# Lightweight RNAseq analysis with BioConductor

A. Lesniewska[1,2]    M.J. Okoniewski[2]

[1]Institute of Computer Science
Poznan University of Technology, Poland

[2]Functional Genomics Center
UNI ETH Zurich, Switzerland

Bioconductor Developer Meeting Europe - 17-18. 11. 2010

# Outline

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Outline

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## State of the technology
### RNA-seq

- The coverage of `SOLID` starts to be enough to run whole transcriptomes RNAseq for higher species.
- 300-900M of reads per run
- Mapping is being constantly improved

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## State of the technology
### RNA-seq

- The coverage of `SOLID` starts to be enough to run whole transcriptomes RNAseq for higher species.
- 300-900M of reads per run
- Mapping is being constantly improved

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## State of the technology
### RNA-seq

- The coverage of `SOLID` starts to be enough to run whole transcriptomes RNAseq for higher species.
- 300-900M of reads per run
- Mapping is being constantly improved

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Our assumptions

- We can use `database storage`

- Recent improvements in DB engines allow fast access: `indexing, partitioning`

- R as the analysis environment – good statistics, comparison to microarrays

- BioConductor library as the way of publishing the analytical API

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Our assumptions

- We can use `database storage`
- Recent improvements in DB engines allow fast access: `indexing, partitioning`
- R as the analysis environment – good statistics, comparison to microarrays
- BioConductor library as the way of publishing the analytical API

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Our assumptions

- We can use `database storage`
- Recent improvements in DB engines allow fast access: `indexing, partitioning`
- R as the analysis environment – good statistics, comparison to microarrays
- BioConductor library as the way of publishing the analytical API

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Our assumptions

- We can use `database storage`
- Recent improvements in DB engines allow fast access: `indexing, partitioning`
- R as the analysis environment – good statistics, comparison to microarrays
- BioConductor library as the way of publishing the analytical API

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Our assumptions

- We can use `database storage`
- Recent improvements in DB engines allow fast access: `indexing, partitioning`
- R as the analysis environment – good statistics, comparison to microarrays
- BioConductor library as the way of publishing the analytical API

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Outline

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Exonmap paradigms

- Database for accessing the annotations
- Gene or a group at a time – not everything
- Translation of `genes<->transcripts<->exons`
- Filtering of interesting genes and exons
- Splicing analyses and plots

It worked => Let's do the same for RNAseq. . .

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Exonmap paradigms

- Database for accessing the annotations

- Gene or a group at a time – not everything

- Translation of `genes<->transcripts<->exons`

- Filtering of interesting genes and exons

- Splicing analyses and plots

It worked => Let's do the same for RNAseq. . .

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Exonmap paradigms

- Database for accessing the annotations
- Gene or a group at a time – not everything
- Translation of `genes<->transcripts<->exons`
- Filtering of interesting genes and exons
- Splicing analyses and plots

It worked => Let's do the same for RNAseq. . .

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Exonmap paradigms

- Database for accessing the annotations
- Gene or a group at a time – not everything
- Translation of `genes<->transcripts<->exons`
- Filtering of interesting genes and exons
- Splicing analyses and plots

It worked => Let's do the same for RNAseq. . .

Motivation
Contribution
Summary and future developments

State of the technology
**Exonmap paradigms**
Data Mining

## Exonmap paradigms

- Database for accessing the annotations
- Gene or a group at a time – not everything
- Translation of `genes<->transcripts<->exons`
- Filtering of interesting genes and exons
- Splicing analyses and plots

It worked => Let's do the same for RNAseq...

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Exonmap paradigms

- Database for accessing the annotations
- Gene or a group at a time – not everything
- Translation of `genes<->transcripts<->exons`
- Filtering of interesting genes and exons
- Splicing analyses and plots

It worked => Let's do the same for RNAseq. . .

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Exonmap paradigms

- Database for accessing the annotations
- Gene or a group at a time – not everything
- Translation of `genes<->transcripts<->exons`
- Filtering of interesting genes and exons
- Splicing analyses and plots

It worked => Let's do the same for RNAseq. . .

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Outline

Motivation
Contribution
Summary and future developments
State of the technology
Exonmap paradigms
Data Mining

# Lindell&Aumann window algorithm



Figure: algorithm & implementation

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Lindell&Aumann window algorithm

- Linear complexity
- Finds irreducible regions
- Applicable directly to coverage on genome data
- Follows biological intuitions
- Biological interpretation of consistent "exonic" region

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Lindell&Aumann window algorithm

- **Linear complexity**
- **Finds irreducible regions**
- Applicable directly to coverage on genome data
- Follows biological intuitions
- Biological interpretation of consistent "exonic" region

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Lindell&Aumann window algorithm

- Linear complexity
- Finds irreducible regions
- Applicable directly to coverage on genome data
- Follows biological intuitions
- Biological interpretation of consistent "exonic" region

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Lindell&Aumann window algorithm

- Linear complexity
- Finds irreducible regions
- Applicable directly to coverage on genome data
- Follows biological intuitions
- Biological interpretation of consistent "exonic" region

Motivation
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

## Lindell&Aumann window algorithm

- Linear complexity
- Finds irreducible regions
- Applicable directly to coverage on genome data
- Follows biological intuitions
- Biological interpretation of consistent "exonic" region

**Motivation**
Contribution
Summary and future developments

State of the technology
Exonmap paradigms
Data Mining

# Irreducible region



Jumping over local „holes"

Minsup=3

Getting rid of narrow peaks

Mincov=2

00222414012260000051

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Outline

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# How it works?



Figure: The flow of RNA seq data processing in the xmapcore
database and the rnaSeqMap library.

Motivation
**Contribution**
Summary and future developments

Schema of the library
**Processing**
Analysis pipelines

# Outline

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Data cleaning and preparation

- Libraries prepared and sequenced
- Raw data files transferred
- Colorspace reads mapped
- Samtools
- AWK script to get the simple, but biiiiig tables
- Import into MySQL

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Data cleaning and preparation

- Libraries prepared and sequenced
- Raw data files transferred
- Colorspace reads mapped
- Samtools
- AWK script to get the simple, but biiiiig tables
- Import into MySQL

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Data cleaning and preparation

- Libraries prepared and sequenced
- Raw data files transferred
- Colorspace reads mapped
- Samtools
- AWK script to get the simple, but biiiiig tables
- Import into MySQL

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Data cleaning and preparation

- Libraries prepared and sequenced
- Raw data files transferred
- Colorspace reads mapped
- Samtools
- AWK script to get the simple, but biiiig tables
- Import into MySQL

Motivation    Schema of the library
Contribution    Processing
Summary and future developments    Analysis pipelines

# Data cleaning and preparation

- Libraries prepared and sequenced
- Raw data files transferred
- Colorspace reads mapped
- Samtools
- AWK script to get the simple, but biiiiig tables
- Import into MySQL

Motivation
Contribution
Summary and future developments
Schema of the library
Processing
Analysis pipelines

## Data cleaning and preparation

- Libraries prepared and sequenced
- Raw data files transferred
- Colorspace reads mapped
- Samtools
- AWK script to get the simple, but biiiiig tables
- Import into MySQL

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

- MySQL >= 5.1
- Xmapcore database (denormalized Ensembl)
- Seq_reads table – with experiment number and genome coordinates of each read
- Indexed
- Partitioned into chromosome
- Average genome range query: 30s laptop, 5s fgcz-s-024

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

- MySQL >= 5.1

- Xmapcore database (denormalized Ensembl)

- Seq_reads table – with experiment number and genome coordinates of each read

- Indexed

- Partitioned into chromosome

- Average genome range query: 30s laptop, 5s fgcz-s-024

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

- MySQL >= 5.1
- Xmapcore database (denormalized Ensembl)
- Seq_reads table – with experiment number and genome coordinates of each read
- Indexed
- Partitioned into chromosome
- Average genome range query: 30s laptop, 5s fgcz-s-024

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

- MySQL >= 5.1
- Xmapcore database (denormalized Ensembl)
- Seq_reads table – with experiment number and genome coordinates of each read
- Indexed
- Partitioned into chromosome
- Average genome range query: 30s laptop, 5s fgcz-s-024

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

- MySQL >= 5.1
- Xmapcore database (denormalized Ensembl)
- Seq_reads table – with experiment number and genome coordinates of each read
- Indexed
- Partitioned into chromosome
- Average genome range query: 30s laptop, 5s fgcz-s-024

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

- MySQL >= 5.1
- Xmapcore database (denormalized Ensembl)
- Seq_reads table – with experiment number and genome coordinates of each read
- Indexed
- Partitioned into chromosome
- Average genome range query: 30s laptop, 5s fgcz-s-024

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

Databases:

- xmapcore
- or basic (3 tables gene,trenscript,exon) in xmapcore-like format
- maybe easily produced from non-Ensembl annotation for rare-species

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

Databases:

- xmapcore
- or basic (3 tables gene,trenscript,exon) in xmapcore-like format
- maybe easily produced from non-Ensembl annotation for rare-species

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Database back-end

Databases:

- xmapcore
- or basic (3 tables gene,trenscript,exon) in xmapcore-like format
- maybe easily produced from non-Ensembl annotation for rare-species

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Stored procedures

- Region reads in given sample
- Gene <-> Transcript <-> Exon <-> reads
- Genes on a chromosome
- Intergenic regions on a chromosome

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Stored procedures

- Region reads in given sample
- `Gene <-> Transcript <-> Exon <-> reads`
- Genes on a chromosome
- Intergenic regions on a chromosome

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Stored procedures

- Region reads in given sample
- `Gene <-> Transcript <-> Exon <-> reads`
- Genes on a chromosome
- Intergenic regions on a chromosome

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Stored procedures

- Region reads in given sample
- `Gene <-> Transcript <-> Exon <-> reads`
- Genes on a chromosome
- Intergenic regions on a chromosome

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Classes in R

- SeqReads – a collection of reads for samples in a given genomic region
- NucleotideDistribution (S3 class) – nucleotide by nucleotide distribution of measured feature
    - Coverage of reads
    - Fold change
    - Splicing Index
    - Significant regions

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Classes in R

- SeqReads – a collection of reads for samples in a given genomic region
- NucleotideDistribution (S3 class) – nucleotide by nucleotide distribution of measured feature
    - Coverage of reads
    - Fold change
    - Splicing Index
    - Significant regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Classes in R

- SeqReads – a collection of reads for samples in a given genomic region
- NucleotideDistribution (S3 class) – nucleotide by nucleotide distribution of measured feature
    - Coverage of reads
    - Fold change
    - Splicing Index
    - Significant regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Classes in R

- SeqReads – a collection of reads for samples in a given genomic region
- NucleotideDistribution (S3 class) – nucleotide by nucleotide distribution of measured feature
    - Coverage of reads
    - Fold change
    - Splicing Index
    - Significant regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Classes in R

- SeqReads – a collection of reads for samples in a given genomic region
- NucleotideDistribution (S3 class) – nucleotide by nucleotide distribution of measured feature
    - Coverage of reads
    - Fold change
    - Splicing Index
    - Significant regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Interesting genes

- Good coverage
- Good coverage of exons
- Interesting splicing index
- Interesting new regions – novel exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Interesting genes

- Good coverage

- Good coverage of exons

- Interesting splicing index

- Interesting new regions – novel exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Interesting genes

- Good coverage
- Good coverage of exons
- Interesting splicing index
- Interesting new regions – novel exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Interesting genes

- Good coverage
- Good coverage of exons
- Interesting splicing index
- Interesting new regions – novel exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Interesting genes

- Good coverage
- Good coverage of exons
- Interesting splicing index
- Interesting new regions – novel exons

More algorithms to establish within the framework!!

Motivation

Contribution

Summary and future developments

Schema of the library

Processing

Analysis pipelines

# Interesting intergenic regions

- Irreducible regions with good coverage

- We treat them as novel genes and run gene-style analysis

- Looking for exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments
Schema of the library
Processing
Analysis pipelines

# Interesting intergenic regions

- Irreducible regions with good coverage
- We treat them as novel genes and run gene-style analysis
- Looking for exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Interesting intergenic regions

- Irreducible regions with good coverage
- We treat them as novel genes and run gene-style analysis
- Looking for exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments
Schema of the library
Processing
Analysis pipelines

## Interesting intergenic regions

- Irreducible regions with good coverage
- We treat them as novel genes and run gene-style analysis
- Looking for exons

More algorithms to establish within the framework!!

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Output

- Iranges objects – for interesting regions
- DESeq object – gene/exon level expression - for the significance analysis with DESeq
- Lists of interesting features

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Output

- Iranges objects – for interesting regions
- DESeq object – gene/exon level expression - for the significance analysis with DESeq
- Lists of interesting features

Motivation
Contribution
Summary and future developments

Schema of the library
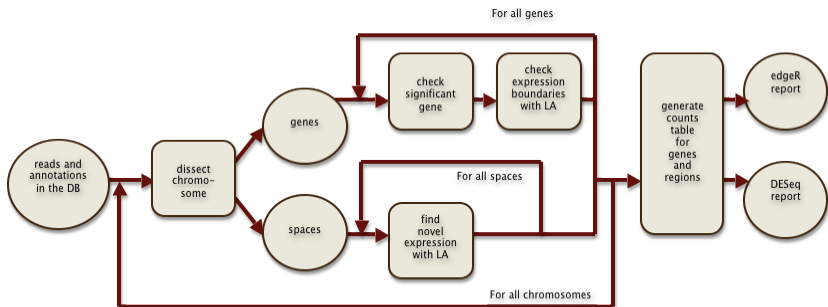Processing
Analysis pipelines

## Output

- Iranges objects – for interesting regions
- DESeq object – gene/exon level expression - for the significance analysis with DESeq
- Lists of interesting features

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
**Analysis pipelines**

# Outline

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
**Analysis pipelines**

# An example of rnaSeqMap analysis pipeline

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Analysis pipelines

- Get all the genes from a chromosome
  - Check for interesting features
  - Check possible gene extensions – expression closely around the gene
- Get all the intergenic regions on chromosome
  - Find novel expressed regions
  - Describe the regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Analysis pipelines

- Get all the genes from a chromosome
    - Check for interesting features
        - Check possible gene extensions – expression closely around the gene
- Get all the intergenic regions on chromosome
    - Find novel expressed regions
    - Describe the regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Analysis pipelines

- Get all the genes from a chromosome
  - Check for interesting features
  - Check possible gene extensions – expression closely around the gene
- Get all the intergenic regions on chromosome
  - Find novel expressed regions
  - Describe the regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Analysis pipelines

- Get all the genes from a chromosome
    - Check for interesting features
    - Check possible gene extensions – expression closely around the gene
- Get all the intergenic regions on chromosome
    - Find novel expressed regions
    - Describe the regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Analysis pipelines

- Get all the genes from a chromosome
    - Check for interesting features
    - Check possible gene extensions – expression closely around the gene
- Get all the intergenic regions on chromosome
    - Find novel expressed regions
    - Describe the regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Analysis pipelines

- Get all the genes from a chromosome
  - Check for interesting features
  - Check possible gene extensions – expression closely around the gene
- Get all the intergenic regions on chromosome
  - Find novel expressed regions
  - Describe the regions

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Analysis pipelines - code

```
test.gene<-function(g,exps,nsums,mi,ms)
{
 rs <- newSeqReadsFromGene(g)
 rs <- addExperimentsToReadset(rs,exps)
 nd.cov <- getCoverageFromRS(rs,exps)
 nd.cov <- normalizeBySum(nd.cov, nsums)
 nd.reg <- findRegionsAsND(nd.cov,as.int(mi),ms=ms)
 ir.reg <- findRegionsAsIR(nd.cov,as.int(mi),ms=ms)
 cat("region search algorithm...\n")
 out <- g
 out <- c(out, apply(distribs(nd.cov),2,max))
 out <- c(out, apply(distribs(nd.cov),2,mean))
 out <- c(out, apply(distribs(nd.reg),2,max))
}
```

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
**Analysis pipelines**

# Analysis pipelines - code

```
test.space<-function(exps,ch,st,en,str,nsums,mi,ms)
{
g.ch <- rnaSeqMap:::.chromosome.number(ch)
rs <- newSeqReads(g.ch,st,en,str)
rs <- addExperimentsToReadset(rs,exps)
nd.cov <- getCoverageFromRS(rs,exps
nd.cov <- normalizeBySum(nd.cov, nsums)
nd.reg <- findRegionsAsND(nd.cov,as.int(mi),ms=ms)
out <- c(ch,st, en, str)
out <- c(out, apply(distribs(nd.cov),2,max))
out <- c(out, apply(distribs(nd.cov),2,mean))
out <- c(out, apply(distribs(nd.reg),2,max))
}
```

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
**Analysis pipelines**

# Analysis pipelines - code

```
my.genes<-geneInChromosome(22, 200000, 204000,1)
my.spaces<-spaceInChromosome(22, 200000, 204000,1)
    interesting.genes <- NULL
for (i in 1:length(my.genes))
{    cat ("Running gene ", i , "----------\n")
    interesting.genes <- rbind(interesting.genes,
      test.gene(my.genes[i], 1:6, nsums))}
   interesting.spaces <- NULL
for (i in 1:(dim(my.spaces))[1])
{    cat ("Running space ", i , "----------\n")
    interesting.spaces <- rbind(interesting.spaces,
      test.space(1:2, 22,my.spaces[i,1],
      my.spaces[i,2],my.spaces[i,3] ))}
```

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Advantages of rnaSeqMap

- Complex analysis of huge data on a small machine - awk, MySQL, R do not have big requirements
- Flexible and fine-grained approach to transcriptomics
  - Not a single nucleotide can hide, if it is expressed
  - Flexible boundaries of expression regions – we rely on Ensembl, but do not have to trust it blindly

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Advantages of rnaSeqMap

- Complex analysis of huge data on a small machine - awk, MySQL, R do not have big requirements
- Flexible and fine-grained approach to transcriptomics
    - Not a single nucleotide can hide, if it is expressed
    - Flexible boundaries of expression regions – we rely on Ensembl, but do not have to trust it blindly

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

# Advantages of rnaSeqMap

- Complex analysis of huge data on a small machine - awk, MySQL, R do not have big requirements
- Flexible and fine-grained approach to transcriptomics
  - Not a single nucleotide can hide, if it is expressed
  - Flexible boundaries of expression regions – we rely on Ensembl, but do not have to trust it blindly

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Challenges

- Size and allocation of RAM memory to run big regions – we have to run one chromosome at a time
- Speed of queries for reads data – not bad now
- Speed of analysis – optimized by rewriting in C
- Installation is not simple – but still simpler than many other systems

Motivation
**Contribution**
Summary and future developments

Schema of the library
Processing
**Analysis pipelines**

## Challenges

- Size and allocation of RAM memory to run big regions – we have to run one chromosome at a time
- Speed of queries for reads data – not bad now
- Speed of analysis – optimized by rewriting in C
- Installation is not simple – but still simpler than many other systems

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Challenges

- Size and allocation of RAM memory to run big regions – we have to run one chromosome at a time
- Speed of queries for reads data – not bad now
- Speed of analysis – optimized by rewriting in C
- Installation is not simple – but still simpler than many other systems

Motivation
Contribution
Summary and future developments

Schema of the library
Processing
Analysis pipelines

## Challenges

- Size and allocation of RAM memory to run big regions – we have to run one chromosome at a time
- Speed of queries for reads data – not bad now
- Speed of analysis – optimized by rewriting in C
- Installation is not simple – but still simpler than many other systems

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Outline

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

## Numeric results

- In total of 38546 genes and pseudogenes, there are:
  - 6863 genes with expression regions >10 for all 6 patients
  - 24172 genes with expression >10 at least for one patient
  - 14375 genes with no irreducible regions >10 in any patient
  - 9912 genes with at least 100 reads mapped in total in 6 samples
  - 5822 genes with no reads at all

Similar to detection on microarrays, however coverage is still too low to detect splicing in most cases. . .

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

## Numeric results

- In total of 38546 genes and pseudogenes, there are:
  - 6863 genes with expression regions >10 for all 6 patients
  - 24172 genes with expression >10 at least for one patient
  - 14375 genes with no irreducible regions >10 in any patient
  - 9912 genes with at least 100 reads mapped in total in 6 samples
  - 5822 genes with no reads at all

Similar to detection on microarrays, however coverage is still too low to detect splicing in most cases. . .

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

## Numeric results

- In total of 38546 genes and pseudogenes, there are:
    - 6863 genes with expression regions >10 for all 6 patients
    - 24172 genes with expression >10 at least for one patient
    - 14375 genes with no irreducible regions >10 in any patient
    - 9912 genes with at least 100 reads mapped in total in 6 samples
    - 5822 genes with no reads at all

Similar to detection on microarrays, however coverage is still too low to detect splicing in most cases. . .

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

## Numeric results

- In total of 38546 genes and pseudogenes, there are:
    - 6863 genes with expression regions >10 for all 6 patients
    - 24172 genes with expression >10 at least for one patient
    - 14375 genes with no irreducible regions >10 in any patient
    - 9912 genes with at least 100 reads mapped in total in 6 samples
    - 5822 genes with no reads at all

Similar to detection on microarrays, however coverage is still too low to detect splicing in most cases. . .

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
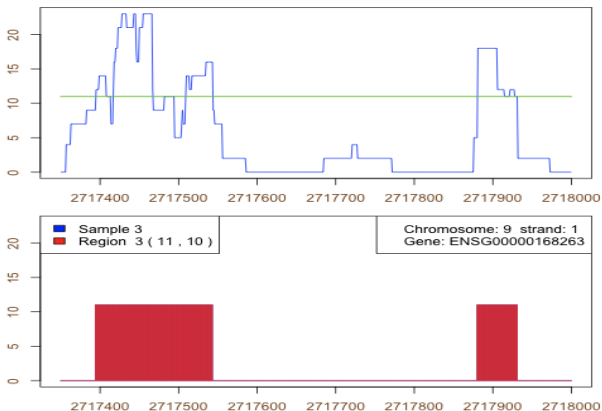Splicing index

## Numeric results

- In total of 38546 genes and pseudogenes, there are:
  - 6863 genes with expression regions >10 for all 6 patients
  - 24172 genes with expression >10 at least for one patient
  - 14375 genes with no irreducible regions >10 in any patient
  - 9912 genes with at least 100 reads mapped in total in 6 samples
  - 5822 genes with no reads at all

Similar to detection on microarrays, however coverage is still too low to detect splicing in most cases. . .
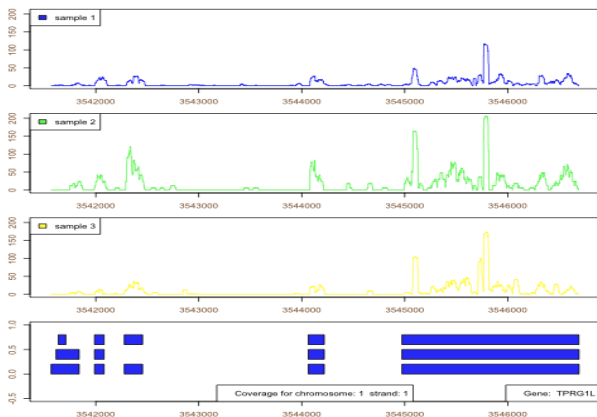
Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Outline

## 1 Motivation

- State of the technology
- Exonmap paradigms
- Data Mining

## 2 Contribution

- Schema of the library
- Processing
- Analysis pipelines

## 3 Summary and future developments

- Numeric results
- Examplary plots
- Splicing index

Motivation
Contribution
Summary and future developments

Numeric results
**Examplary plots**
Splicing index

# Irreducible regions of coverage

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Examplary plot

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Examplary plot

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Outline

1 **Motivation**
   - State of the technology
   - Exonmap paradigms
   - Data Mining

2 **Contribution**
   - Schema of the library
   - Processing
   - Analysis pipelines

3 **Summary and future developments**
   - Numeric results
   - Examplary plots
   - Splicing index
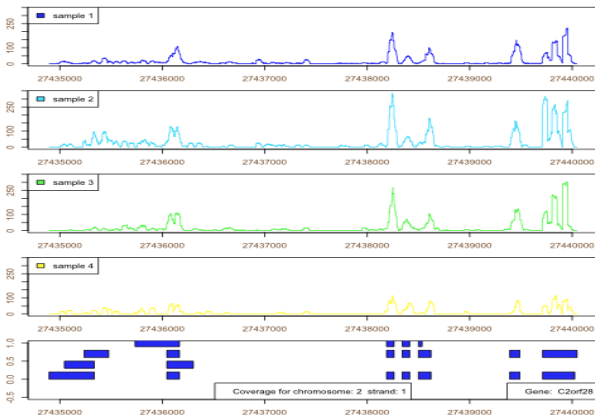
Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Splicing indeks

- Similar to original in Gardina et al.
- Normalized to +/- 1
- Calculated on each nucleotide

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Splicing indeks

- Similar to original in Gardina et al.
- Normalized to +/- 1
- Calculated on each nucleotide

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

## Splicing indeks

- Similar to original in Gardina et al.
- Normalized to +/- 1
- Calculated on each nucleotide

Motivation
Contribution
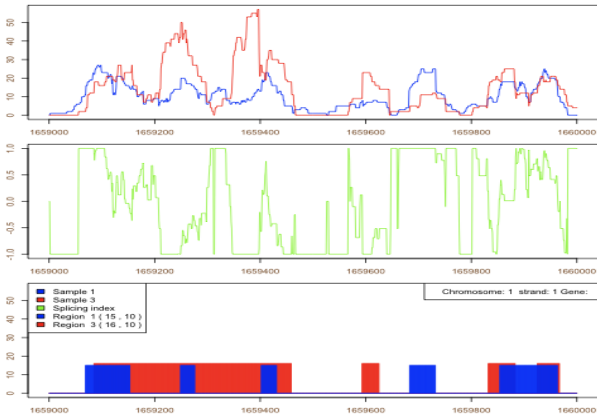Summary and future developments

Numeric results
Examplary plots
Splicing index

## Splicing index

$$
SI(n) = \begin{cases}
0, \text{ if} & (E_{1n} = 0 \land E_{2n} = 0) \\
1, \text{ if} & (E_{1n} = 0 \land E_{2n} = 0) \quad \lor \left( \frac{E_{1n}}{G_{1n}} \cdot \frac{E_{2n}}{G_{2n}} > 2 \right) \\
-1, \text{ if} & (E_{1n} = 0 \land E_{2n} = 0) \quad \lor \left( \frac{E_{1n}}{G_{1n}} \cdot \frac{E_{2n}}{G_{2n}} < 0.5 \right) \\
log_2 \left( \frac{E_{1n}}{G_{1n}} \cdot \frac{E_{2n}}{G_{2n}} \right) & \text{in all other cases}
\end{cases}
$$

Where $E_{1n}$ and $E_{2n}$ are the coverage values for a given nucleotide, while $G_{1n}$ and $G_{2n}$ are the counts of reads in the region or gene.

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Splicing index

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

## Future developments

- exon/isoform discovery
- paired end reads
- new splicing index forms
- parallel execution with snow, multicore,. . .
- . . . etc

Motivation
Contribution
Summary and future developments

Numeric results
Examplary plots
Splicing index

# Summary

- The library rnaSeqMap in Bioconductor 2.7

- ...

- Have fun!!!

# For Further Reading I

📄 Aumann Y, Lindell Y:
J. Intell. Inf. Syst. 2003, **20**(3):255–283.

📄 Gardina et al.:
BMC Genomics 2006, 7:325.

📄 Yates T, Okoniewski MJ, Miller CJ
Nucleic Acids Research 2008, **36(suppl 1)**:D780–D786.

📄 Okoniewski M, Yates T, Dibben S, Miller C
Genome Biology 2007, **8(5)**:R79.

# Acknowledgements

- Functional Genomic Center Zurich:
  Marzanna Künzli-Gontarczyk
  Sirisha Aluri
  Weihong Qi
  Hubert Rehrauer
  Tanguy Le Carrour
  Remy Bruggmann
  Hansruedi Baetschmann
- Kinderspital Zürich:
  Beat Scheaffer
  Marco Wachtel
- Institute of Molecular Systems Biology, ETH:
  Lucia Bautista Borrego
- PICR Manchester:
  Tim Yates