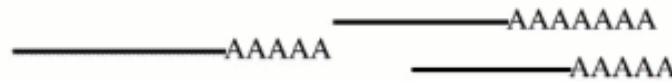# Methods for Analysis of RNA-Seq data.
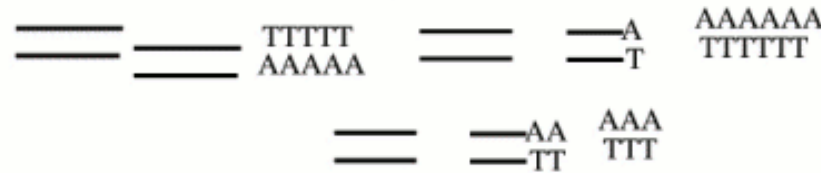
Hugues Richard

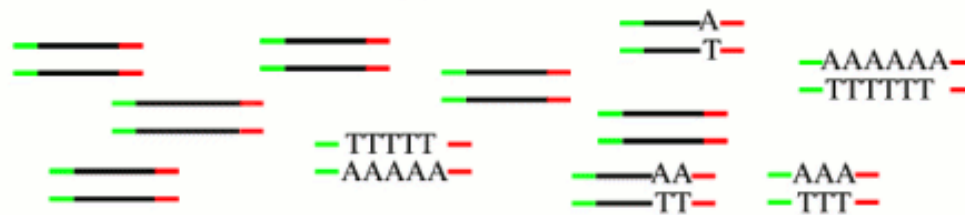# RNA-Seq protocol (no strand information)



extraction of poly-A RNAs

conversion into ds-cDNA and shearing
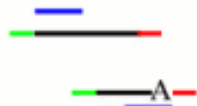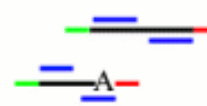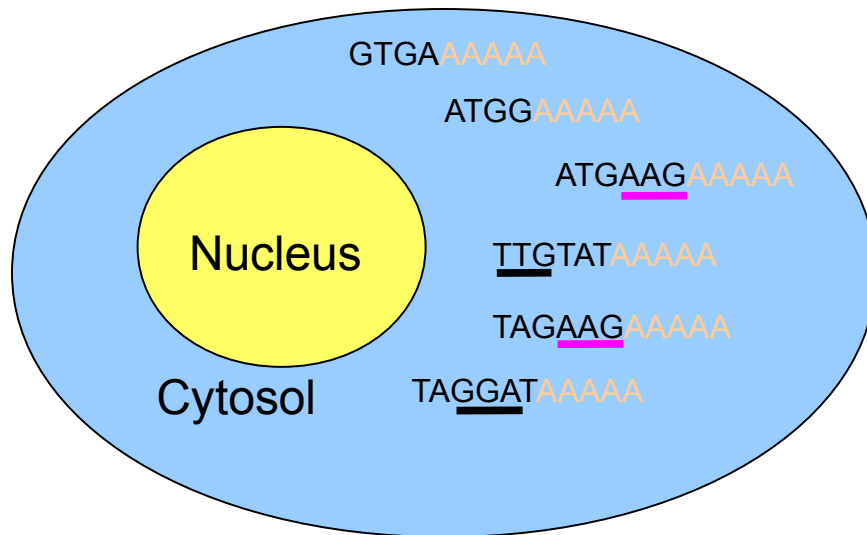
amplification and adapter ligation

sequencing

single end (SET)

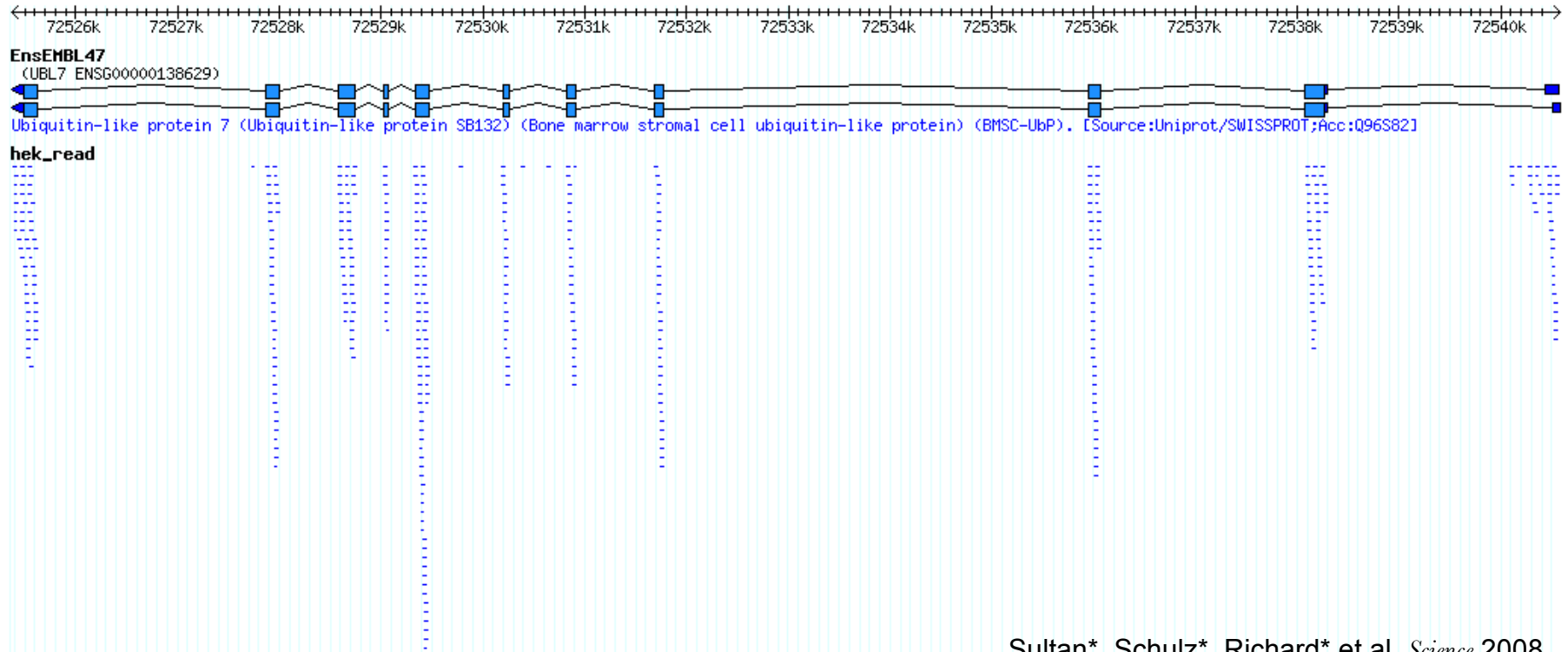paired-end (PET)

M. Schulz©

"bag of transcript positions"

TTG
GGA
AAG

Variability of the counts (sampling)
influenced by:
 _ region length
 _ copy number

# Mapping the reads



Aggregate counts on exons/genes

Normalize by number of possible hits (RPKM)

Sultan*, Schulz*, Richard* et al. *Science* 2008

Mortazavi*, Williams* et al. *Nat Methods* 2008

Wang*, Sandberg* et al. *Nature* 2008

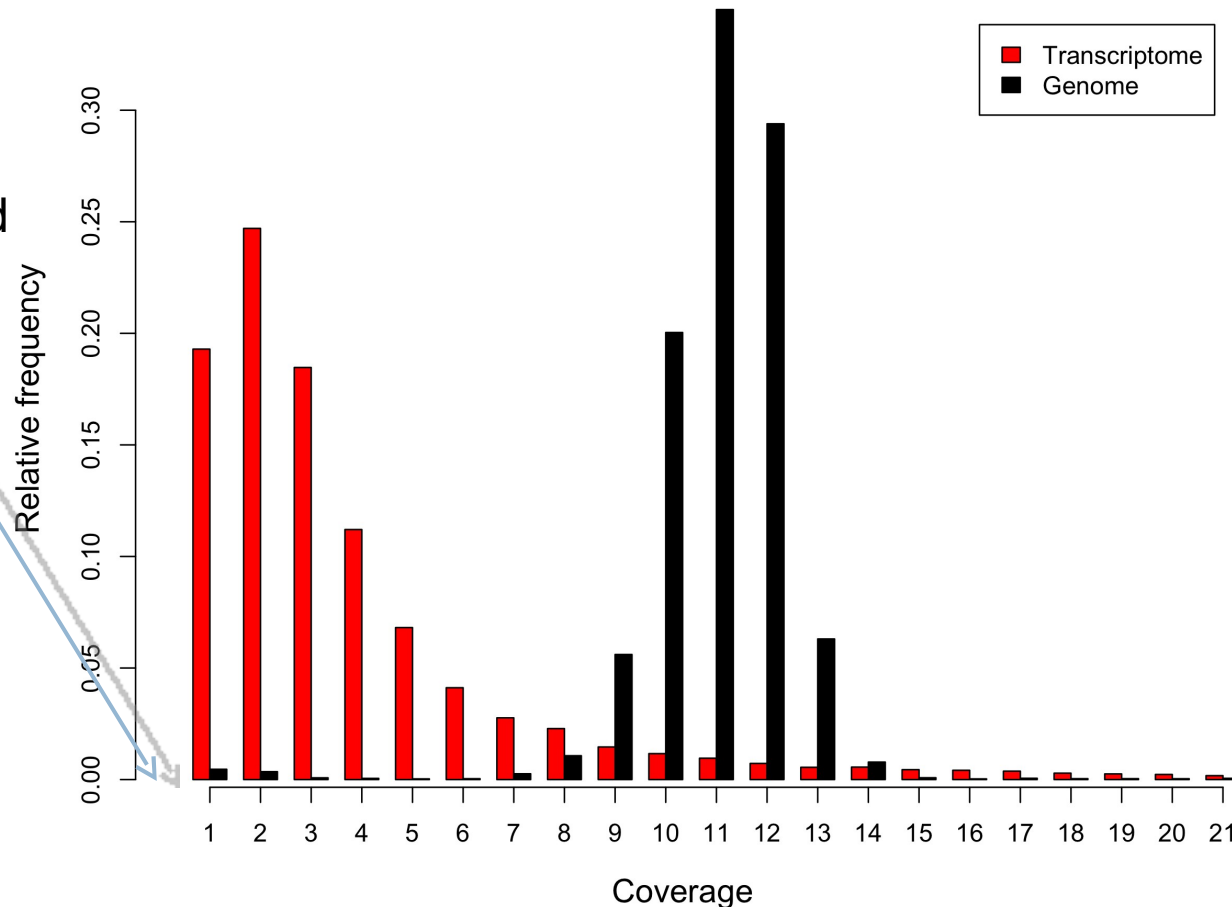Cloonan*, Forest*, Kolle* et al. *Nat Methods* 2008

# Outline

- 𝕰𝖘𝖙𝖎𝖒𝖆𝖙𝖎𝖓𝖌 𝖉𝖊𝖕𝖙𝖍 𝖔𝖋 𝖘𝖊𝖖𝖚𝖊𝖓𝖈𝖎𝖓𝖌.
  - Do we see all the transcriptional units ?

- Infering new events:
  - Detection of new transcriptional units.
  - Detection of Alternative Splicing Events.

- RGASP competition:
  - Transcriptome assembly with Oases (Schulz/Zerbino)
  - Reads remapping with RazerS (Weese)

MAX-PLANCK-GESELLSCHAFT

UPMC PARISUNIVERSITAS

# Transcriptome vs genome assembly

- RNA-Seq reads are distributed according to transcript expression levels.

# of
non observed
genes ?

# What is current coverage ?

- Total number of transcripts :

  - Count of a transcriptional unit $i$:

  $$X_i \sim \mathcal{P}(\lambda_i)$$

  - Unspecified counts:

  $$f(c; \mathbf{g}) = \int_0^\infty \exp^{-\lambda} \frac{\lambda^c}{c!} \, d\mathbf{g}(\lambda)$$
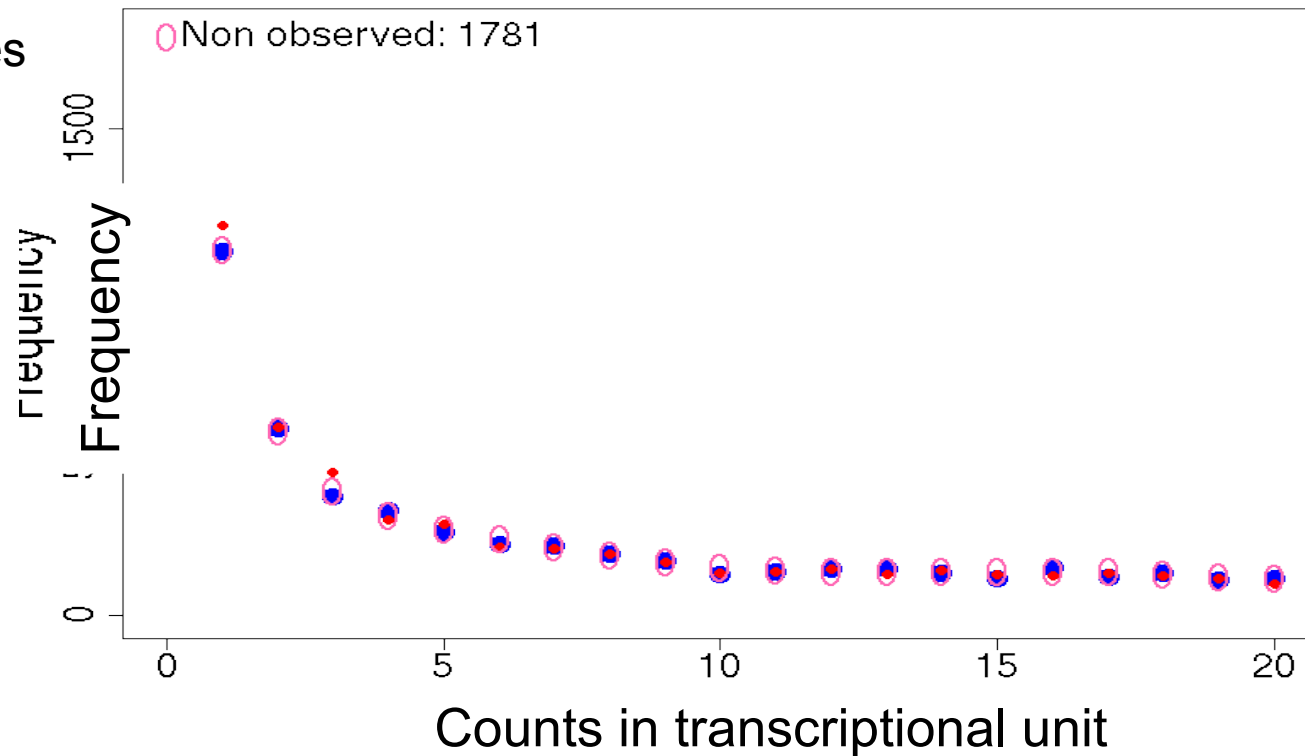
  - If we estimate $g$, then:

  $$\hat{N} = \left\langle \frac{\#\{\text{observed}\}}{1 - f(0, \hat{g})} \right\rangle$$

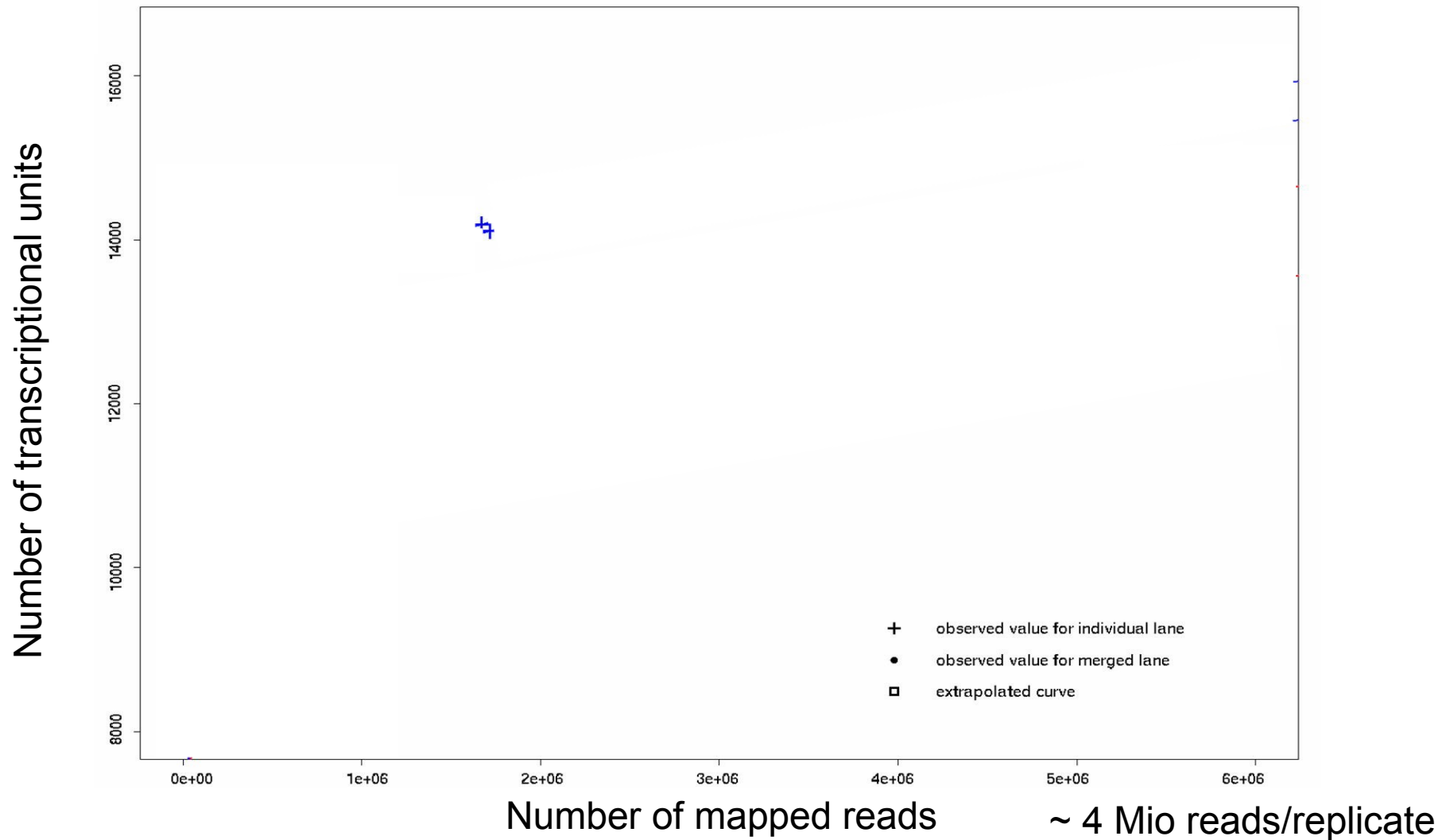- Expected number of new discoveries with more lanes

UPMC PARIS UNIVERSITAS

# Estimating count frequency law

2 lanes



Penalized Non Parametric Maximum Likelihood method   (Wang & Lindsay 05)

# Dynamic range



Number of transcriptional units

Number of mapped reads

~ 4 Mio reads/replicate

+ observed value for individual lane
• observed value for merged lane
□ extrapolated curve

# Outline

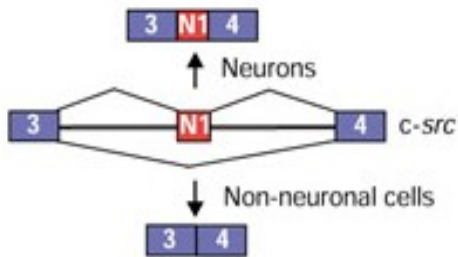- Estimating depth of sequencing.
  - do we see all the genes ?

- Infering new events:
  - Detection of new transcriptional units.
  - Alternative Splicing events detection.

- RGASP competition:
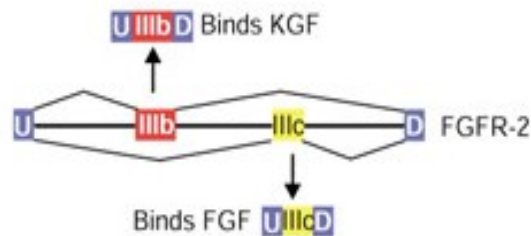  - Transcriptome assembly with Oases (Schulz/Zerbino)
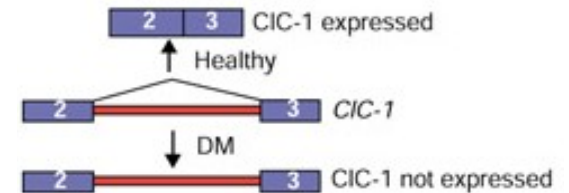  - Mapping with RazerS (Weese)

UPMC PARISUNIVERSITAS

# Alternative Exon Events (AEEs)



Ladd and Cooper 2002

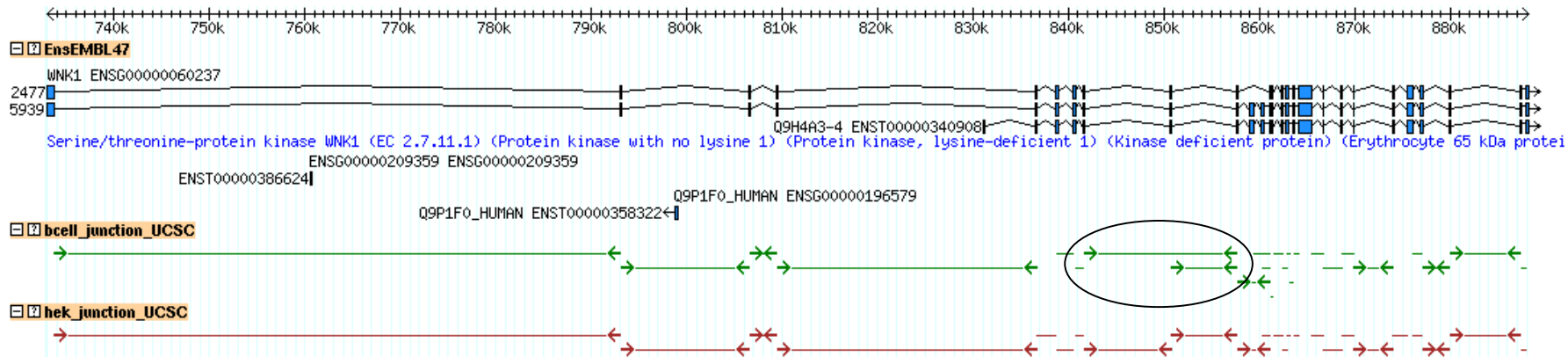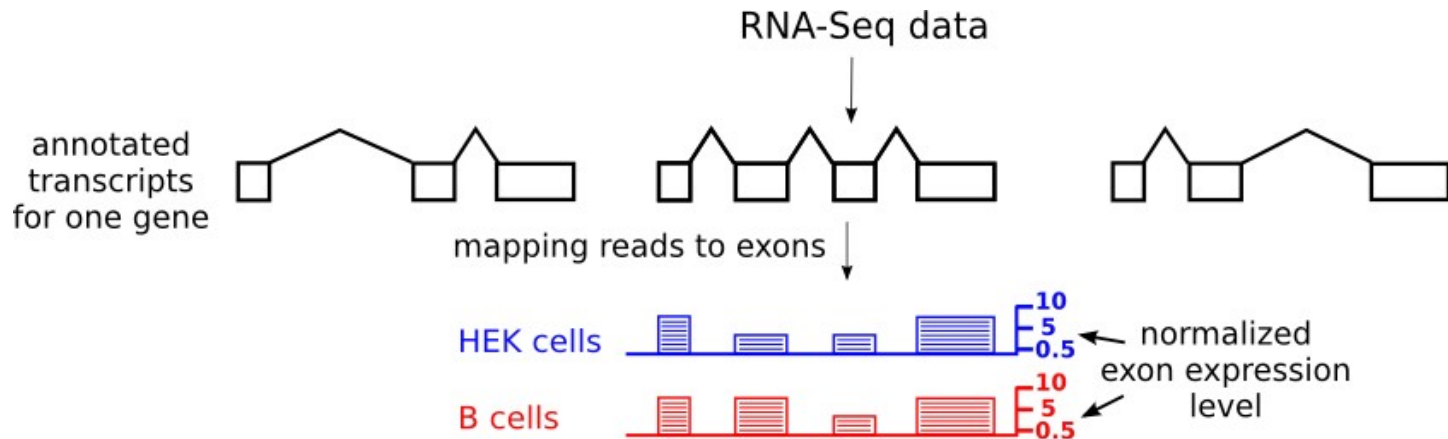# Splice junctions



Align unmatched reads to a set of artificially created junctions
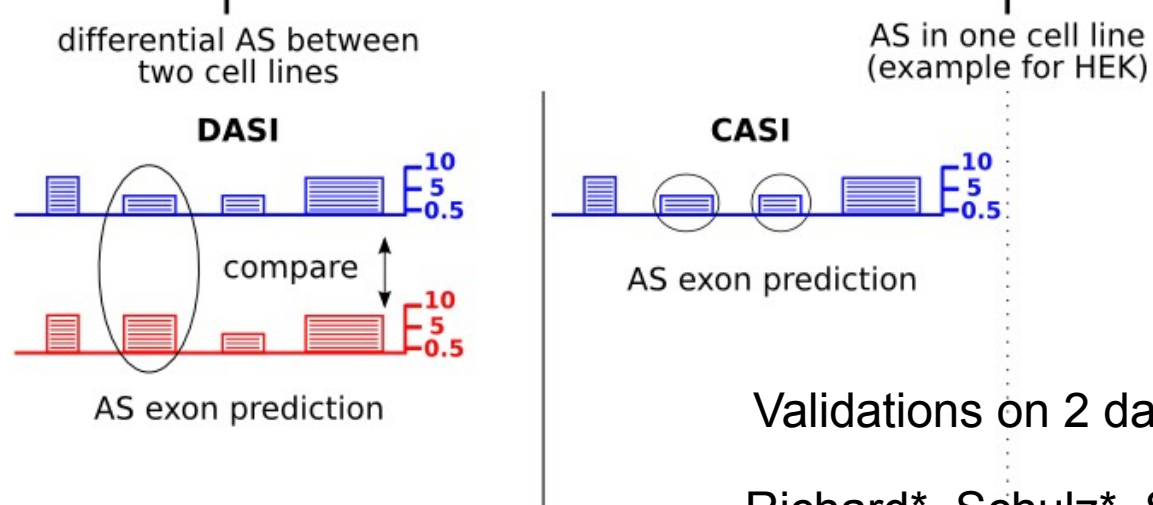
Use a splice junction aligner (Tophat, QPalma, GEMM)

# Detecting AEE from Exons Expression

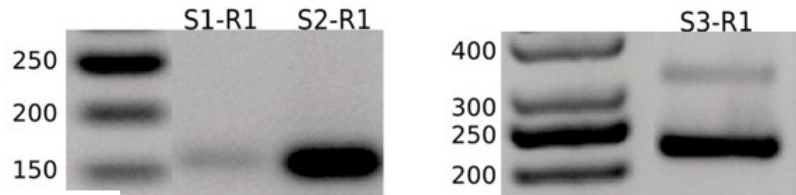

Validations on 2 datasets (HEK and B cells)

Richard*, Schulz*, Sultan* et al. *NAR* 2010

# Within cell AS : CASI



61 events tested
35 validated

# Robustness (simulation)



1000 simulations of alternative exon events
of one gene with 6 exons and 300 reads in total

# Robustness (bootstrap)



(500 repetitions)

Bootstrap:

_ Randomly alter exon boundaries

_ Monitor changes in prediction

# Alternative Polyadenylation in HIP2[+]



[+]Sandberg et al. Science, 2008

# Detecting AEE from Exons Expression
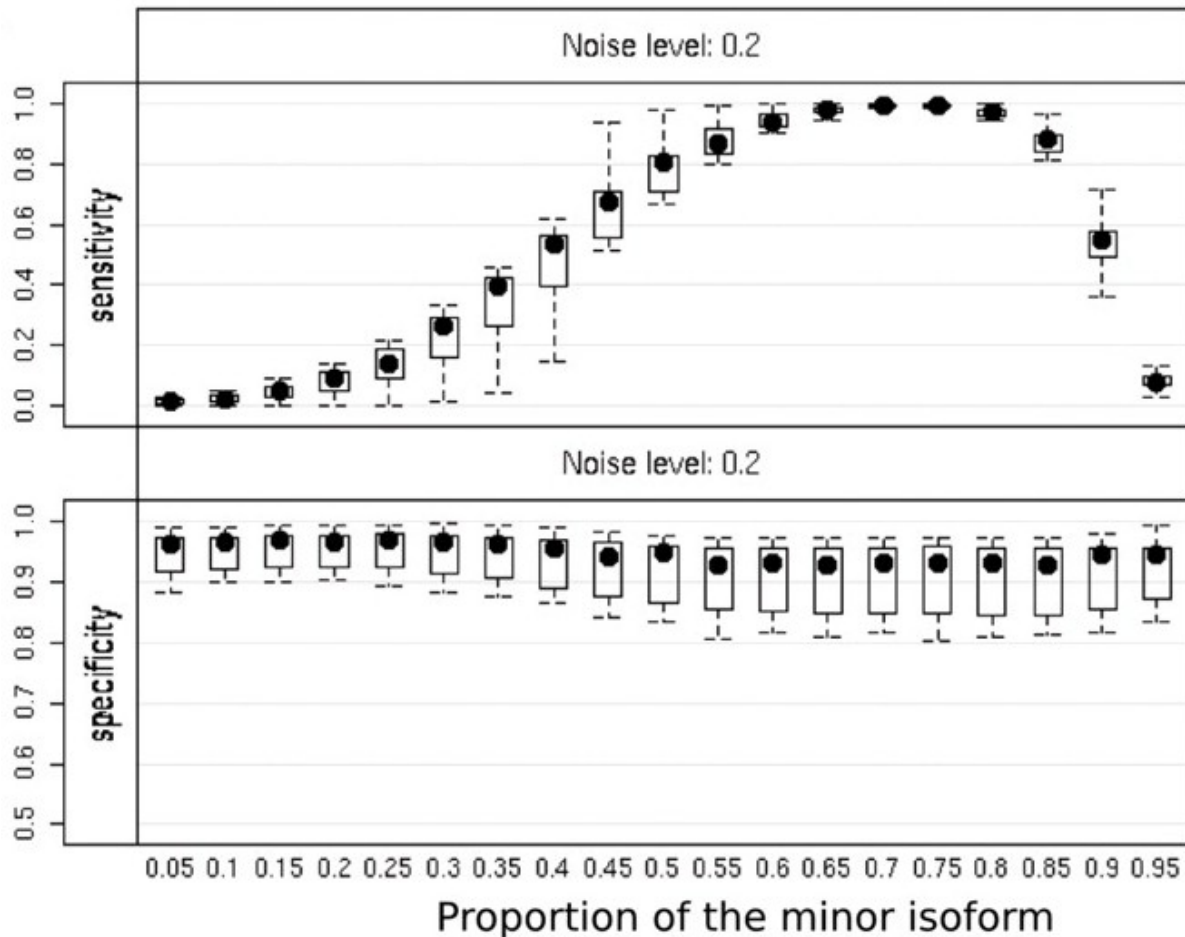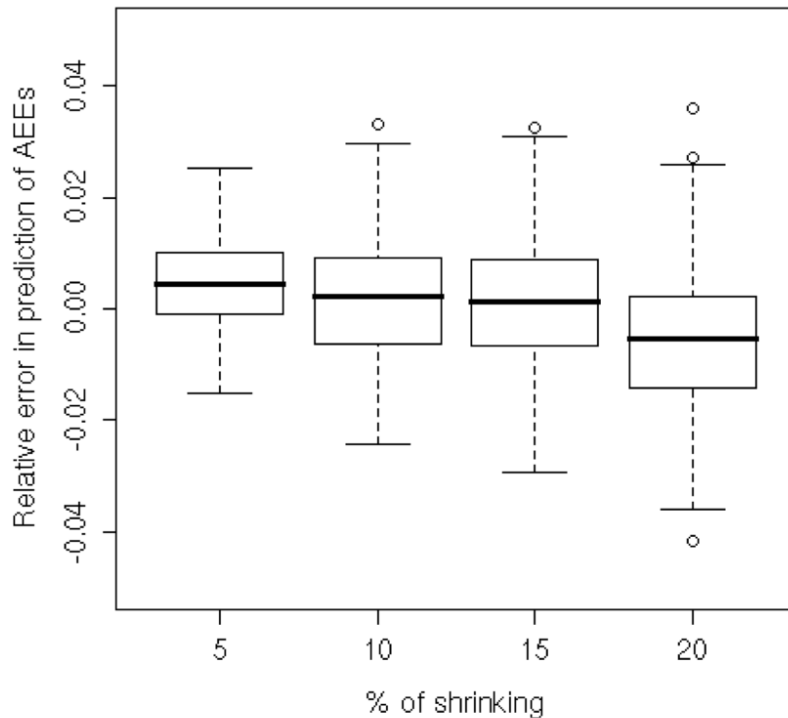
# Isoform quantification: POEM



Count for each isoform j

$$Z_j \sim \mathcal{P}(\lambda_j)$$

$$\lambda_j = \lambda \cdot s_j \cdot p_j$$

binning within every exon at random :

$$q_{i,j} = \frac{l_i}{s_j} \cdot I_{i,j}$$

We observe only the counts falling within each subexon

# Isoform quantification: POEM



Comparison to:
- qPCR               (47 events)
- Estimate from junction counts
                     (267 events)

grey - POEM
black - qPCR

# Validation



Junction reads

qRT-PCR

PCC:                    0.65                                    0.81

# Comparison to Exon arrays

_ produced 4 replicates Exon Array hybridizations for each cell line
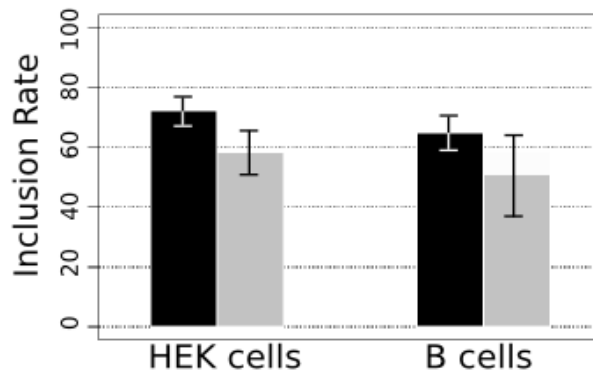_ based on ENSEMBL 25% more exona are detected by RNA-Seq



DASI
FDR = 0.05
613 genes

overlap
93 genes

MIDAS
FDR = 0.2
934 genes

|DASI| > 2

|SI| > 0.5

365 genes
968 exons

10 genes
16 exons

390 genes
512 exons

UPMC
PARIS UNIVERSITAS

# Outline

- Coverage for count based problems.
  - do we see all the genes ?

- Infering new events:
  - Detection of new transcriptional units.
  - Alternative Splicing events detection.

- RGASP competition:
  - Transcriptome assembly with Oases (Schulz/Zerbino)
  - Mapping with RazerS (Weese)

# Transcripts assembly

Imagine two transcripts:
TAGTCGAG   GCTT        TAGAGACAG
TAGTCGAG   TCCGA        TAGAGACAG

```
AGTCGAG CTTTAGA  CGATGAG CTTTAGA
  GTCGAGG  TTAGATC  ATGAGGC      GAGACAG
    GAGGCTC   GTCCGAT AGGCTTT GAGACAG
AGTCGAG     TAGATCC ATGAGGC  TAGAGAA
TAGTCGA  CTTTAGA CCGATGA      TTAGAGA
   CGAGGCT  AGATCCG TGAGGCT  AGAGACA
TAGTCGA GCTTTAG TCCGATG  GCTTTAG
  TCGATTGC     GATCCGA GAGGCTT AGAGACA
TAGTCGA     TTAGATC GATGAGG TTTAGAG
  GTCGAGG TCTAGAT    ATGAGGC  TAGAGAC
    AGGCTTT  GTCCGAT AGGCTTT GAGACAG
AGTCGAG    TTAGATA ATGAGGC     AGAGACA
    GGCTTTA  TCCGATG     TTTAGAG
    CGAGGCT TAGATCC  TGAGGCT     GAGACAG
AGTCGAG   TTTAGATC  ATGAGGC TTAGAGA
```
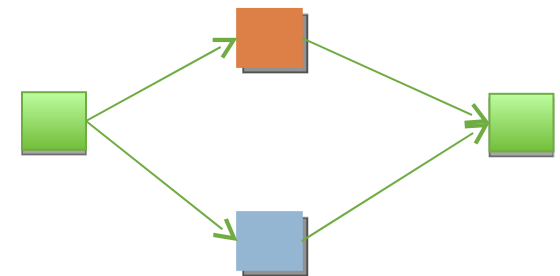
Assemble reads
into contigs



Oases:
  _De Bruijn graph
  _Velvet framework

Schulz, Zerbino et al (in preparation)
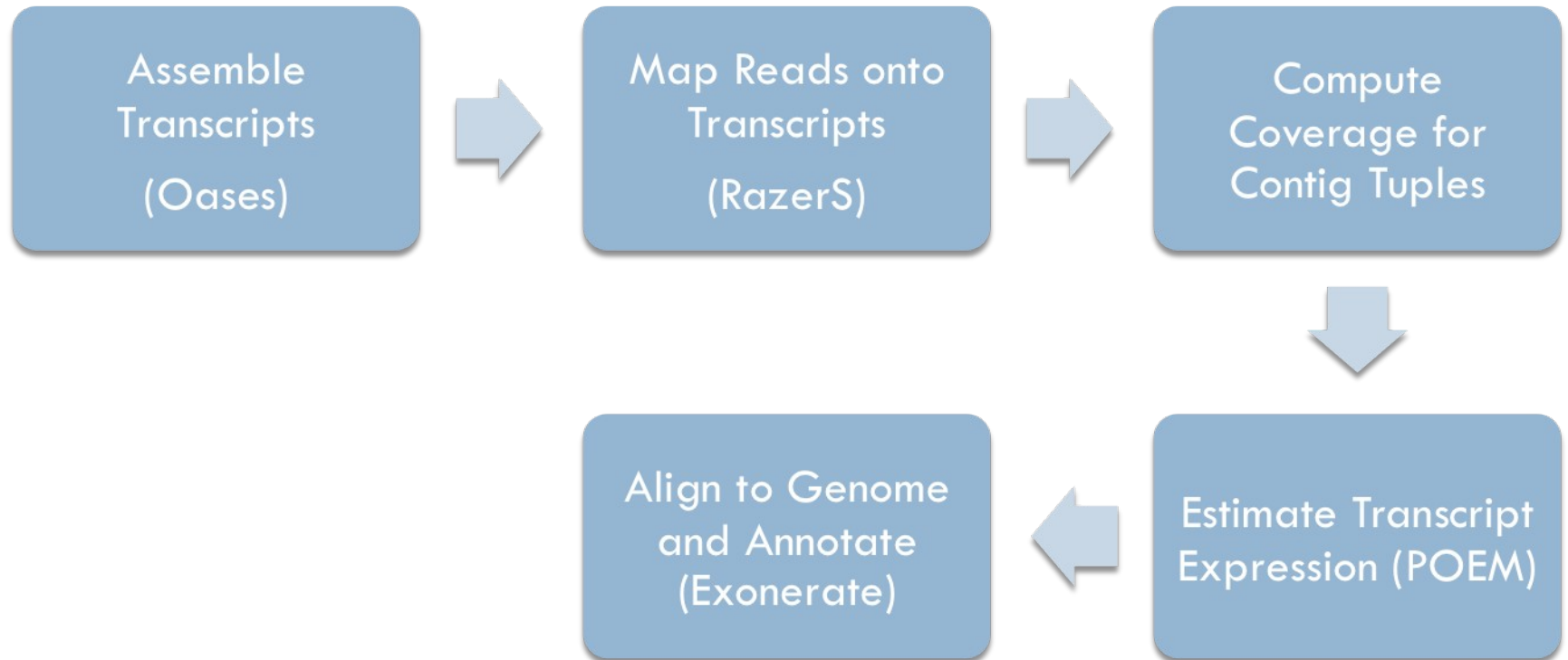
# De Novo transcripts assembly

- Advantages:
  - No reference or bad quality genome
  - Cancer transcriptomics (genes fusion)
  - Micro exons

- Assembly specific challenges:
  - sequencing errors are hard to rescue
  - differentiation of (post-)transcriptional modifications like alt. splicing, alt. polyadenylation, alt. first exon or trans-splicing
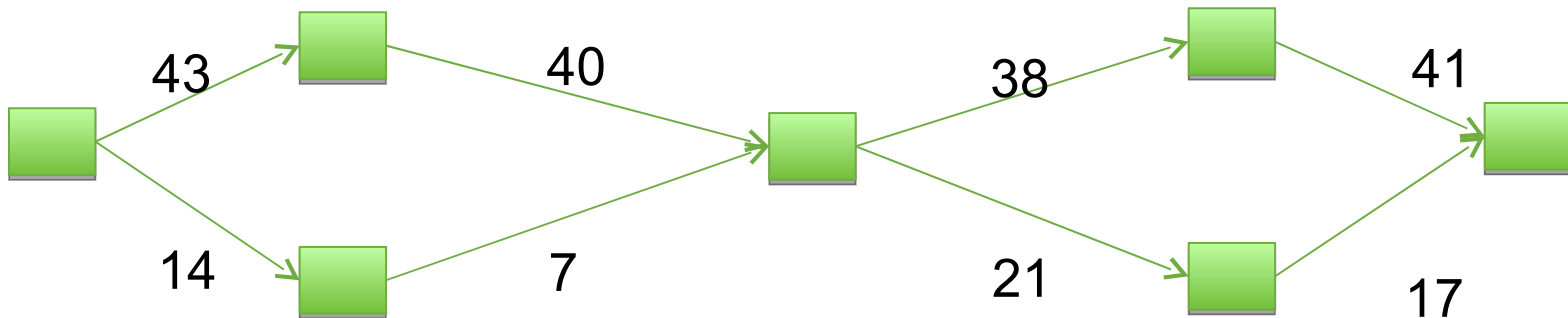  - judgement of assembly quality without reference
- paralogous domain genes

# Workflow for RGASP

# From contigs to transcripts

- weighted contig graph

  - reconstruction hard (Lacroix et al. WABI 2008)

- iterative maximum likelihood reconstruction method (heuristic) (Lee 2003)

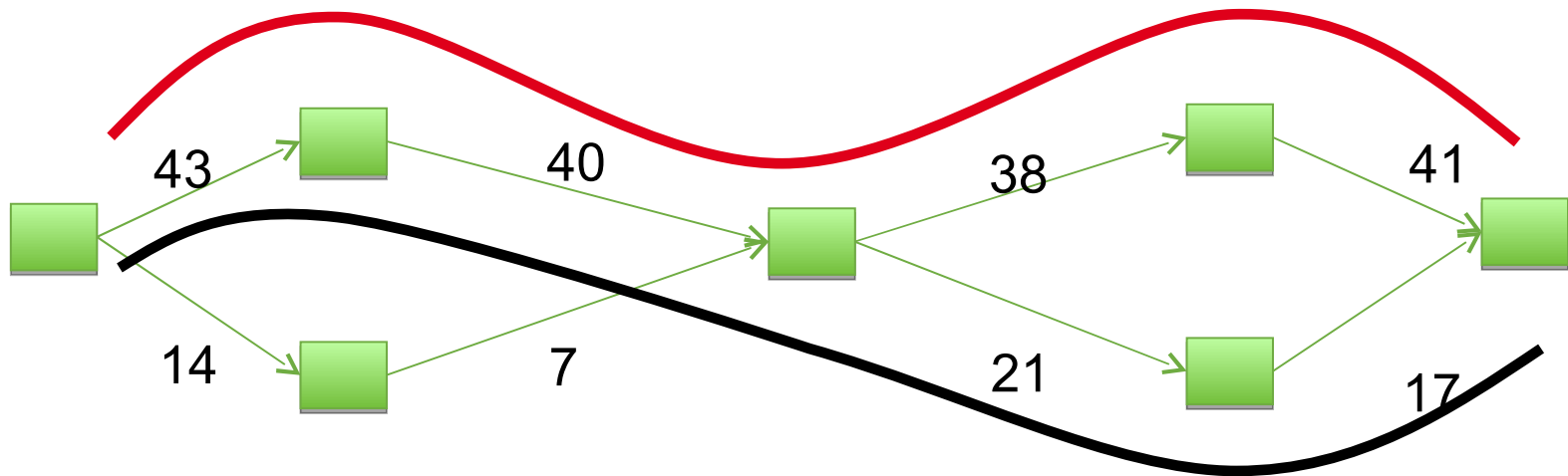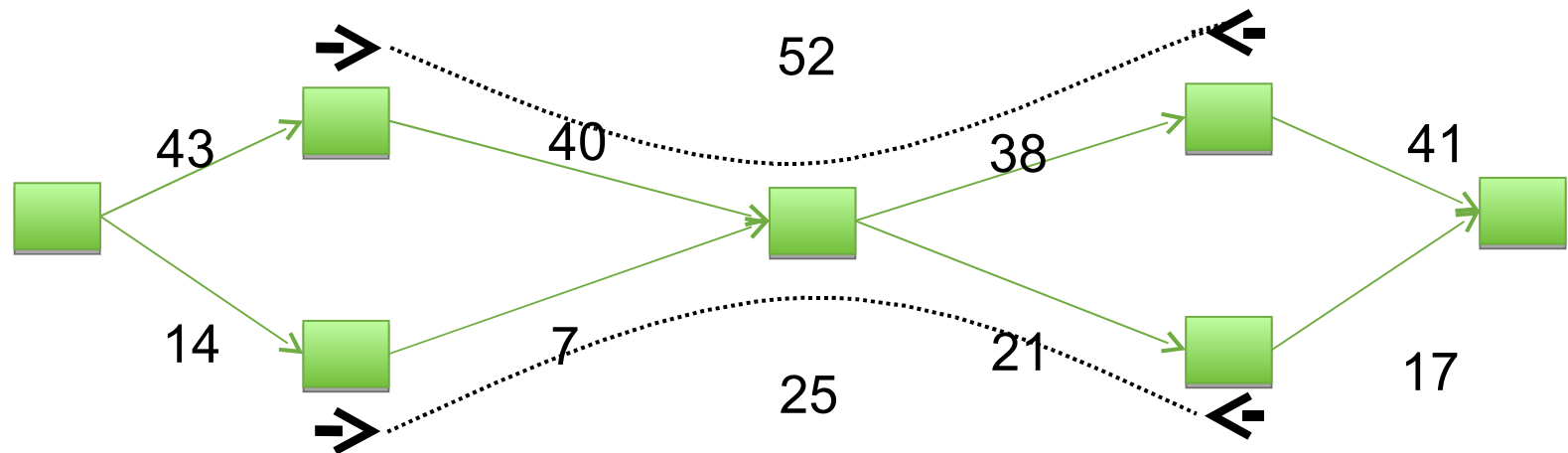# From contigs to transcripts

- weighted contig graph

  - reconstruction hard (Lacroix et al. WABI 2008)

- iterative maximum likelihood reconstruction method (heuristic) (Lee 2003)

# From contigs to transcripts

- weighted contig graph

  - reconstruction hard (Lacroix et al. WABI 2008)

- incorporation of paired-end reads (transitive reduction)

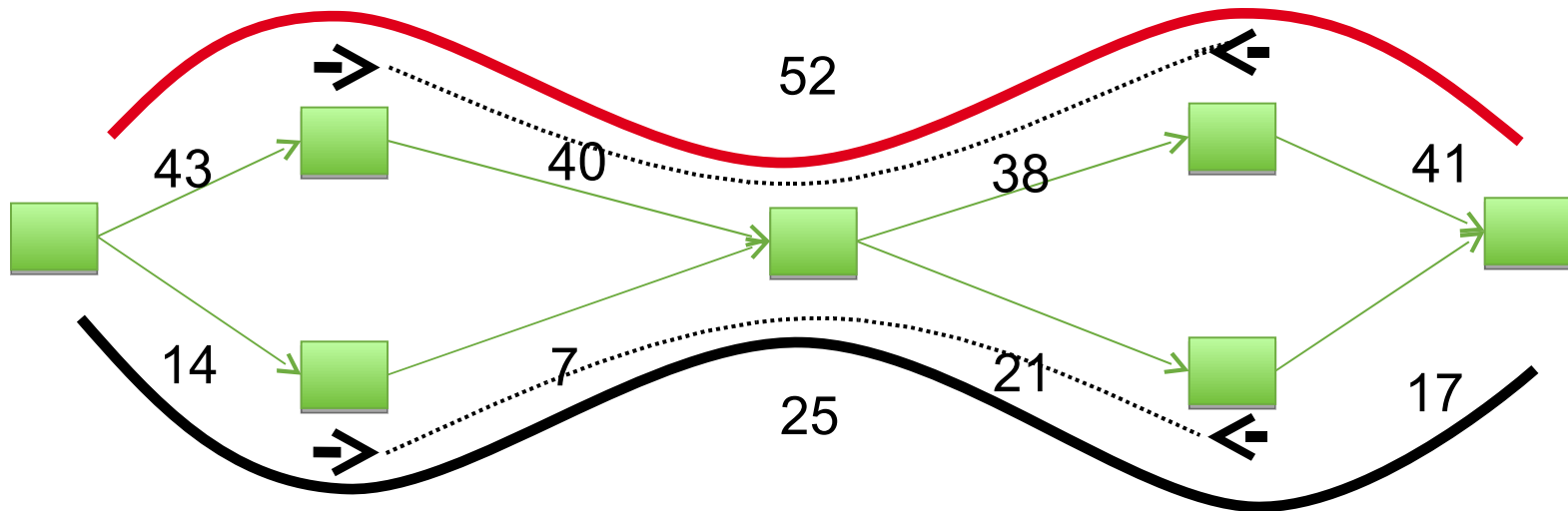# From contigs to transcripts

- weighted contig graph

  - reconstruction hard (Lacroix et al. WABI 2008)

- incorporation of paired-end reads (transitive reduction)
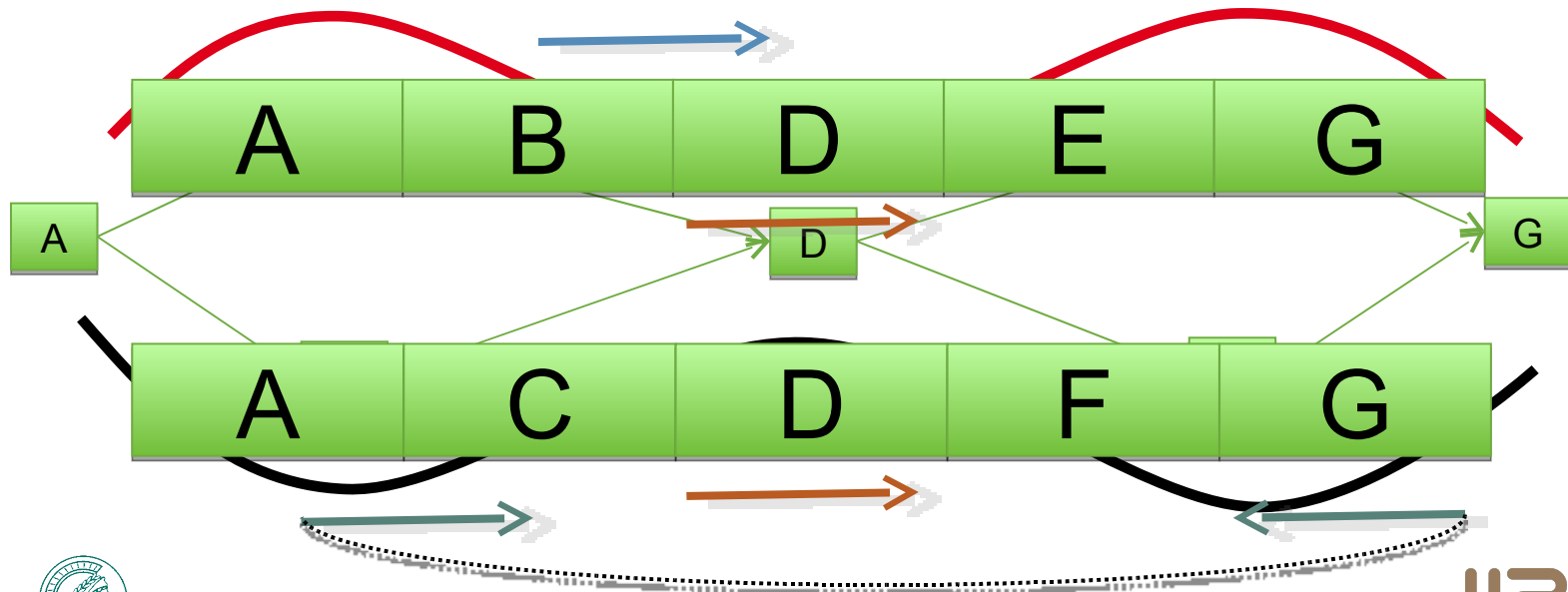
# Map Reads onto Transcripts

- map all reads onto assembled transcripts

- allow for mismatches, indels

- record all (multiple) matches

- RazerS - Fast Read Mapping with Sensitivity Control (Weese et al. 2009)

# Compute Coverage for Contig Tuples

- for every read match find tuples of covered contigs

- count for every tuple the number of covering reads
  - reweight multiple and paired-end matches

# Perspectives

- Estimating the depth of sequencing
  - Estimates for the total number of *transcripts*

- Infering new events
  - Automatic correction of experimental biases

- Stay tuned for RGASP results

# Acknowledgements

**MPIMG (Berlin)**

*Marcel Schulz*

*Marc Sultan*

Alon Magen
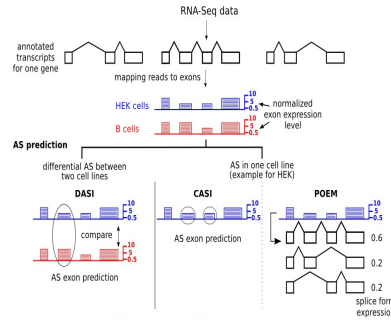Tatjana Borodina
Aleksey Soldatov
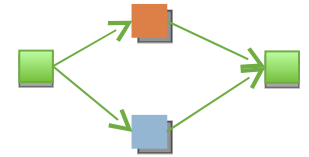Dmitri Parkhomchuk
Stefan Haas
Martin Vingron
Hans Lehrach
Marie-Laure Yaspo



**RGASP pipeline:**

*Marcel Schulz*
*Daniel Zerbino* (UCSC, former EBI)
*David Weese* (FU Berlin)

Ewan Birney (EBI)
Knut Reinert (FU Berlin)
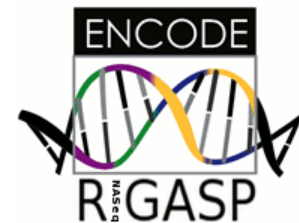Martin Vingron (MPIMG)

Now at Cordeliers
Come and visit !

# Results on Drosophila Paired-end reads

read length=36
k=21
fragment length=200

| #S2-DRSC Paired-reads | # loci | # transcripts | median length | mapped to genome | Overlapping ENSEMBL transcripts | #identified exons |
|---|---|---|---|---|---|---|
| 43,836,085* | 32,905 | 53,299 | 214 | 89% | 61% | 72,892 |

human, fly, and worm predictions
submitted to RGASP competition

*high quality reads

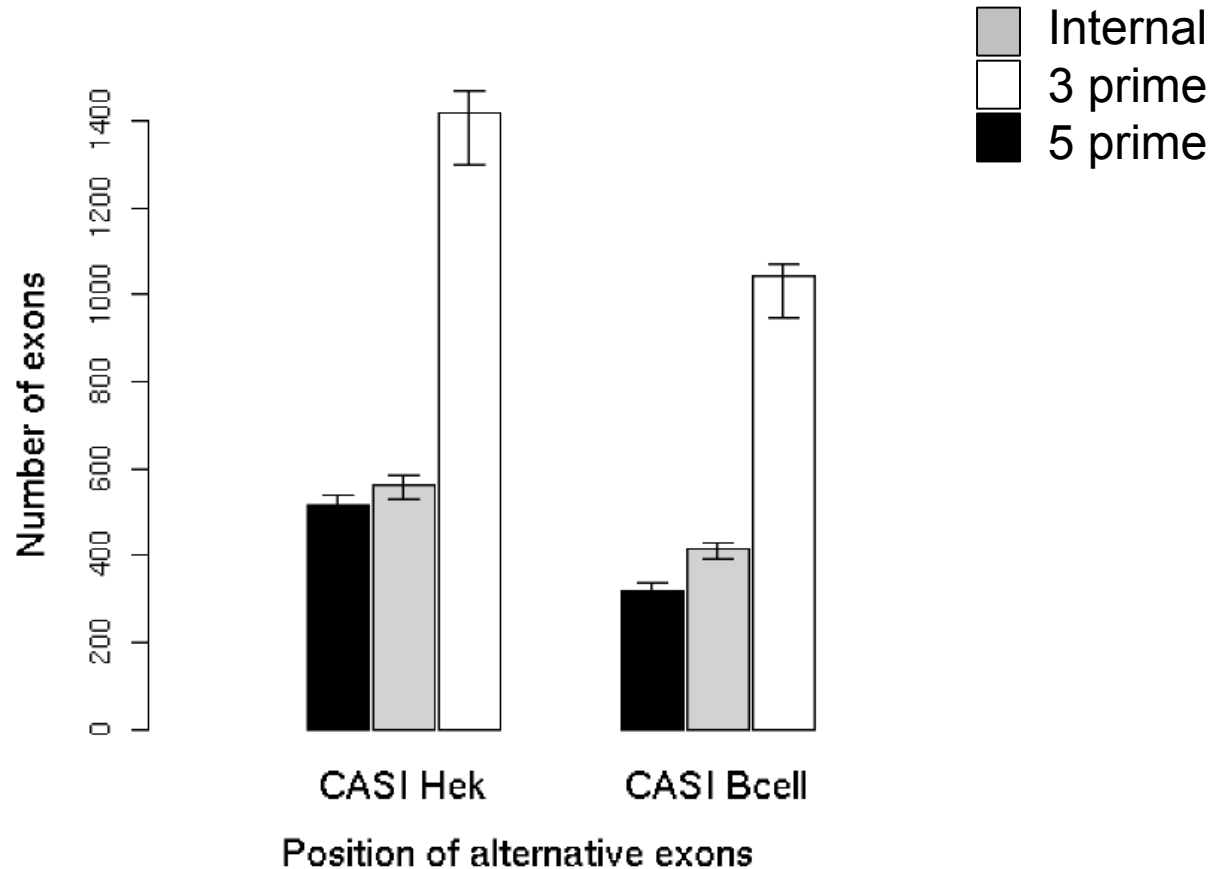# Statistical tests and assumptions



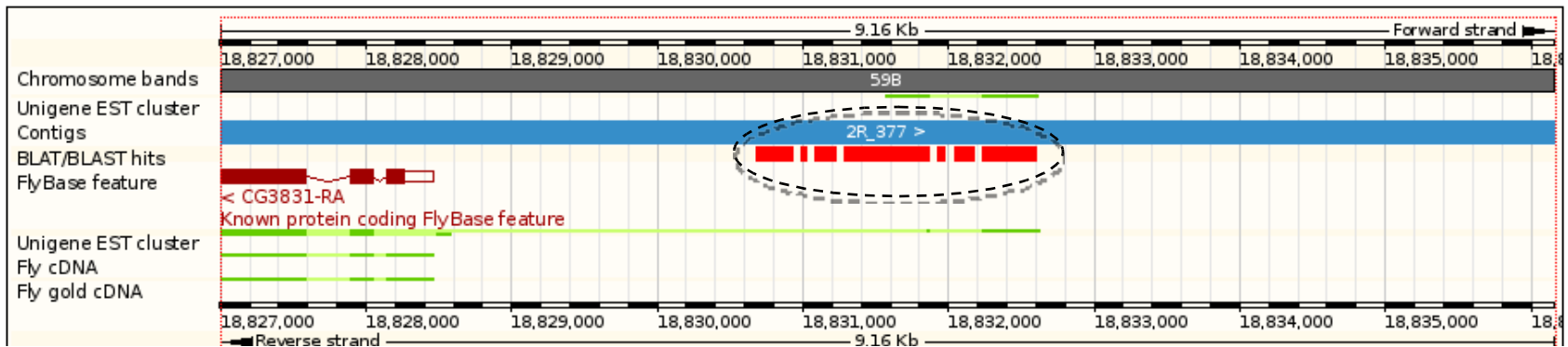exon based index is a robust z-score estimated with median and median absolut deviation
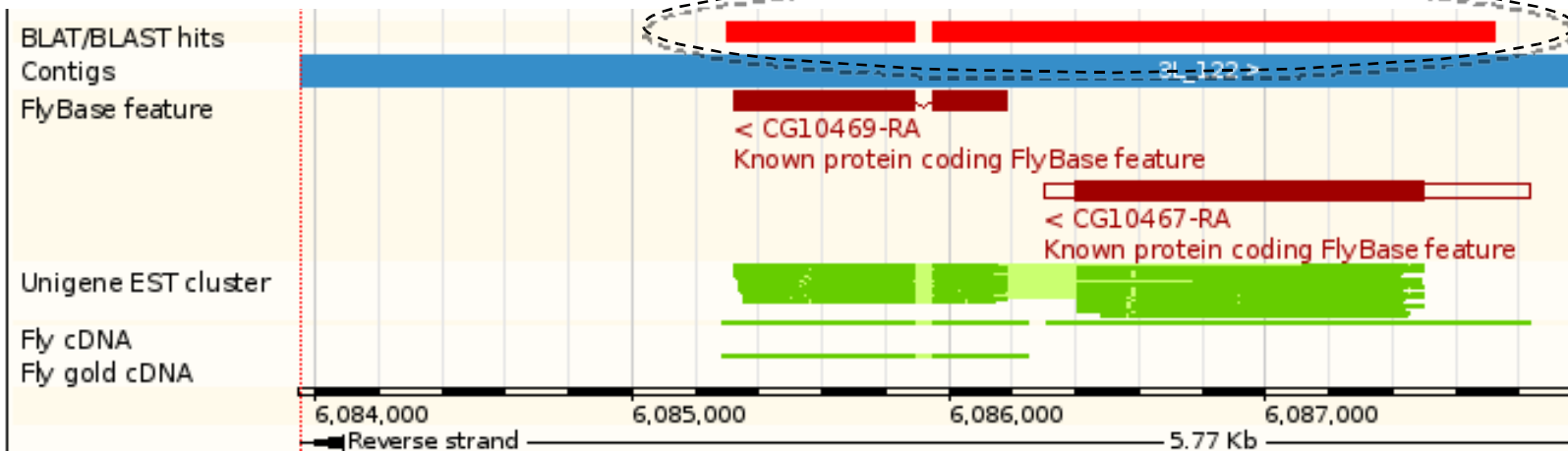
# CASI complements splice junctions

# Examples

unannotated gene with 7 exons (locus 136)

assembled transcript



transcript connecting 2 genes (locus 2589)

# Exons Arrays vs RNA-Seq