



# KNIME & Next Generation Sequencing at the Pasteur Institute

Bernd Jagla, Ph.D.  
*Institut Pasteur*  
*Génopole – Plate-forme 2 Puces à ADN*

# Thanks

- **KNIME**
  - Karol Kozak
  - KNIME community
- **PF2:**
  - Odile Sismeiro
  - Marie-Agnes Dillies
  - Jean-Yves Coppee
  - Caroline Proux
  - Guillaume Soubigou
- **GBrowse:**
  - Lincoln Stein and the rest of the GBrowse/BioPerl community
- **Mobylye:**
  - Hervé Ménager, Bertrand Néron

# Content

- Introduction
- NGS data processing
- How KNIME fits in this

# Introduction

- ***Plate-forme 2 Puces à ADN – what we do***
  - *Tiling arrays, Microarrays, NGS*
- ***Putting it into context – Introduction to small RNAs***
- ***State-of-the-art technologies: 2<sup>nd</sup> generation sequencing***

# *Plate-forme 2 Puces à ADN*

- Tiling arrays
  - Scan a known genome with equally spaced probes
- Micro arrays
  - Scan a known genome for regions of interest (e.g. genes/mRNA)
- Next generation sequencing
  - Unbiased sequencing of unknown RNA or DNA molecules

# Small RNA

- size ranges from 20 to 27 nucleotides
- miRNA, siRNA and piRNA
- play essential roles in the all eukaryotes kingdoms (protists, animals, plants, and fungi)
- involved in different phenomena essential for
  - genome stability, development, adaptive responses to biotic and abiotic pressures, control of transposable elements, antiviral defence...
- effects of small RNAs on gene expression are generally inhibitory

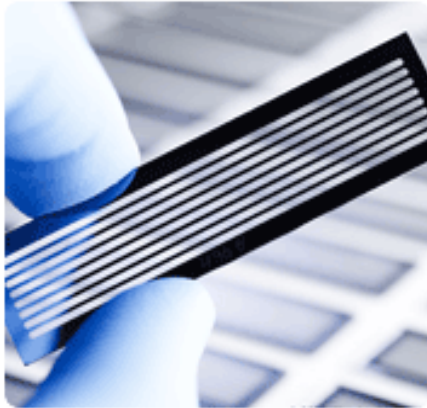
# Illumina Solexa technology

- Ultra High Throughput sequencing:
  - One machine = 150 M sequences of 36 - ~150 nucleotides per week
  - One human genome in two weeks
  - Sanger has about 40 machines
  - Beijing Genomic Institute has eq. >300 machines with the throughput described here

# Equipment

---

Flow cell



Cluster Station



Genome Analyser



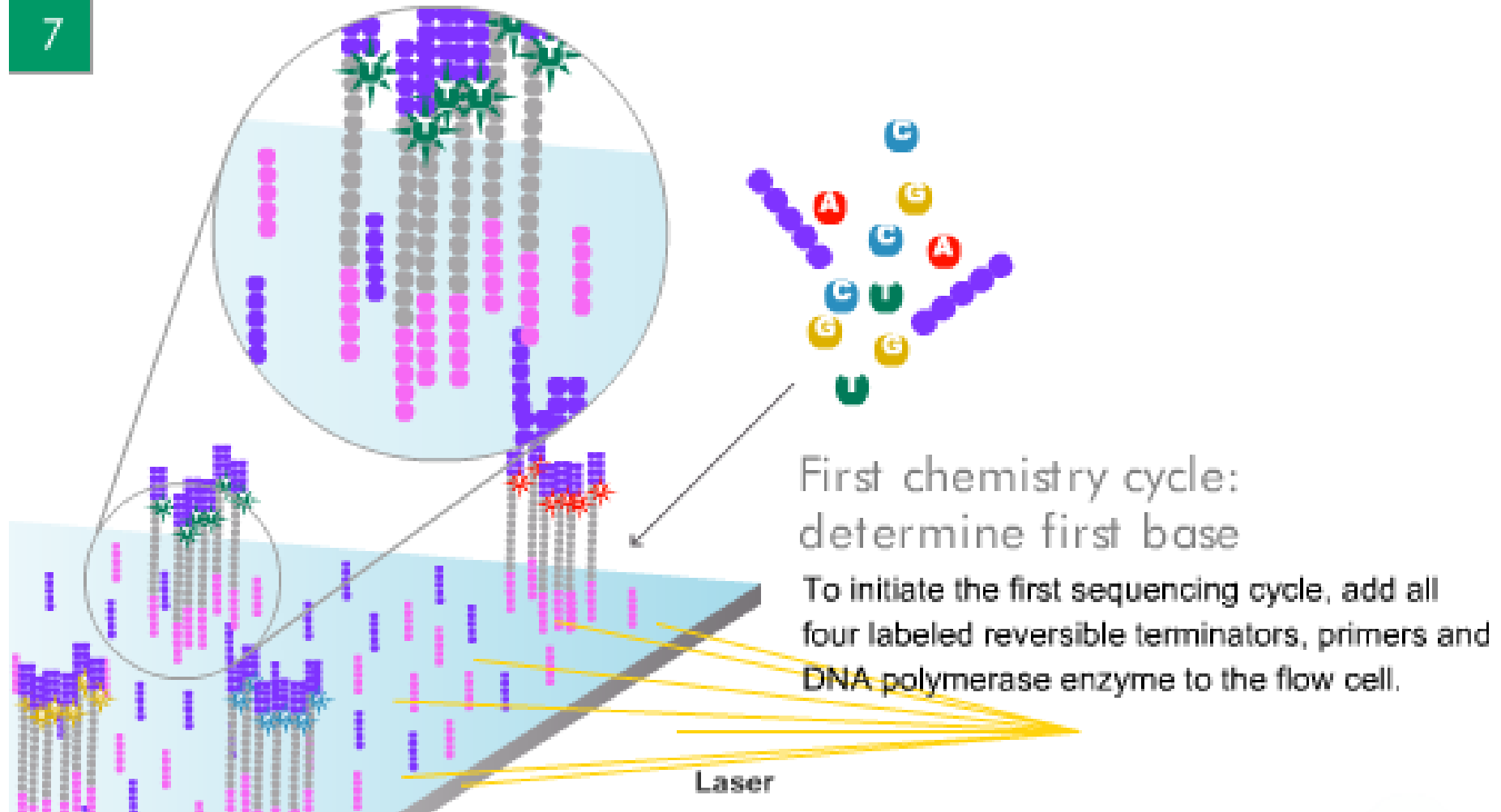
Paired-end module





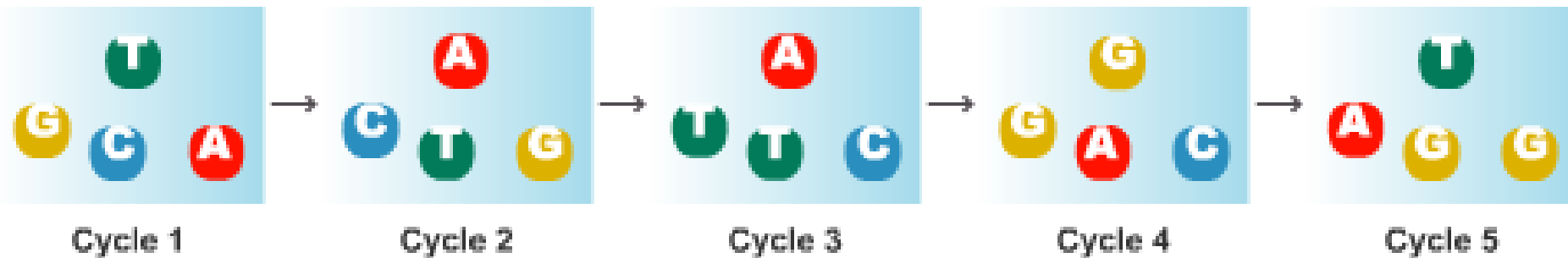
# Principle: sequencing by synthesis

7



# Sequence by synthesis

11

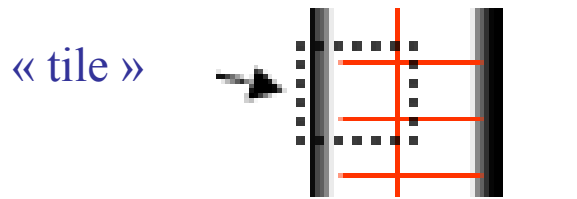
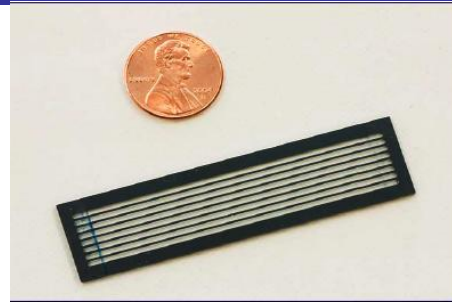


GCTGA....

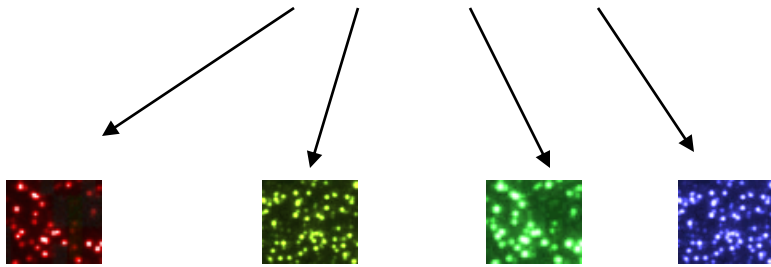
Sequence read over multiple chemistry cycles

Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

# 1 Cycle



genome analyser



7 Mb + 7 Mb + 7 Mb + 7 Mb = 28 Mb  
1 + 1 + 1 + 1 = 4 files

x 100

x 100



x 8

x 8

2,8 Gb

400 files

**22,4 Gb**

3200 files

\* Each file contains 50 000 - 150 000 lines (clusters)

# 1 run

1 cycle	22.4 Gb – 3200 files
---------	----------------------

36 cycles	806 Gb – 115.200 files
-----------	------------------------

72 cycles	1.6 Tb – 230.400 files
-----------	------------------------

72 cycles	3.2 Tb – 460.800 files
-----------	------------------------

Paired end	
------------	--

# NGS data processing

- ***NGS Data primary data analysis***
  - *Image analysis*
  - *Base calling*
  - *Alignment*
  - *Quality control*
- ***NGS Data secondary data analysis***
  - *Functional annotation*
  - *Visualization*
  - *Classification*
  - *Quantitative/Qualitative analysis*

# NGS – primary data analysis

Images

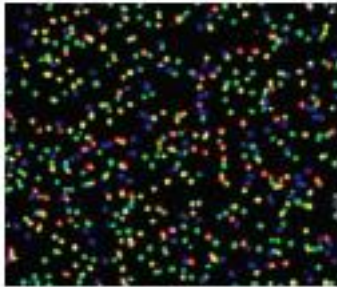


Image Analysis

Line No	X	Y	Ch1-FCST	Ch2-FCST
5	117	584	1000	495.1 589.9 628.7 255.6
5	117	573	593	56.5 613.8 250.5 458.8
5	117	488	788	1583.2 705.1 48.8 87.4
5	117	688	888	1343.8 788.0 88.8 138.8
5	117	1107	1207	58.8 818 857.5 888.2
5	117	1054	488	284.7 884.4 87.2 85.5
5	117	887	588	785.1 888.8 83.2 388.8
5	117	945	1788	83.2 541.9 884.7 885.7
5	117	888	774	885.5 555.2 83.2 885.8
5	117	888	1188	334.8 814.8 18.7 78.5
5	117	547	1782	843.8 784.9 888.8 888.5
5	117	887	1774	81.8 81.8 828.3 1388.8

Base Calling

```

ATGGCCTGGGCTAGTTTCGATTTACGAT
CCTGGGCTAGTTTCGATTTACGATCGAT
GCTAGTTTCGATTTACGATCGATCGTTG
ATCGATCGTTGCATGCTGGGGTAGTGC
TTCGATTTACGATCGATCGTTGCATGCT
TCGATTTACGATCGATCGTTGCATGCTG
ICTAGTTTCGATTTACGATCGATCGTTGC
TCGATTTACGATCGATCGTTGCATGCTG
TACGATCGATCGTTGCATGCTGGGGTA
TCGATCGTTGCATGCTGGGGTAGTGCCT
TCGATTTACGATCGATCGTTGCATGCTG
CGATTTACGATCGATCGTTGCATGCTGG
TAGTTTCGATTTACGATCGATCGTTGCA
TGATTTACGATCGATCGTTGCATGCTGG
ACGATCGATCGTTGCATGCTGGGGTAG
    
```

Aligned Reads

```

TCGCTAAGGCTAAGTTTCATGCTAAGGTTTCGAA
A-GCCTAAGGCTAAGTTTCATGCTAAGGTTTCGAA
AT-GCTAAGGCTAAGTTTCATGCTAAGGTTTCGAA
ATG-GTAAGGCTAAGTTTCATGCTAAGGTTTCGAA
ATGC-TAAGGCTAAGTTTCATGCTAAGGTTTCGAA
ATGCG-AAGGCTAAGTTTCATGCTAAGGTTTCGAA
ATGCGT-AGGCTAAGTTTCATGCTAAGGTTTCGAA
ATGCGTA-GCTAAGTTTCATGCTAAGGTTTCGAA
ATGCGTAA-CTAAGTTTCATGCTAAGGTTTCGAA
    
```

- Transfert of image files
- Identification of clusters

- Cluster position
  - Clusters intensities
  - Signal to noise analysis
- => One file per cycle

- Intensity correction
- Quality assesment
- Calculate sequence for each cluster

- Align sequences to reference genome
- Sequence files
- Statistics
- Graphics

Number of files

360.000+

5.000

10.000

36.000

Disc space used

TBs

~100 GB

~ 80 GB

~ 25 GB

# NGS – *secondary data analysis*

- Making sense of the alignments
  - Visualization (Genome Browser)
  - Statistical interpretation

# KNIME - Content

- *Why KNIME*
- *Mobylye*
- *GBrowse*
- *Implemented workflows*
- *Under development:*
  - *Upload to GBrowse*
  - *Mobylye*
  - *SSH*
  - *Cytoscape*
- *Availability*
- *Missing Features*



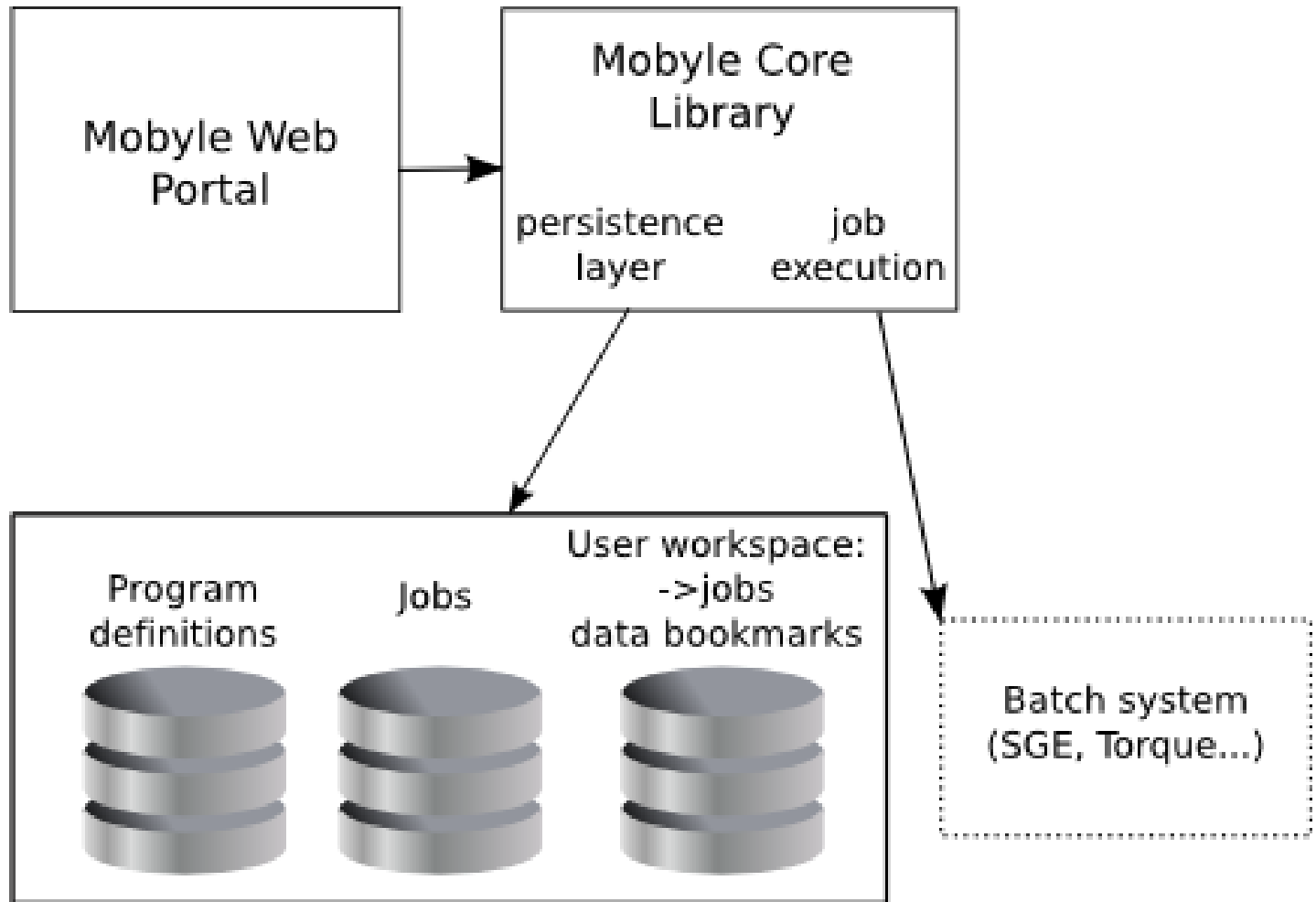
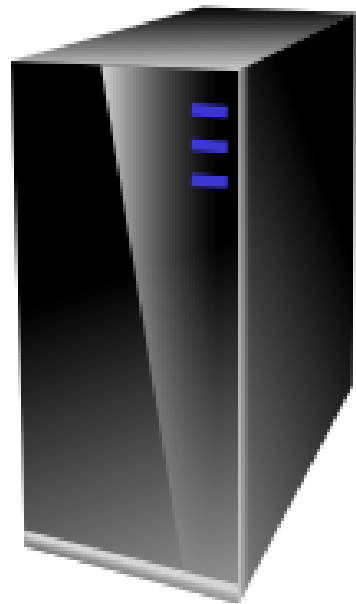
# Why KNIME

- Replace command line scripts
- Graphical
- Tool integration
- Can deal with large number of rows
- Open-source
- Flexible
- Workflows are executable from command line
- Interact with data (Highlight)
- **Very good support**

# Using KNIME

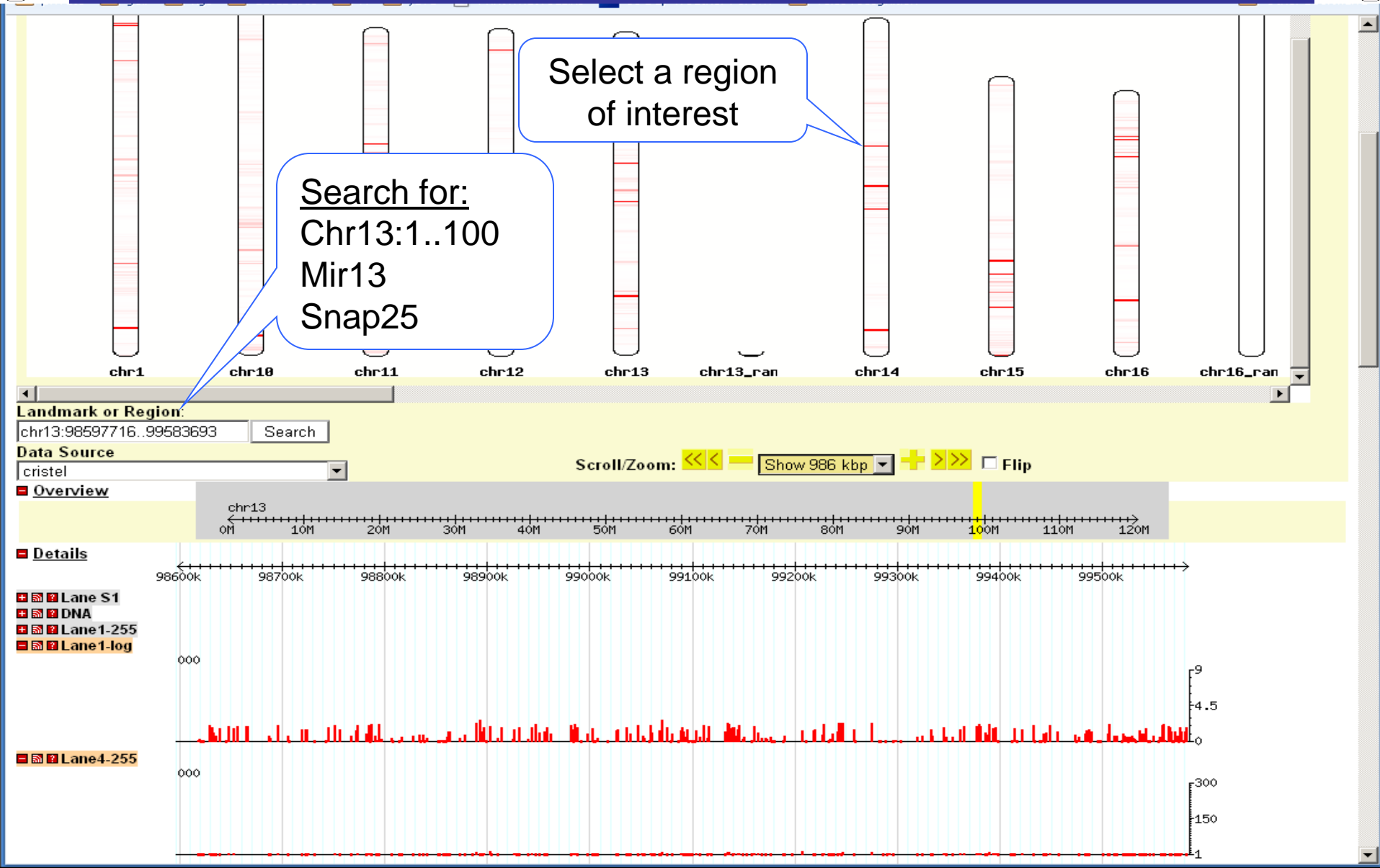
1. Develop workflow using sample data
2. Execute through command line via
  - Mobylye, Web-interface, directly...
3. Examine results using KNIME
4. Goto 1.

# Mobyle



> 250 programs integrating, mostly dealing with sequence data

# Visualization: GBrowse



Details

10412.6k 10412.7k 10412.8k 10412.9k 10413k 10413.1k 10413.2k 10413.3k 10413.4k 10413.5k 10413.6k 10413.7k 10413.8k

- + Lane S1
- + DNA
- Lane 1-log

- + Lane 4-log
- Lane 3-log

new miRNA?

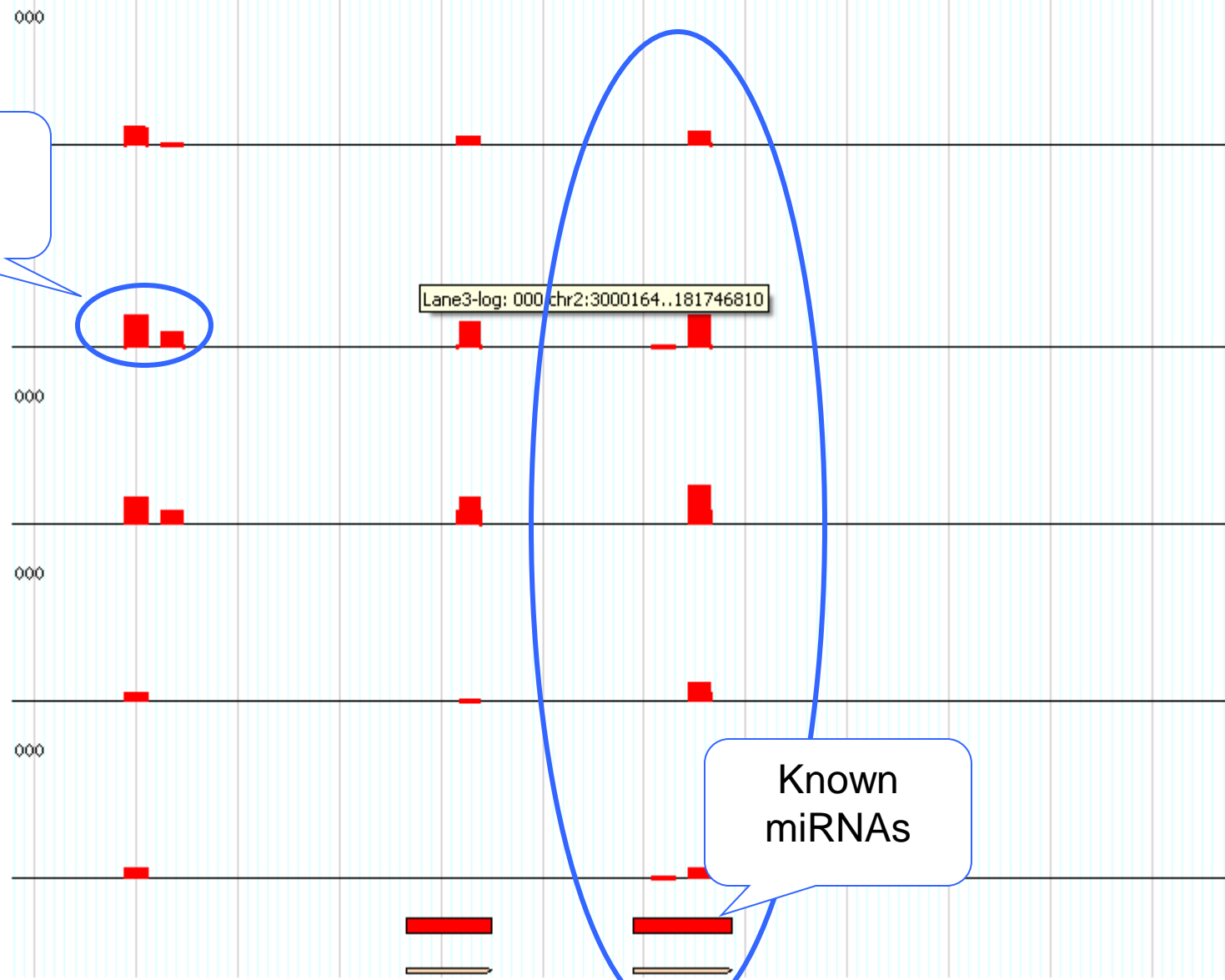
- Lane 2-log

- Lane 7-log

- Lane 6-log

- smRNA

- miRNA



Landmark or Region:

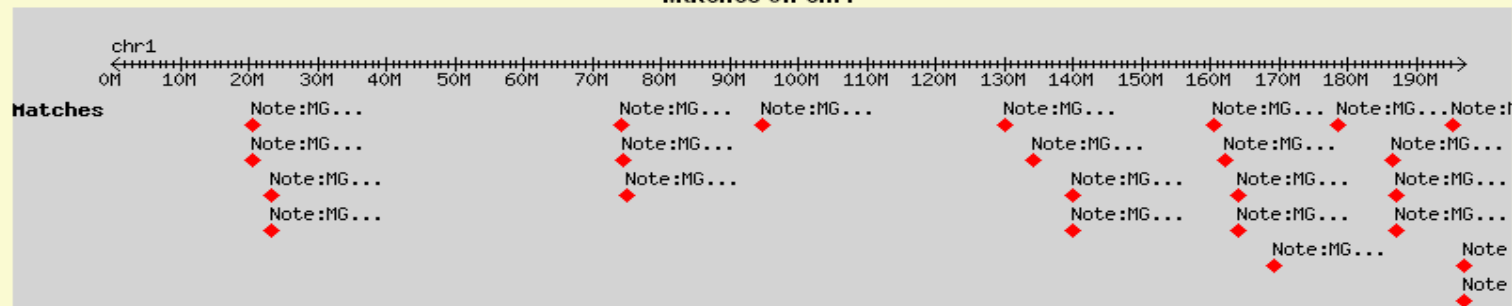
mir

Data Source

crispr

The following 555 regions match your request.

Matches on chr1



<a href="#">Note:MGI:2676881</a>	<a href="#">Mir206</a>	<a href="#">chr1:20.67..20.67 Mbp (73 bp)</a>	score=10
<a href="#">Note:MGI:3618720</a>	<a href="#">Mir133b</a>	<a href="#">chr1:20.67..20.67 Mbp (119 bp)</a>	score=10
<a href="#">Note:MGI:2676907</a>	<a href="#">Mir30a</a>	<a href="#">chr1:23.28..23.28 Mbp (71 bp)</a>	score=10
<a href="#">Note:MGI:3619048</a>	<a href="#">Mir30c-2</a>	<a href="#">chr1:23.3..23.3 Mbp (84 bp)</a>	score=10
<a href="#">Note:MGI:3836960</a>	<a href="#">Mir192b</a>	<a href="#">chr1:74.25..74.25 Mbp (67 bp)</a>	score=10
<a href="#">Note:MGI:2676901</a>	<a href="#">Mir26b</a>	<a href="#">chr1:74.44..74.44 Mbp (85 bp)</a>	score=10
<a href="#">Note:MGI:3619376</a>	<a href="#">Mir375</a>	<a href="#">chr1:74.95..74.95 Mbp (64 bp)</a>	score=10
<a href="#">Note:MGI:2676834</a>	<a href="#">Mir149</a>	<a href="#">chr1:94.75..94.75 Mbp (66 bp)</a>	score=10
<a href="#">Note:MGI:2676813</a>	<a href="#">Mir128-1</a>	<a href="#">chr1:130.1..130.1 Mbp (70 bp)</a>	score=10
<a href="#">Note:MGI:3618732</a>	<a href="#">Mir135b</a>	<a href="#">chr1:134.1..134.1 Mbp (97 bp)</a>	score=10
<a href="#">Note:MGI:3629589</a>	<a href="#">Mir181a-1</a>	<a href="#">chr1:139.9..139.9 Mbp (87 bp)</a>	score=10
<a href="#">Note:MGI:3618735</a>	<a href="#">Mir181b-1</a>	<a href="#">chr1:139.9..139.9 Mbp (80 bp)</a>	score=10
<a href="#">Note:MGI:3629597</a>	<a href="#">Mir488</a>	<a href="#">chr1:160.4..160.4 Mbp (109 bp)</a>	score=10
<a href="#">Note:MGI:3836959</a>	<a href="#">Mir1927</a>	<a href="#">chr1:162.2..162.2 Mbp (109 bp)</a>	score=10
<a href="#">Note:MGI:3618742</a>	<a href="#">Mir199a-2</a>	<a href="#">chr1:164.1..164.1 Mbp (110 bp)</a>	score=10
<a href="#">Note:MGI:2676890</a>	<a href="#">Mir214</a>	<a href="#">chr1:164.2..164.2 Mbp (110 bp)</a>	score=10
<a href="#">Note:MGI:3629598</a>	<a href="#">Mir689-1</a>	<a href="#">chr1:169.3..169.3 Mbp (109 bp)</a>	score=10
<a href="#">Note:MGI:3619366</a>	<a href="#">Mir350</a>	<a href="#">chr1:178.7..178.7 Mbp (99 bp)</a>	score=10
<a href="#">Note:MGI:3837225</a>	<a href="#">Mir1981</a>	<a href="#">chr1:186.6..186.6 Mbp (82 bp)</a>	score=10
<a href="#">Note:MGI:2676858</a>	<a href="#">Mir194-1</a>	<a href="#">chr1:187.1..187.1 Mbp (67 bp)</a>	score=10
<a href="#">Note:MGI:2676891</a>	<a href="#">Mir215</a>	<a href="#">chr1:187.1..187.1 Mbp (112 bp)</a>	score=10
<a href="#">Note:MGI:2676880</a>	<a href="#">Mir205</a>	<a href="#">chr1:195.3..195.3 Mbp (68 bp)</a>	score=10
<a href="#">Note:MGI:3619047</a>	<a href="#">Mir29b-2</a>	<a href="#">chr1:196.9..196.9 Mbp (81 bp)</a>	score=10
<a href="#">Note:MGI:2676906</a>	<a href="#">Mir29c</a>	<a href="#">chr1:196.9..196.9 Mbp (88 bp)</a>	score=10

# Availability

- Currently through SVN for registered users of Pasteur project web-site ([projets.pasteur.fr](http://projets.pasteur.fr))
- Planed: open source of workflows and nodes (sourceforge???)

# DEMO

- Tiling array analysis
- NGS data analysis