecutive reads into blocks, using blockbuster

4. Discard clusters with <2 blocks and/or <10 reads (small information content)
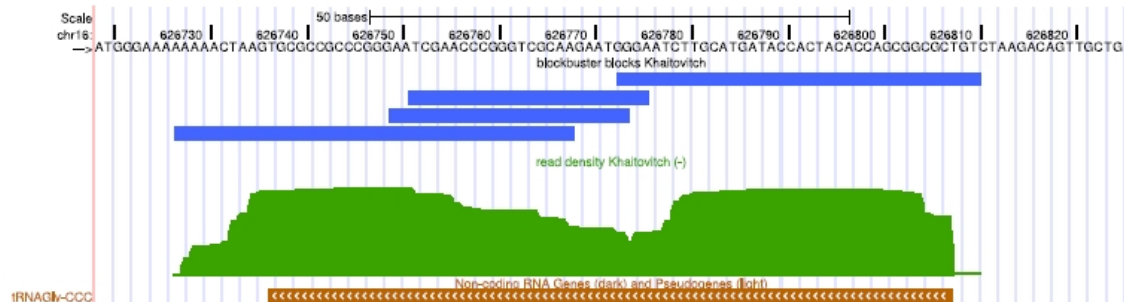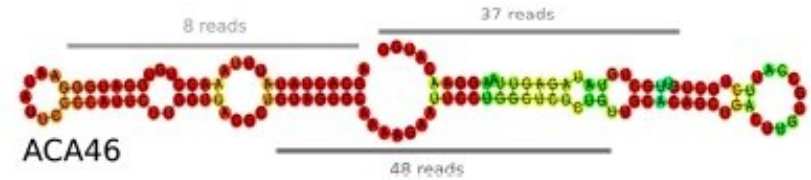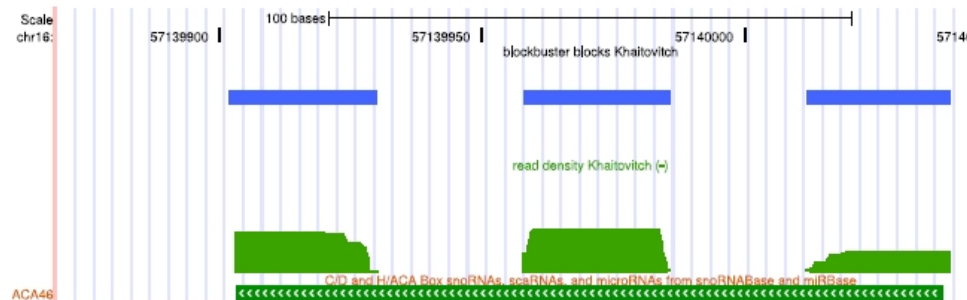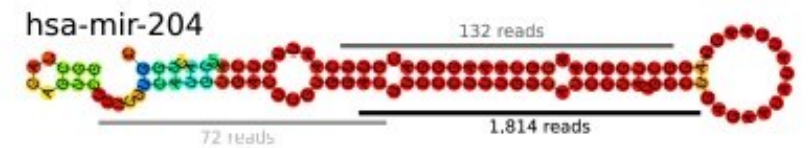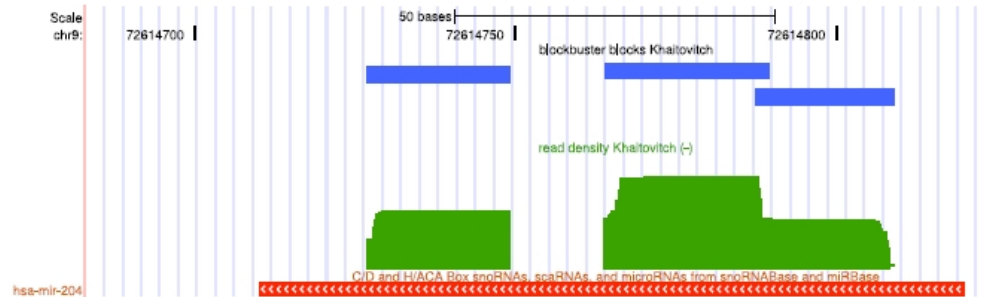
# Data preparation

# Data preparation

# Different ncRNAs have different read patterns

# Classifier (random forest)

- Trained on human reads mapping to human ncRNAs (khaitovitch – brain data – 454 sequenzer)
- Trainingset:
  - 243 miRNAs
  - 19 snoRNAs (H/ACA)
  - 116 snoRNAs (C/D)
  - 336 tRNAs
  - 273 other ncRNAs

### Confusion Matrix (10-fold cross validation)

| miRNA | snoRNA (H/ACA) | snoRNA (C/D) | tRNA | other | <- classified as |
|---|---|---|---|---|---|
| **232 (95%)** | 3 | 2 | 4 | 2 | miRNA |
| 8 | **3 (2%)** | 1 | 3 | 4 | snoRNA (H/ACA) |
| 2 | 2 | **71 (61%)** | 16 | 25 | snoRNA (C/D) |
| 3 | 0 | 5 | **291 (87%)** | 37 | tRNA |
| 9 | 1 | 15 | 53 | **195 (71%)** | other |

# Classification of worm ncRNAs

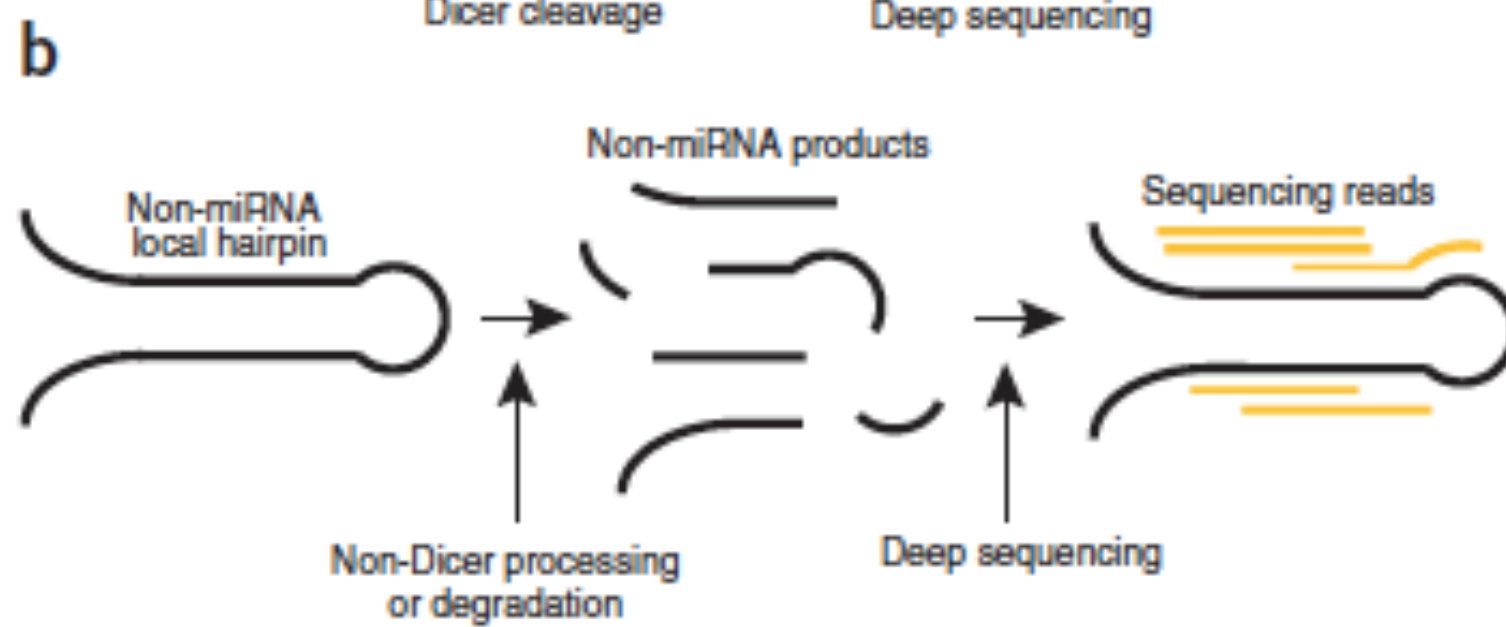|  | Recall | | | | PPV |
|---|---|---|---|---|---|
|  | miRNA | snoRNA | tRNA | | |
| **MDC set** | **61%**<br>(49 of 80) | **0%**<br>(0 of 28) | **76%**<br>(438 of 573) | | 94% |
| **Bartel set** | **85%**<br>(61 of 71) | **5%**<br>(1 of 8) | **76%**<br>(438 of 573) | | 95% |
| **Berezikov set** | **79%**<br>(26 of 33) | **0%**<br>(0 of 0) | **31%**<br>(52 of 169) | | 96% |

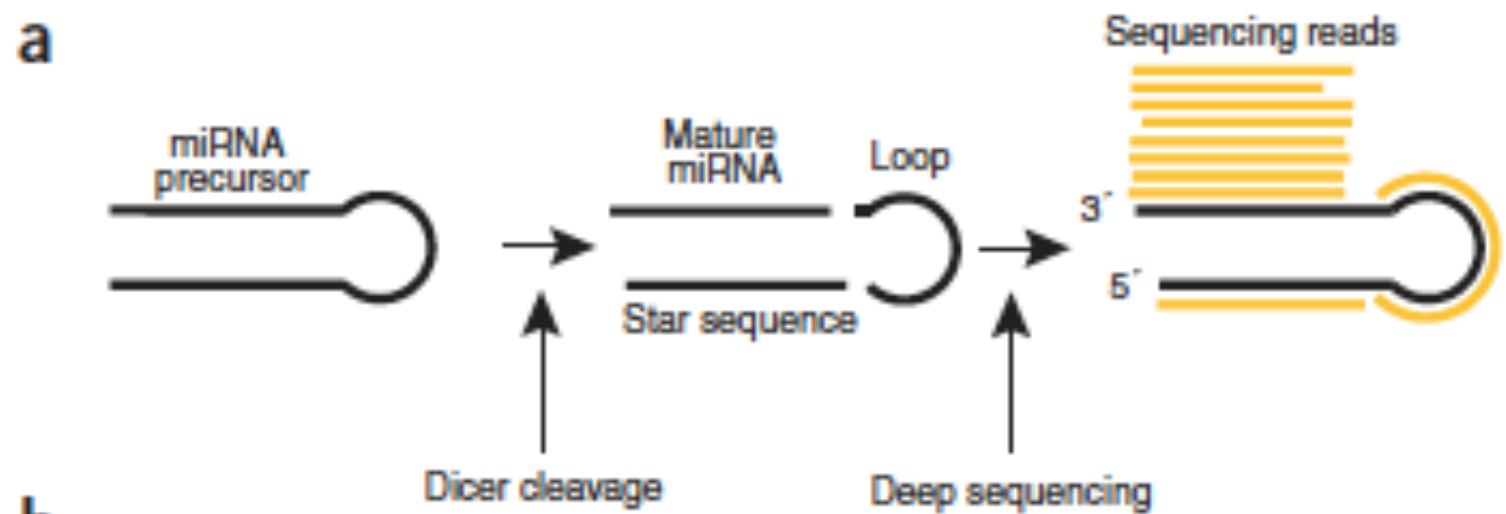# miRDeep

# Discovering microRNAs from deep sequencing data using miRDeep

Marc R Friedländer[1], Wei Chen[2], Catherine Adamidi[1], Jonas Maaskola[1], Ralf Einspanier[3], Signe Knespel[1] & Nikolaus Rajewsky[1]

**The capacity of highly parallel sequencing technologies to detect small RNAs at unprecedented depth suggests their value in systematically identifying microRNAs (miRNAs). However, the identification of miRNAs from the large pool of sequenced transcripts from a single deep sequencing run remains a major challenge. Here, we present an algorithm, miRDeep, which uses a probabilistic model of miRNA**

and 454 Life Sciences/Roche, can sequence DNA orders of magnitude faster and at lower cost than Sanger sequencing and are evolving so rapidly that increases in sequencing speed by at least another order of magnitude seem likely over the next few years. Although the Solexa/Illumina system can produce ~32 million sequencing reads in one run, read length is currently limited to 35 bp. In contrast, the current 454 platform yields reads up to 200 bases each, although the number of reads

**a**

miRNA precursor

Mature miRNA    Loop

Star sequence

Dicer cleavage

Sequencing reads

3´
5´

Deep sequencing

**b**

Non-miRNA local hairpin

Non-miRNA products

Non-Dicer processing or degradation

Sequencing reads

Deep sequencing

# miRanalyzer

# miRanalyzer: a microRNA detection and analysis tool for next-generation sequencing experiments

Michael Hackenberg[1], Martin Sturm[2], David Langenberger[3,4], Juan Manuel Falcón-Pérez[5] and Ana M. Aransay[1,*]

[1]Functional Genomics Unit, CIC bioGUNE, CIBERehd, Technology Park of Bizkaia, 48160 Derio, Bizkaia, Spain, [2]Institute for Bioinformatics and Systems Biology, German Research Center for Environmental Health, Ingolstädter Landstrasse 1, D-85764 Neuherberg, [3]Department of Genome-Oriented Bioinformatics, Wissenschaftszentrum Weihenstephan, Technische Universitat München, 85350 Freising, [4]Bioinformatics Group, Department of Computer Science, University of Leipzig, Haertelstr. 16-18, D-04107 Leipzig, Germany and [5]Metabolomics Unit, CIC bioGUNE, CIBERehd, Technology Park of Bizkaia, 48160 Derio, Bizkaia, Spain

## ABSTRACT

Next-generation sequencing allows now the sequencing of small RNA molecules and the estimation of their expression levels. Consequently, there will be a high demand of bioinformatics tools to cope with the several gigabytes of sequence data generated in each single deep-sequencing experiment. Given this scene, we developed

## INTRODUCTION

The recent years witnessed a profound change in our understanding of the regulation of gene expression. Small non-coding RNA especially came into focus as it became clear that they are key players in many cellular processes by post-transcriptionally regulating gene expression via either degradation, translational repression, or both (1,2). MicroRNAs, belonging to the family of small non-coding RNAs, are endogenous in many animal and
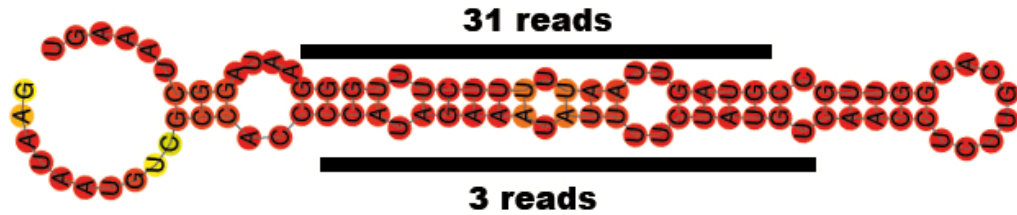
# Results

# miRNA prediction in worm

| MDC | miRDeep | miRanalyzer | Classifier |
|---|---|---|---|
| miRDeep | **335** | 50 | 3 |
| miRanalyzer | 50 | **650** | 4 |
| Classifier | 3 | 4 | **5** |

| Bartel | miRanalyzer | Classifier |
|---|---|---|
| miRanalyzer | **67** | 1 |
| Classifier | 1 | **7** |

| Berezikov | miRanalyzer | Classifier |
|---|---|---|
| miRanalyzer | **27** | 0 |
| Classifier | 0 | **1** |

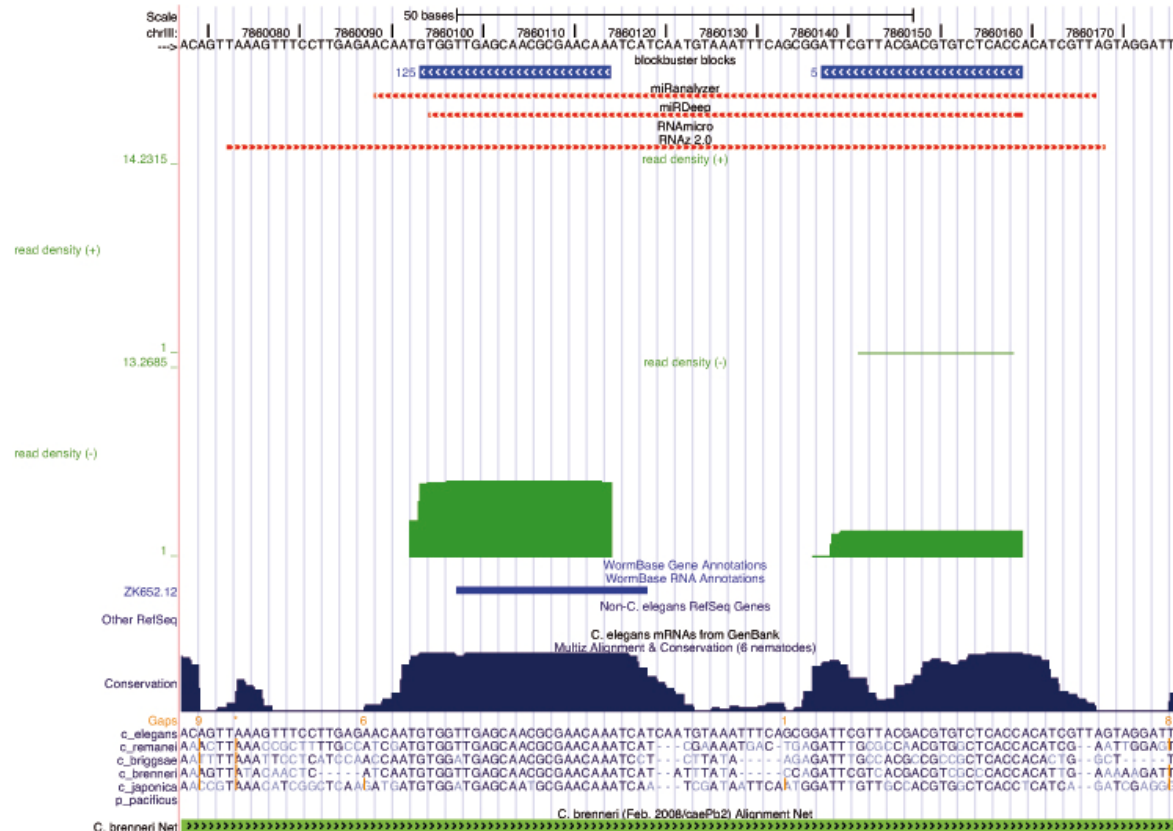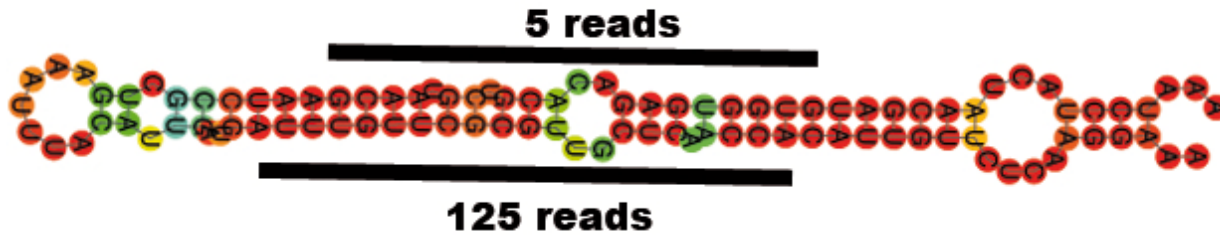# 3 candidates found by the classifier, miRanalyser and miRDeep (MDC – set)



chrX:2,329,018-2,329,107

# 3 candidates found by the classifier, miRanalyser and miRDeep (all sets)

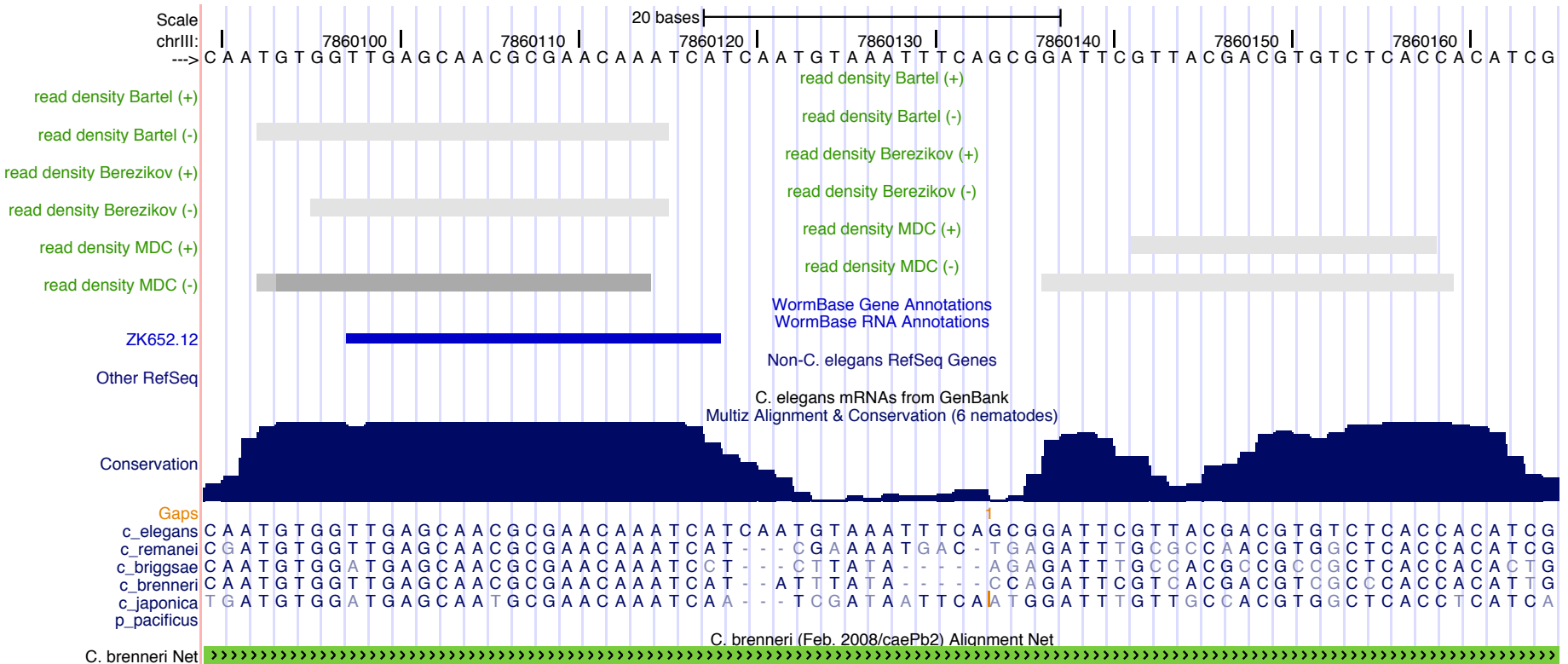# 3 candidates found by the classifier, miRanalyser and miRDeep (MDC – set)



**chrIII:7,860,027-7,860,227**
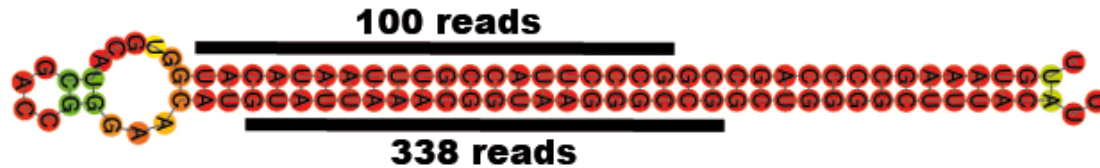
Antisense to cel-miR-356

RNAz support

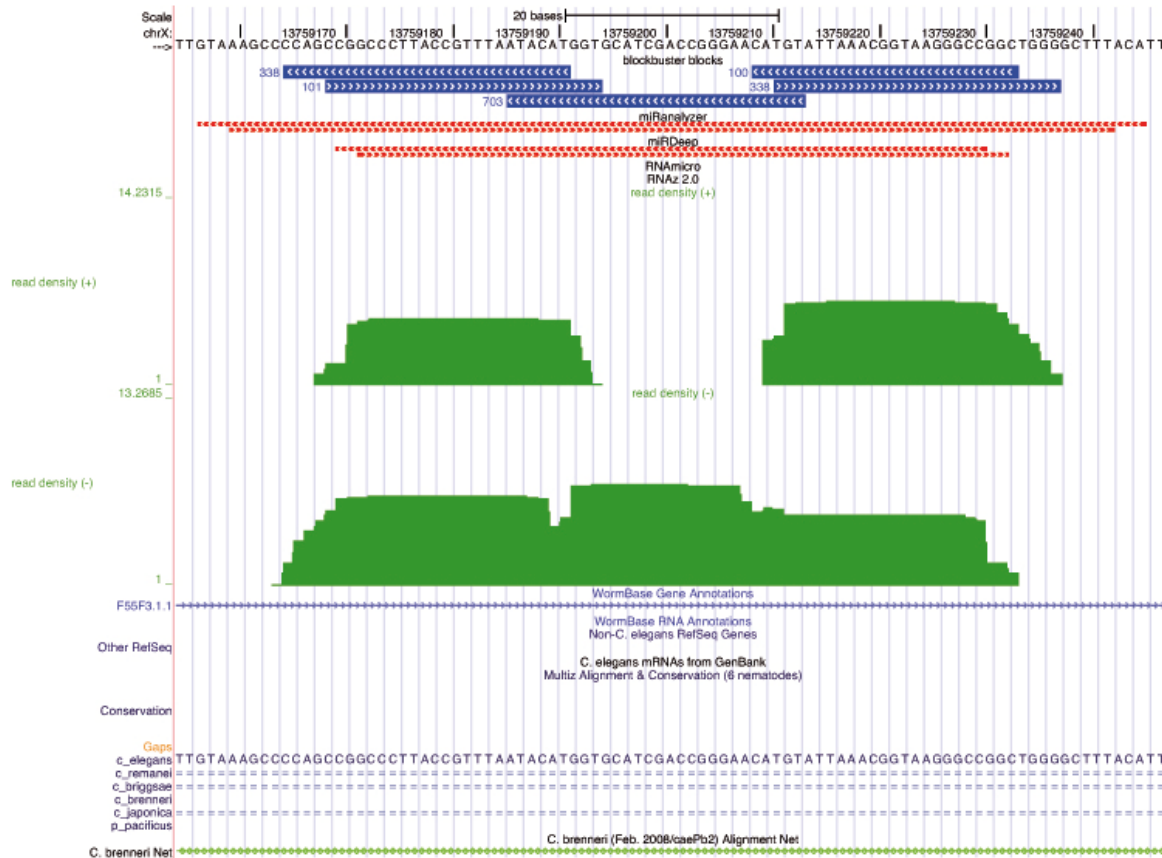# 3 candidates found by the classifier, miRanalyser and miRDeep (all sets)



- miRanalyzer predictions for all datasets
- No classification (only one block in other datasets)

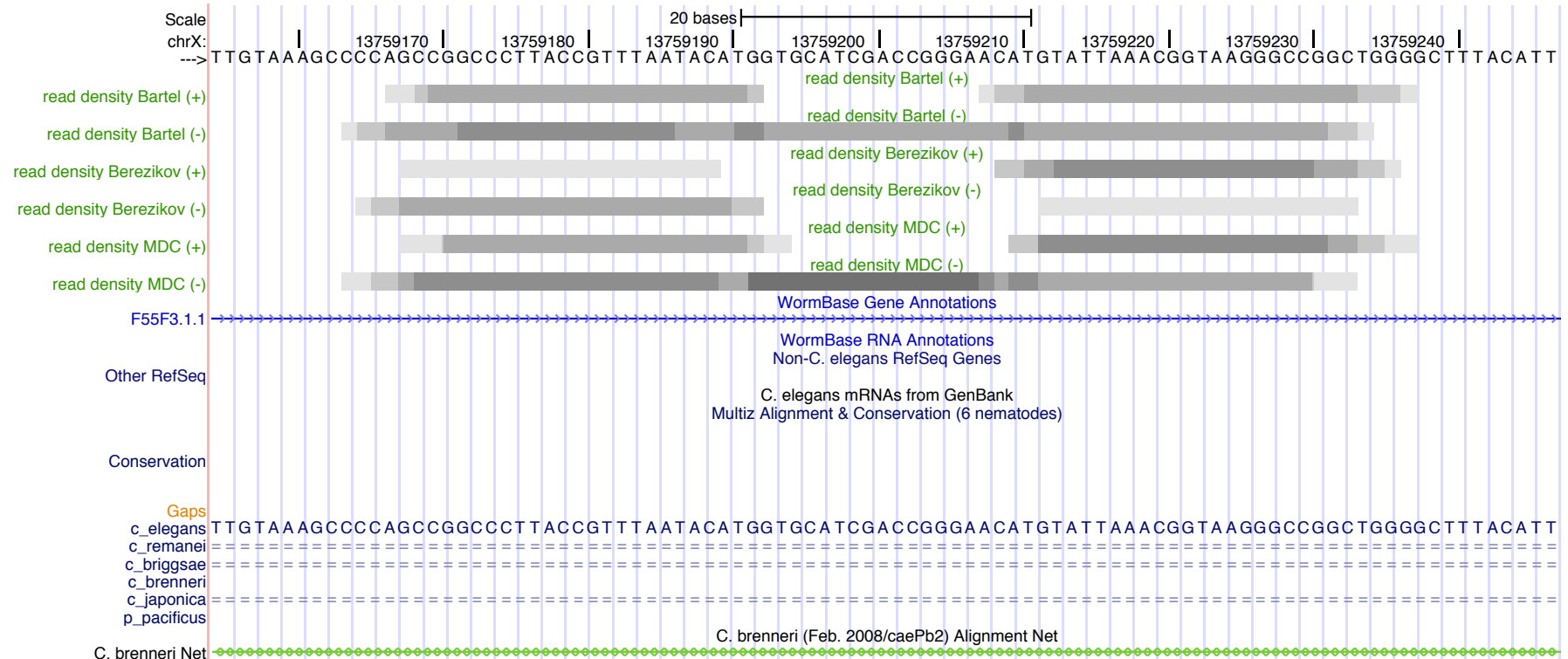# 3 candidates found by the classifier, miRanalyser and miRDeep (MDC – set)



chrX:13,759,155-13,759,247

Sense miRNA detected by all three methods

Antisense miRNA only detected by miRanalyzer and miRDeep

# 3 candidates found by the classifier, miRanalyser and miRDeep (all sets)



- miRanalyzer predictions for all datasets
- classification for sets from Bartel and MDC (no classification for Berezikov)

# Conclusion

- There are a lot of not annotated ncRNAs (amongst others microRNAs) waiting to be found by using deep sequencing data

- Classifier can be used between species (patterns seem to be conserved)

- Different experiments (454, Solid, Solexa) generate comparable patterns

# Acknowledgements