

Hugues Richard

Technology Transforming Biology

Until late 20th Century



Hypothesis Generation
and Validation

21th Century and Beyond



**Mathematics &
Computation**

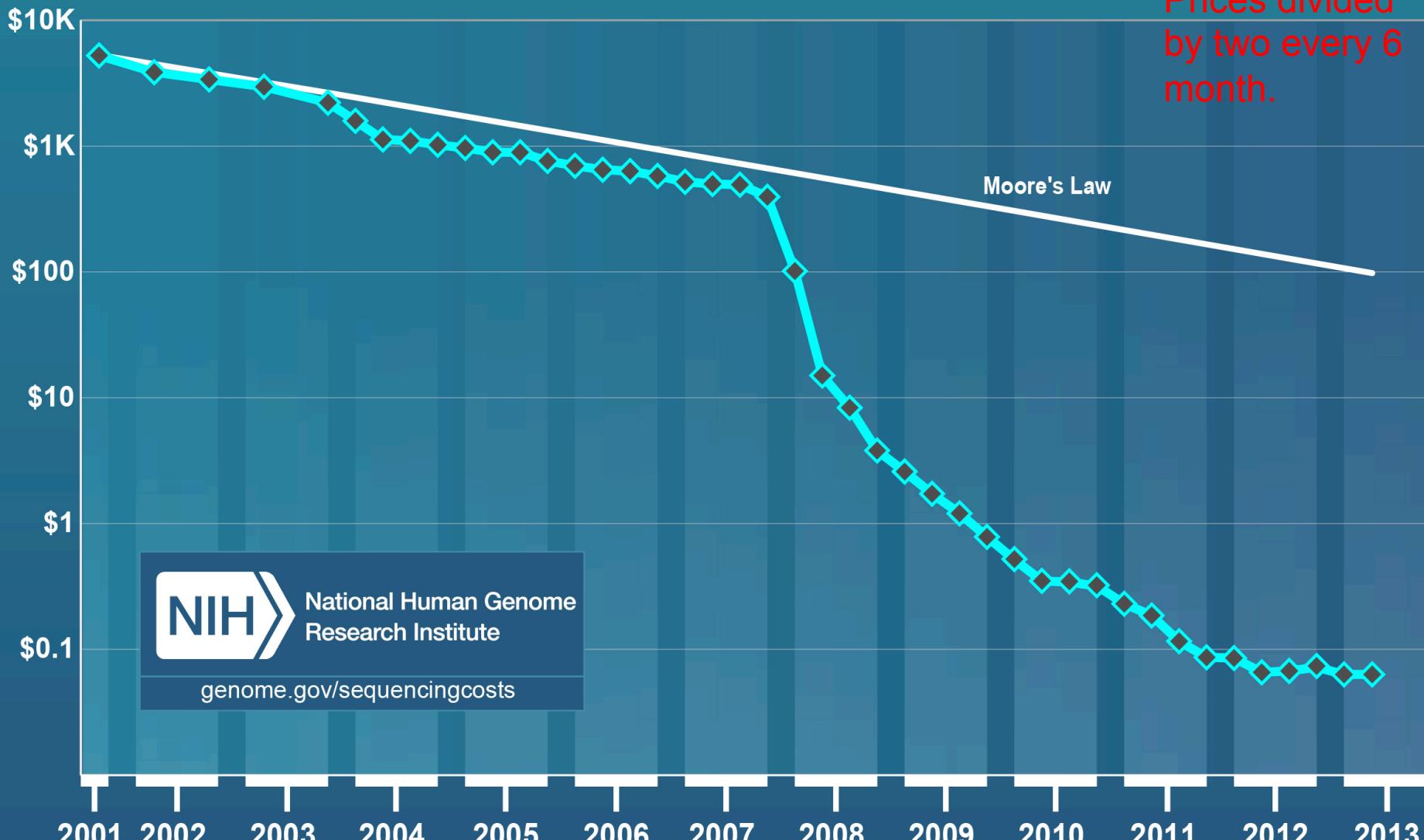


Hypothesis Generation
and Validation

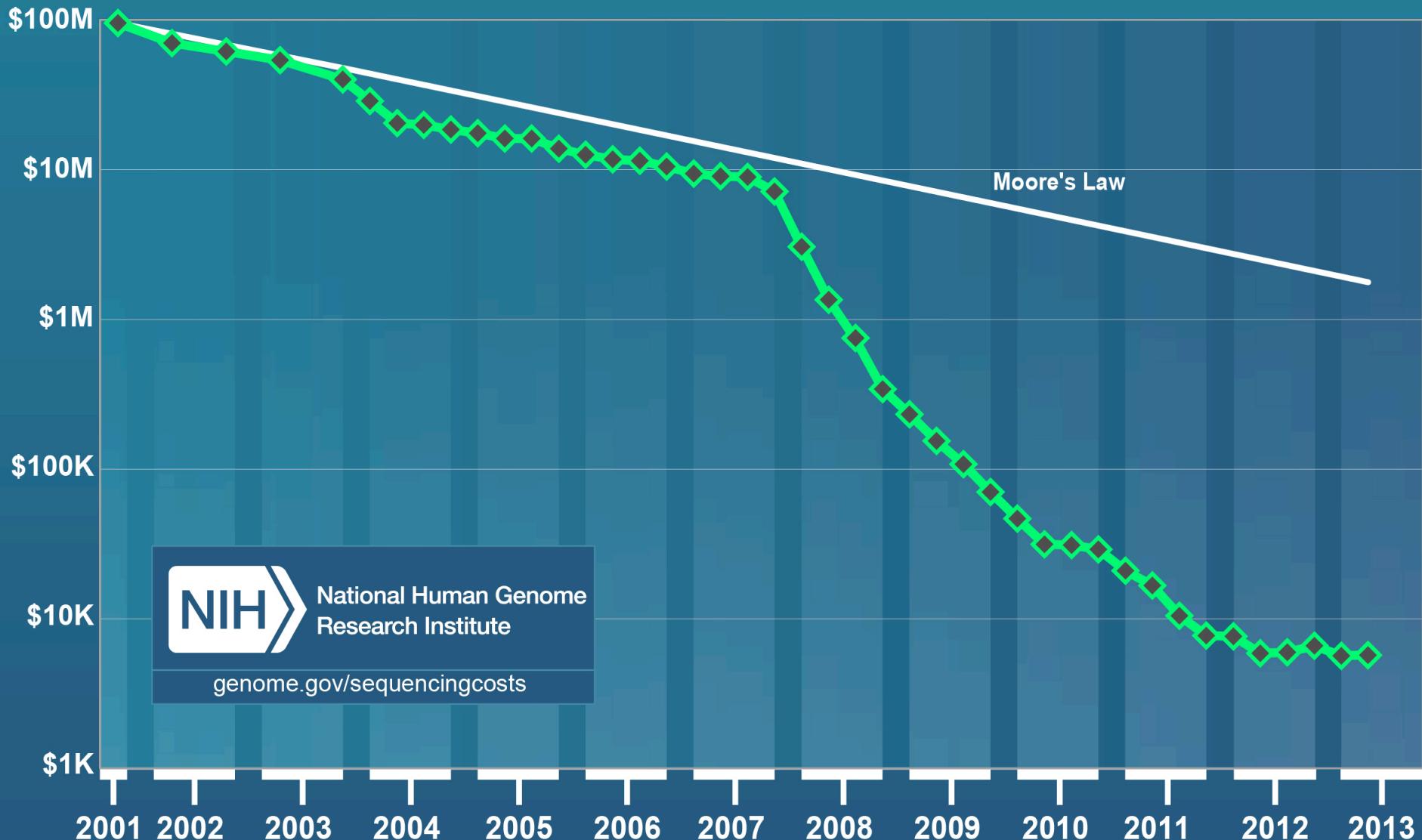
High throughput technologies

Cost per Raw Megabase of DNA Sequence

Prices divided
by two every 6
month.



Cost per Genome



Analysis of NGS data

- Lectures:
 - 12/01: High-throughput sequencing, sequence alignment and error correction
 - 19/01: Genome sequencing and resequencing
 - genome assembly
 - Annotation of sequence polymorphism (A. Gillet)
 - 26/01: Transcriptome analysis
- Tutorials:
 - 12/01, 26/01: Pan & Core Genome (A. Ugarte)
 - 19/01: detection of structural variants (A. Gillet)



1. Genome (Re)Sequencing

Sequence alignment / Analysis of variants

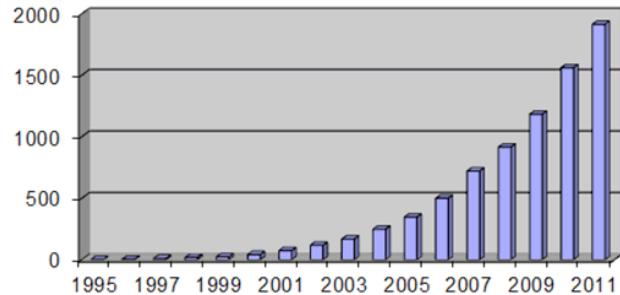
Outline

- Brief historical background on sequencing
 - the many applications of High throughput sequencing
- A typical resequencing workflow
 - Reads alignment
 - Basic statistics after alignment
 - Detection of sequence polymorphisms
 - Genome Assembly



Timeline of genomes

- 19 years of sequencing
 - 1995: first bacterial genome *H. Influenza*
 - 1996: first eukaryotic genome *S. Cerevisiae*
 - 2001: Completion of the Draft of the Human genome.
 - 10 Billion Dollar and 10 years to complete.
 - Now, almost 2k genomes
(free download)
- 2005 on: Massive development of new technologies
 - One human genome < 5,000\$, less than 1 week.

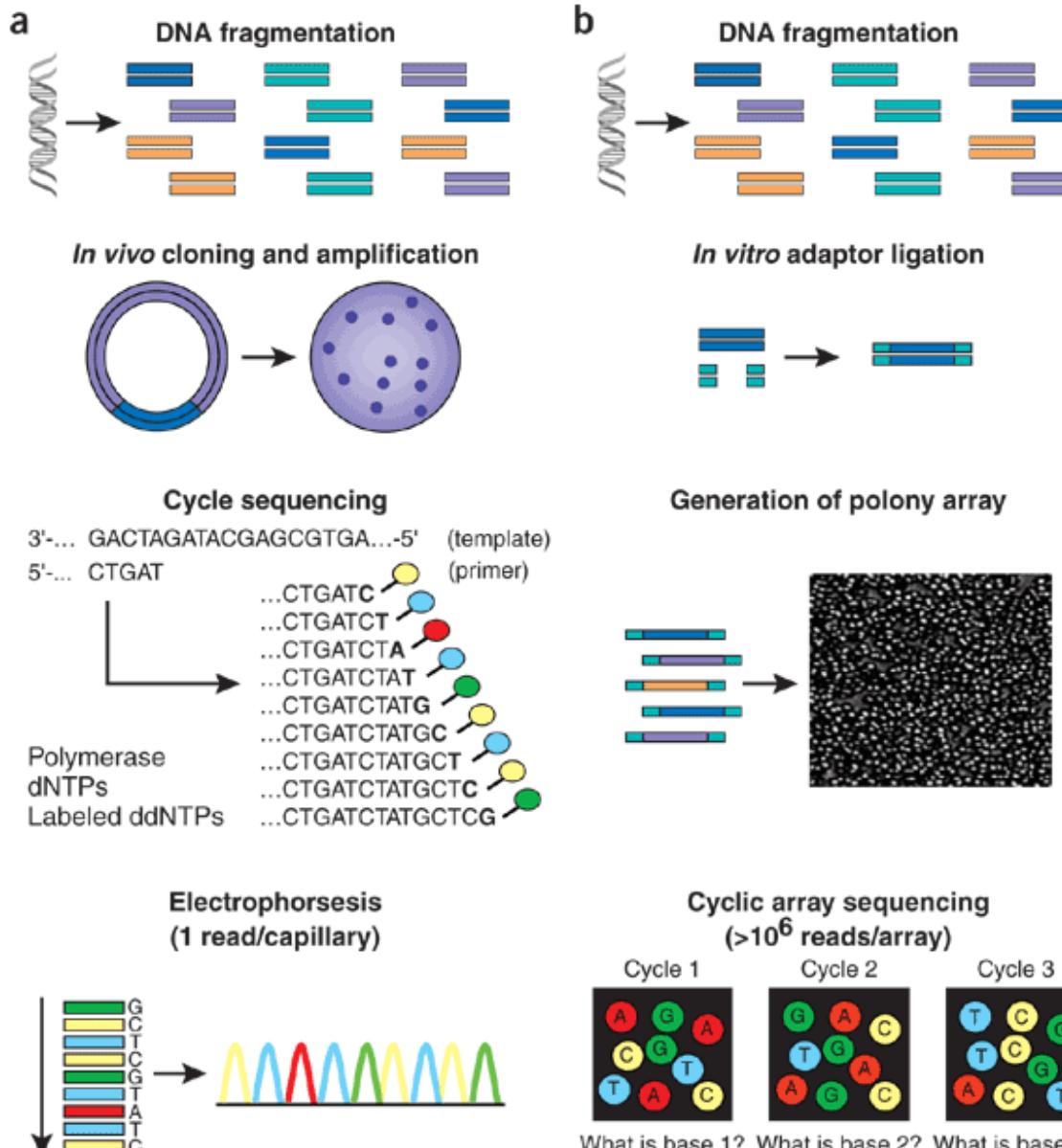




Changes in instrument capacity over the past decade, and the timing of major sequencing projects

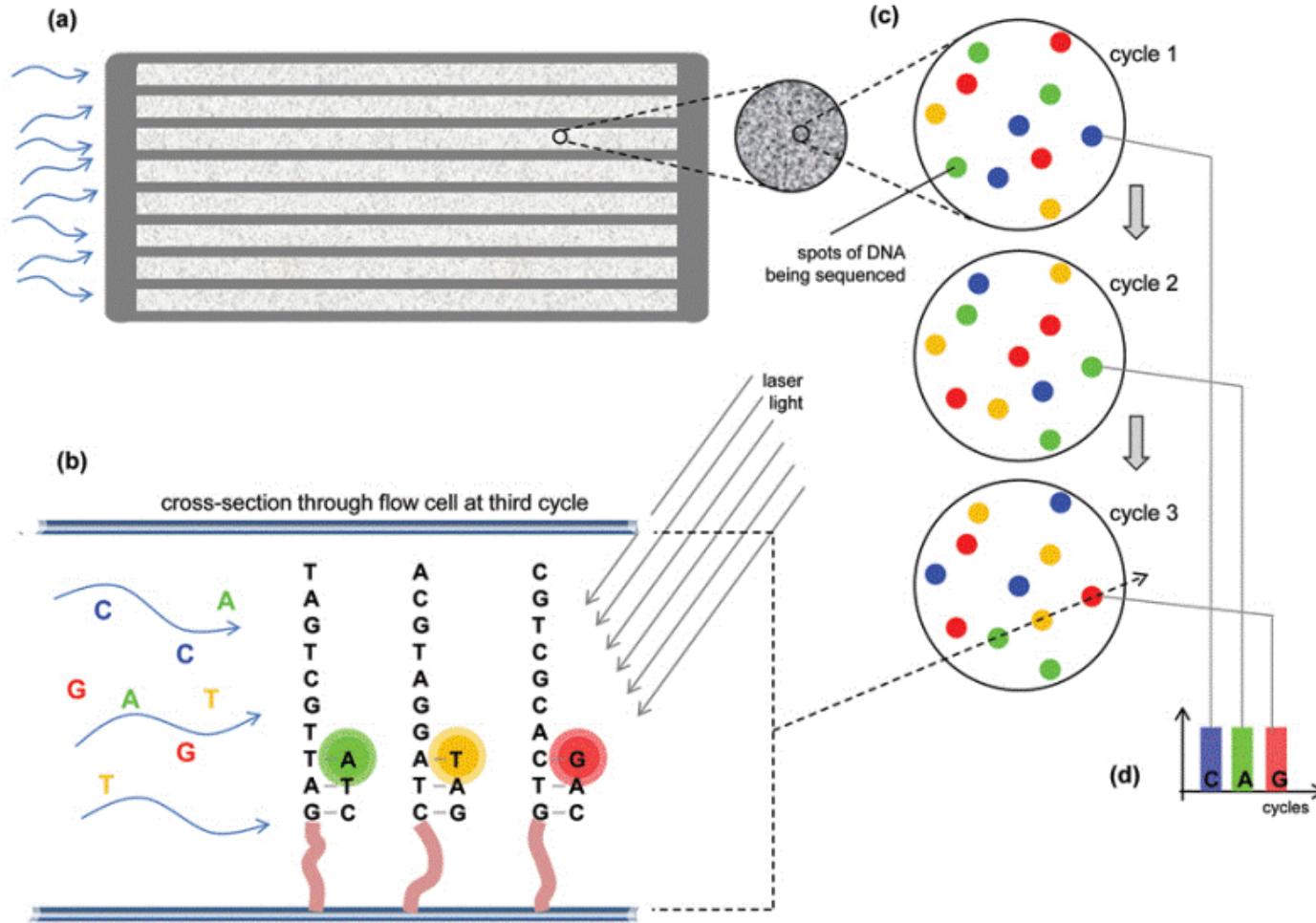
nature
ER Mardis. *Nature* 470, 198-203 (2011)

Specificity of 2nd generation technologies

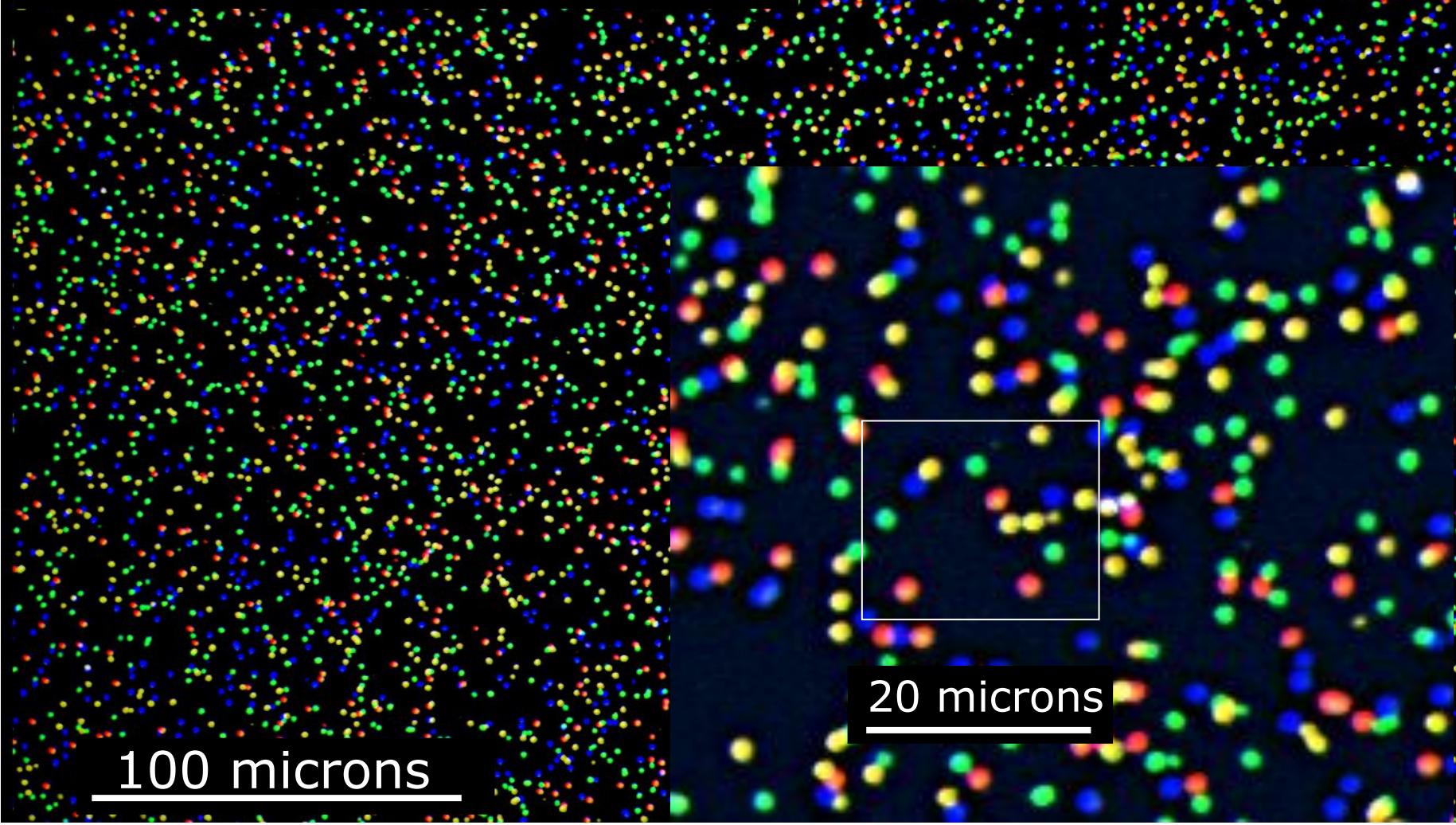


- a. Sanger sequencing
- b. Illumina sequencing

Illumina Sequencing



1 cycle: ~3000 images
~40 million clusters



150 light-years in diameter
about 10 million stars



Copyright: Martin Pugh

Next | 2nd Generation Sequencing

Machine	454 FLX [Roche]	HiSeq 2000 [Illumina]	SOLiD 4 [ABI]
Sequencing Approach	Pyrophosphate Release	Bridge Amplification	Ligation
Read Lengths	≈ 400 2x140 bps	2x100 bp	50 50+35 bp
Paired Ends	Yes	Yes	Yes
Output per Run	≈ 400 Mb in 2-4 d	100 200 Gb in 4 8 d	100 Gb in 6 d
Accuracy depends on	Homopolymer Length (>6 problematic)	Position in the Read	Position in the Read
Costs per Run	\$9,000	\$10,000	\$6,000



Numbers as of 2011.08

Next Generation Sequencing Platforms



Jusqu'à
1,8 Tb

HiSeqXTen



HiSeq
Jusqu'à
1000Gb

Jusqu'à
100 Gb

Solid5500

Ion Proton



Jusqu'à
10 Gb

PGM



De 10 Mb à
1,5 Gb

NextGen500
Jusqu'à
120 Gb

MiSeq
Jusqu'à
15 Gb



Jusqu'à
700 Mb

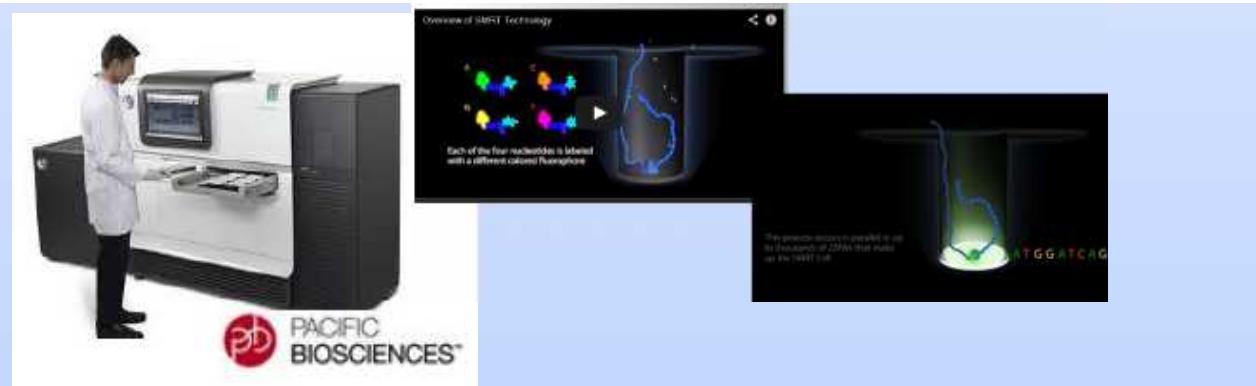
Junior

Jusqu'à
35 Mb

Recherche

Clinique,
Diagnostic

Single molecule sequencing

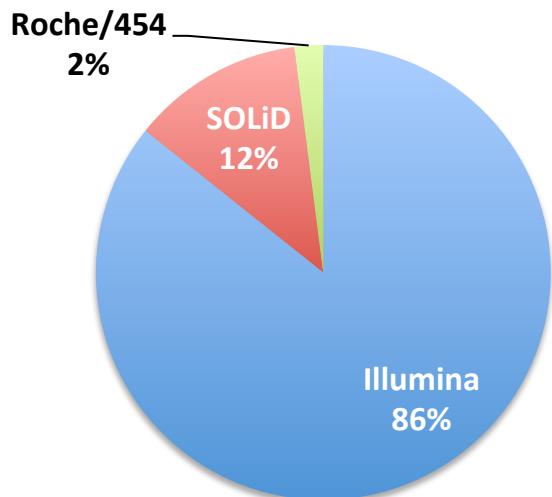


Oxford Nanopore Minion



High-Throughput Sequencing

technology	instrument	run time	yield [Mb/run]	read length [bp]	costs [\$/Mb]	error rate [%]
Sanger	3730xl (capillary)	2 h	0.06	650	1500	0.1-1
Illumina	HiSeq 2000	8 d	200,000	2 × 100	0.10	≥ 0.1
SOLiD	SOLiD 4	12 d	71,400	50 + 35	0.11	> 0.06
Roche/454	FLX Titanium	10 h	500	400	12.4	1
SMRT™	PacBio RS	0.5-2 h	5-10	860-1100	11-180	16
HeliScope™	Helicos	N/A	28,000	35	N/A	N/A



(Glenn, 2011)

More than 100 Terabases publicly available
in the Short Read Archive (Kodama et al., 2012)

Big data

- One sequencing experiment:
 - > 100 Gb of sequence
Human genome 3.3 Gb

LARGE HADRON COLLIDER
13 PETABYTES (2010)
The proton collider, near Geneva, Switzerland, generates about 15 petabytes of data per year — even after rejecting 99.9995% of collisions.

HEAVY-DUTY DATA

The computer-storage space required to support projects in the digital humanities is now starting to rival that of big-science projects.

BIG SCIENCE

SLOAN DIGITAL SKY SURVEY

50 TERABYTES
The survey, begun in 1998 using a 2.5-metre telescope in New Mexico, has discovered nearly half-a-billion asteroids, stars, galaxies and quasars.

GENBANK

530 GIGABYTES
This database, which stores publicly available sequenced DNA, included 127 billion bases at the latest count.

BIG HUMANITIES

CULTUROMICS N-GRAMS VIEWER

300 GIGABYTES (English only)
The string of letters in this corpus of 5 million books is 1,000 times longer than the human genome.

YEAR OF SPEECH

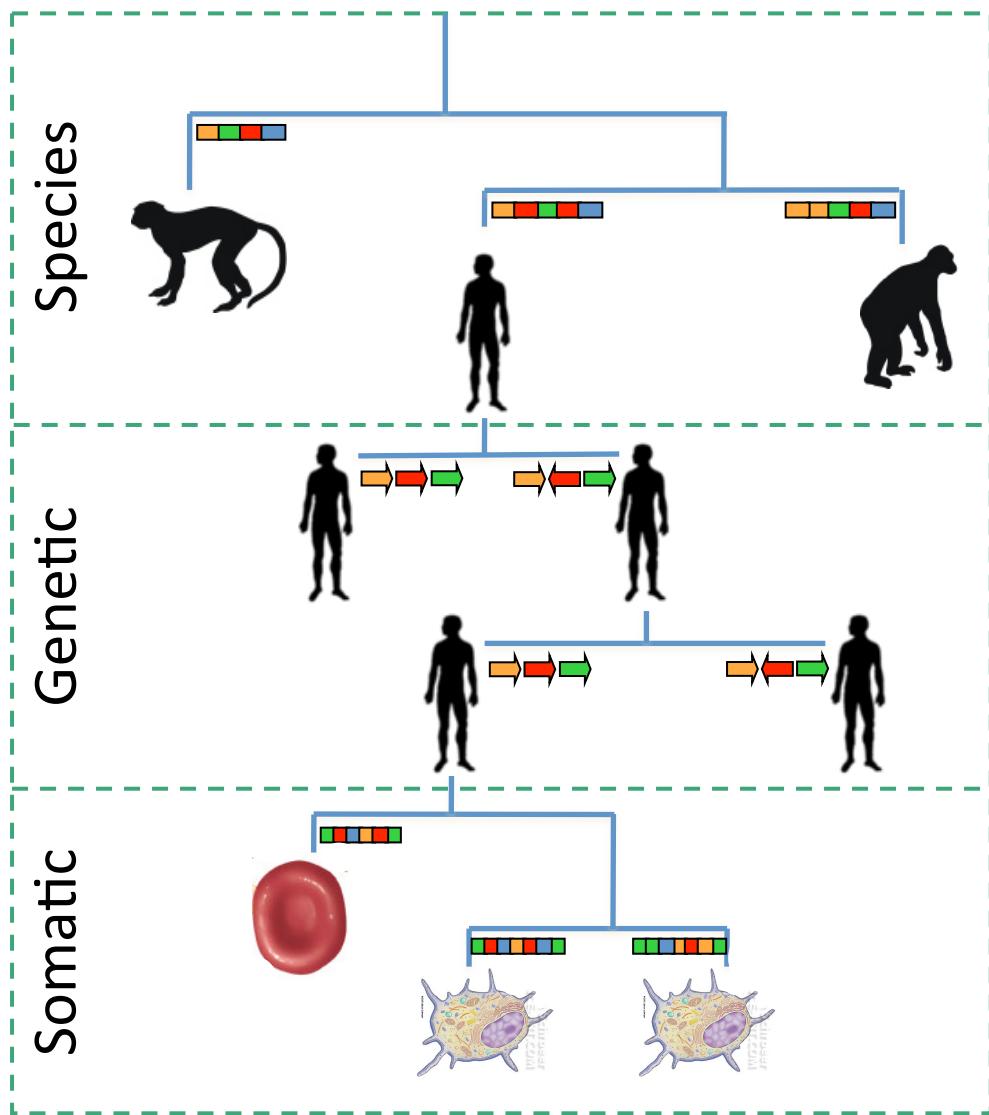
1 TERABYTE
This database includes recordings from telephone conversations, broadcast news, talk shows and US Supreme Court arguments.

UNIVERSITY OF SOUTHERN CALIFORNIA SHOAH ARCHIVE

200 TERABYTES
This archive stores 52,000 videotaped interviews with Holocaust survivors from 56 countries.

1 petabyte = 1,024 terabytes = 1,048,576 gigabytes

Genome Sequencing and Comparison

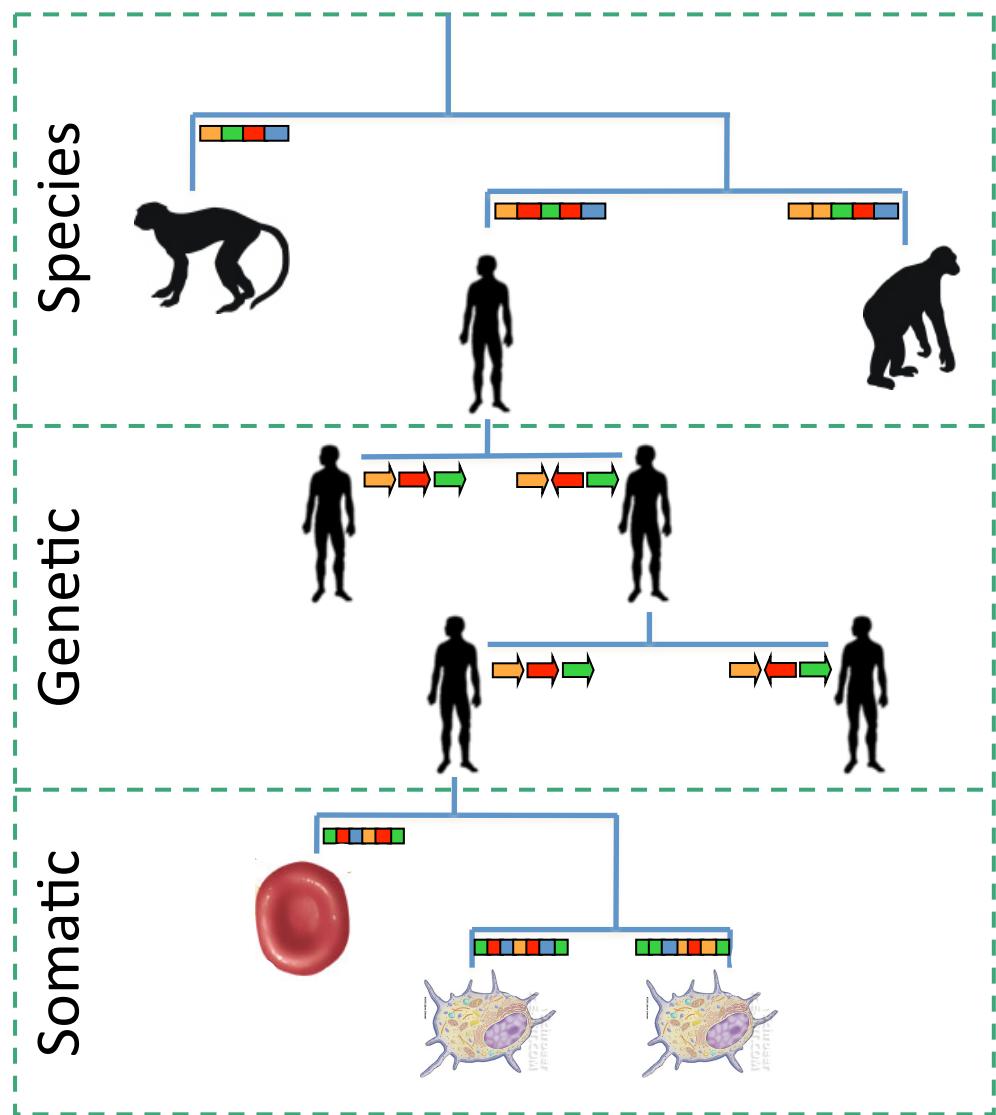


Comparative genomics
Differences between species?

Personal genomics
Genetic basis for traits of individuals?

Cancer genomics
Which somatic mutations lead to cancer?

Petabytes of Genomes



1000 Genomes
A Deep Catalog of Human Genetic Variation



 **Autism
Genome
10K**

THE CANCER GENOME ATLAS 



**International
Cancer Genome
Consortium**

Genome (re)-sequencing

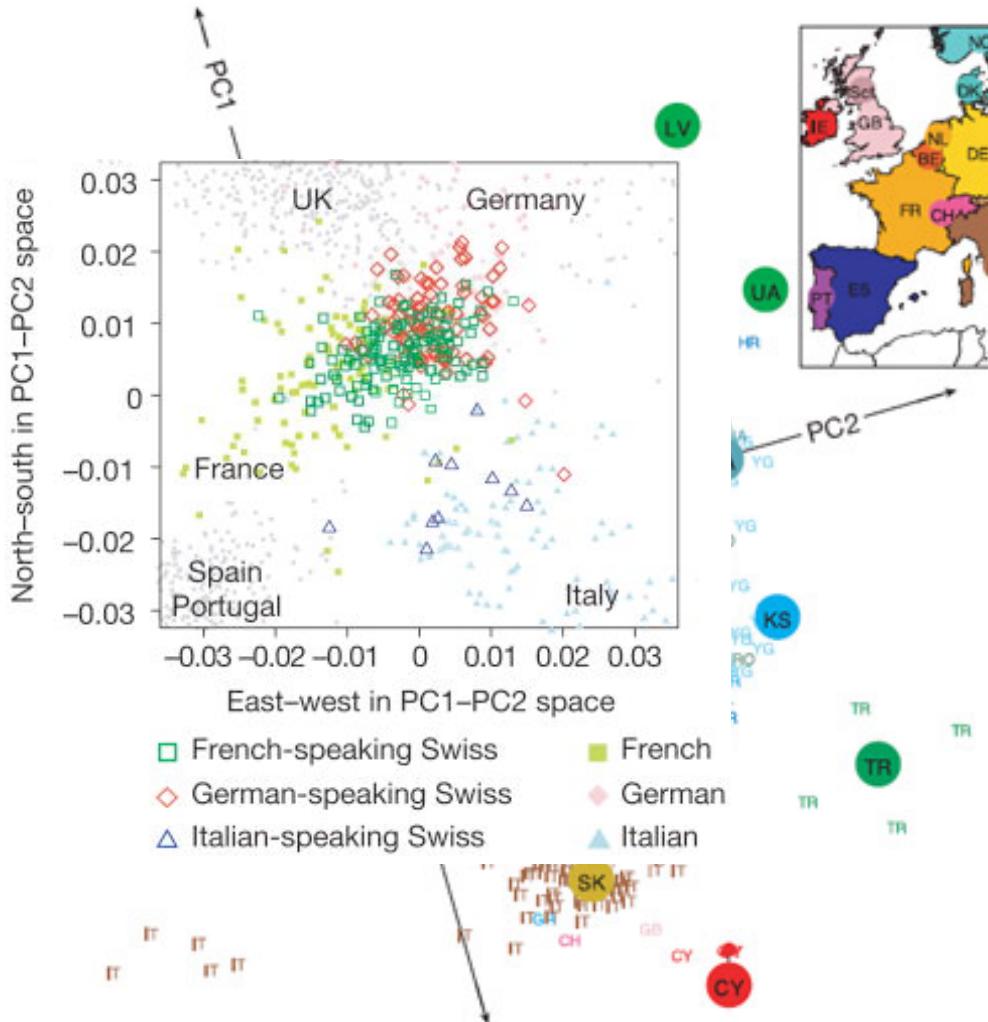
- Evolutionary studies:
 - 10k genomes on vertebrate species
- Gene Wide association studies
 - Autism genome 10k
- Population genomics:
 - 1001 genomes (A. Thaliana)
- Personal genomics
- Tumor cell profiling (cancer genomics)



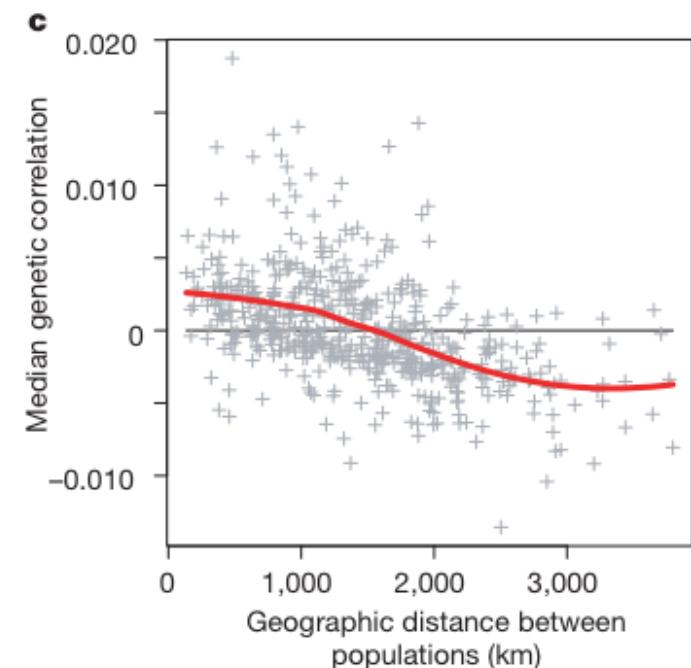
The Cancer Genome Atlas  *Understanding genomics
to improve cancer care*



Genotype variation



PCA on 3,192 individual from all Europe.
Data are alleles at a subset of genomic loci (~ 200k loci)



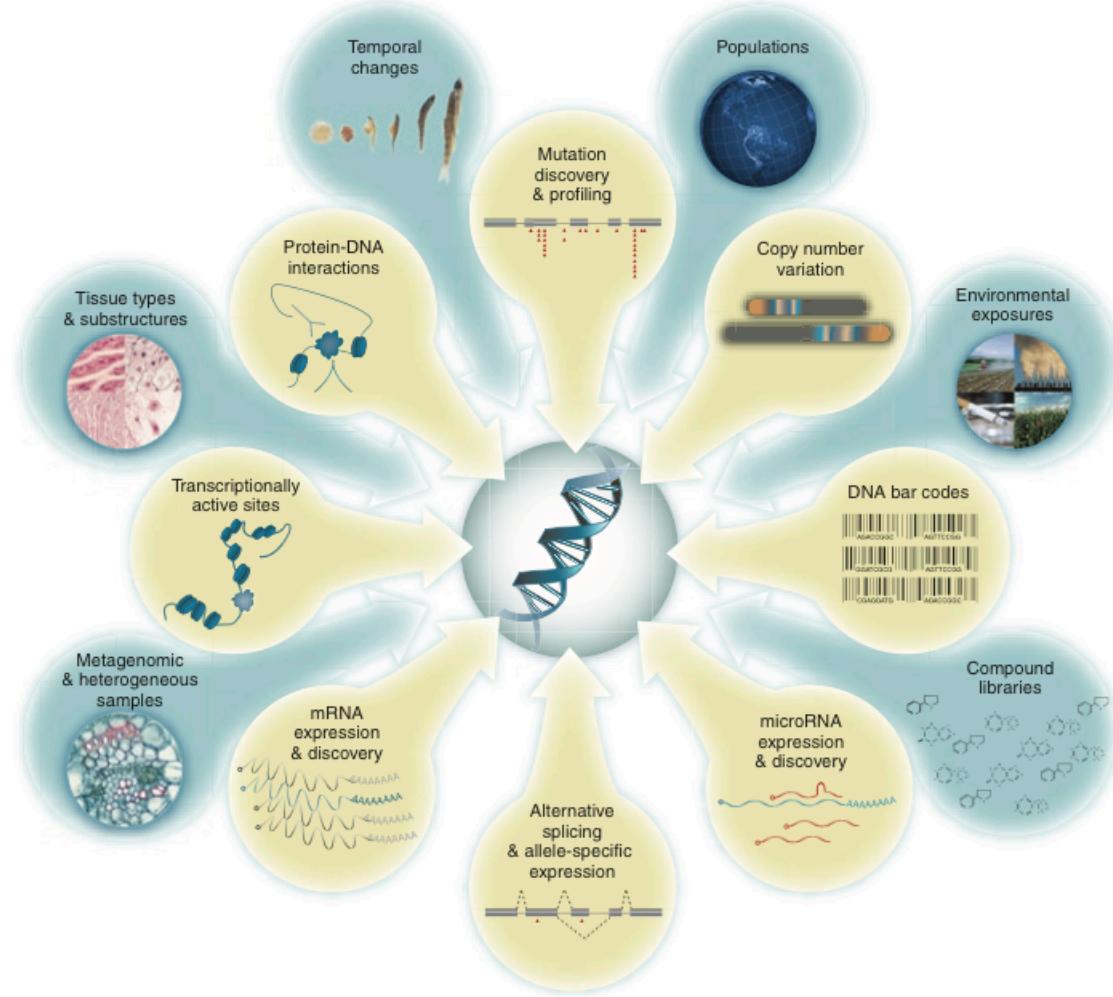
Sum-up

- Increasing availability of genomes will have a profound impact.
 - Personal Computers in the 80s, now Personal Genomes
- Large amounts of data:
 - Need efficient algorithms
 - Many (Mio) sequences + errors → Statistical reasoning
- Main questions for primary analysis in genomics:
 - Assign each sequence to its originating position
 - Annotate sequence variation w.r. to the reference



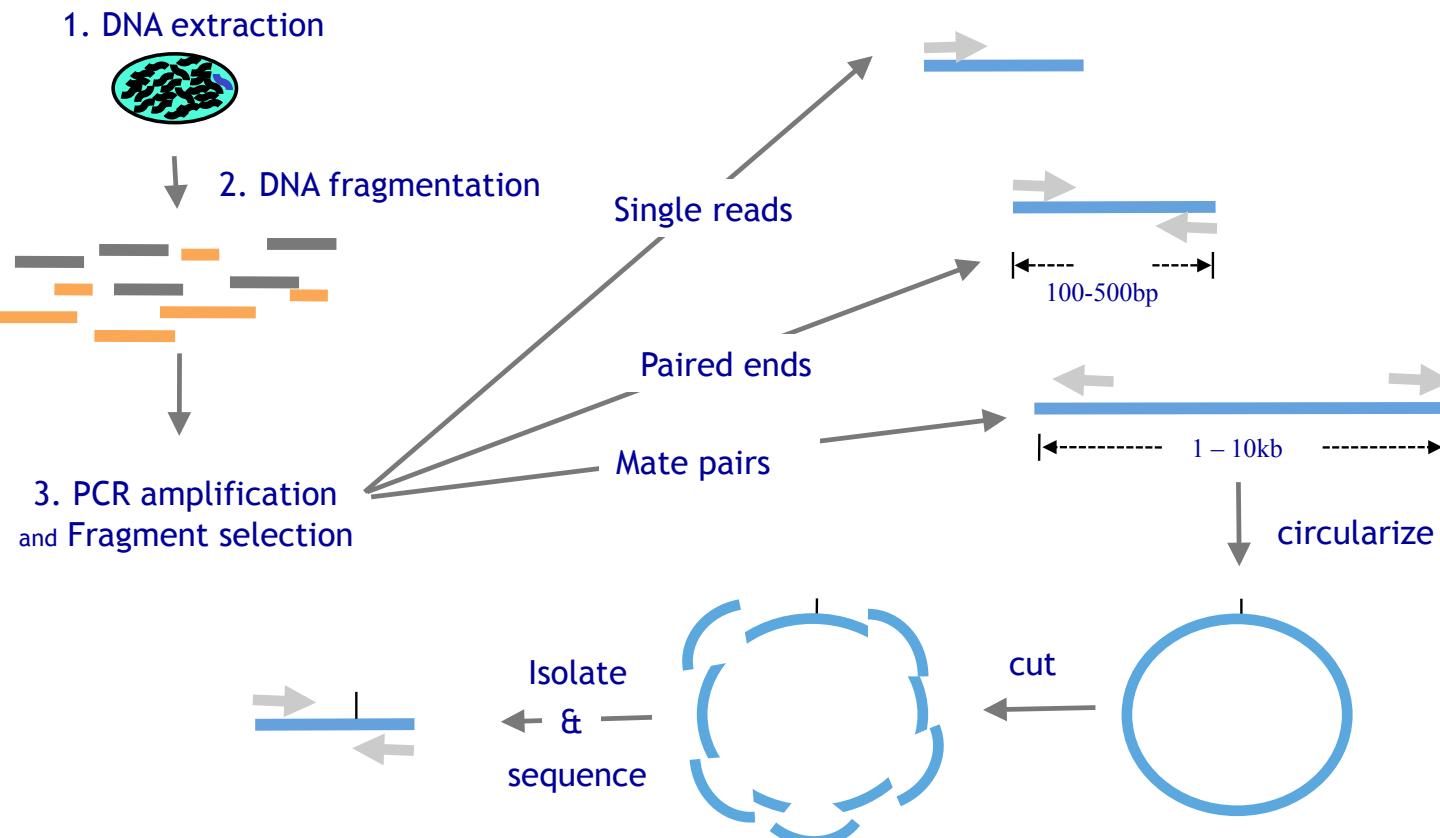
The many applications of NGS

- Genomic
 - (re)sequencing
 - metagenomic (tutorial)
 - targeted sequencing
- Functional
 - Protein-DNA (chIP-Seq)
 - mRNA (RNA-seq, CAGE)
 - small RNA (μ RNA-Seq)
 - methylome (meDIP-Seq)
 - Chromosomal structure
(Hi-Seq, 3C)



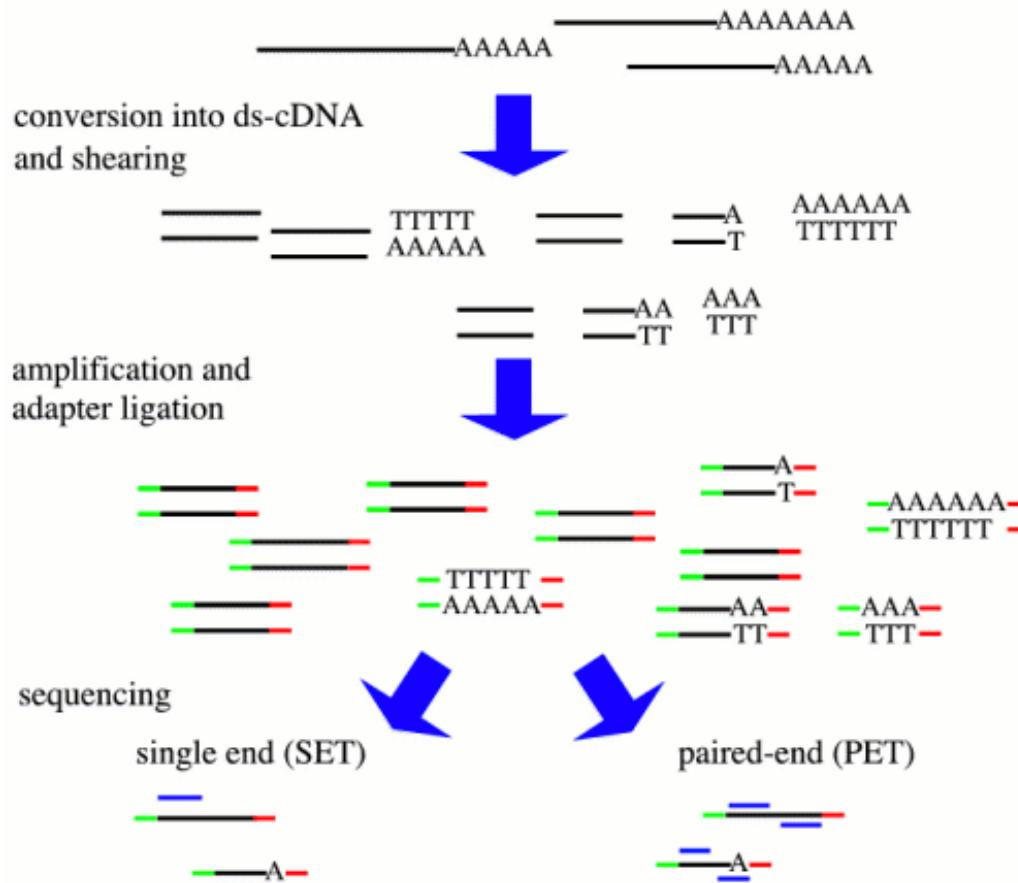
The sequences have to be aligned or assembled

Genome sequencing



RNA-Seq protocol

extraction of poly-A RNAs

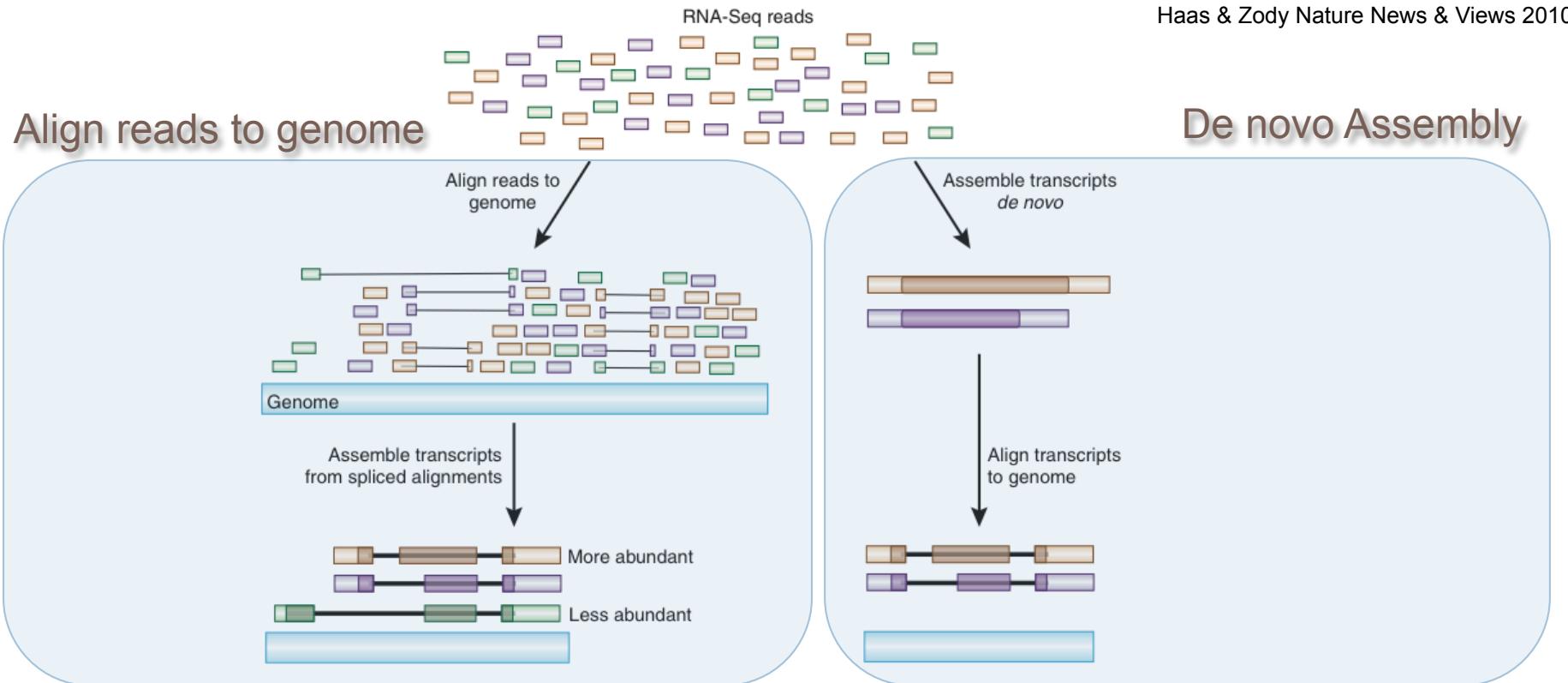


Whole Transcriptome Shotgun Sequencing

M. Schulz[©]



From reads to transcripts



- (Re)-annotation of gene models (transcripts)
- Transcripts quantification
- Transcriptional profiles

3C-based technologies

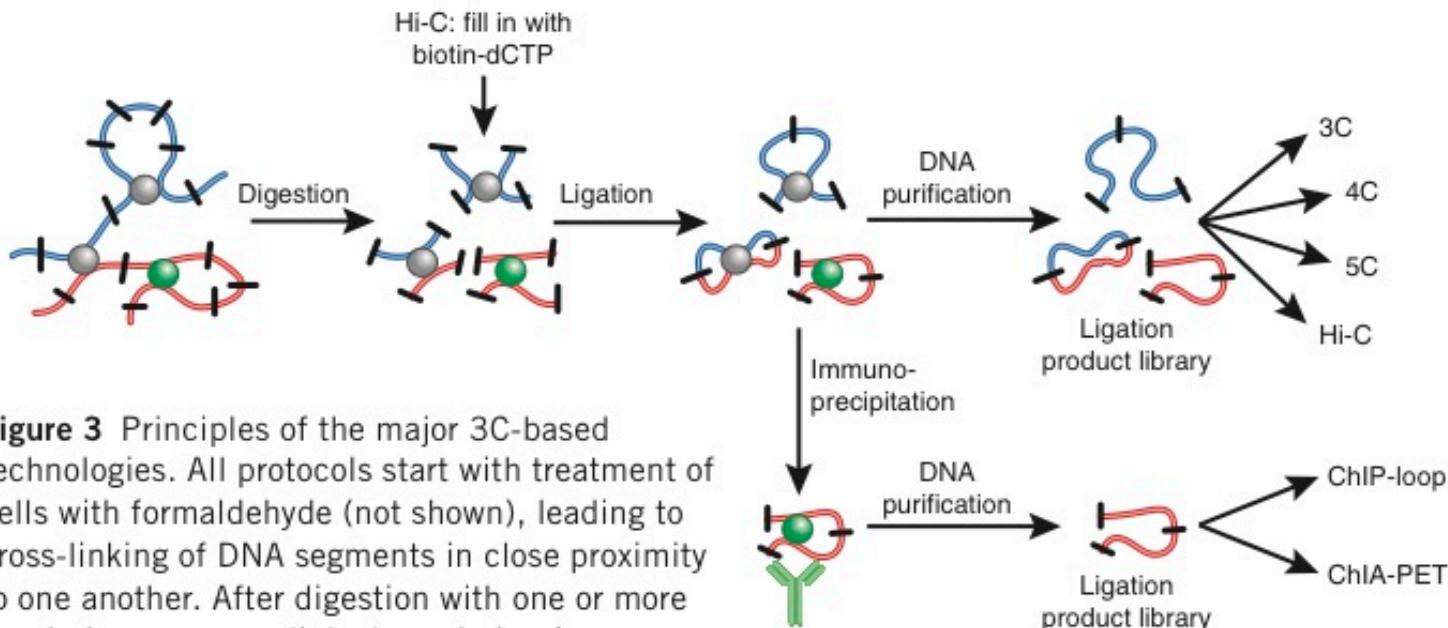


Figure 3 Principles of the major 3C-based technologies. All protocols start with treatment of cells with formaldehyde (not shown), leading to cross-linking of DNA segments in close proximity to one another. After digestion with one or more restriction enzymes, linked restriction fragments are intramolecularly ligated. In the case of Hi-C, the ends of the restriction fragments are first filled in with biotinylated dNTPs before ligation to facilitate purification of ligation junctions using streptavidin-coated beads. Single or multiple ligation events are detected directly (using 3C, 4C, 5C and Hi-C), or immunoprecipitation is first used to enrich for DNA associated with a protein of interest (using ChIP-loop and ChIA-PET). See **Table 2** for overview of different detection strategies and their scope.



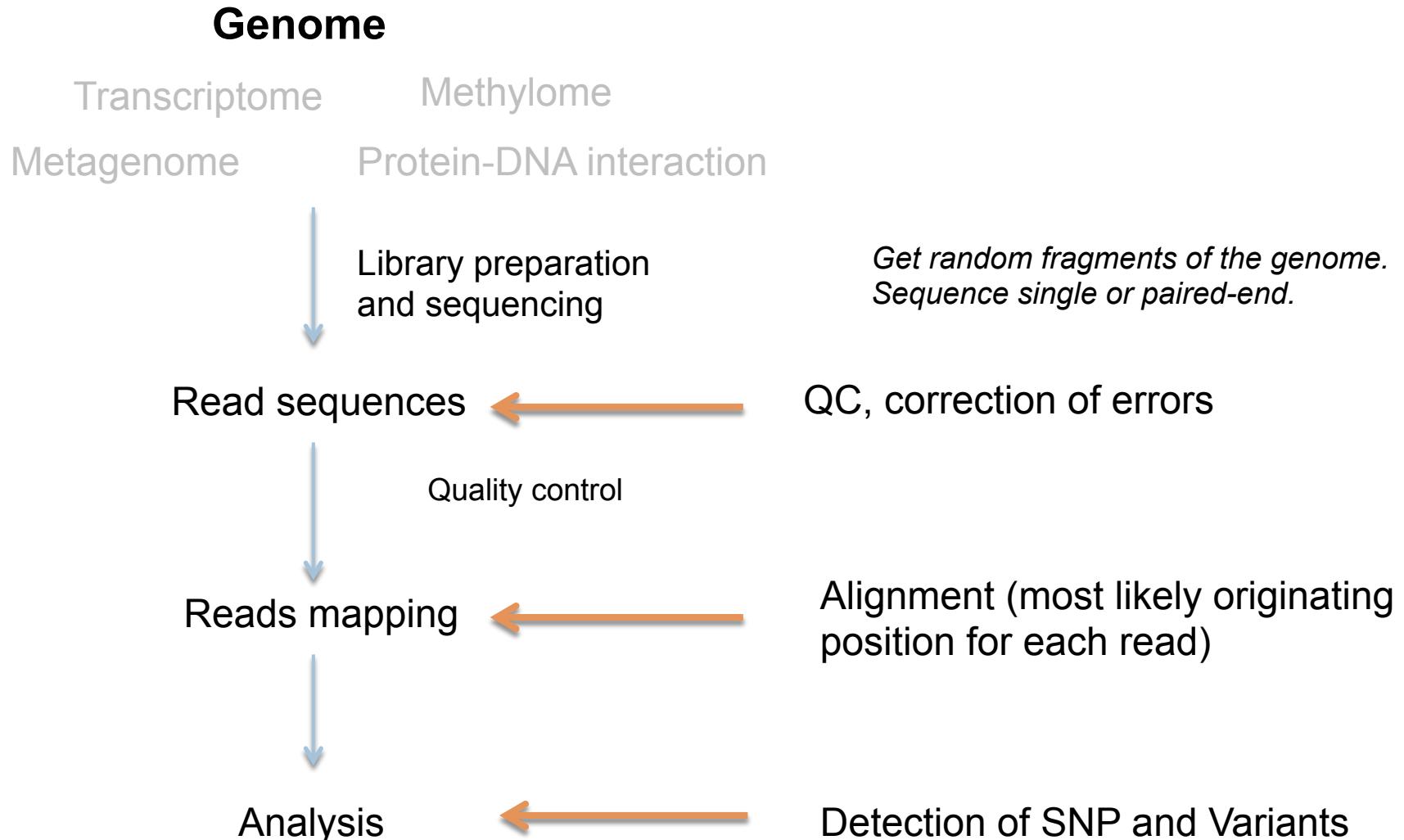
Outline

- Brief historical background on sequencing

- A typical resequencing workflow
 - Reads alignment
 - Basic statistics after alignment
 - Automatic correction of sequencing errors
 - Detection of sequence polymorphisms
 - Genome Assembly

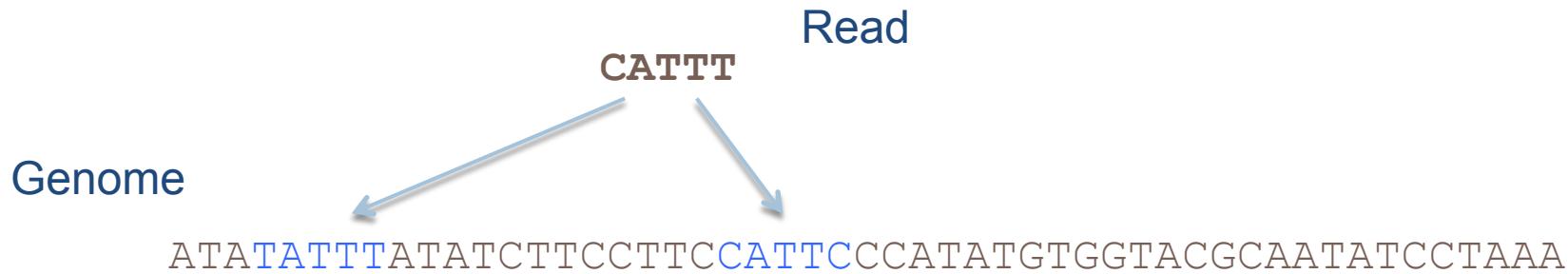


NGS workflow



Question 1: read alignment

- We sequenced a read. Where should we place it on the genome ?



- Genome is available
 - Recall pairwise alignment ?
 - Has to be done in short time.



Time is money

- Evaluating by dynamic programming would cost $k \cdot n$ operations for each read.
 - Nowadays need to align ~500Mio reads / exp.
- Need to do cost saving cuts by either:
 - Index the genome with an appropriate data structure
 - Hash table of length k
 - Suffix tree of the genome
 - For each read, filter the set of possible starting positions with the data structure.



High throughput read alignment

- Usually semi-global alignment
 - Sometimes allow soft trimming at one/both end
- Many aligners exist now.

Program	Algorithm	SOLiD	Long ^a	Gapped	PE ^b	Q ^c
Bfast	hashing ref.	Yes	No	Yes	Yes	No
Bowtie	FM-index	Yes	No	No	Yes	Yes
BWA	FM-index	Yes ^d	Yes ^e	Yes	Yes	No
MAQ	hashing reads	Yes	No	Yes ^f	Yes	Yes
Mosaik	hashing ref.	Yes	Yes	Yes	Yes	No
Novoalign ^g	hashing ref.	No	No	Yes	Yes	Yes

^aWork well for Sanger and 454 reads, allowing gaps and clipping.

^bPaired end mapping. ^cMake use of base quality in alignment. ^dBWA

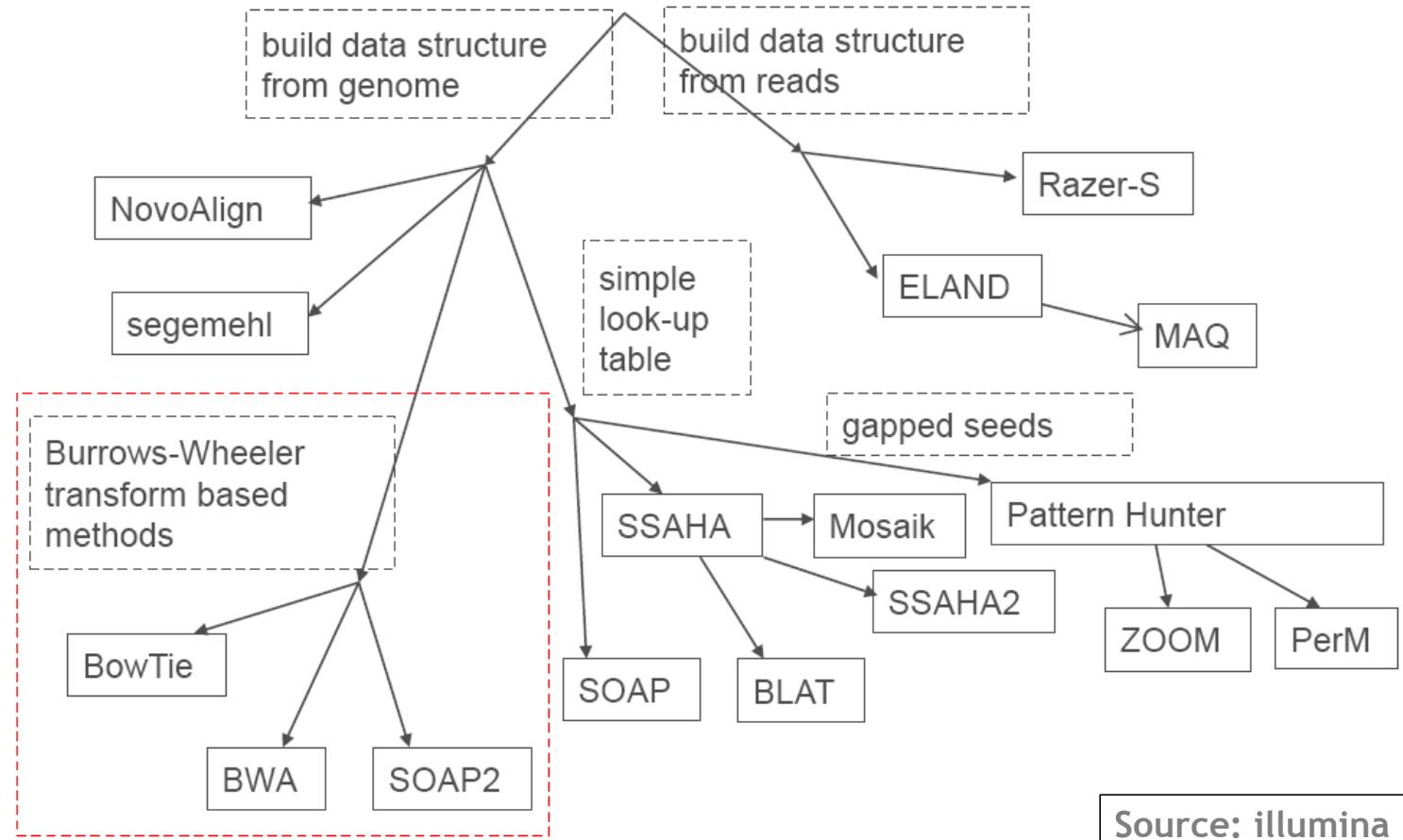
trims the primer base and the first color for a color read. ^eLong-read alignment implemented in the BWA-SW module. ^fMAQ only does

gapped alignment for Illumina paired-end reads. ^gFree executable for

non-profit projects only.



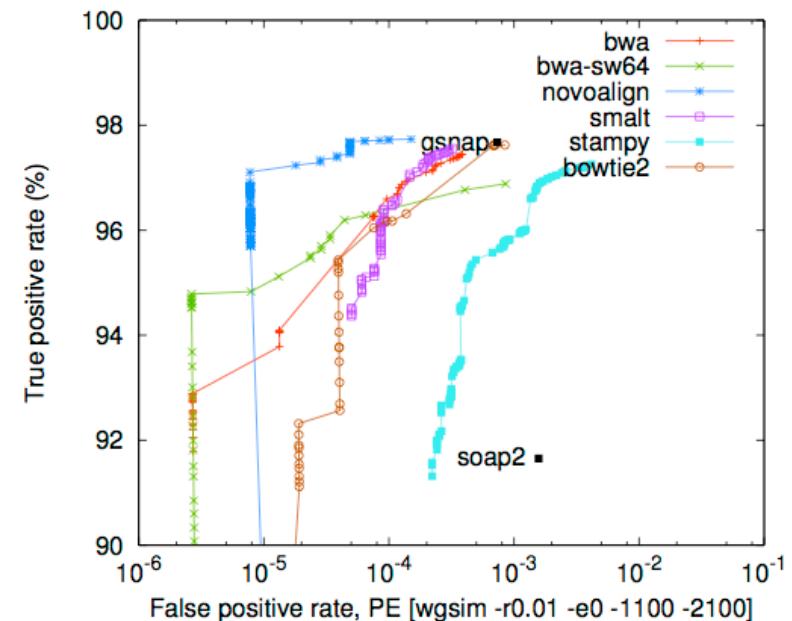
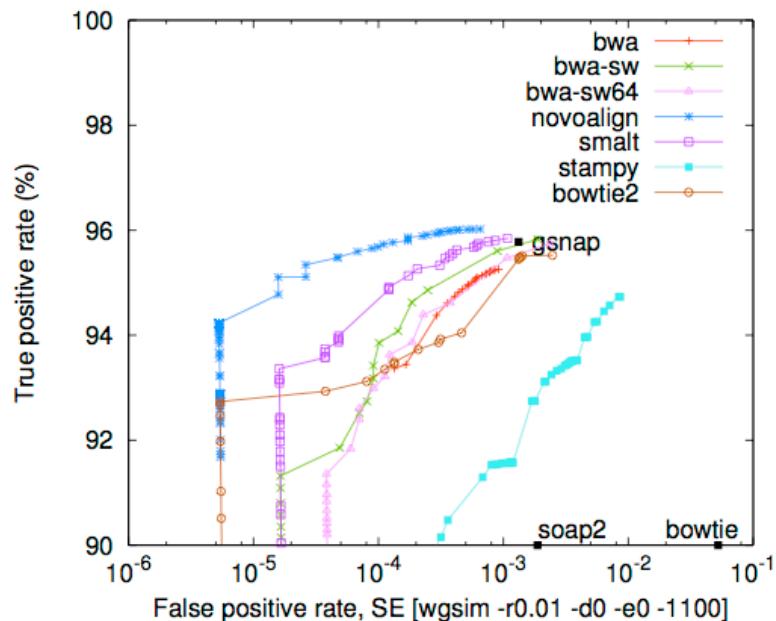
Read mapping tools



Last generation of mappers combine genome and read indexing (Massai, OLeggo)



Comparison read mappers



from Heng Li (BWA) benchmark page: <http://lh3lh3.users.sourceforge.net/alnROC.shtml>
 100k reads simulated with wgsim **without** sequencing errors

An alignment result

Sequence coverage
(# of reads covering a position)

SNPs

```

ATCCTGATTGGTGAACGTTATCGACGATCCGATCGA
ATCCTGATTGGTGAACGTTATCGACGATCCGATCGA
    CGGTGAACGTTATCGACGATCCGATCGAACTGTCAGC
    GGTGAACGTTATCGACGTTCCGATCGAACTGTCAGCG
    TGAACGTTATCGACGTTCCGATCGAACTGTCAGCGGC
    TGAACGTTATCGACGTTCCGATCGAACTGTCAGCGGC
    TGAACGTTATCGACGTTCCGATCGAACTGTCAGCGGC
        GTTATCGACGATCCGATCGAACTGTCAGCGCAAGCT
        TTATCGACGATCCGATCGAACTGTCAGCGCAAGCT

```

```
ATCCTGATTGGTGAACGTTATCGACGATCCGATCGAACTGTCAGCGCAAGCTGATCGATCGATGCTAGTG
```

reference genome

```

TTATCGACGATCCGATCGAACTGTCAGCGCAAGCT
    TCGACGATCCGATCGAACTGTCAGCGCAAGCTGATCG
        ATCCGATCGAACTGTCAGCGCAAGCTGATCG    CGAT
        TCCGATCGAACTGTCAGCGCAAGCTGATCG    CGATC
        TCCGATCGAACTGTCAGCGCAAGCTGATCGATCGA
            GATCGAACTGTCAGCGCAAGCTGATCG    CGATCGA
            AACTGTCAGCGCAAGCTGATCG    CGATCGATGCTA
            TGTCAGCGCAAGCTGATCGATCGATCGATGCTAG
            TCAGCGCAAGCTGATCGATCGATCGATGCTAGTG

```

INDELS



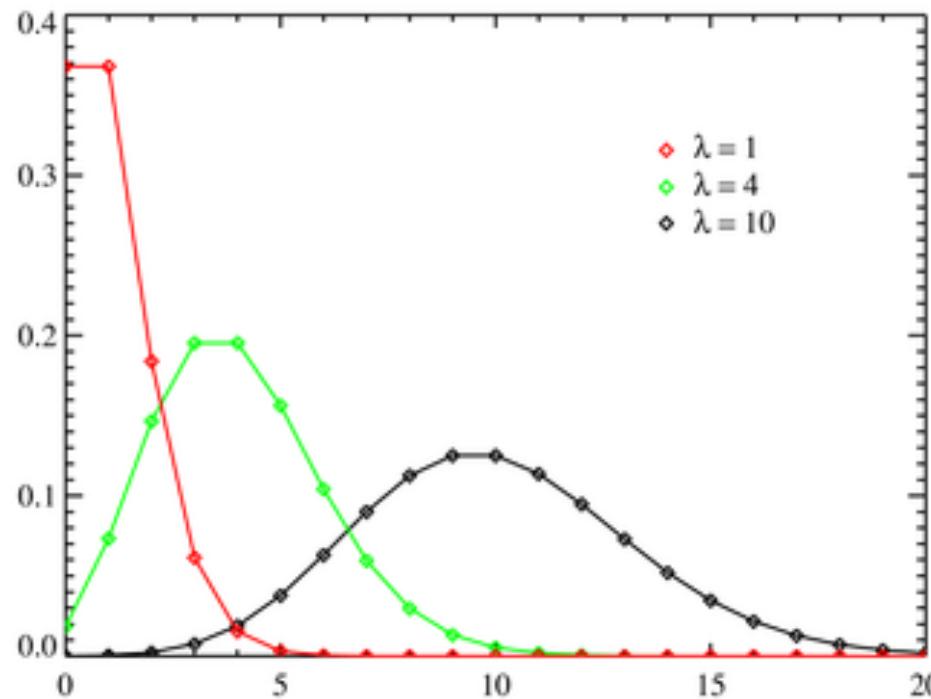
Question 2: Sequence Coverage ?

- Every read sample a position along the genome.
 - Coverage: Number of reads overlapping a position
 - What is the expected distribution ?
- Suppose uniform sampling.



Poisson Distribution

- Only one parameter, the expected coverage λ
- Variance and expectation are equal

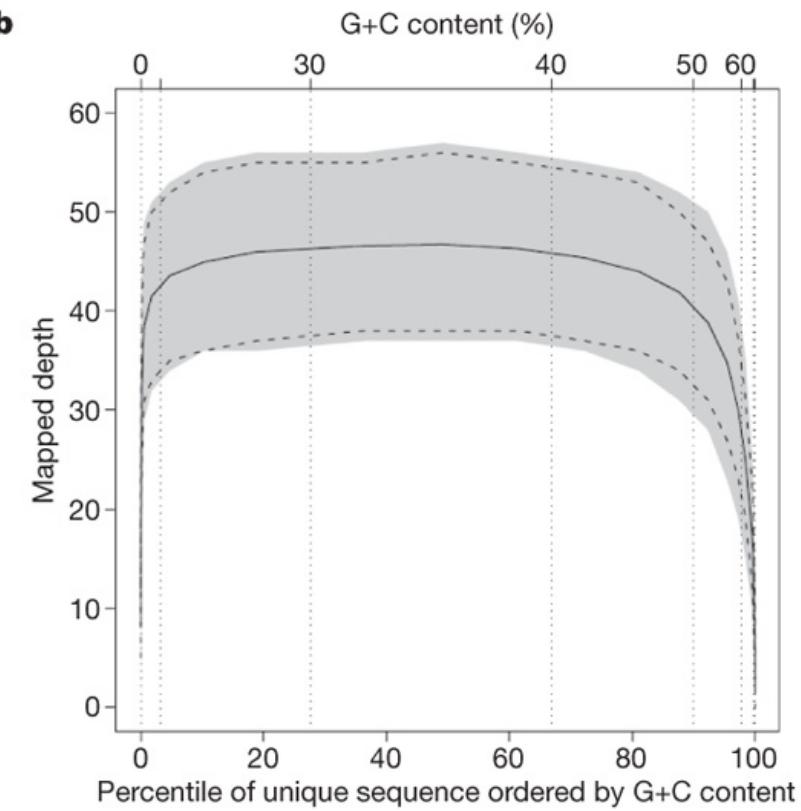
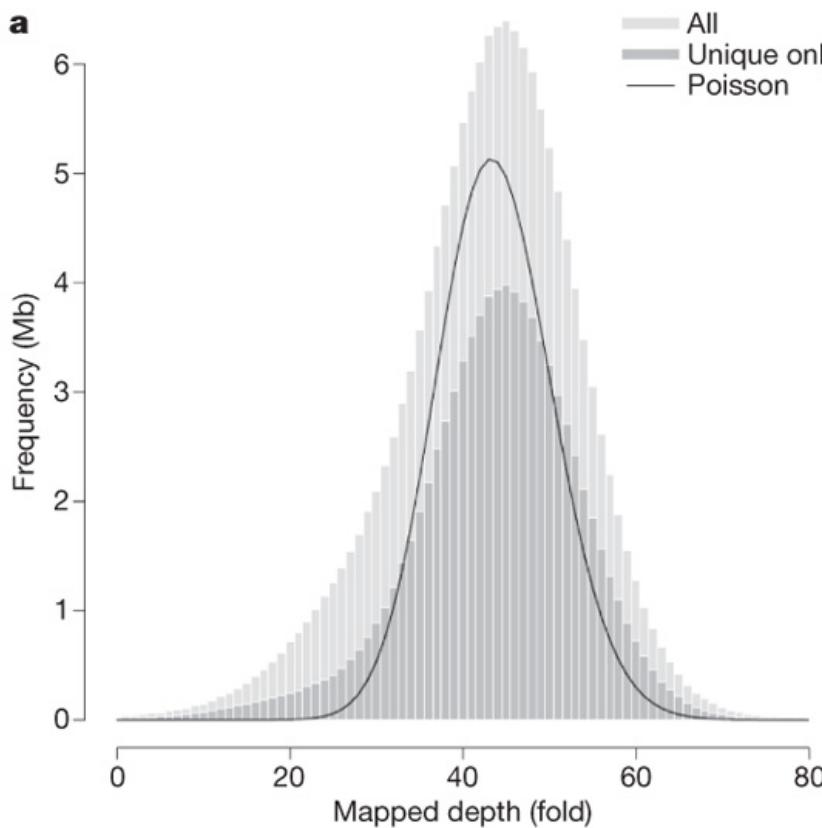


Probability distribution function for various parameters of λ



In practice

- Poisson assumes uniform sampling, but multiple artifacts in the data (also: PCR duplicates)



Bentley et al. Nature 2008

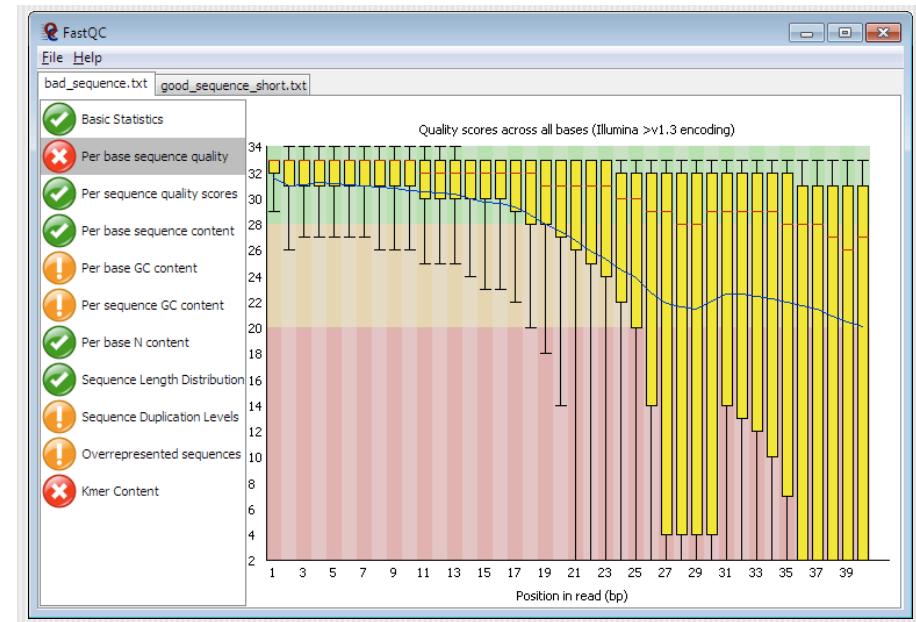


Quality evaluation

- Sequence Quality score
 - reported by the sequencer

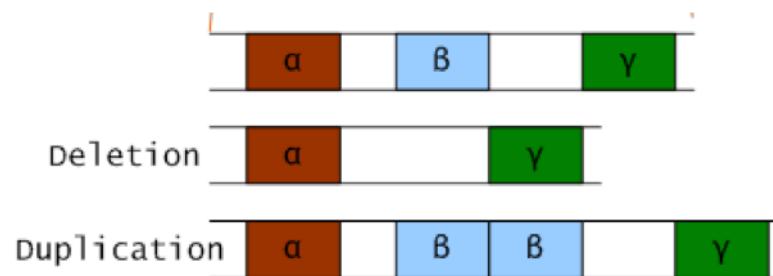
Phred score: $-10 \cdot \log_{10}(P_{\text{err}})$

- % mapped reads
- coverage distribution



Variant Calling with a reference

- Single Nucleotide Variants (SNV)
 - 1 bp difference to the reference
- Short insertions/deletions (indels)
 - <50 bp variation to the reference
- Structural variants



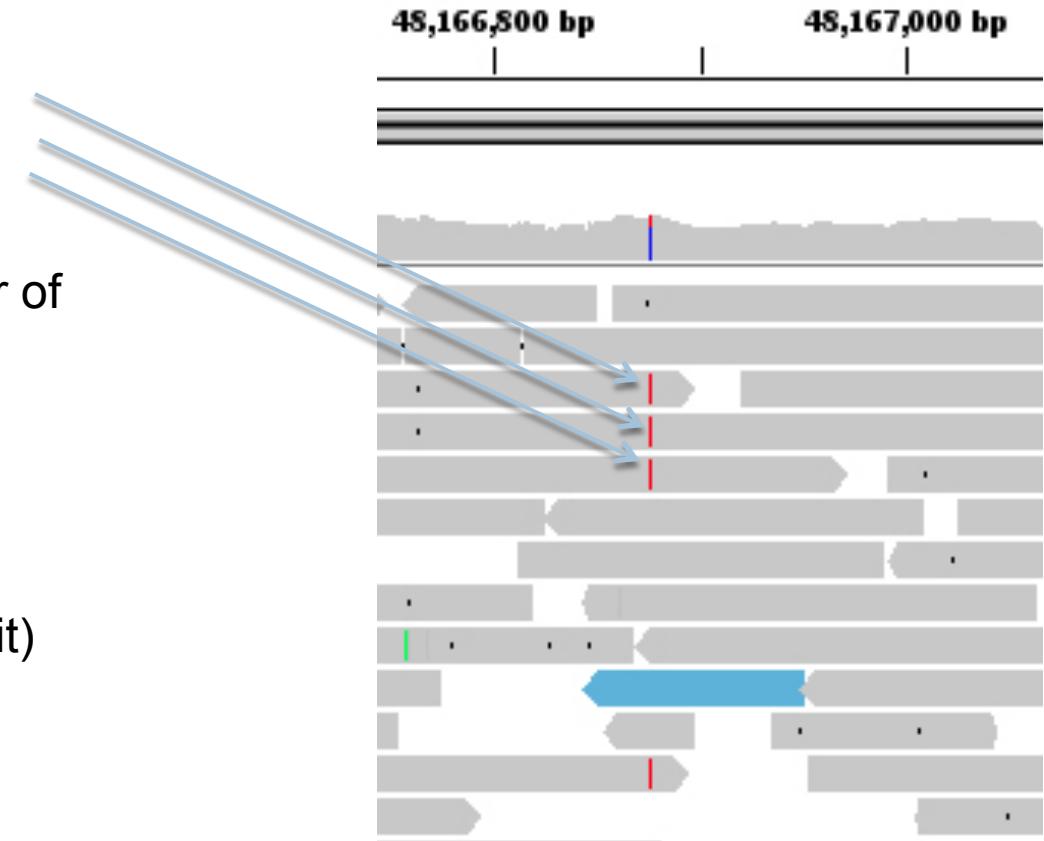
Detection of SNV

Are these genuine or due to sequencing errors ?

Hypothesis testing on the number of reads supporting the SNV vs expected error rate.

Tools:

- GATK (Genome analysis toolkit)
- SAMtools mpileup
- VarScan



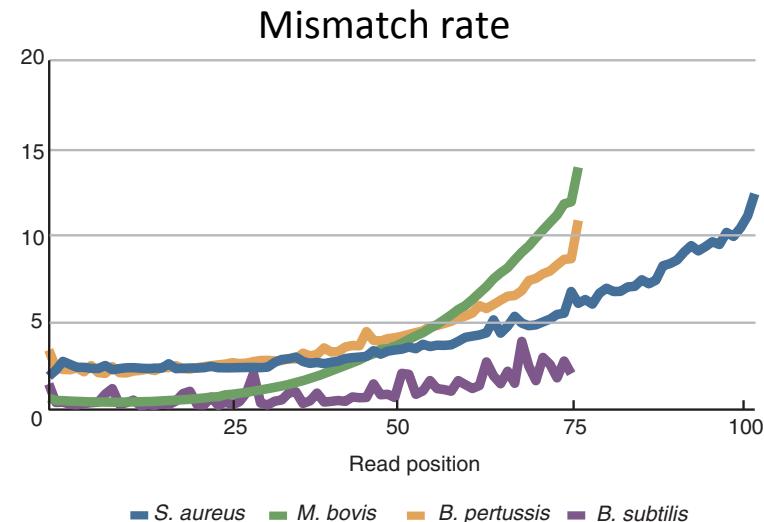
SNV annotation

- Need to be annotated w.r. to genomic information
 - Homo/heterozygous
 - coding/non-coding SNVs
 - Known/unknown (dbSNP, 1k genome)
- Functional implication of a SNV
 - Silent/Non-silent mutation (within CDS)
 - Predicting possible impact of an amino acid substitution on the protein stability (Sift, Polyphen2...)



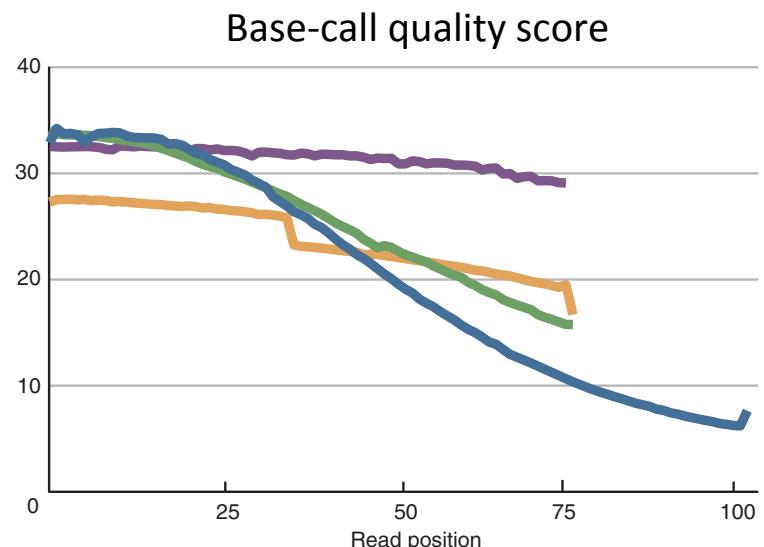
Illumina error profile

- Biases found by Dohm *et al.* (2008):
 - **Mismatches** are dominant errors
 - Positional error profile with an increase towards the 3' end
 - non-uniform substitution matrix
 - 8 times more A→C errors compared to A→G or A→T errors



- Base-call qualities give a hint about errors (Nakamura *et al.*, 2011)

Phred score: $-10 \cdot \log_{10}(P_{\text{err}})$



Errors in other Technologies

- Other platforms are dominated by indel errors due to:
 - Indistinguishable signals for homopolymers (>6bp)
 - Undetected low-intensity signals (single molecule sequencing)

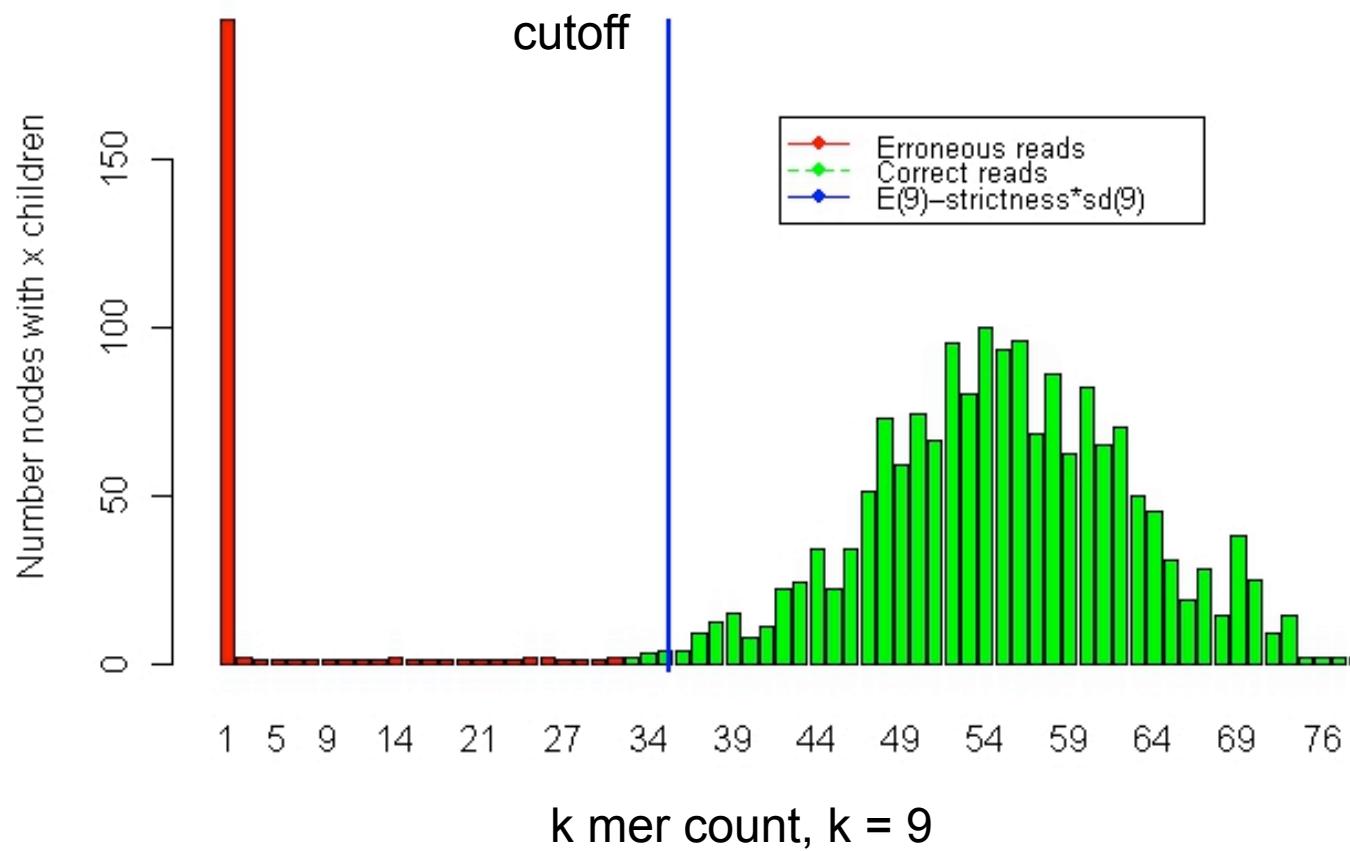
Table I: Characteristics of major NGS platforms as of February 2012

Company	Platform	Read length (bp)	Throughput & time per run	Technique	Dominant error type
Illumina	HiSeq 2000 ^a	36 50 100	105–600 Gb 2–11 days	Reversible terminator	Substitution
Applied Biosystems	5500 SOLiD TM System ^b	356 075	7–9 Gb/day	Sequencing by ligation	–
Complete Genomics		35	–	Ligation based	–
Helicos BioSciences	HeliScope SMS ^c	25–55	21–35 Gb	Single molecule sequencing	Insertion Deletion
454 Life Sciences	GS FLX Titanium XL ^d	≤1000	700 Mb 23 h	Sequencing by synthesis	Insertion Deletion
Ion Torrent	Ion PGM Sequencer 318 ^e	>200	>1 Gb 2 h	Ion semiconductor sequencing	Insertion Deletion
Pacific Biosciences	PacBio RS	1 k–10 k	–	Single molecule sequencing	Insertion Deletion

(Yang *et al.*, 2012)

Detecting Potential Errors

- reads with errors will have a low coverage.



- Automatic correction of sequencing errors without alignment to the genome
→ cf other slides

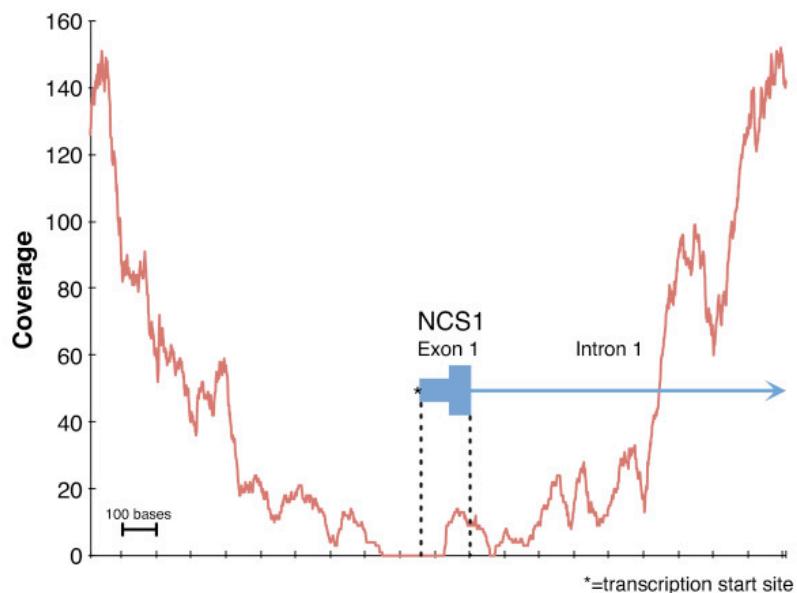


Performance of Genome resequencing

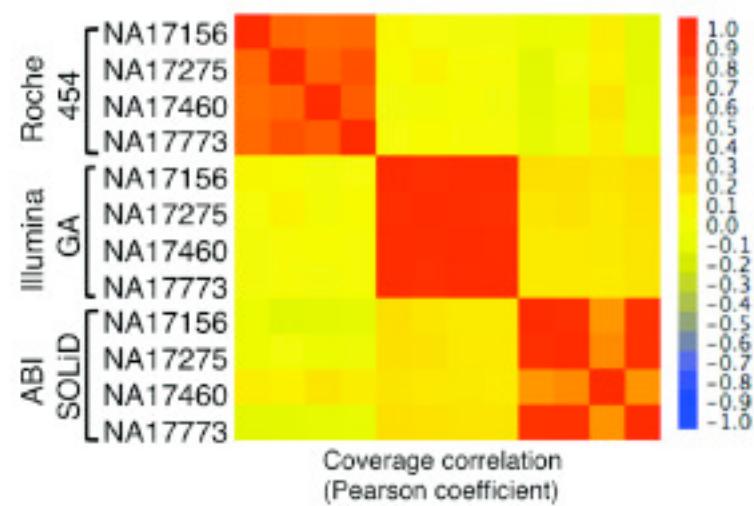
- Depends on artifacts/biases:
 - Each technology/preparation has its own
- Two main artifacts:
 - Sequencing bias (quality score, error correction)
 - Coverage bias
 - Can be related to GC composition
 - Some regions with consistent lower coverage cannot be genotyped
- Assess by comparing multiple technologies with Sanger sequencing



Ross et al. Genome biology 2013

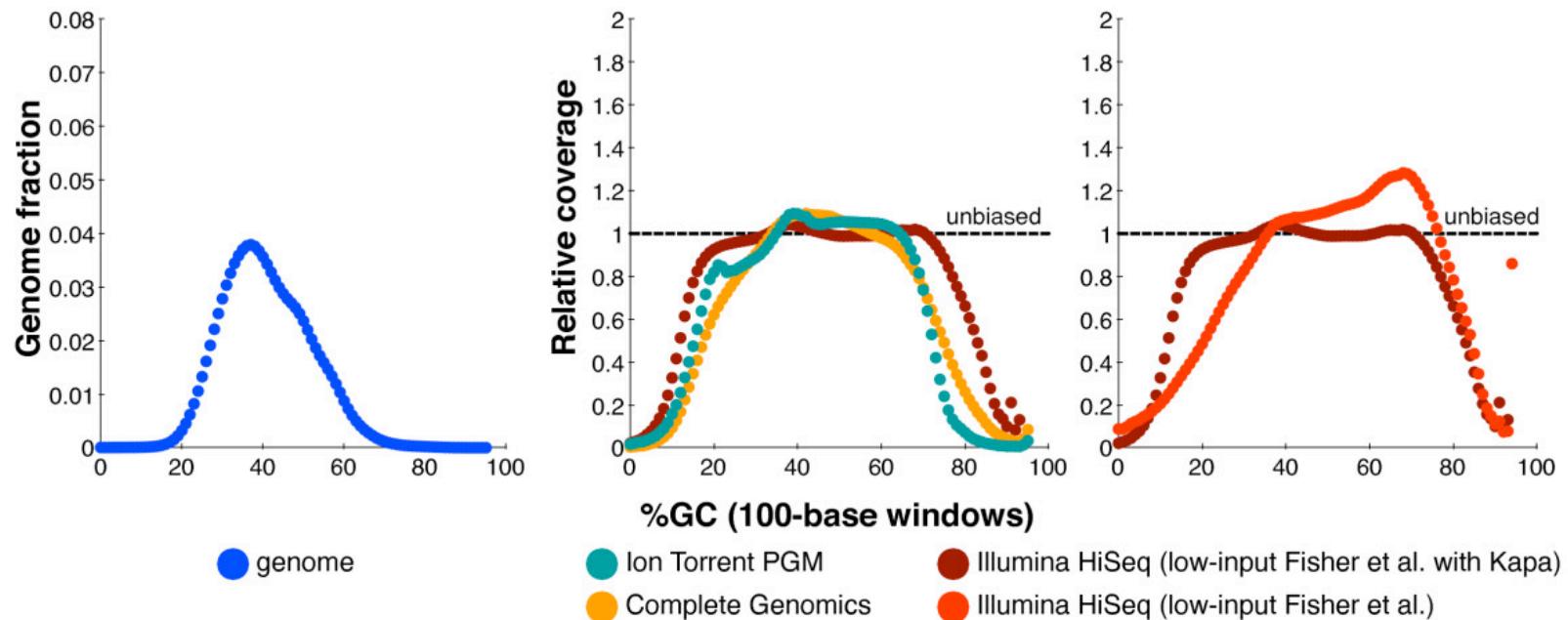


Harismendy et al. Genome biology 2009



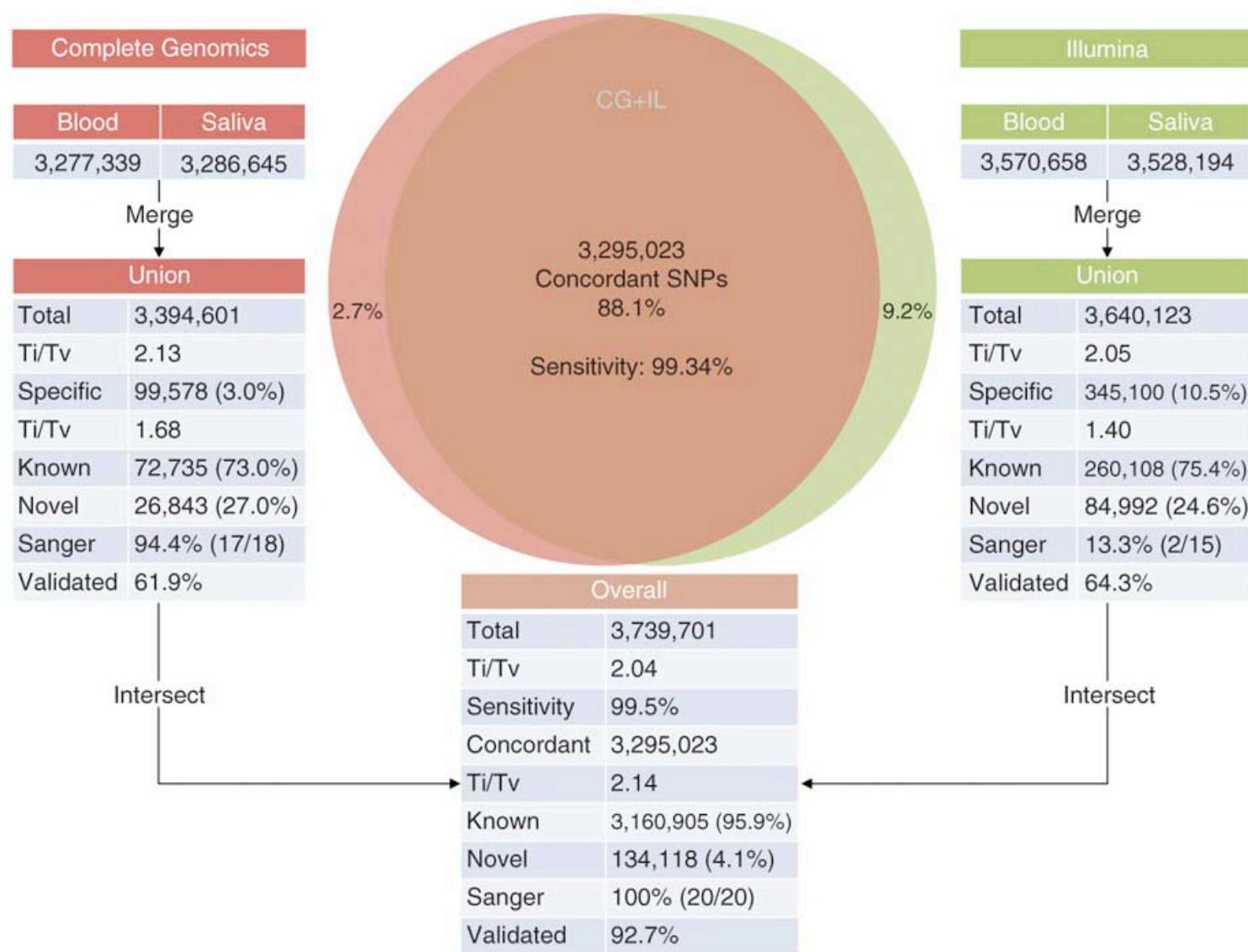
- Some region will never be covered enough for SNV calling
- Frequently linked to high/low GC

- Coverage tend to be correlated with technology (here targeted resequencing)

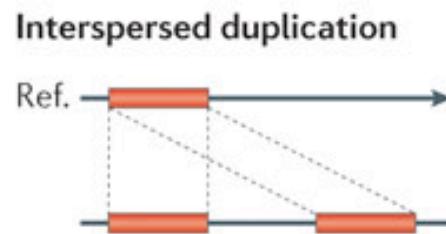
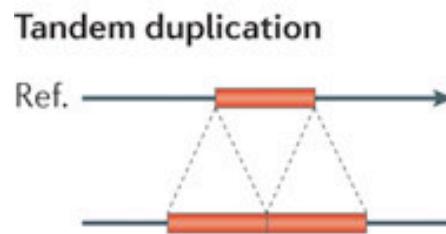
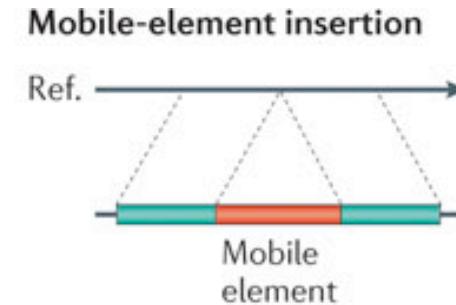
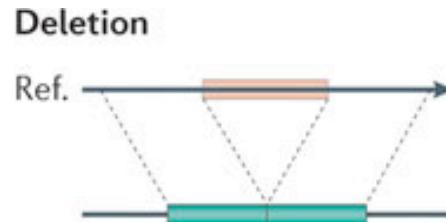


Ross et al. Genome biology 2013

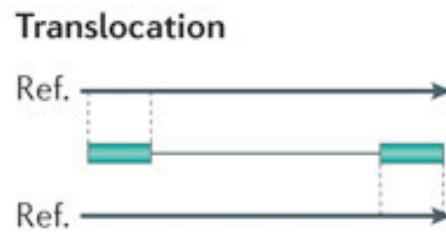
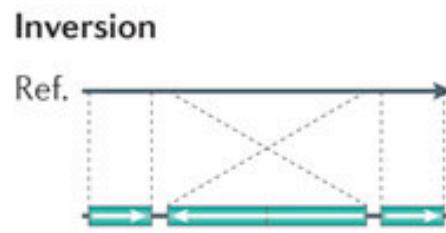


a

Annotation of structural variants (SV)



1k GP: 14k deletions
(38 Mio SNV total)



Also annotation of cancer specific SV

Nature Reviews | Genetics

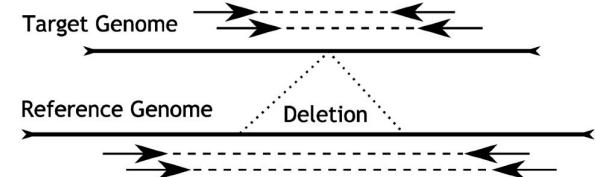
Alkan et al 2011



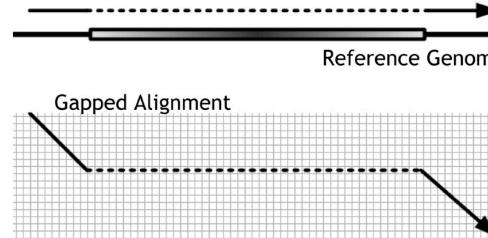
SV detection methods

- Different information can be used (or merged)

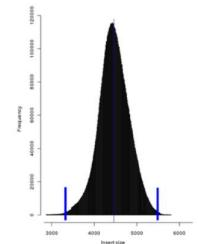
- Paired-end Mapping



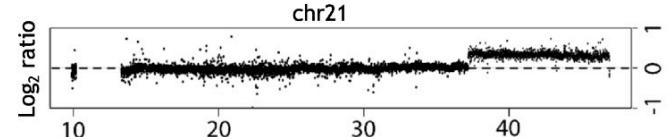
- Split-read alignment



Insert Size Distribution



- Read depth analysis



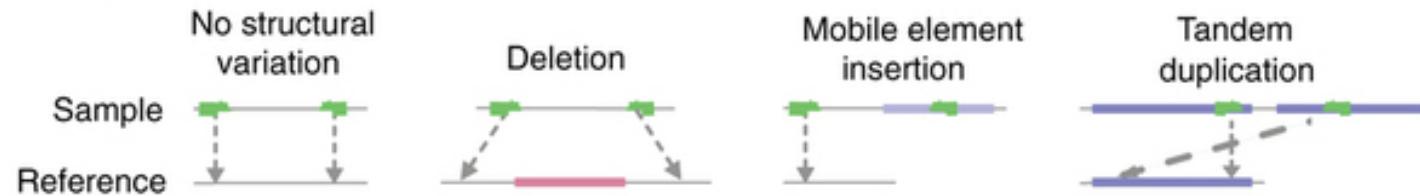
$$\log_2 \frac{\# \text{Reads}_{\text{Diseased}}}{\# \text{Reads}_{\text{Normal}}}$$

- Assembly and mapping

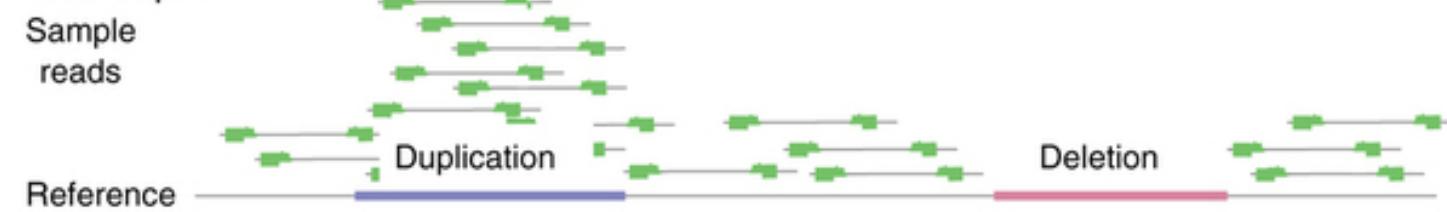


Detection of variants (with Paired end)

Read pairs



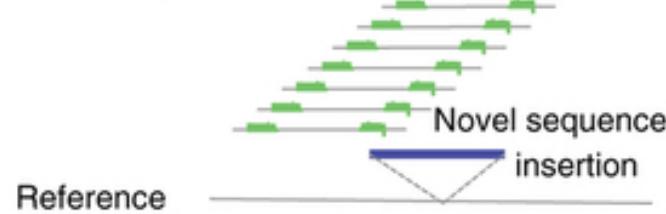
Read depth

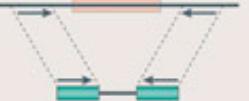
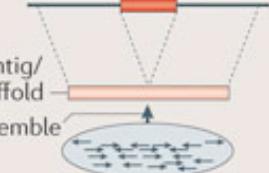
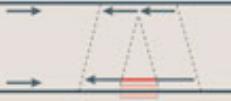
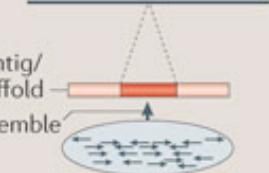
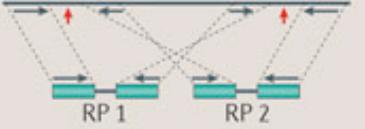
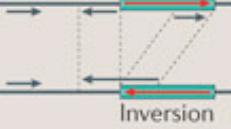
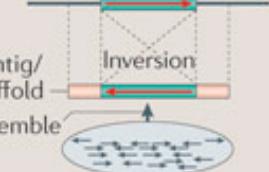
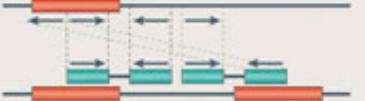
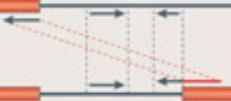
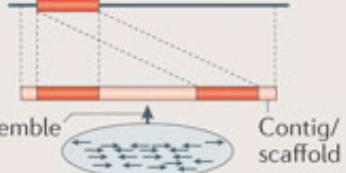
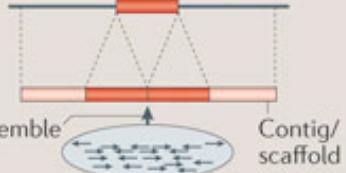


Split reads

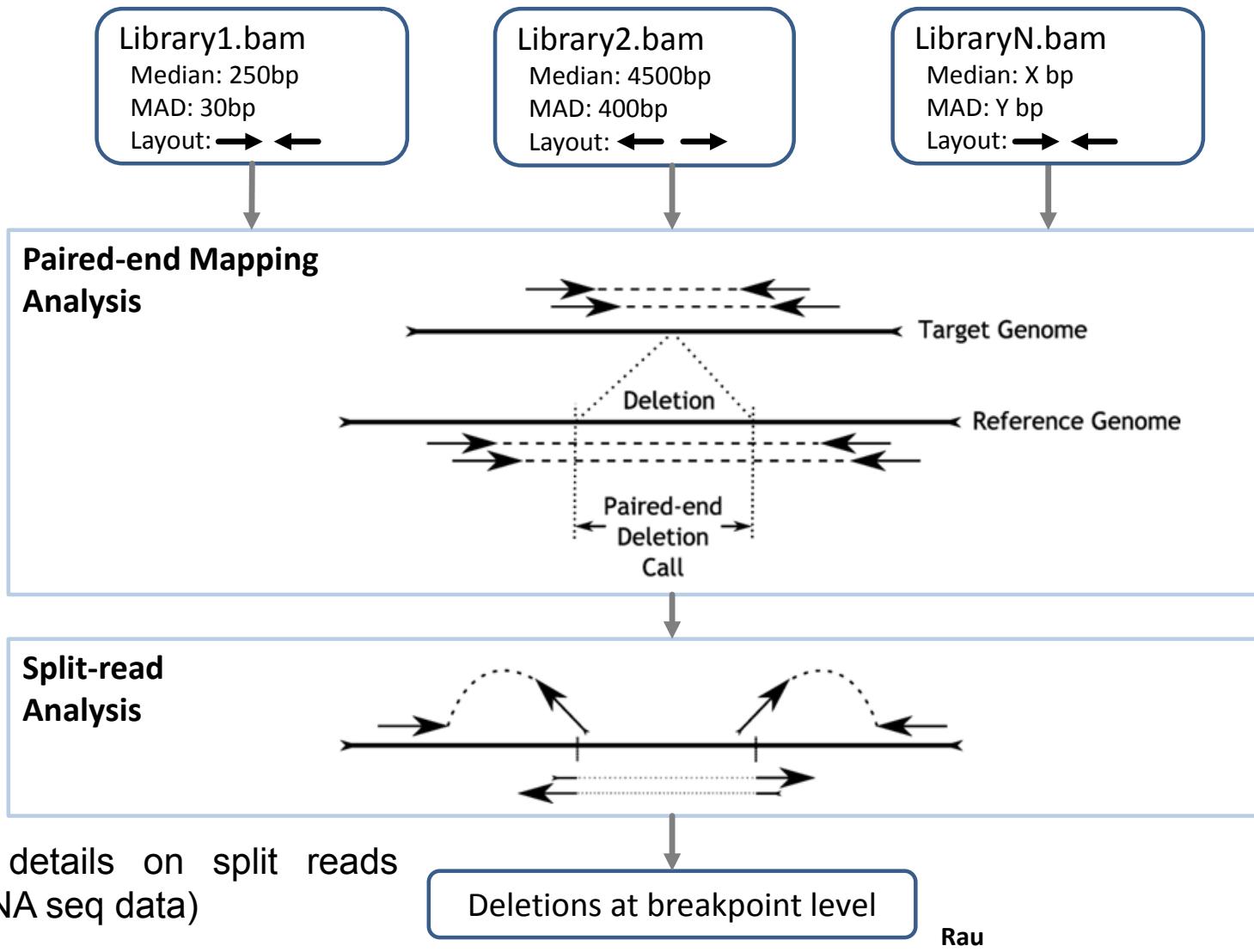


Assembly



SV classes	Read pair	Read depth	Split read	Assembly
Deletion				
Novel sequence insertion		Not applicable		
Inversion		Not applicable		
Interspersed duplication				
Tandem duplication				

Delly: Combine PE and split read analysis



Outline

- Brief historical background on sequencing
- A typical resequencing workflow
 - Reads alignment
 - Basic statistics after alignment
 - Automatic correction of sequencing errors
 - Detection of sequence polymorphisms (A. Gillet)
 - **Genome Assembly**



Sequence Assembly

- Assemble:

- reconstruct the large sequence (the genome) by merging and ordering the reads from the sequencing experiments
- need to summarize overlap between reads in a data structure to reconstruct contiguous sequences (*contigs*)

- Assembly is:

- trivial if repeats are shorter than read length
- computationally intractable with too short reads
(sequencing errors play a role)

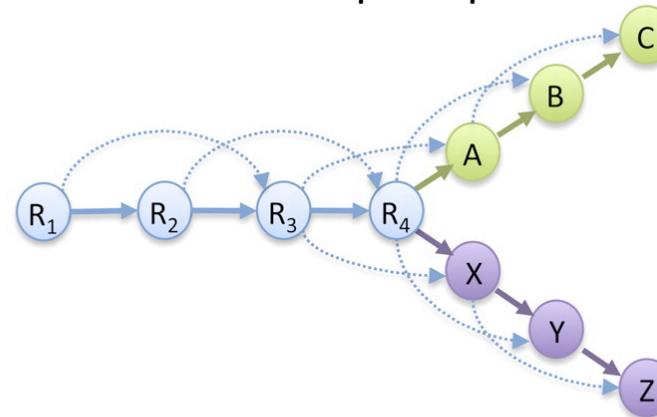


Assembly paradigms

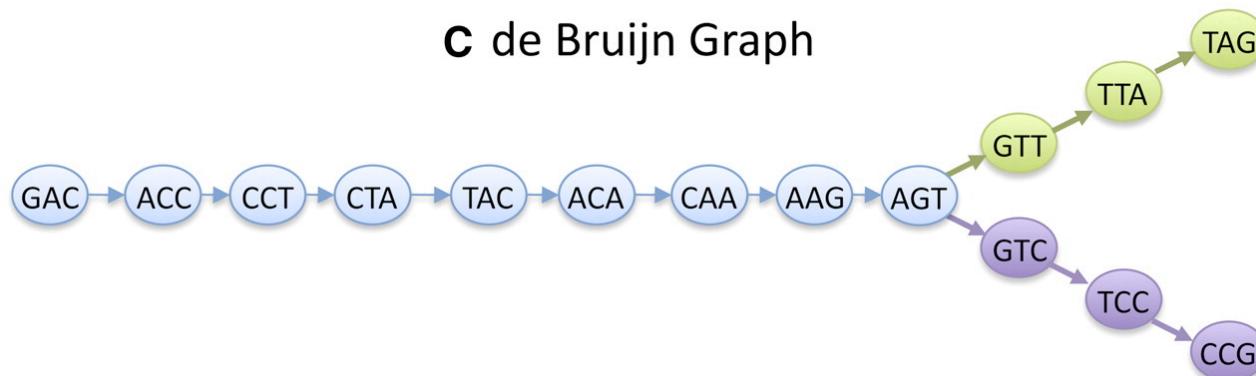
A Read Layout

R ₁ :	GACCTACA
R ₂ :	ACCTACAA
R ₃ :	CCTACAAG
R ₄ :	CTACAAGT
A:	TACAAGTT
B:	ACAAGTTA
C:	CAAGTTAG
X:	TACAAGTC
Y:	ACAAGTCC
Z:	CAAGTCCG

B Overlap Graph

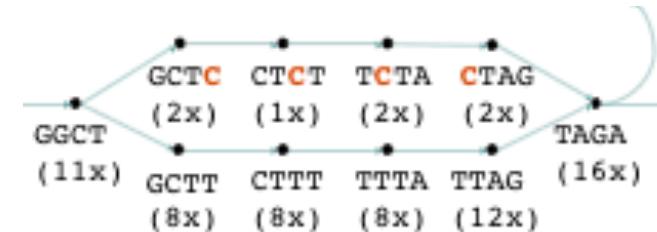


C de Bruijn Graph

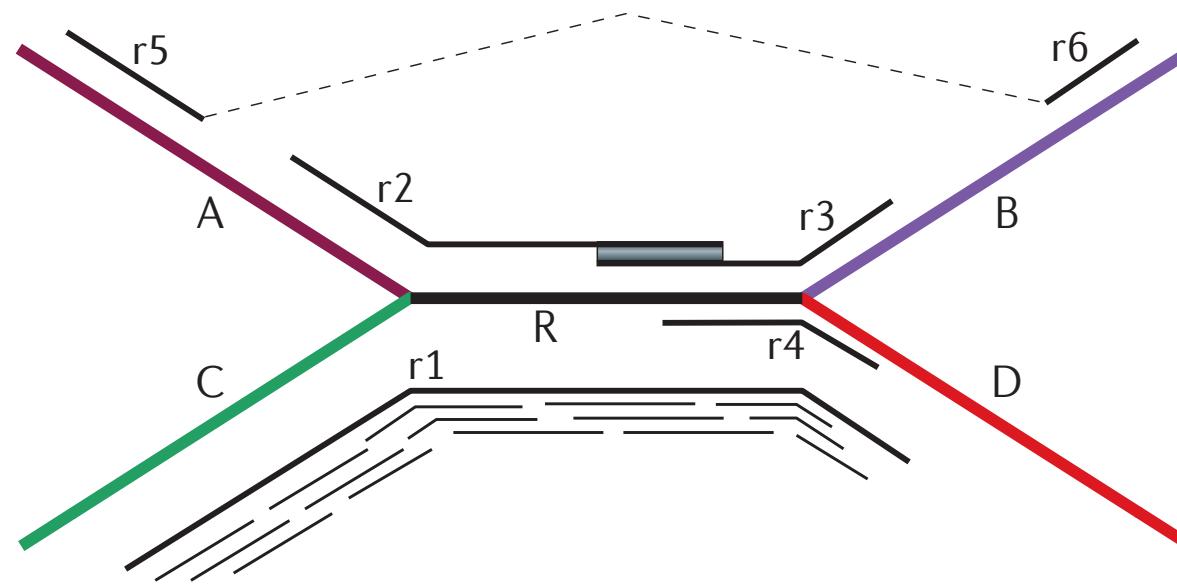


Genome assembly paradigms

- Overlap Layout consensus and string graph
 - repeats are “perfectly” summarized
 - Computationally costly to do all pairwise alignments (optimized with FM-index/burrows wheeler, SGA assembler recently)
- De Bruijn Graph on the set of sequences
 - easy to do a hash map of all reads
 - repeats needs to be transitively reduced
 - errors impact the graph structure



Effect of sequencing parameters



- Possible reconstructions: ARB+CRD or ARD+CRB
- some reads are longer than the repeat region R



Difficulties with assembly

- **Repeats** create cycles/branching
 - need longer reads
 - Use paired-end and mate pair information
- **Sequencing errors**
 - increase the size of the DB graph (Tips and bubbles)
 - need to be lower than intra genomic variations
- Need to **experimentally design** the sequencing experiment (sequence recipes from ALLPATH-LG) according to the genome characteristics and sequencing experiment.



TAGTCGAGGCTTTAGATCCGATGAGGCTTTAGAGACAG
 AGTCGAG CTTTAGA CGATGAG CTTTAGA
 GTCGAGG TTAGATC ATGAGGC GAGACAG
 GAGGCT**C** ATCCGAT AGGCTTT GAGACAG
 AGTCGAG TAGATCC ATGAGGC TAGAGAA
 TAGTCGA CTTTAGA CCGATGA TTAGAGA
 CGAGGCT AGATCCG TGAGGCT AGAGACA
 TAGTCGA GCTTTAG TCCGATG GCT**C**TAG
 TCGAC**CGC** GATCCGA GAGGCTT AGAGACA
 TAGTCGA TTAGATC GATGAGG TTTAGAG
 GTCGAGG **T**CTAGAT ATGAGGC TAGAGAC
 AGGCTTT ATCCGAT AGGCTTT GAGACAG
 AGTCGAG TTAGATT ATGAGGC AGAGACA
 GGCTTTA TCCGATG TTTAGAG
 CGAGGCT TAGATCC TGAGGCT GAGACAG
 AGTCGAG TTTAGATC ATGAGGC TTAGAGA
 GAGGCTT GATCCGA GAGGCTT GAGACAG

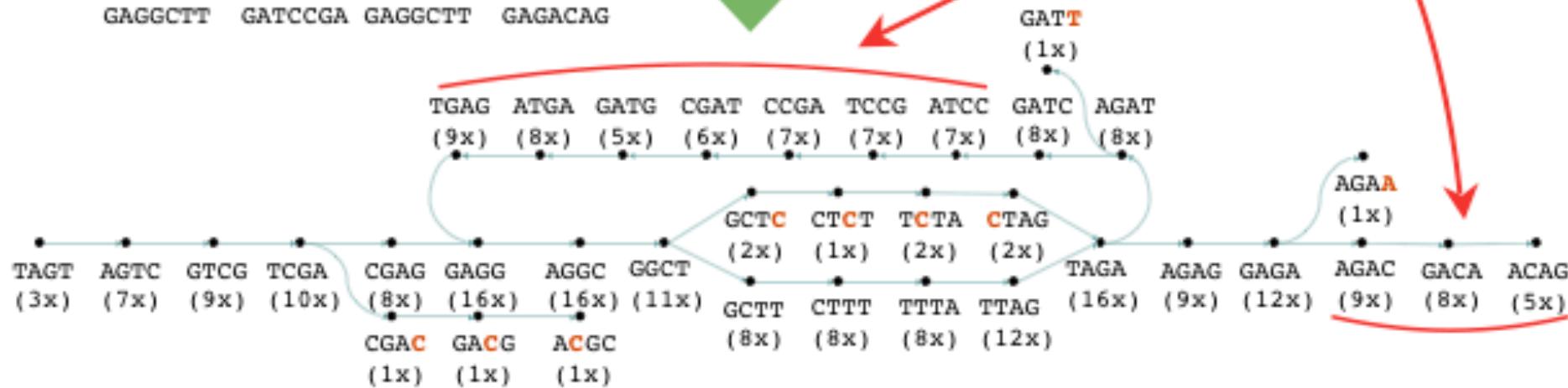


1. Sequencing
(e.g. Solexa, 454...))



2. Hashing

Linear stretches



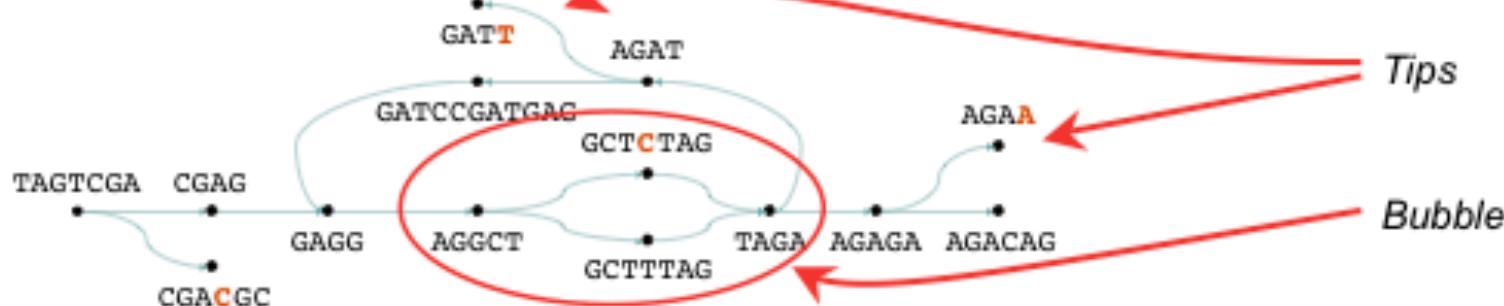
Velvet strategy (Zerbino et al. 2008)

GGCTTTA TCCGATG TTTAGAG
 CGAGGCT TAGATCC TGAGGCT GAGACAG
 AGTCGAG TTTAGATC ATGAGGC TTAGAGA
 GAGGCTT GATCCGA GAGGCTT GAGACAG

2. Hashing



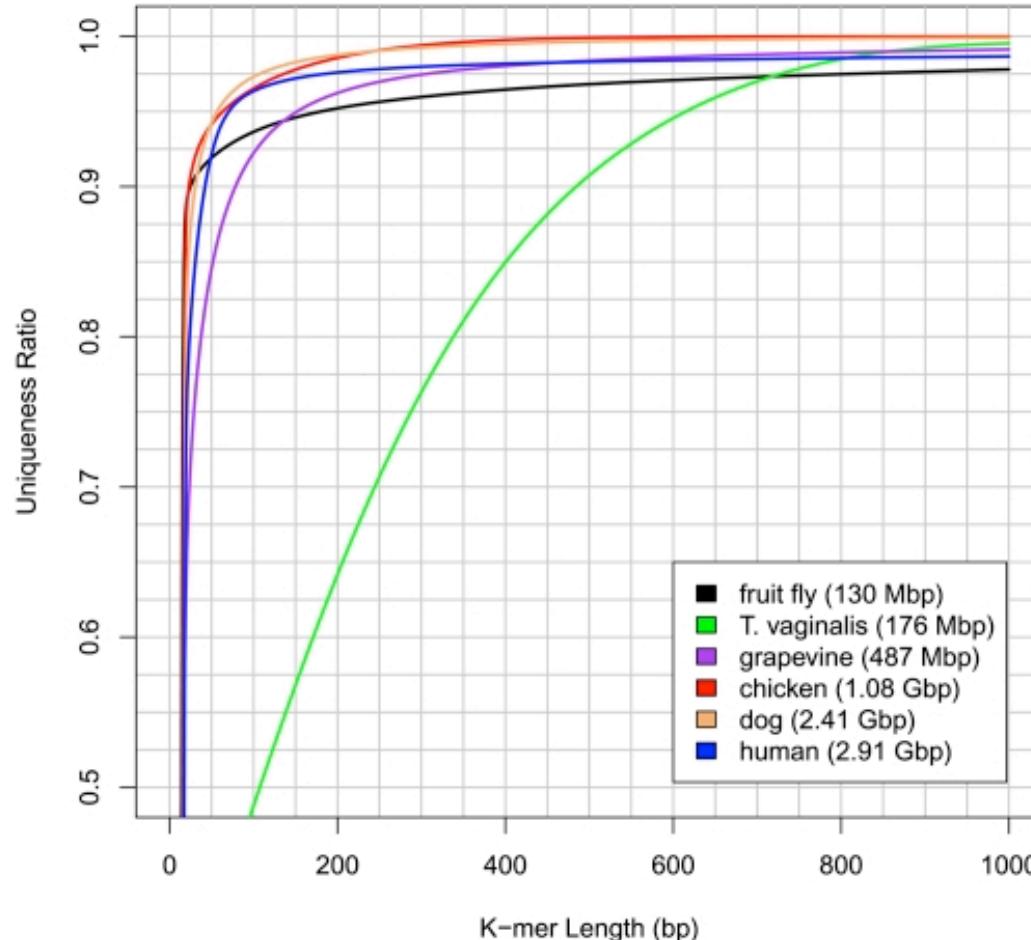
3. Simplification of linear stretches



4. Error removal



Repeat content

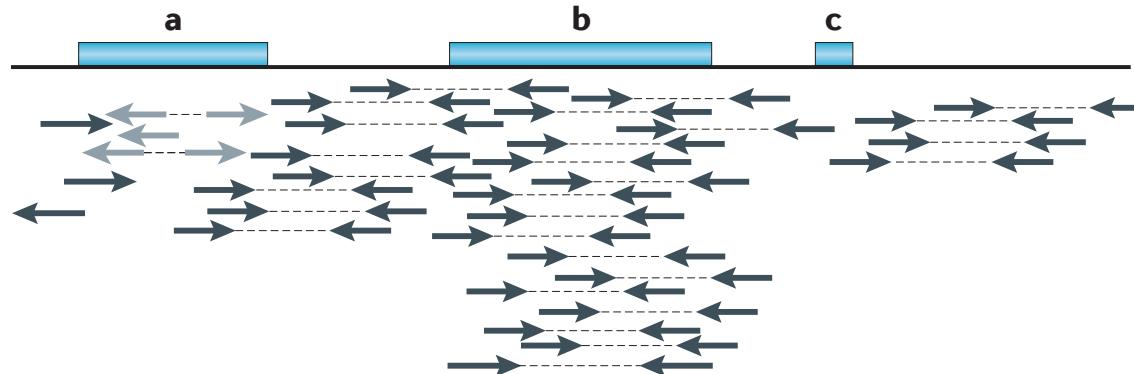


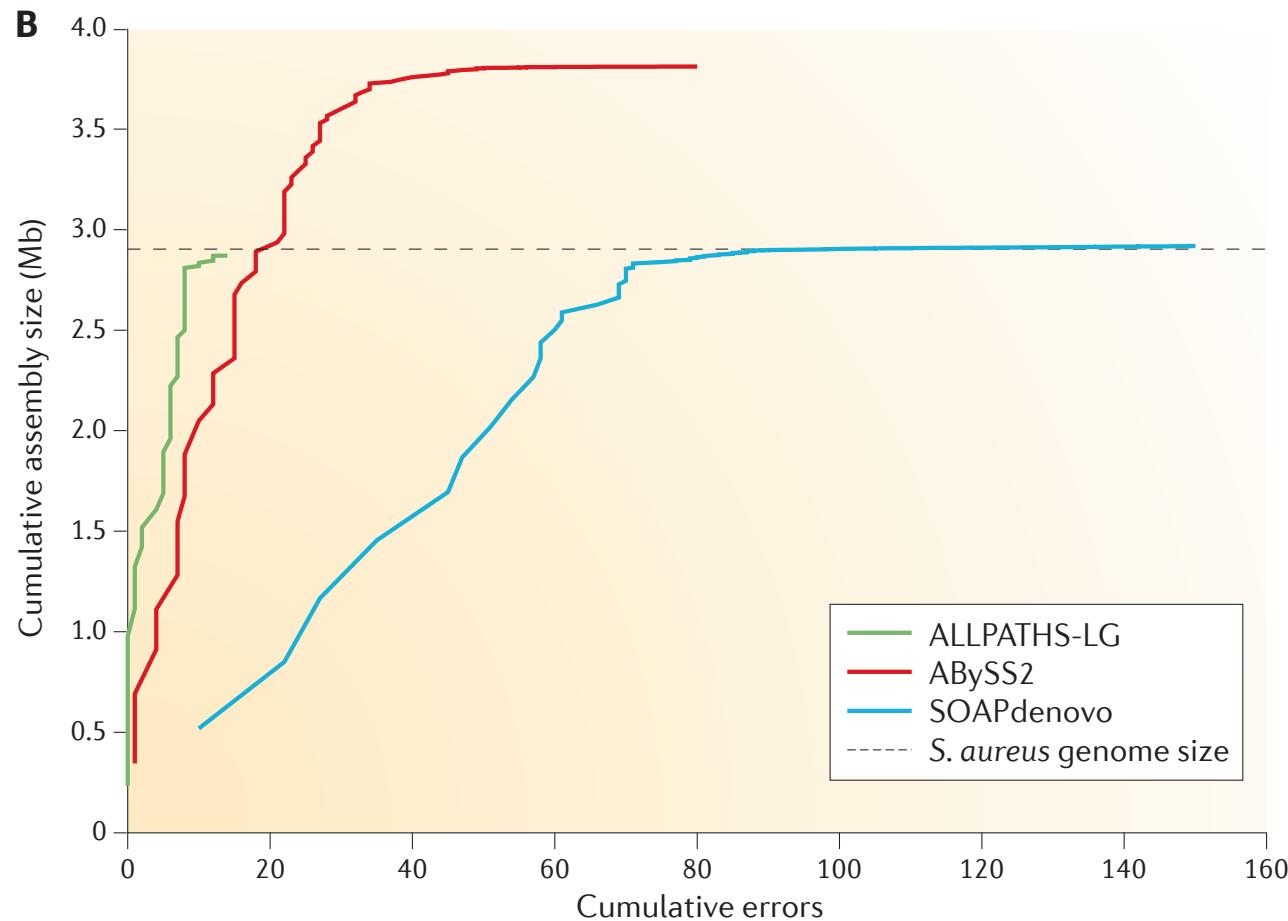
Uniqueness ratio: percentage of the genome that is covered by unique sequences of length k or longer



Assembly validation

- Measure contiguity of assembled sequences
 - total size and number of contigs
 - N50: weighted median contig size
- Evaluate assembly errors
 - compare to BAC sequences, transcriptome, closely related organism
 - consistency after aligning the reads





THE ASSEMBLATHON

- Competitive assessment of *de novo* assembly
 - Assemblathon3 plans to work with limited ressource
 - 20k\$ to each team
 - Use money to buy sequencing
- Assemblathon1:
 - Evaluate on a simulated diploid genome (from chr13)
 - Setting up metrics is not an easy task:
N50-NG50-CPNG50-CC50
- Evaluation of classical tools:
 - Allpath-LG, Velvet, Euler-SR, Ray, SOAPdenovo, Abyss...



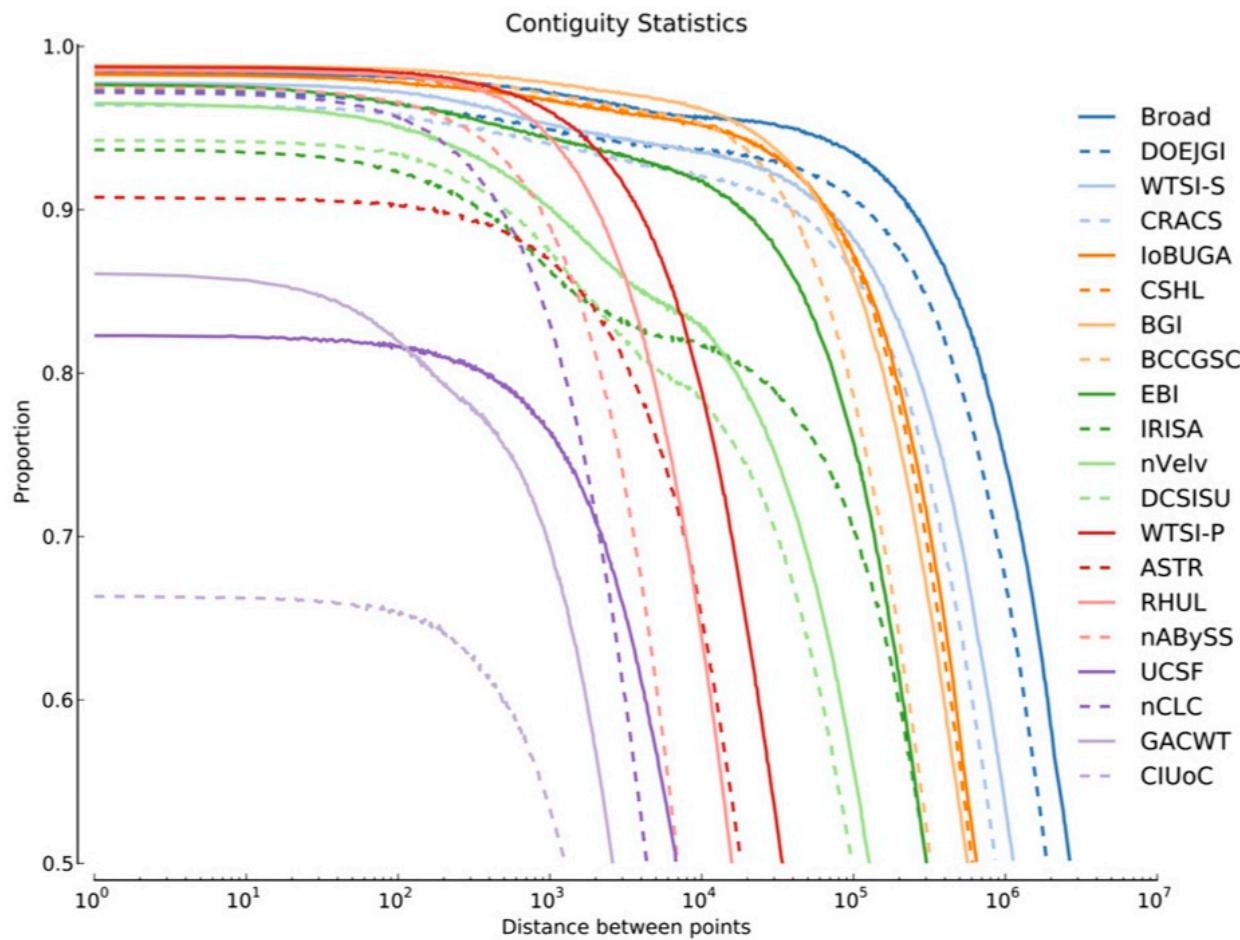


Figure 5. The proportion of correctly contiguous pairs as a function of their separation distance. Each line represents the top assembly from each team. Correctly contiguous 50 (CC50) values are the lowest point of each line. The legend is ordered *top to bottom* in descending order of CC50. Proportions were calculated by taking 100,000,000 random samples and binning them into 2000 bins, equally spaced along a \log_{10} scale, so that an approximately equal number of samples fell in each bin.

Other Applications of Assembly

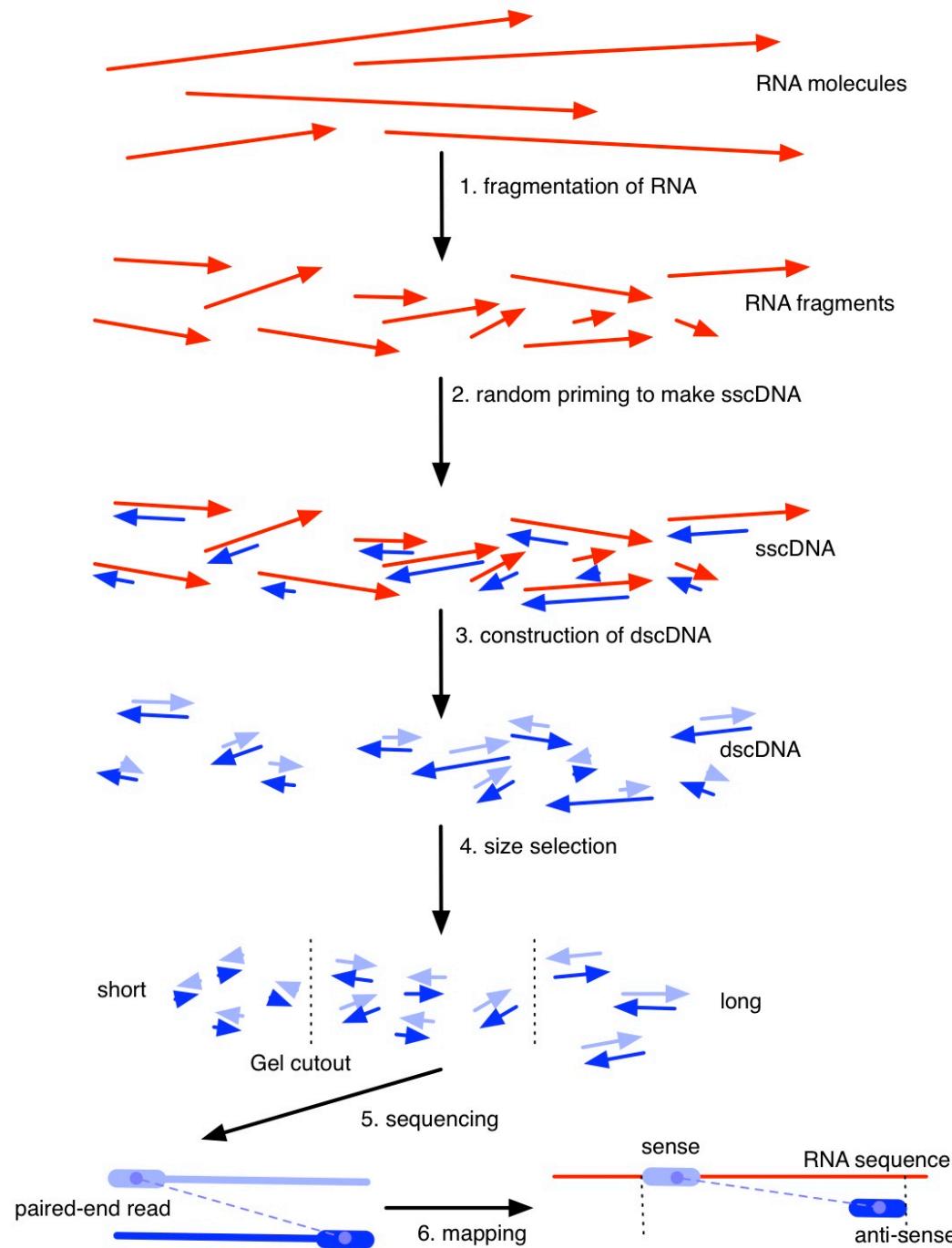
- Whole genome sequencing of isolate genome
 - plenty of material, uniform coverage
- Single cell whole genome sequencing
 - uneven coverage due to amplification
- Transcriptome assembly
 - many transcript contigs with uneven abundances
 - sequence similarity between isoforms
- Metagenomics



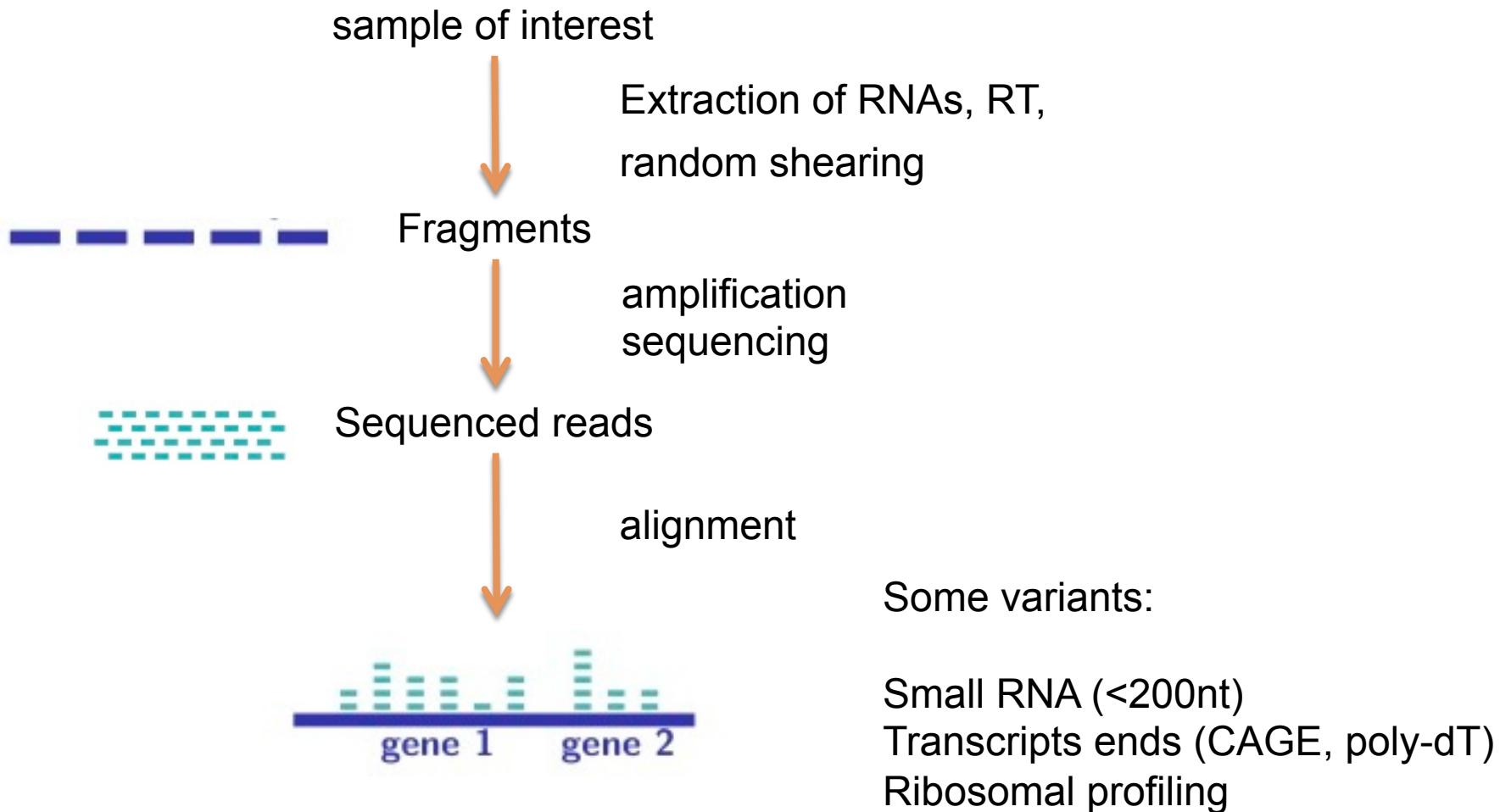
Assemblers	Technology	Availability	Notes
Genome assemblers			
ALLPATHS-LG	Illumina, Pacific Biosciences	ftp://ftp.broadinstitute.org/pub/crd/ALLPATHS/Release-LG	Requires a specific sequencing recipe (BOX 3)
SOAPdenovo	Illumina	http://soap.genomics.org.cn/soapdenovo.html	Also used for transcriptome and metagenome assembly
Velvet	Illumina, SOLiD, 454, Sanger	http://www.ebi.ac.uk/~zerbino/velvet	May have substantial memory requirements for large genomes
ABySS	Illumina, SOLiD, 454, Sanger	http://www.bcgsc.ca/platform/bioinfo/software/abyss	Also used for transcriptome assembly
Metagenome assemblers			
Genovo	454	http://cs.stanford.edu/group/genovo	Uses a probabilistic model for assembly
MetaVelvet	Illumina, SOLiD, 454, Sanger	http://metavelvet.dna.bio.keio.ac.jp	Based on Velvet
Meta-IDBA	Illumina	http://i.cs.hku.hk/~alse/hkubrg/projects/metaidba	Based on IDBA
Transcriptome assemblers			
Trinity	Illumina, 454	http://trinityrnaseq.sourceforge.net	Tailored to reconstruct full-length transcripts; may require substantial computational time
Oases	Illumina, SOLiD, 454, Sanger	http://www.ebi.ac.uk/~zerbino/oases	Based on Velvet
Single-cell assemblers			
SPAdes	Illumina	http://bioinf.spbau.ru/en/spades	
IDBA-UD	Illumina	http://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud	Based on IDBA

2. Analysis of the Transcriptome

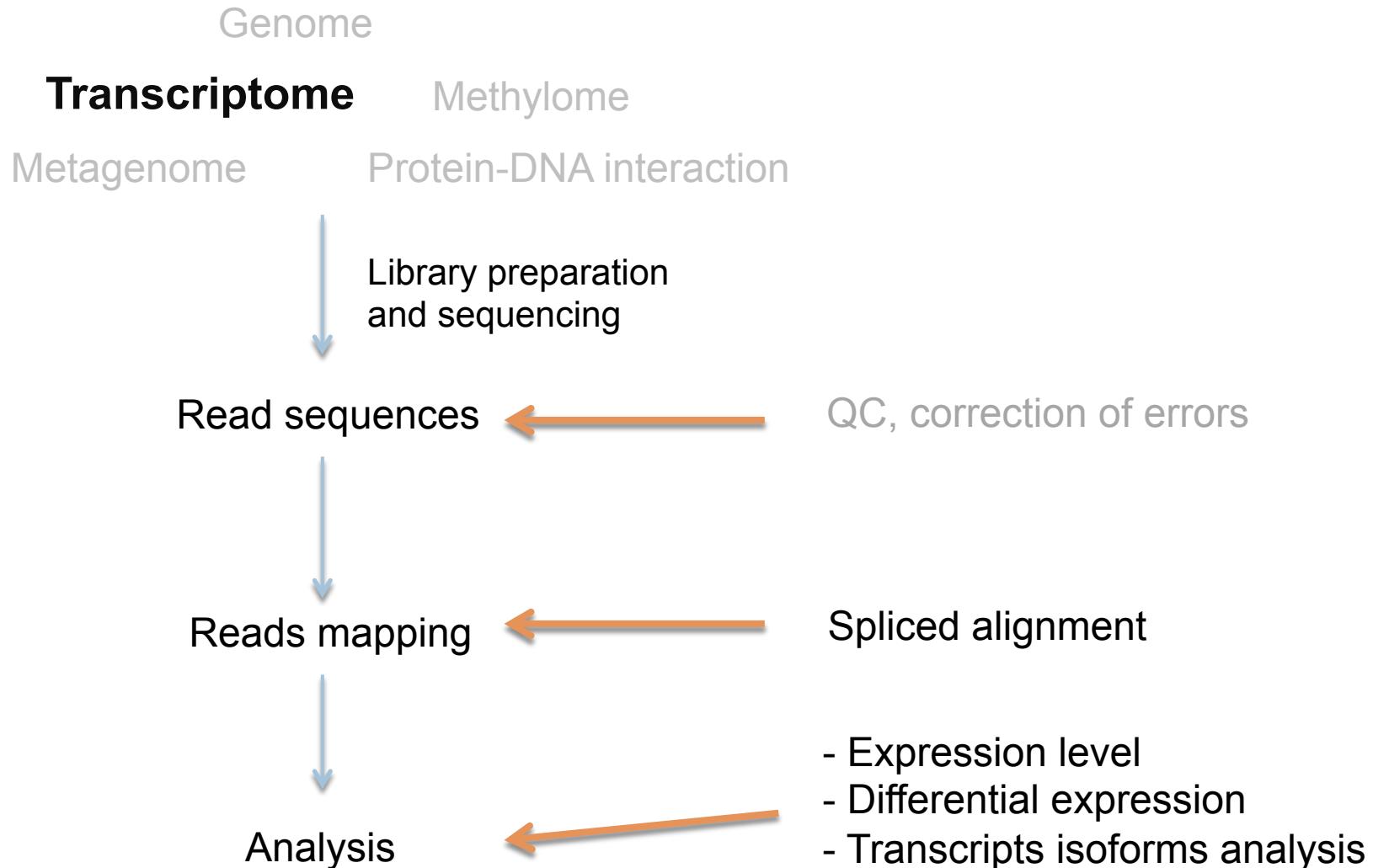
RNA-Seq data analysis



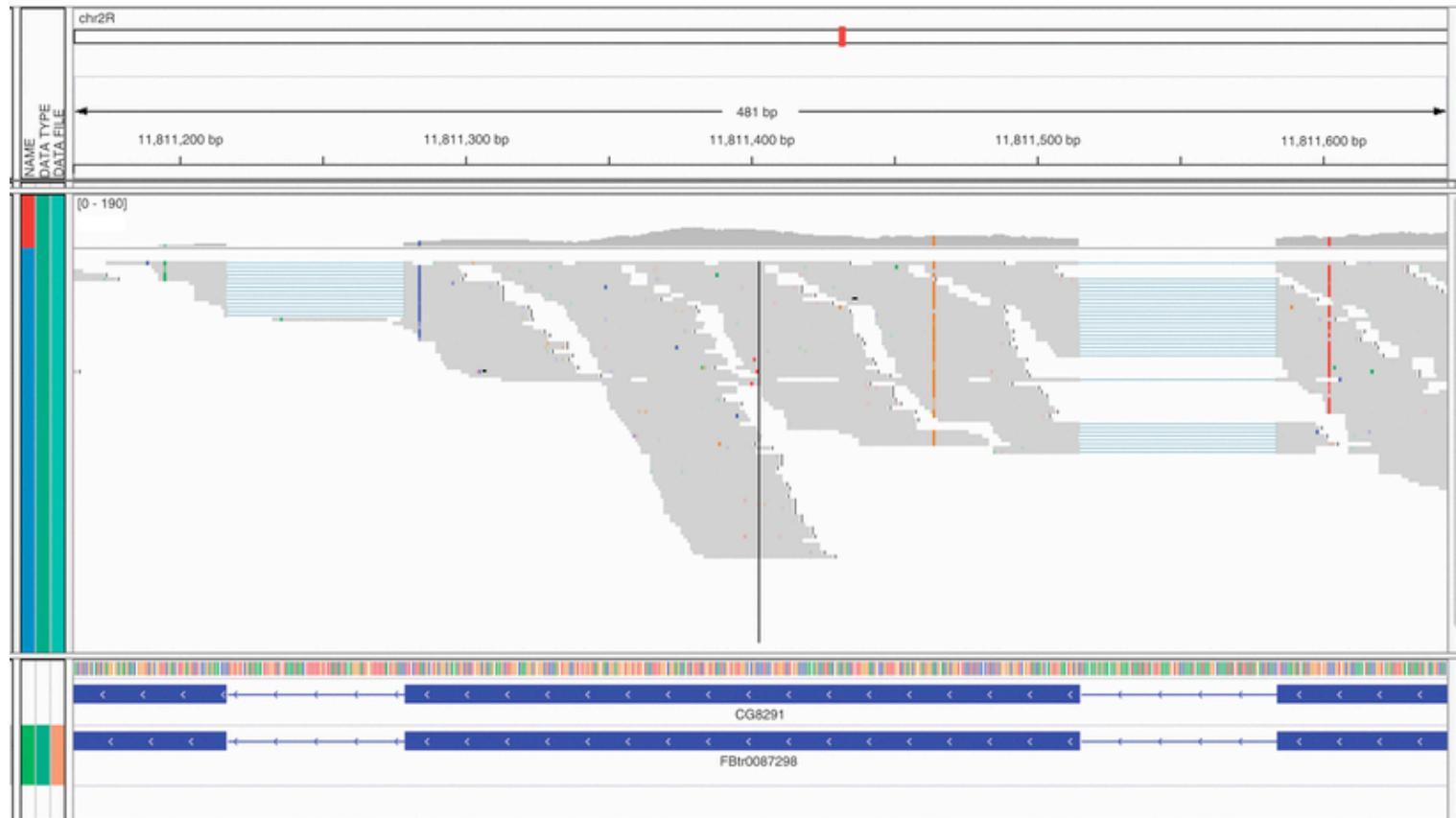
RNA-Seq protocol (2)



RNA-Seq workflow



RNA-Seq data



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

RNA-Seq

- Snapshot of transcripts activity
 - Direct identification (sequencing cDNAs)
 - Digital counting (large dynamic range)
- Questions
 - Determine transcript structure (boundaries, splicing...)
 - post-transcriptional modifications
 - Quantify expression levels/differential expression
- Challenges
 - large quantities of data (1 Gb/Mio reads)
 - still many biases/artifacts
 - transcript reconstruction is non trivial



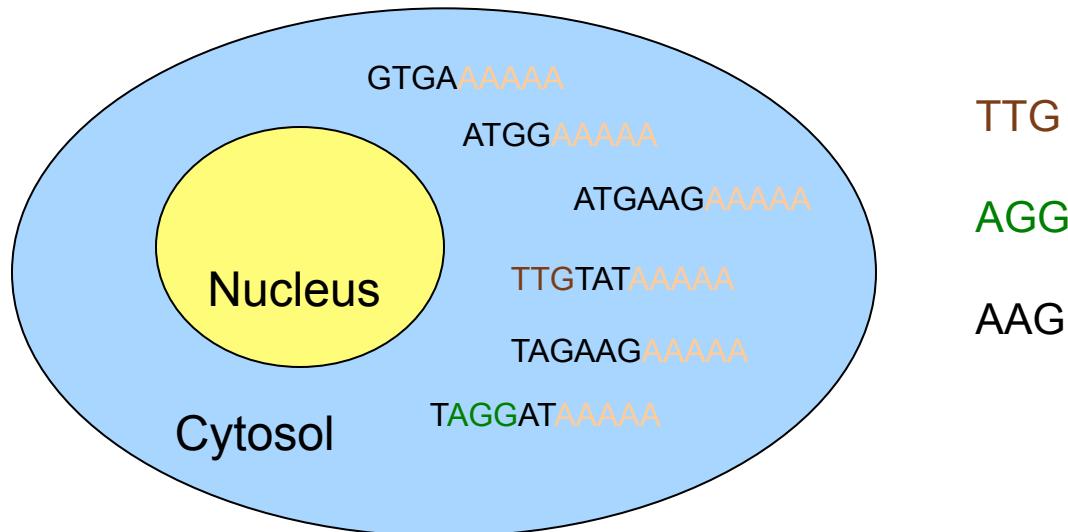
Outline

- Gene/Exon Expression level quantification
 - Sampling model and setup
 - Testing for differential expression
- Analysis of Alternative Splicing/Transcripts isoforms
 - Spliced alignment
 - Transcripts isoforms analysis
 - Transcripts isoforms quantification
 - Transcriptome reconstruction/*de novo* assembly



Counting reads

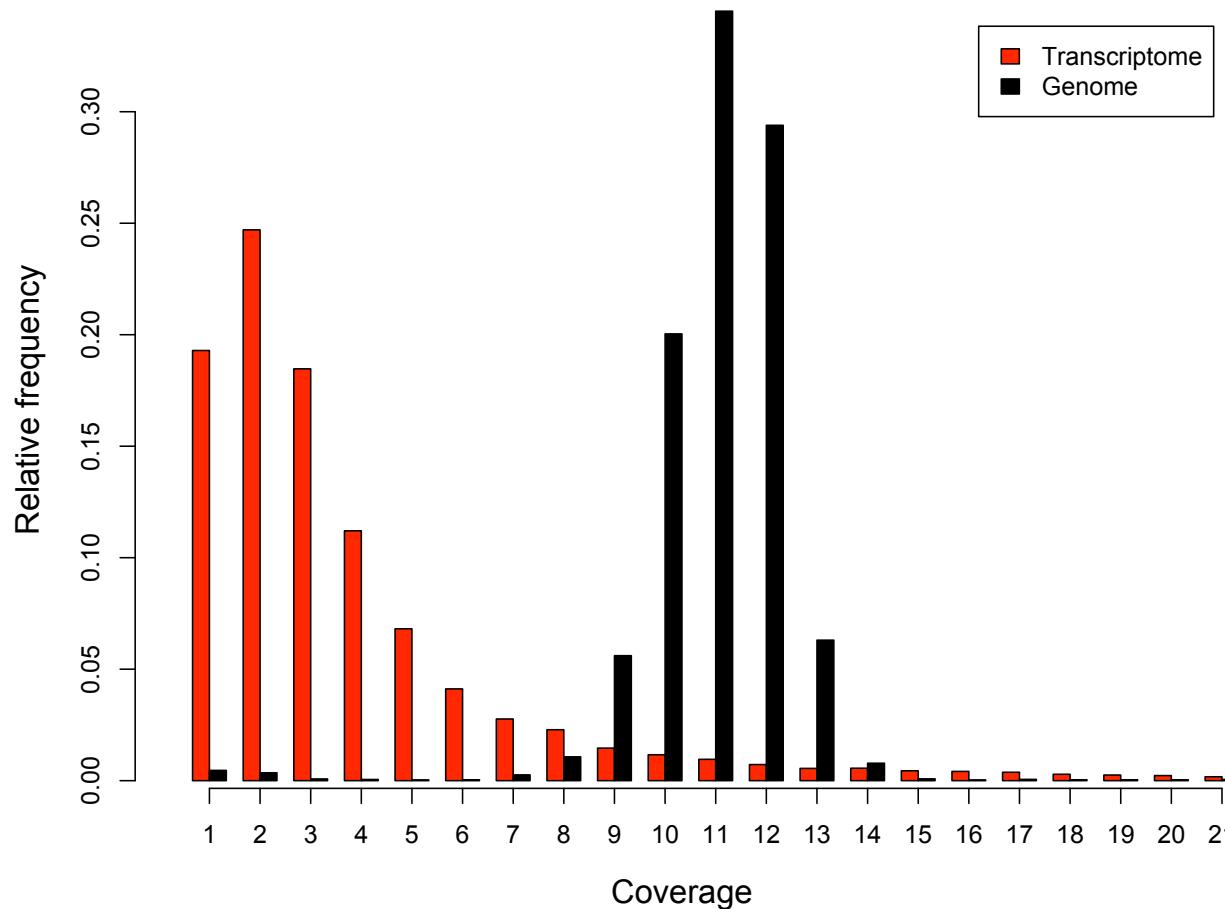
“bag of transcripts positions”



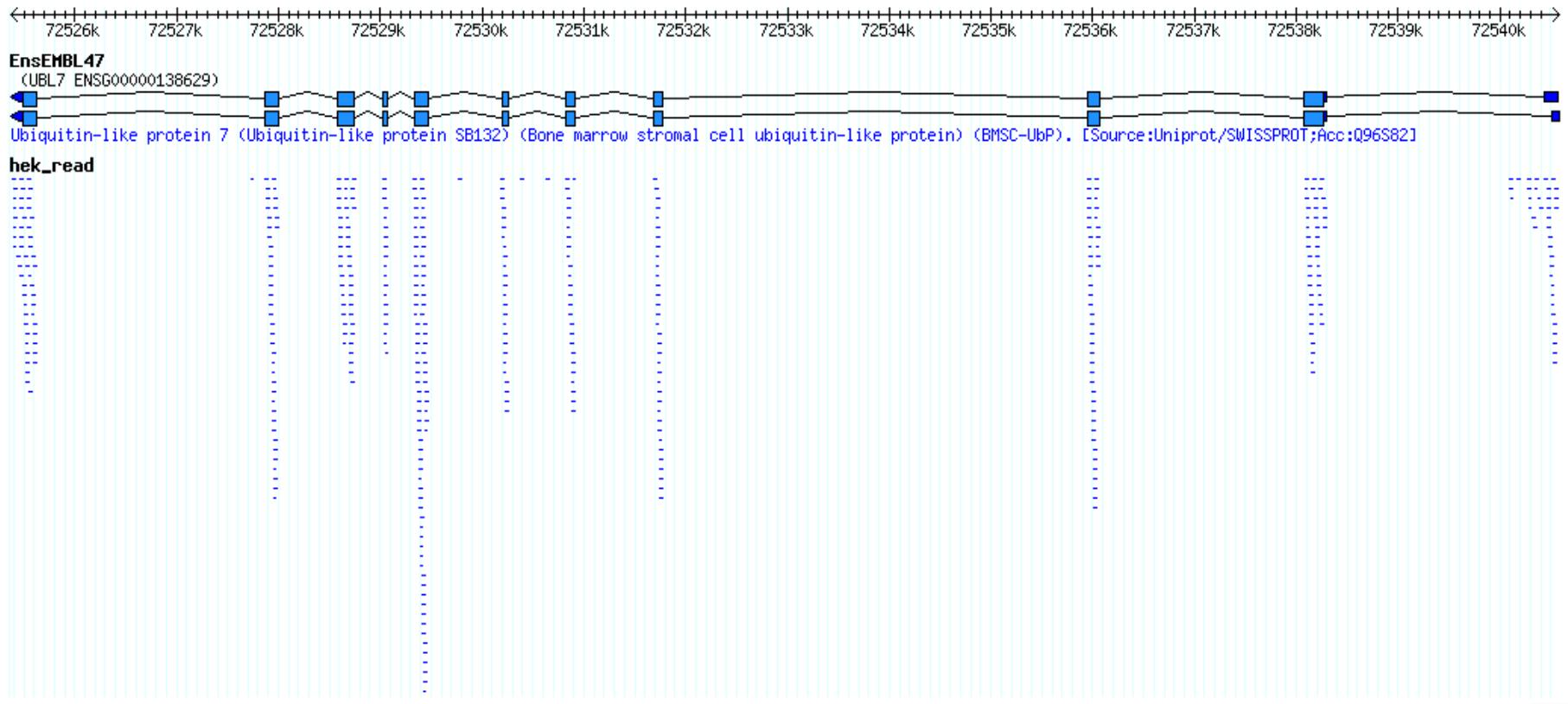
Counts ideally follow a Poisson distribution,
with an influence for transcript length

Transcriptome vs genome coverage

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



Mapping to annotated exons

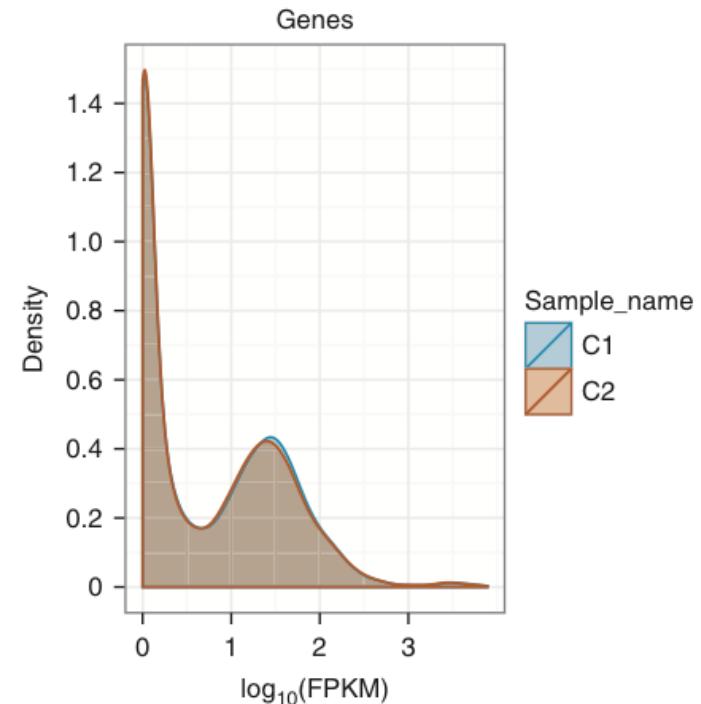


Questions:

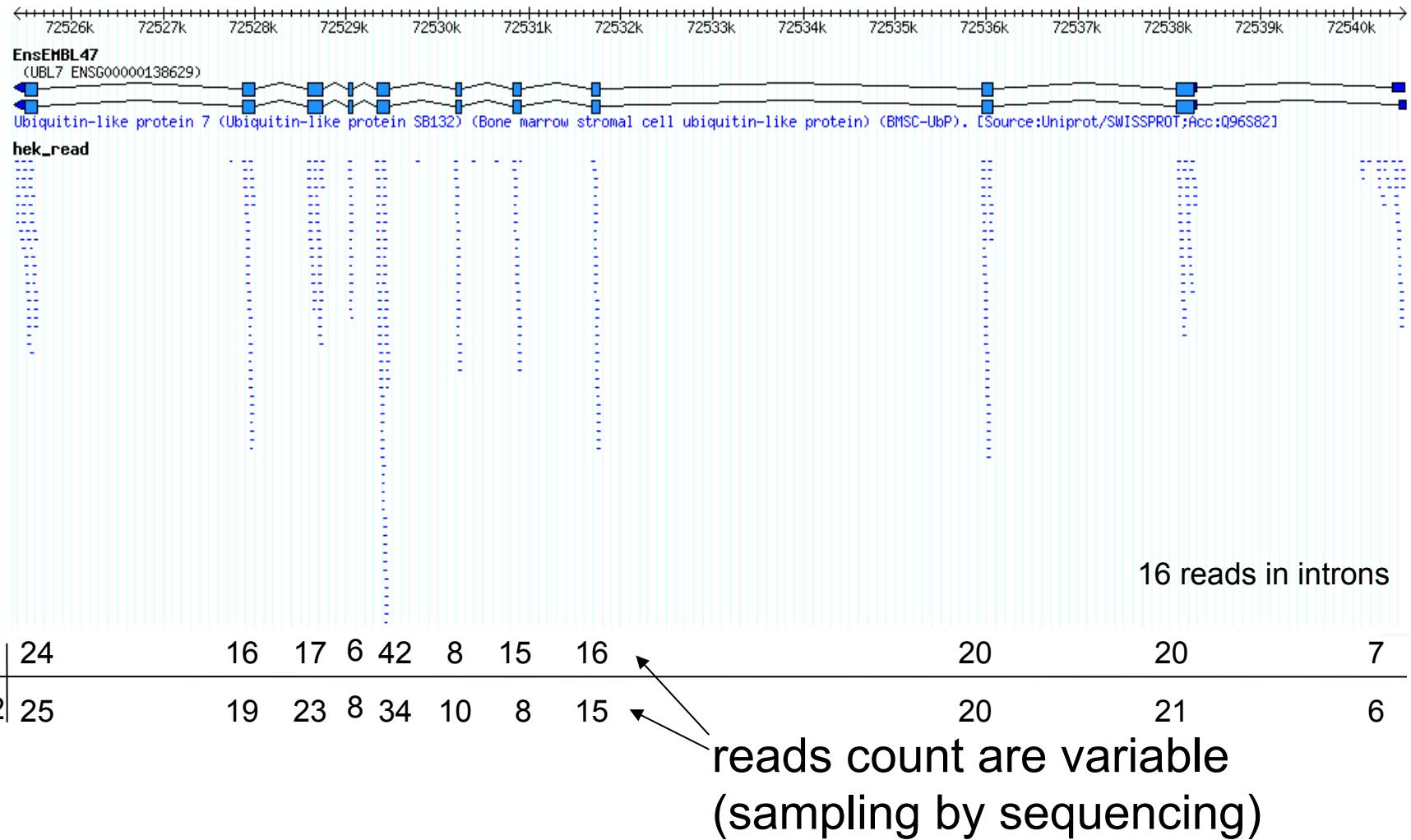
- Is a gene/region expressed?
- How to quantify gene expression level ?

Normalizing expression

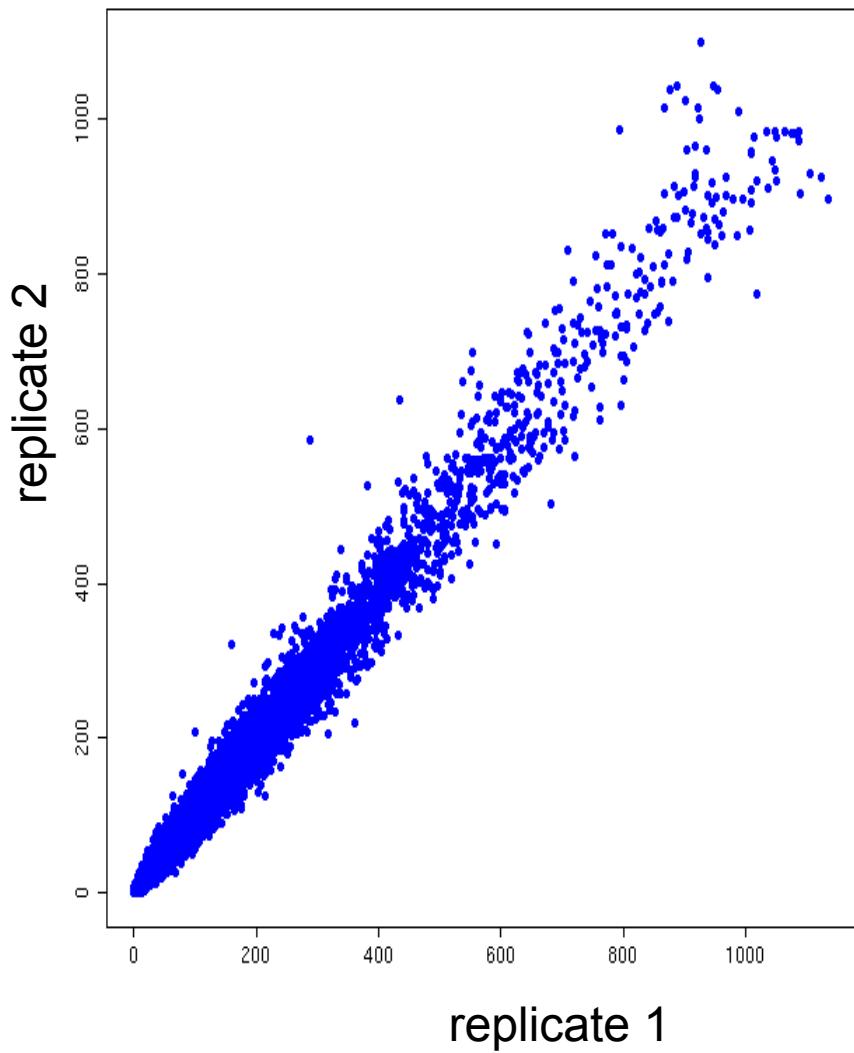
- Count mapped reads on regions (exon, gene...)
 - proportional to sequencing depth and region length.
- RPKM: Reads Per Kilobase per Million (reads mapped) (Mortazavi et al 2008)
 - FragmentPKM for mate pairs.
 - Expressed region >1-10 RPKM (depending on complexity).
 - ! Normalization not robust (a few genes account for most of the reads)



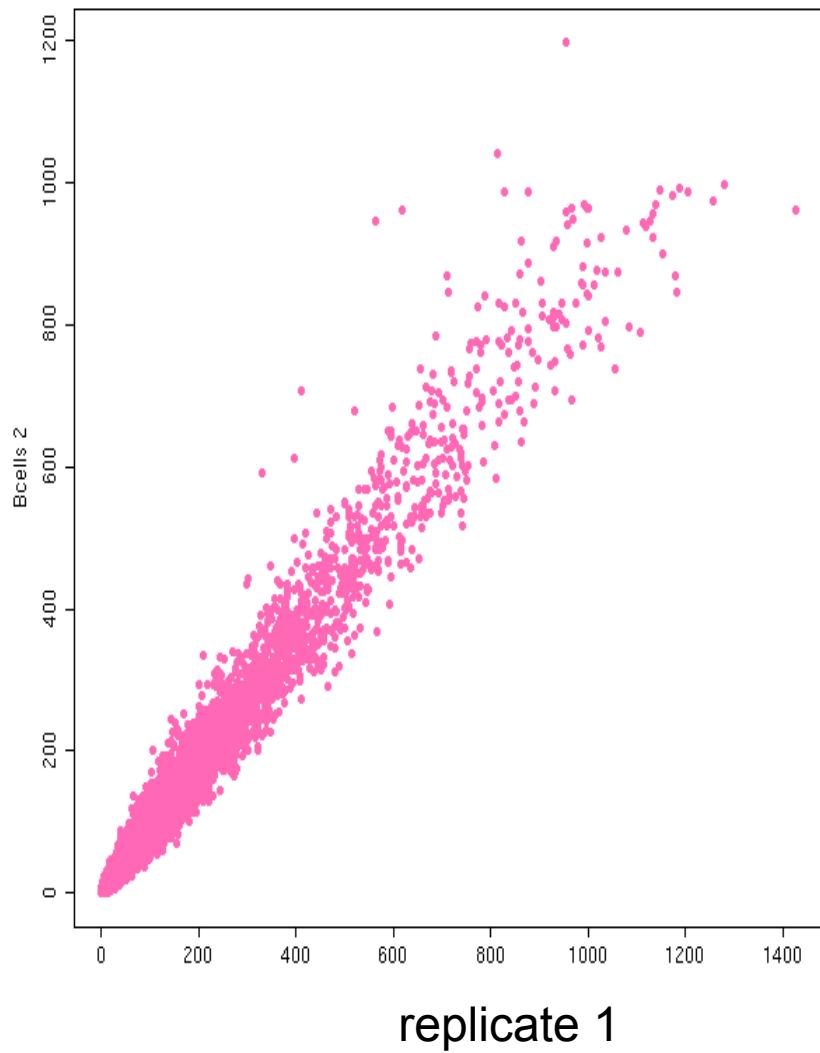
Variability on 2 technical replicates



HEK cells



B cells



Count Model For Expression

- Read counts summarized over a region are variable:
 - sampling of sequence fragments.
 - (library preparation effect)
 - (biological variability)
- Observed counts $Y_{i,j}$ (only with sampling variability):

$$Y_{i,j} \sim \mathcal{P}(\mu_{i,j}) \quad E(Y_{i,j}) = \text{Var}(Y_{i,j}) = \mu_{i,j}$$

- Discrete distribution (different from Gaussian)
- Compare counts from same region
 - no length normalization, only sample size

Testing differential expression

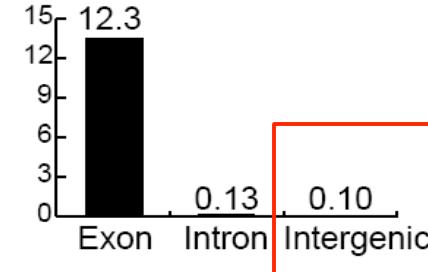
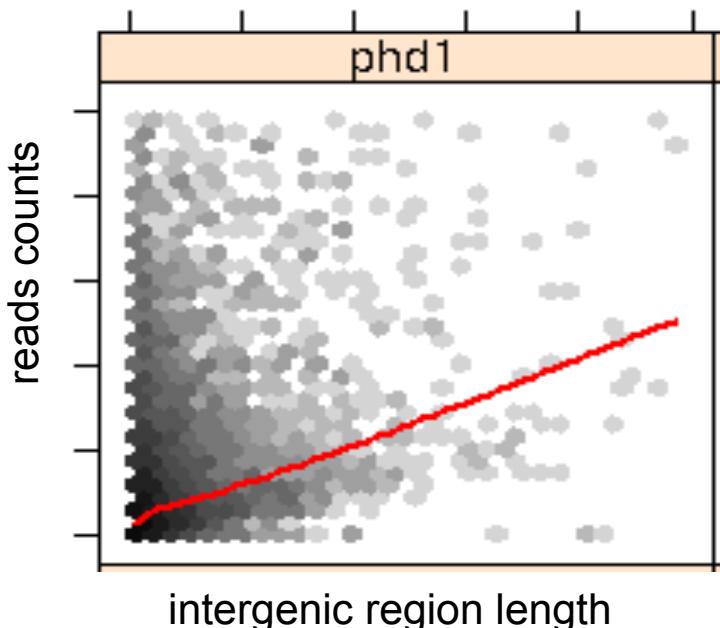
- Are differences in read counts higher than the variability expected from the experiment ?
 - Is the count on a region higher than the noise from experiment?
 - Is the expression of a gene/exon different between two experimental conditions?
- Do statistical hypothesis testing:
 - How unlikely is the observed difference in counts?
(how significant)
 - Need a null model and a test statistics...



Test expression above noise

- Assume a proportion $p_{0,j}$ for noise.

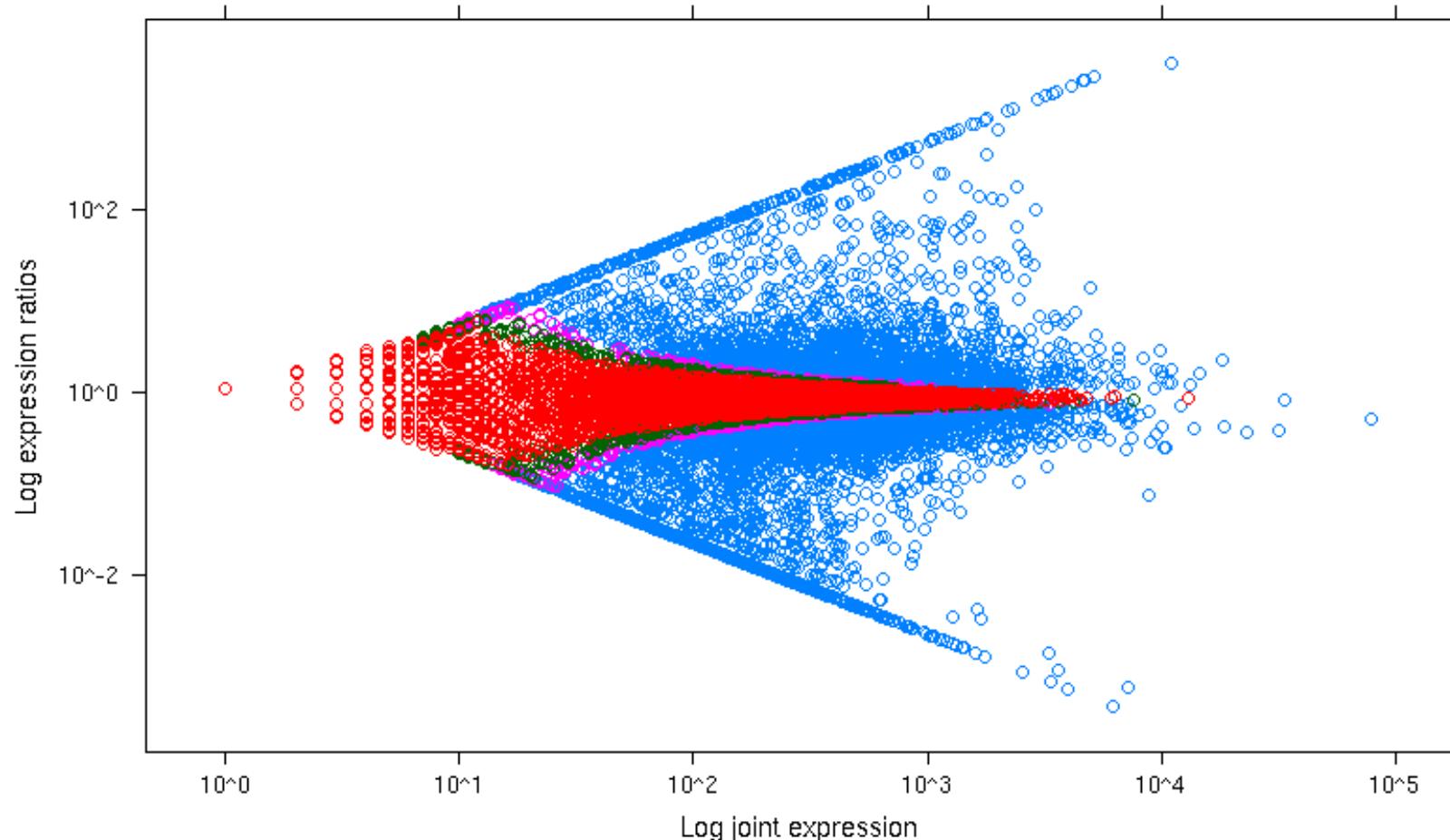
$H_0 = \{ \text{region present in proportion } p_{0,j} \}$
 $p_{0,j}$ estimated on the set of intergenic regions



Wang ET et al, *Nature* (2008)

Noise ~ 0.4 RPKM (Zebrafinch forebrain RNA-Seq)

Higher variability with lower counts



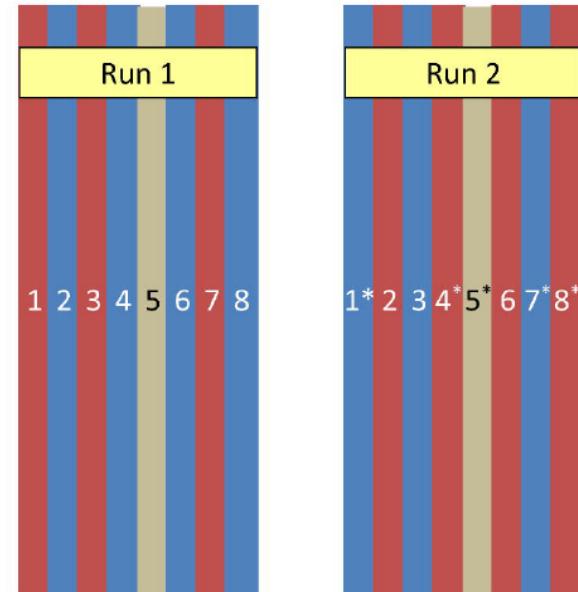
Colors refer to the probability of differential expression

Verifying Poisson assumption

(on technical replicates)

- Assess technical variability in 2 conditions.
 - should obey Poisson distribution
- Only 0.5% of genes show departure from the Poisson distribution.

Illumina study design



Kidney
Liver

* Sequenced at a concentration of 1.5 pM

Marioni et al, *Genome Research*, 2008

Extra variability

- But there are other terms than sampling variability:
 - library preparation
 - biological variability (at gene and sample level)

Var(Expr) = Across Group Variability + Measurement Error + Biological Variability

KD Hansen et al. 2011

- Negative Binomial is usually used to account for overdispersion on counts

Poisson

$$\nu = \mu$$

Poisson + constant CV

$$\nu = \mu + \alpha \mu^2$$

(knowing mean and variance is sufficient)



Estimate variance

- Estimate on each gene:
 - originally developed for SAGE data (Baggerley et al 2005)
 - but limited power due to small number of replicates
- Estimate a term common to most genes using mean variance relationship:
 - EdgeR (Robinson et al 2010, Mc Carthy et al 2012)
 - DESeq (Anders et al 2010, 2012)



DESeq

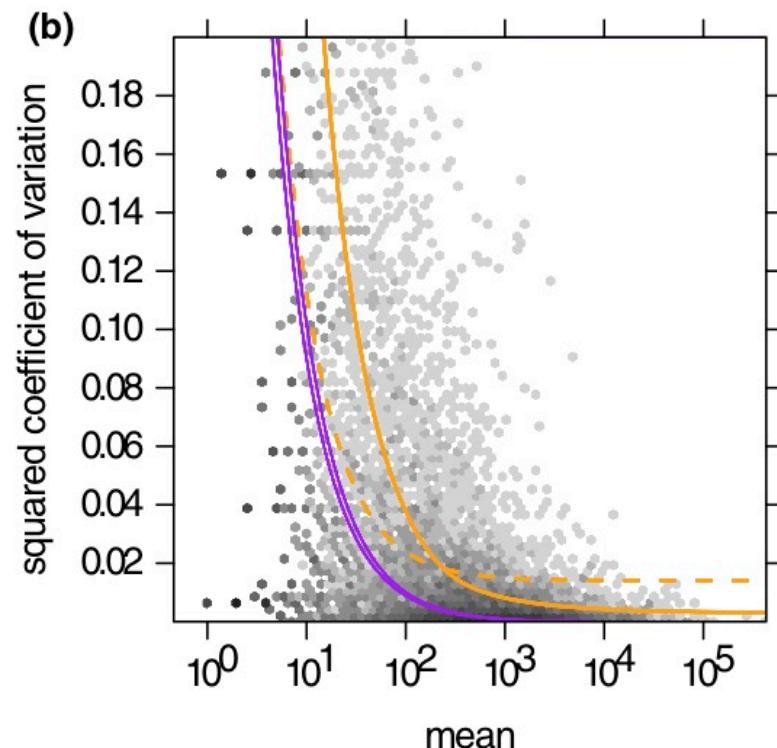
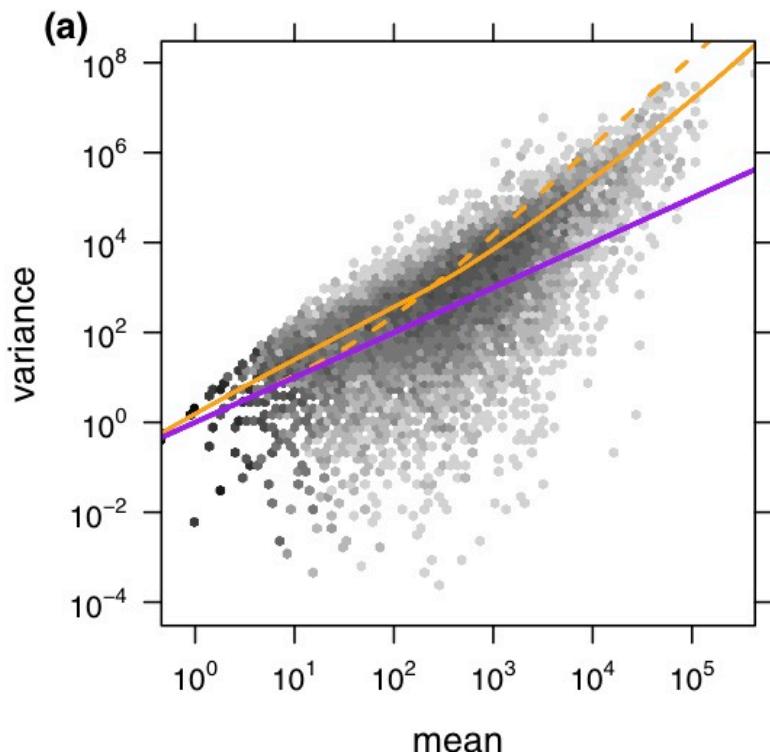
- Estimates the mean vs variance relationship with a local regression on replicates

$$\nu = \mu + f(\mu^2)$$

(also scaling factor from sequencing depth)

- Test for differential expression for the pooled counts of all replicates in each conditions c_A and c_B .
 - probability of observing c_A+c_B or more extreme cases, given the sum is c_A+c_B
 - c_A and c_B distribution are sums of NB and approximated as NB.





Poisson

$$\nu = \mu$$

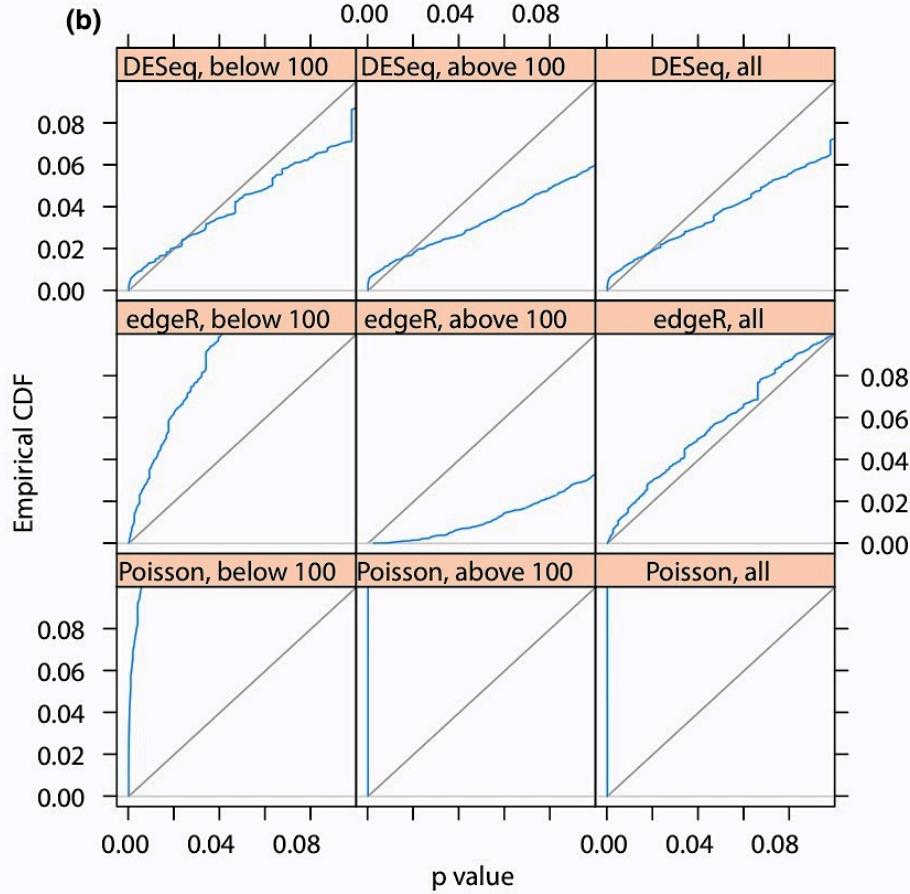
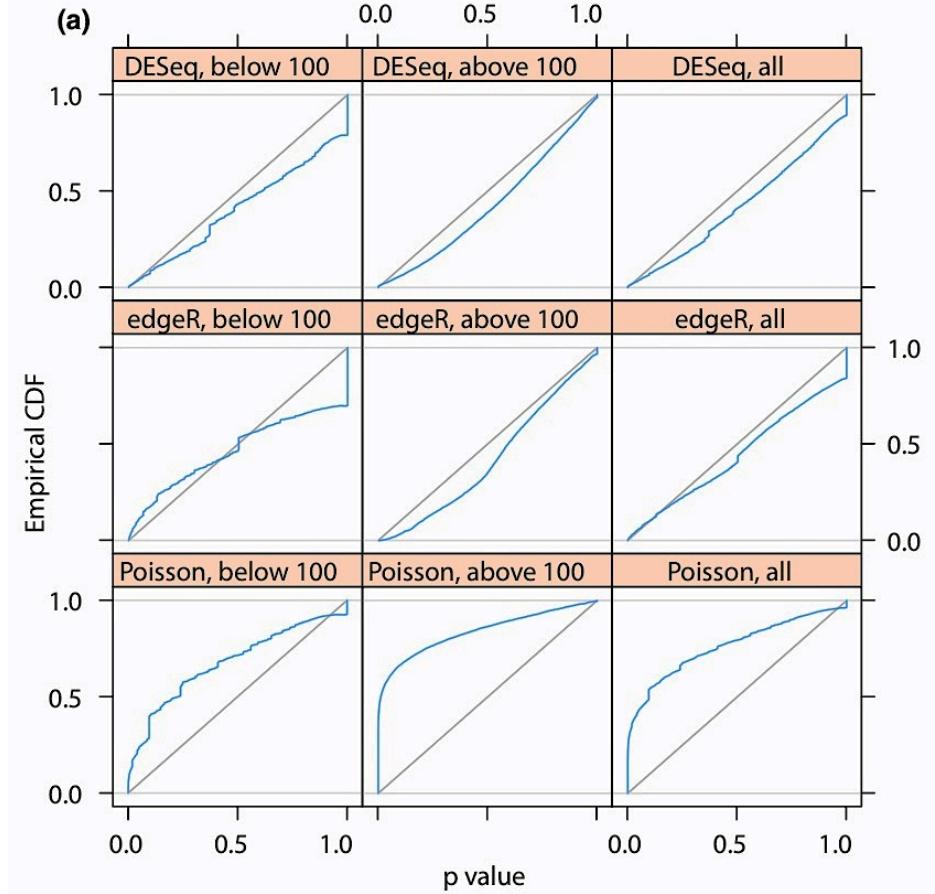
Poisson + constant CV

$$\nu = \mu + \alpha \mu^2$$

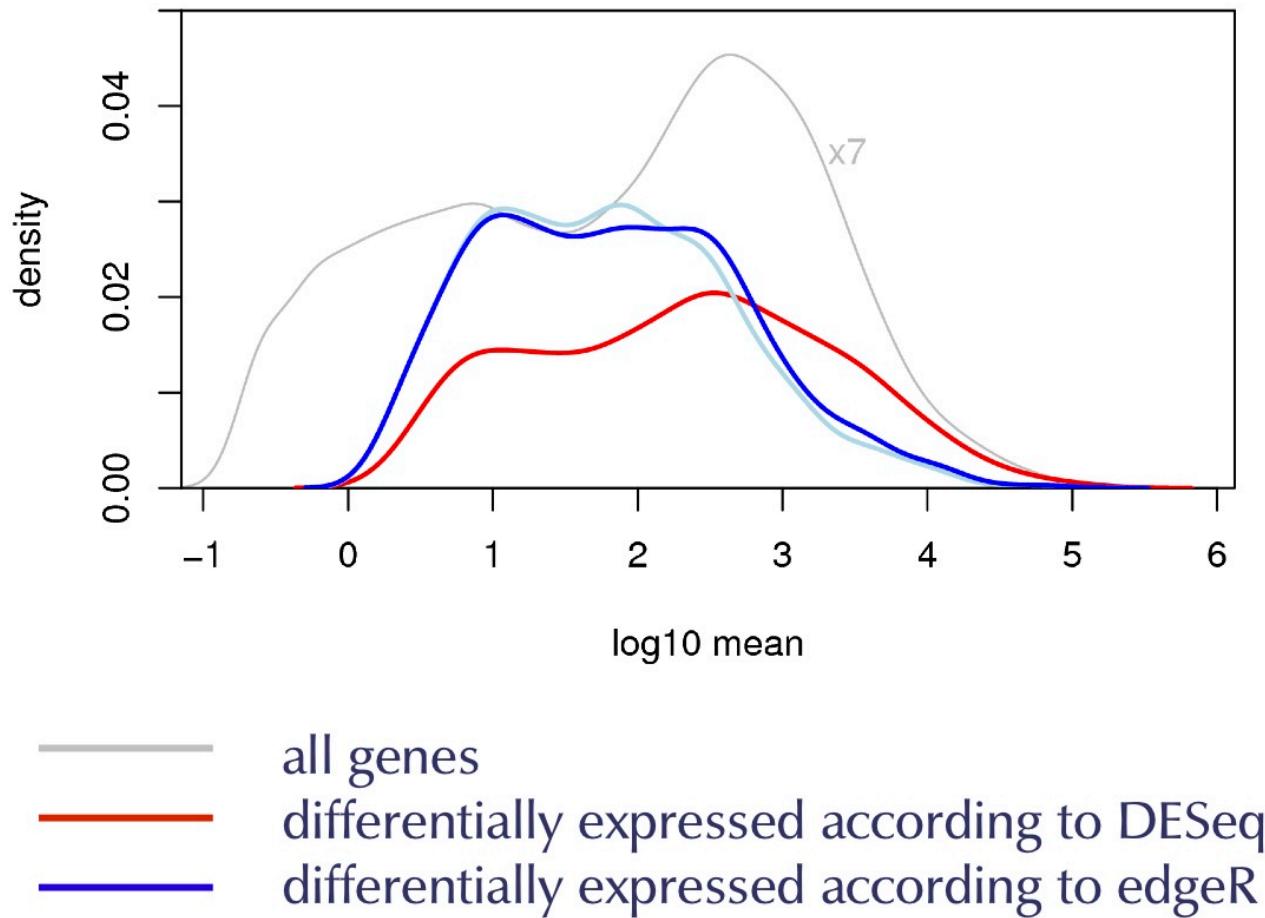
Poisson + local regression

$$\nu = \mu + f(\mu^2)$$





evaluating pvalues distribution on one replicate after estimating variance.
edgeR more sensitive on low counts, but more conservative for higher counts



Alternative to count model

- Transform data to approximate normality
 - Usually not accurate for low counts
 - can use log, square root, arcsin...
- Limma-VOOM: log transformation to counts per million (log-cpm):
 - Can leverage linear model techniques after capturing the mean vs variance relationship on log-cpm
 - Computation of significance levels is faster and more flexible (normal distribution)

Law et al 2013

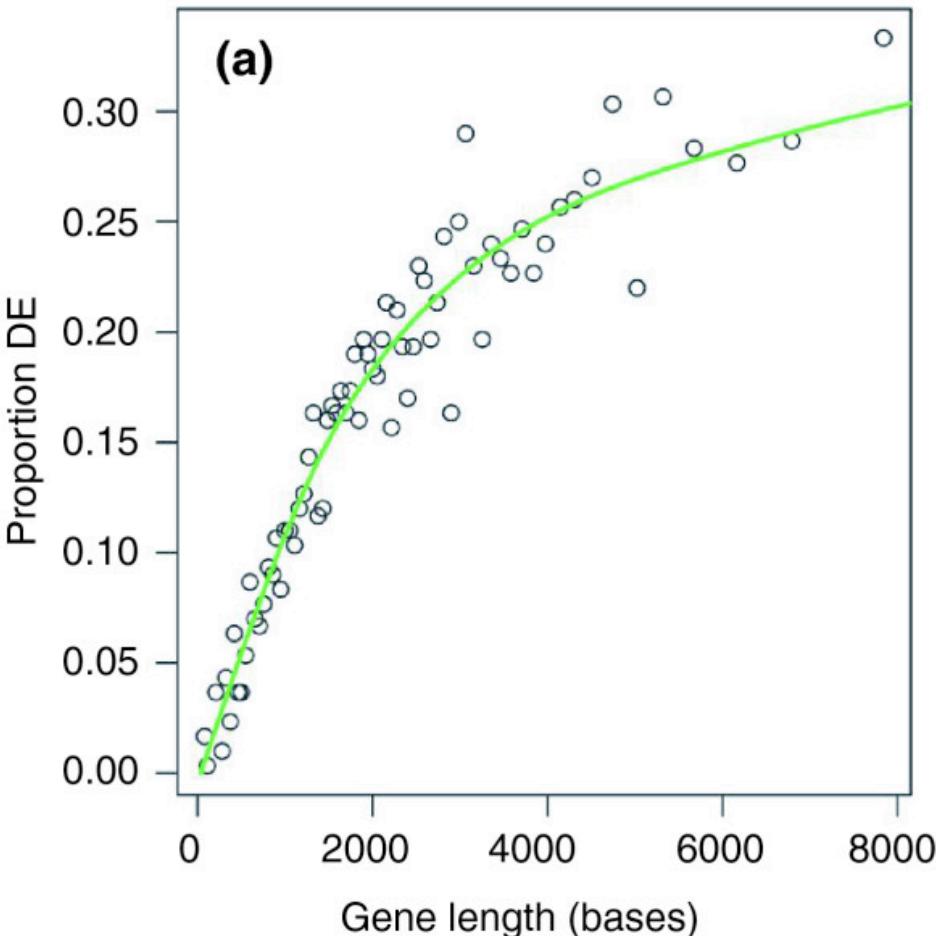


Summary DE

- 2 main terms of variability
 - Variability from sampling
 - mainly affects low counts (<100)
 - Increase power with more sequences
 - Group/biological Variability
 - more present on high counts
 - Increase power with more replicates
- Usually limited number of replicates
 - Distributional assumptions permit a tradeoff by pooling variance estimates on regions with same expression level.



Influence of the region length



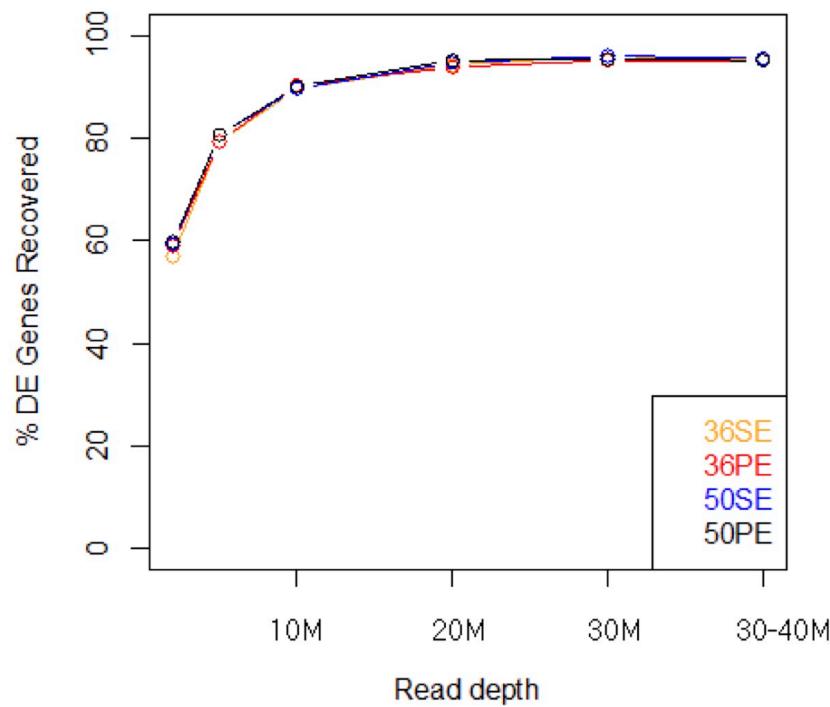
- longer/highly expressed genes more likely to be DE
- biases downstream GO analysis
- reweight genes' pvalues based on their length

Experimental design

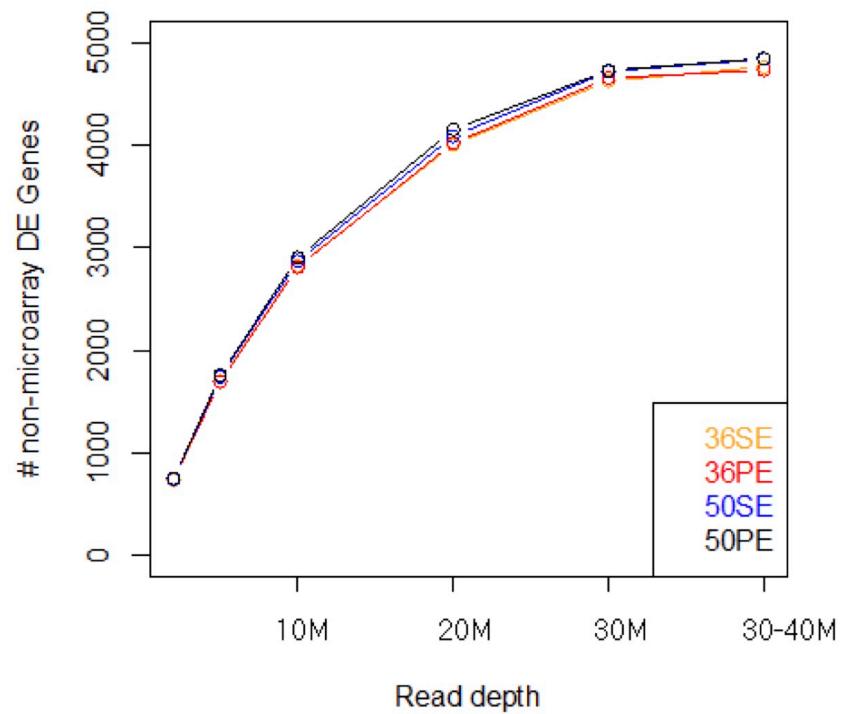
- Need to optimize parameters of the experiment according to budget.
- Sequencing parameters
 - read length (36 to 120 bp)
 - sequencing depth (2Mio to 100 Mio)
 - Single-end or paired-end
- Number of biological replicates



Recovery of Microarray DE Genes



DE Genes unique to RNA-seq

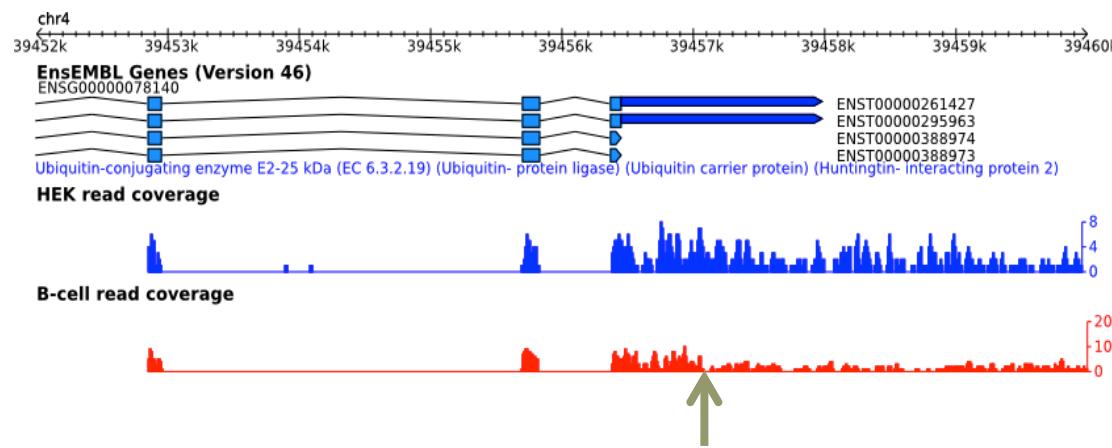


From Suraj Menon, comparison microarray vs RNA-Seq
6 replicates, tamoxifen treatment on MCF7 cell line



Read counts along the sequence

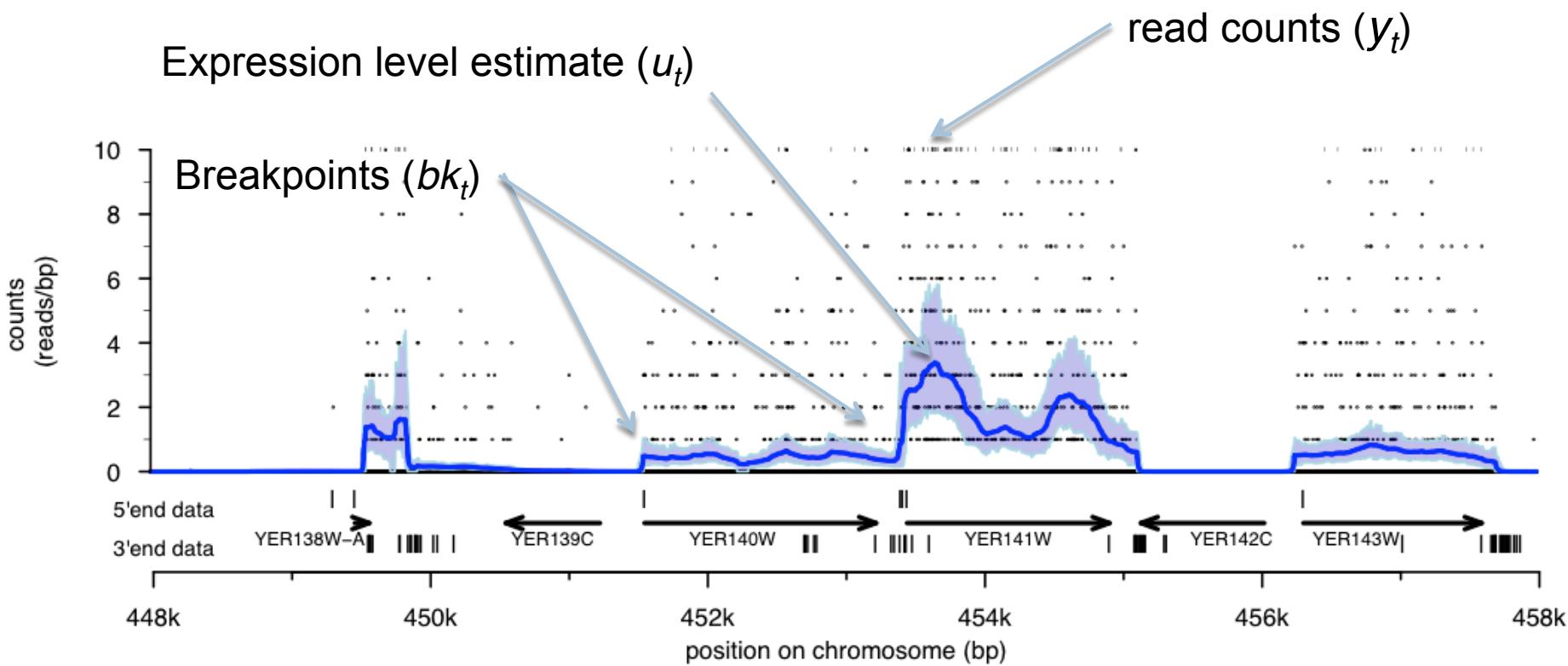
- Detect shift in read density without prior information



- Estimate expression level with confidence intervals



Read counts



Parseq

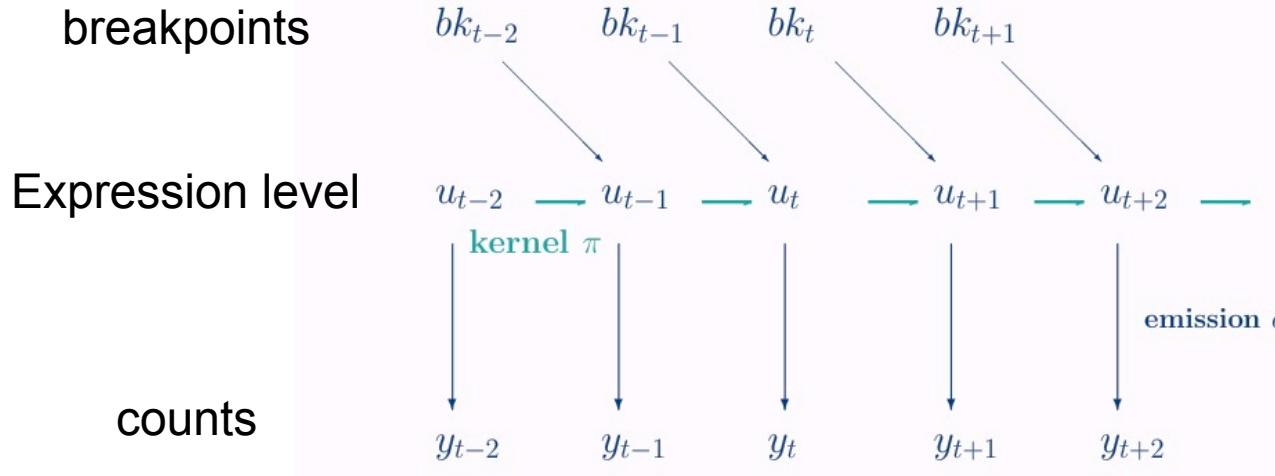


Joint work with B. Mirauta and P. Nicolas.

□ Reconstruct U_t trajectory at a bp resolution:

- Without prior annotation stranded RNA-Seq
- Not necessarily homogeneous

□ Segmentation with Hidden Markov models:



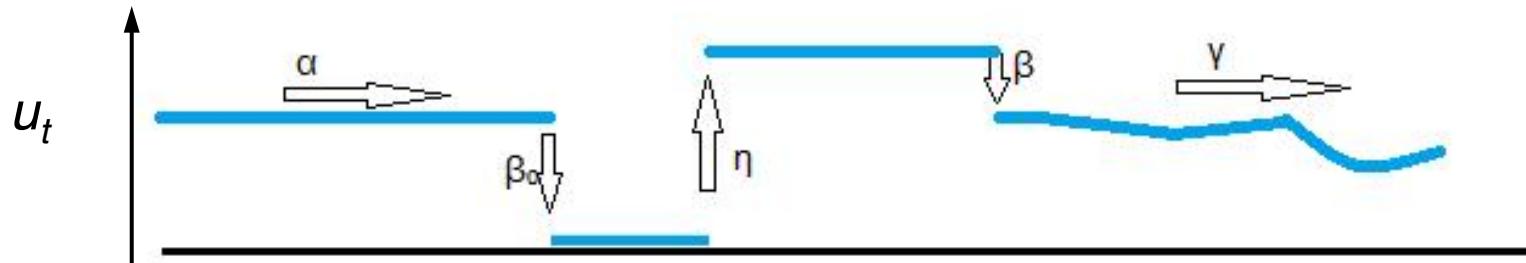
Ideally:

$$y_t \sim \mathcal{P}(u_t)$$

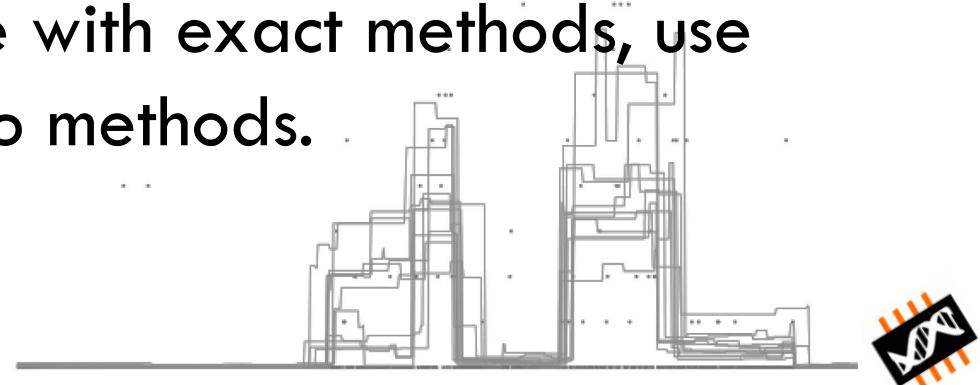


Expression level changes

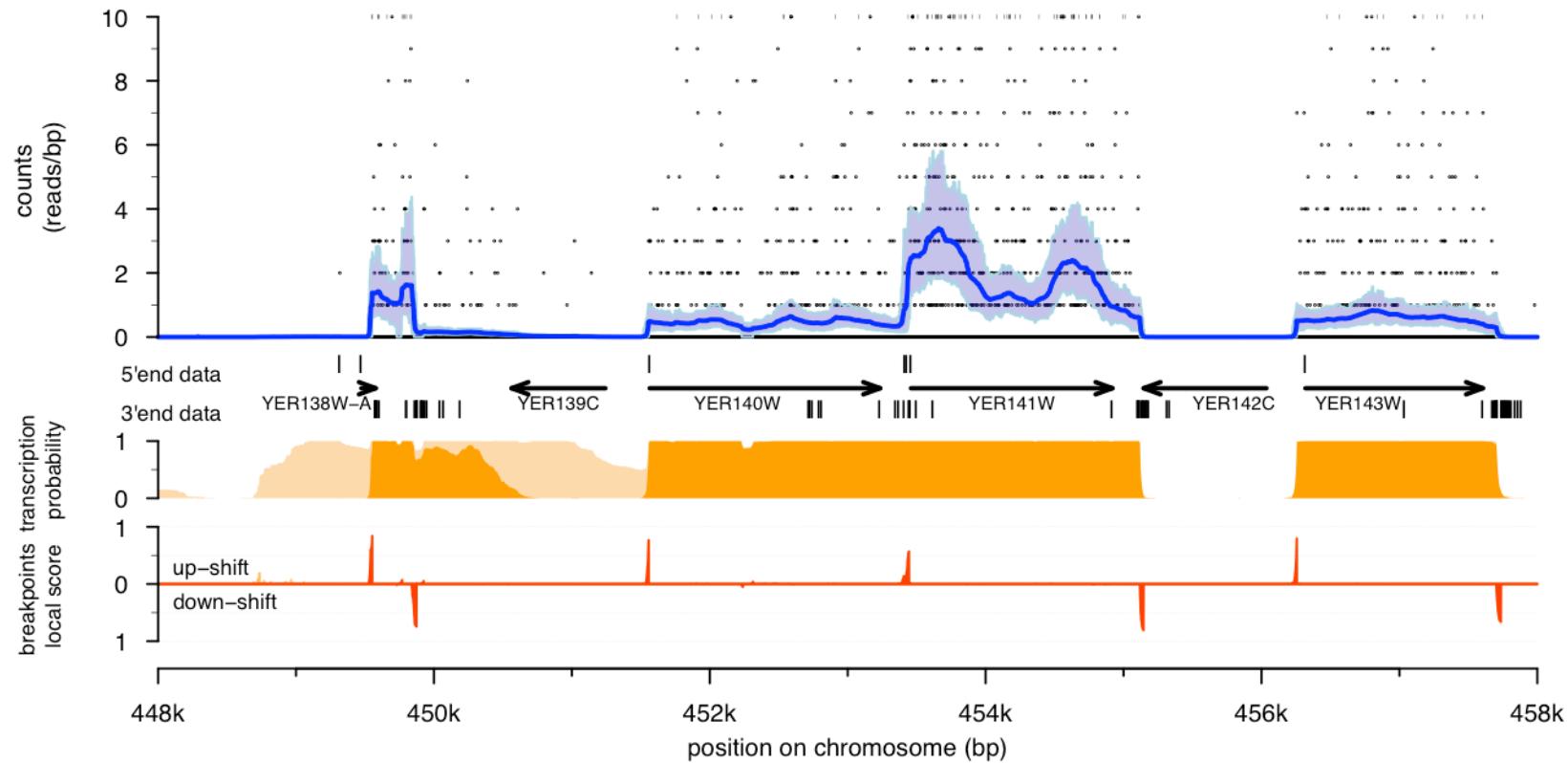
- Various types of change on transition kernel



- Read counts distribution accounts for two types of variability: local and correlated
- Estimation not possible with exact methods, use *Sequential Monte Carlo* methods.



A Parseq output



<http://www.lcqb.upmc.fr/parseq/>



Outline

- Gene/Exon Expression level quantification
 - Sampling model and setup
 - Testing for differential expression
- Analysis of Alternative Splicing/Transcripts isoforms
 - Spliced alignment
 - Transcripts isoforms analysis
 - Transcripts isoforms quantification
 - Transcriptome reconstruction/*de novo* assembly



Introducing Alternative splicing

1 pre mRNA → multiple transcripts (also called isoforms)

Alternative Splicing

The same pre-mRNA...

exon order is ALWAYS preserved!



...can make many different mRNA transcripts



...and many different proteins

Importance of AS

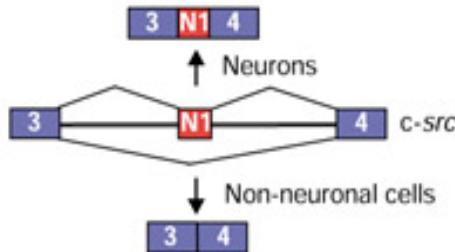
- >90% of human multi-exon genes undergo alternative splicing.
- Important in regulation of gene function.
- Aberrant splicing is a major cause of human diseases.
- An important mechanism for acquisition of evolutionary novelties (e.g. between human and other primates)

1. Xing and Lee, *Nature Reviews Genetics*, 2006, 7: 499-510.
2. Xing and Lee, *PNAS*, 2005, 102(38): 13526 - 13531.
3. Calarco*, Xing*, Caceres*, et al, *Genes & Dev*, 2007, 21:2963-2975.

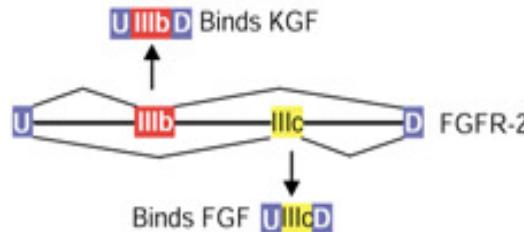


Types of splicing events

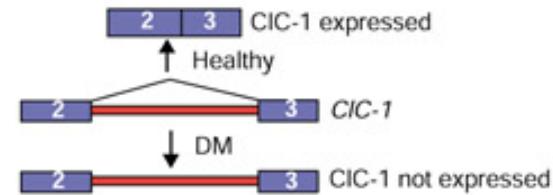
Cassette exon



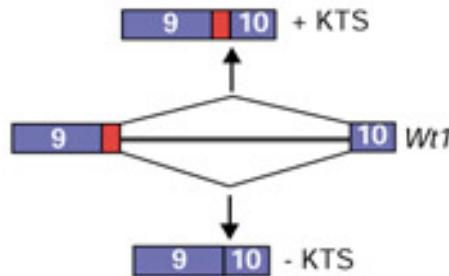
Mutually exclusive exons



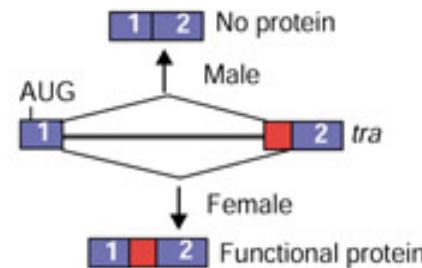
Intron retention



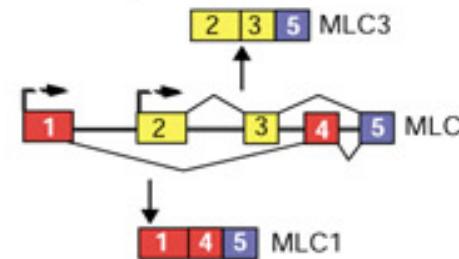
Alternative 5' splice site



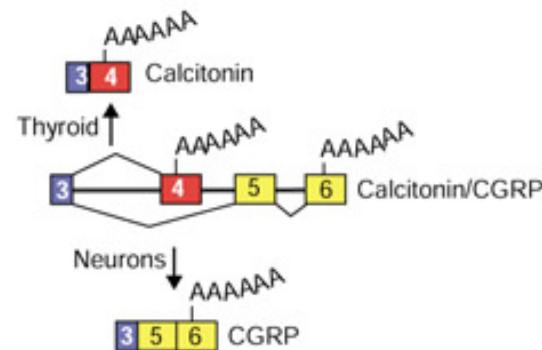
Alternative 3' splice site



Alternative promoter/ first exon



Alternative terminal exon

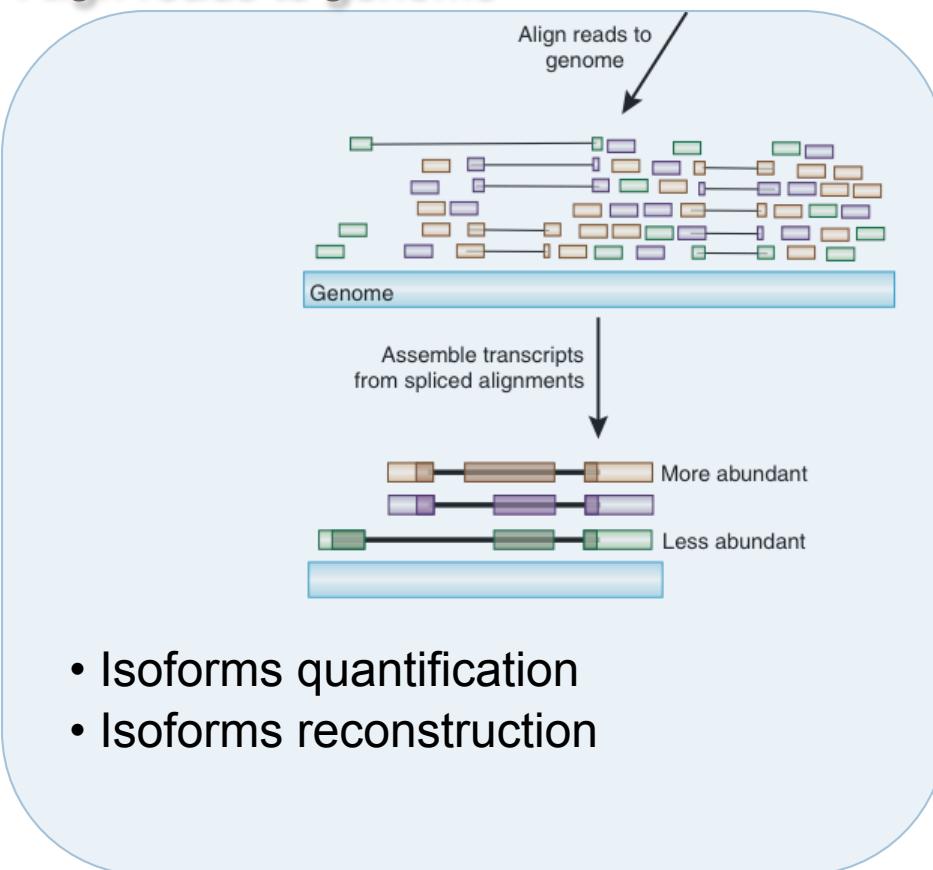


Ladd and Cooper 2002



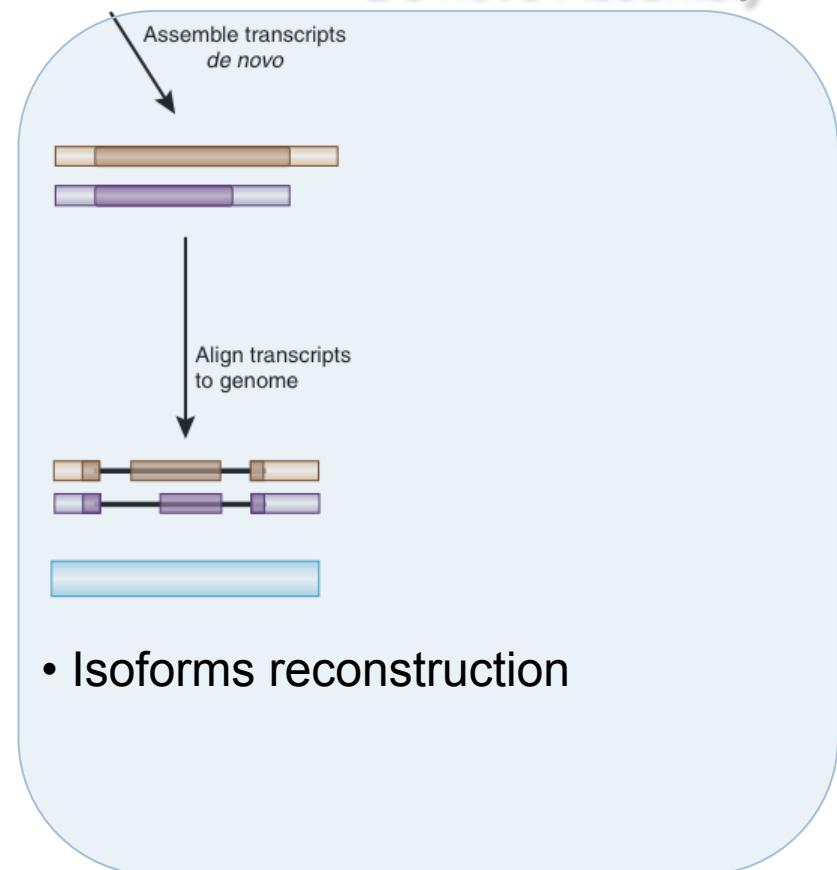
From reads to transcripts

Align reads to genome



Haas & Zody Nature News & Views 2010

De novo Assembly



- Isoforms quantification
- Isoforms reconstruction

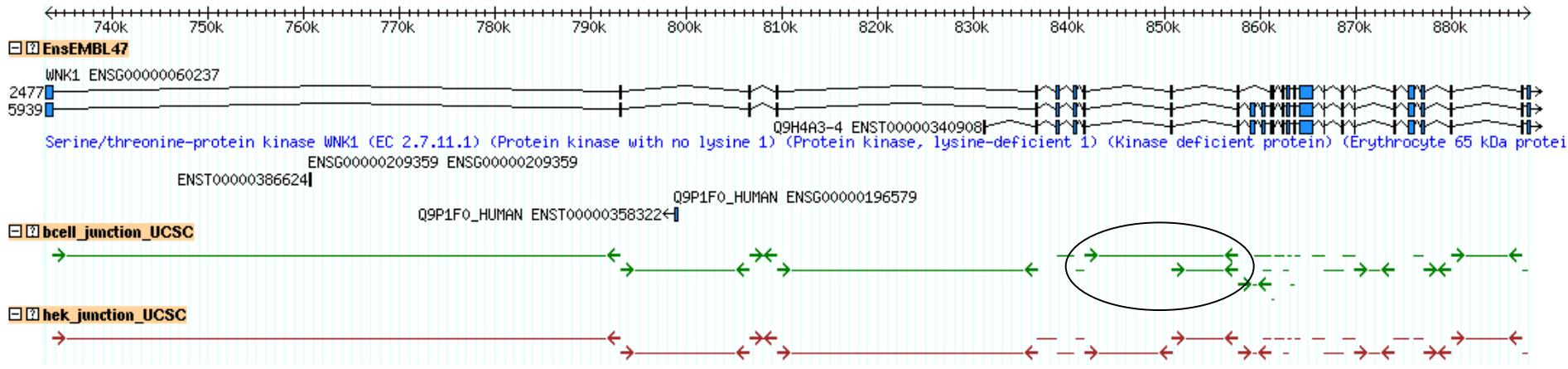
- Isoforms reconstruction

Outline

- Gene/Exon Expression level quantification
 - Sampling model and setup
 - Testing for differential expression
- Analysis of Alternative Splicing/Transcripts isoforms
 - **Spliced alignment**
 - Transcripts isoforms analysis
 - Transcripts isoforms quantification
 - Transcriptome reconstruction/*de novo* assembly



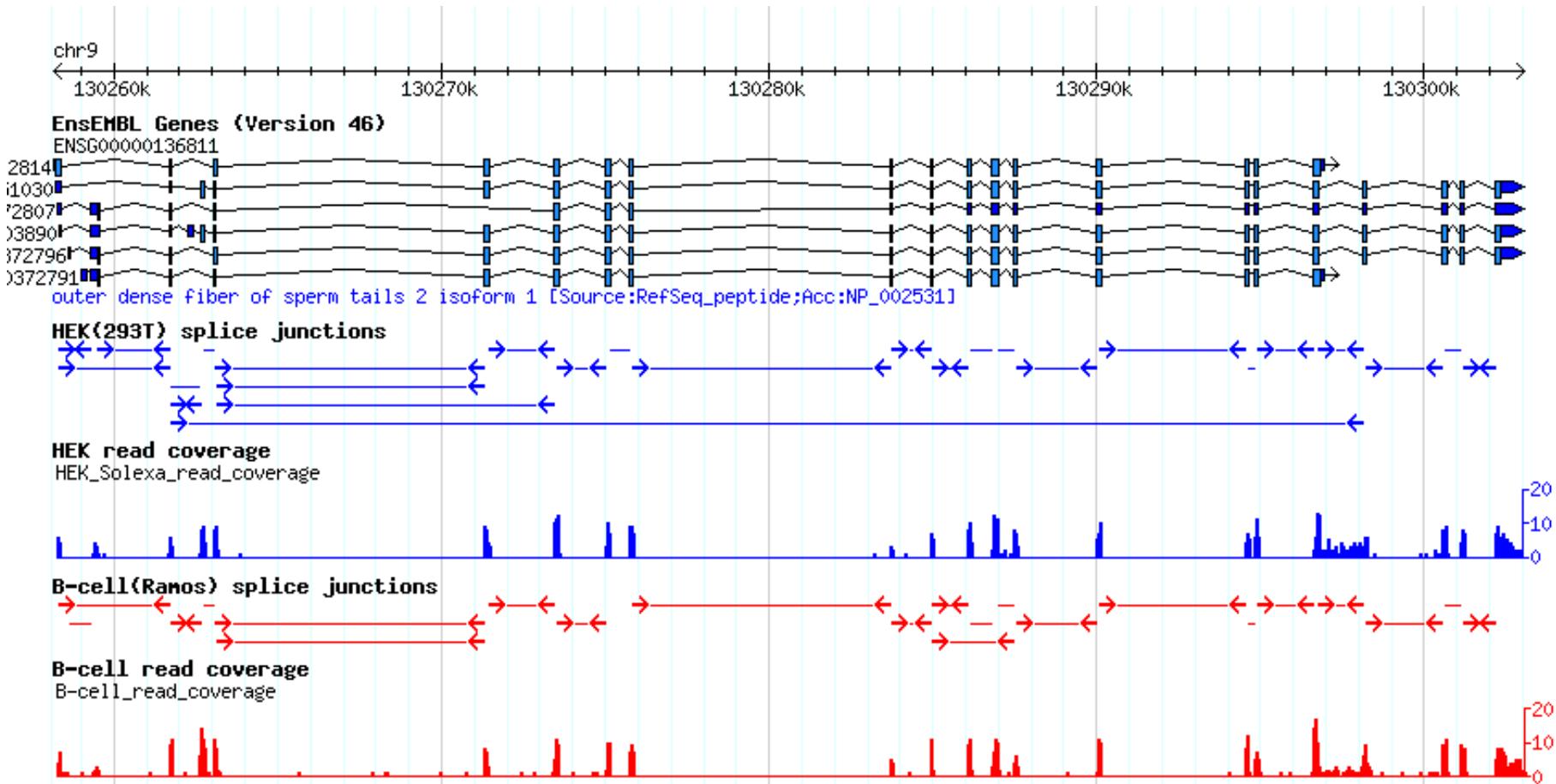
Identify from Splice junctions



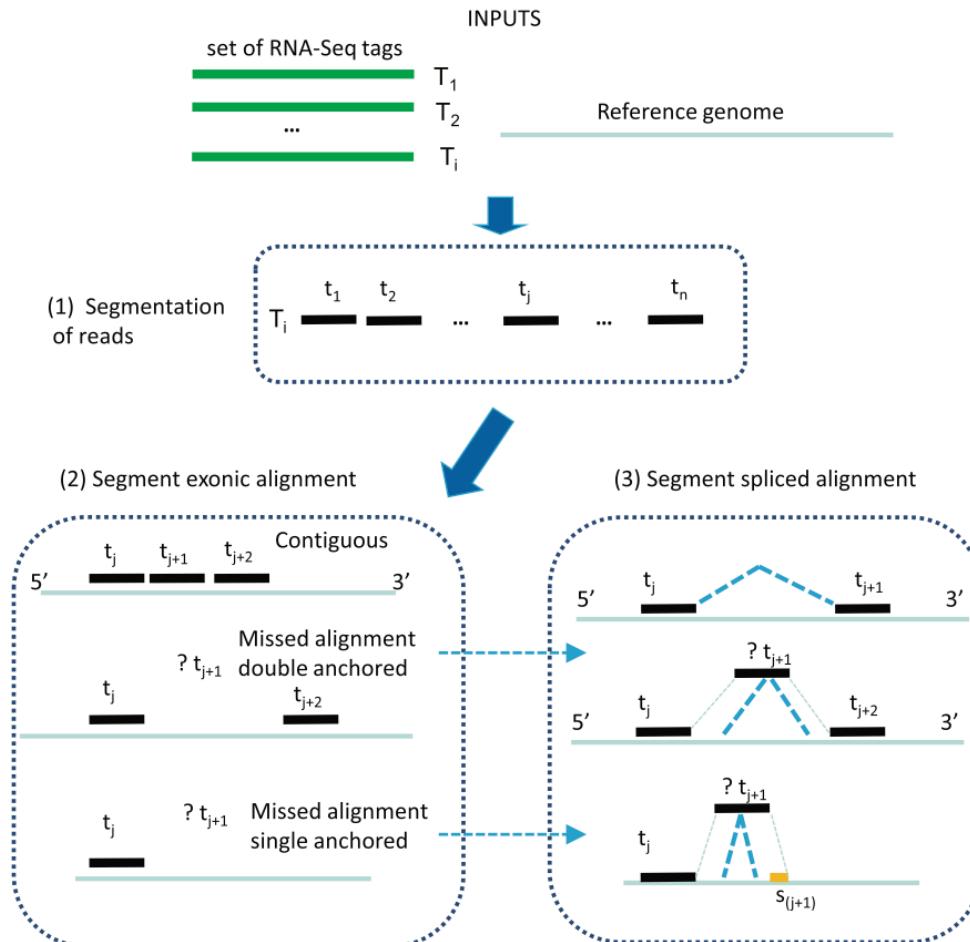
Map to a base of exons-exons junctions:

- Derived from genome annotation.
- Using the reads on exons as a guideline (TopHat).
- Perform alignment with arbitrary gap length, using pair end information, reads cut in half (SpliceMap, MapSplice).
- Use information from splice site in score function

Example junction

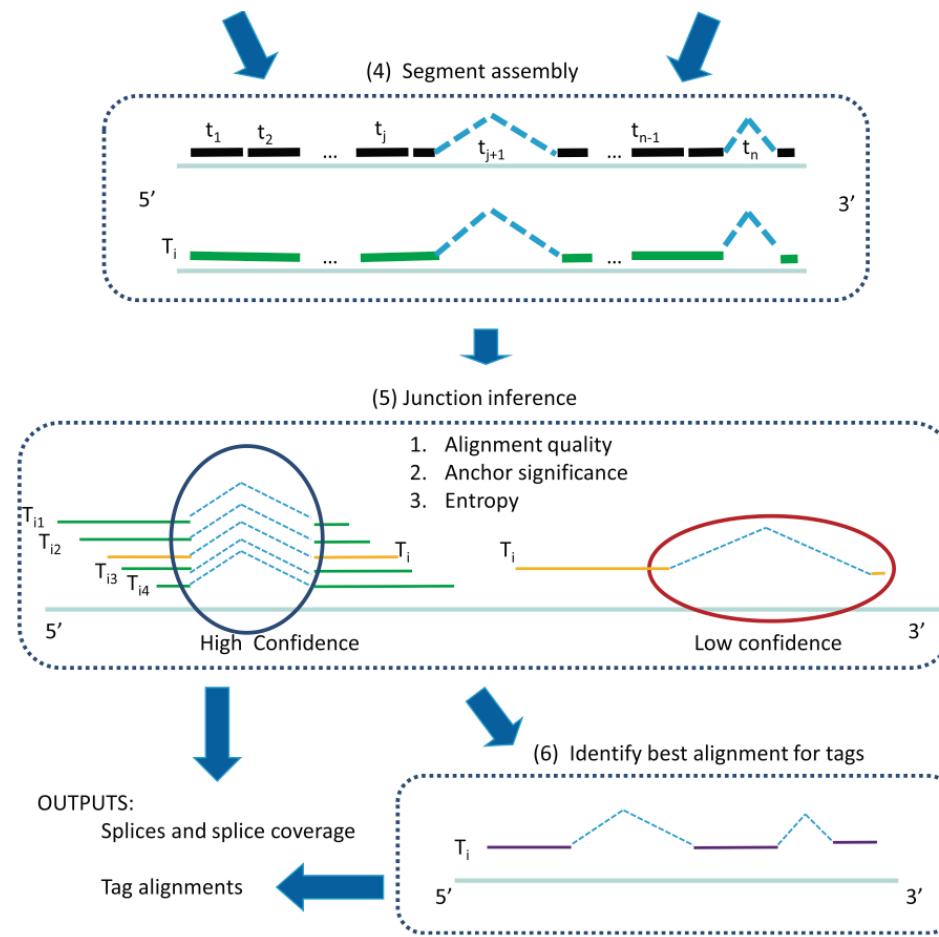


Aligning splice junctions (splicemap)



(2) Segment exonic alignment

(3) Segment spliced alignment



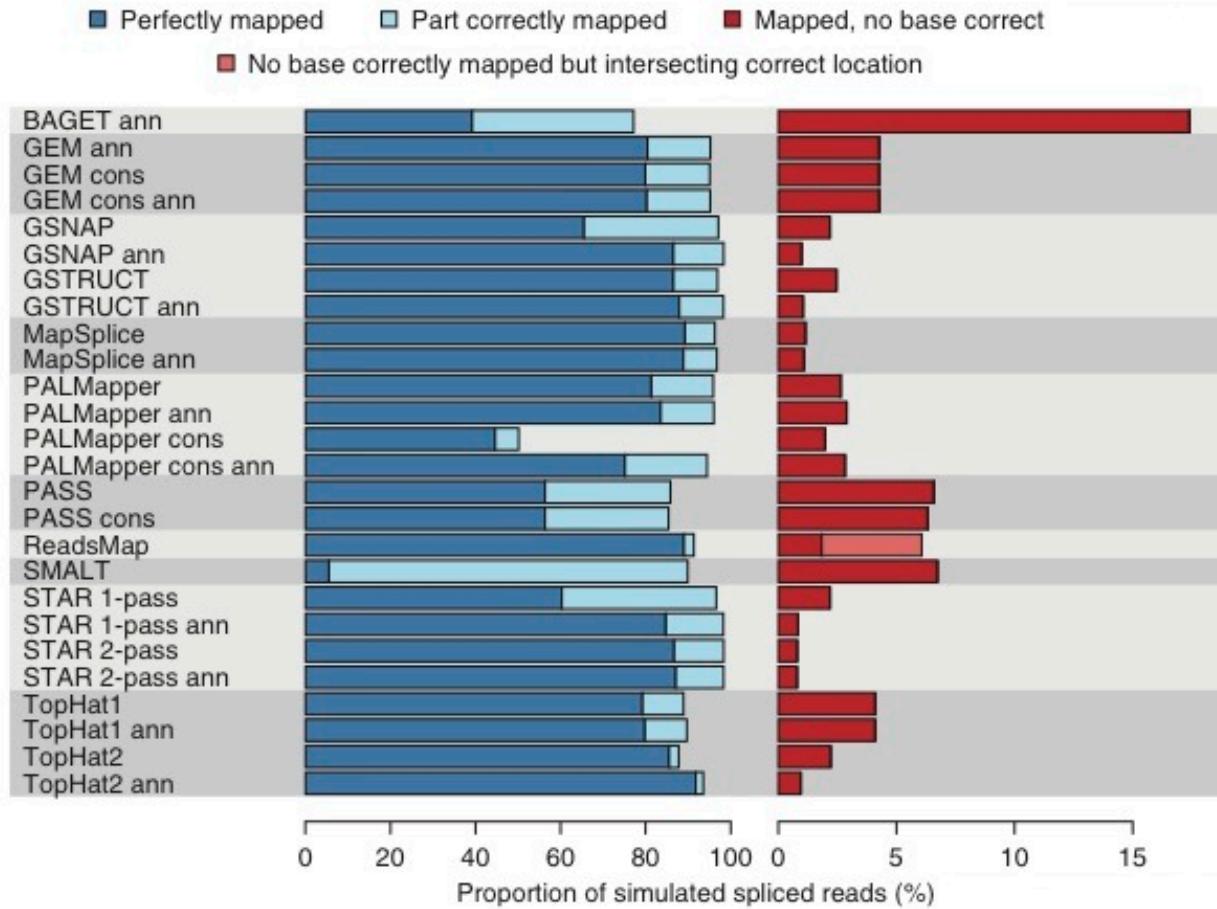
RGASP-Benchmark splice aligners

- RNA-Seq Genome Annotation Assessment Project
 - 9 teams submitted 26 alignment protocols
 - 2 real (human cell line, mouse)
 - Two simulated datasets (1 easy and 1 complex)
- Softwares:
 - GSNAP, MapSplice, PALMapper, ReadsMap, STAR, TopHat.
 - GEM, PASS, GSTRUCT, BAGET.
 - SMALT (contiguous aligner not tuned for splice junctions)

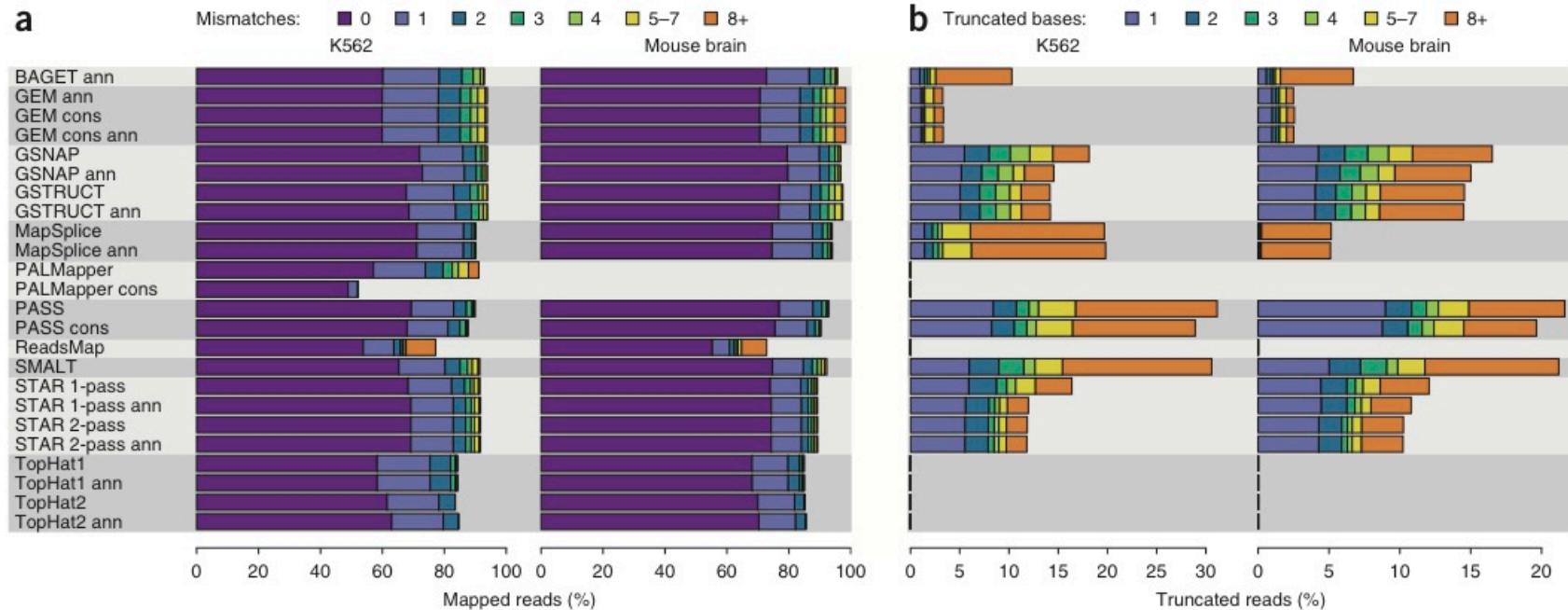
Engström et al 2013



Read placement

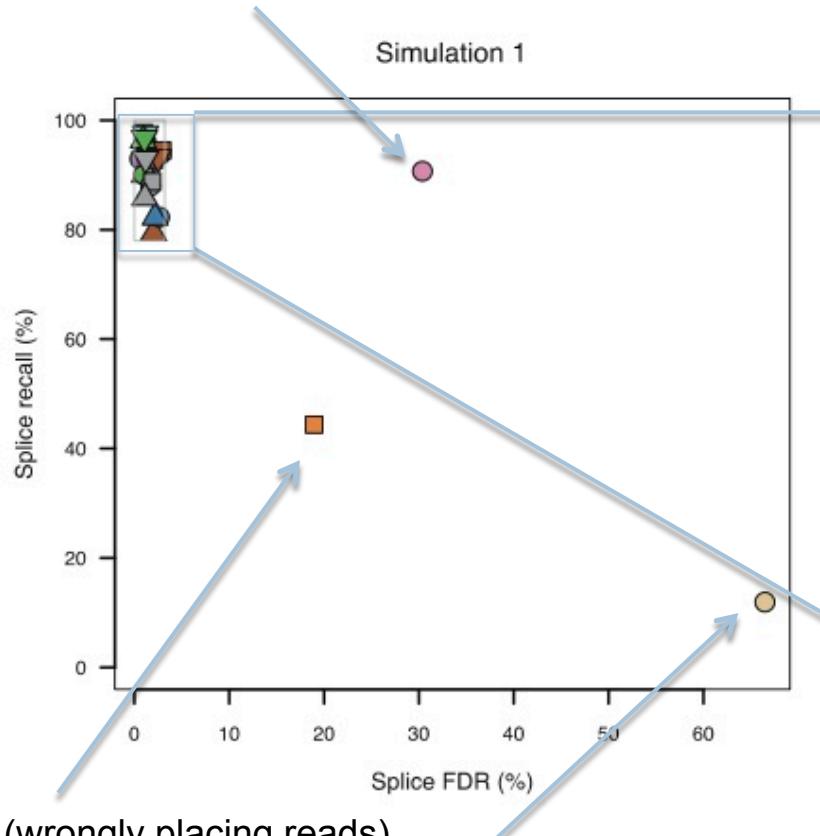


Mismatches



Evaluation

READSmap (overestimating indels)



GSNAP, GSTRUCT

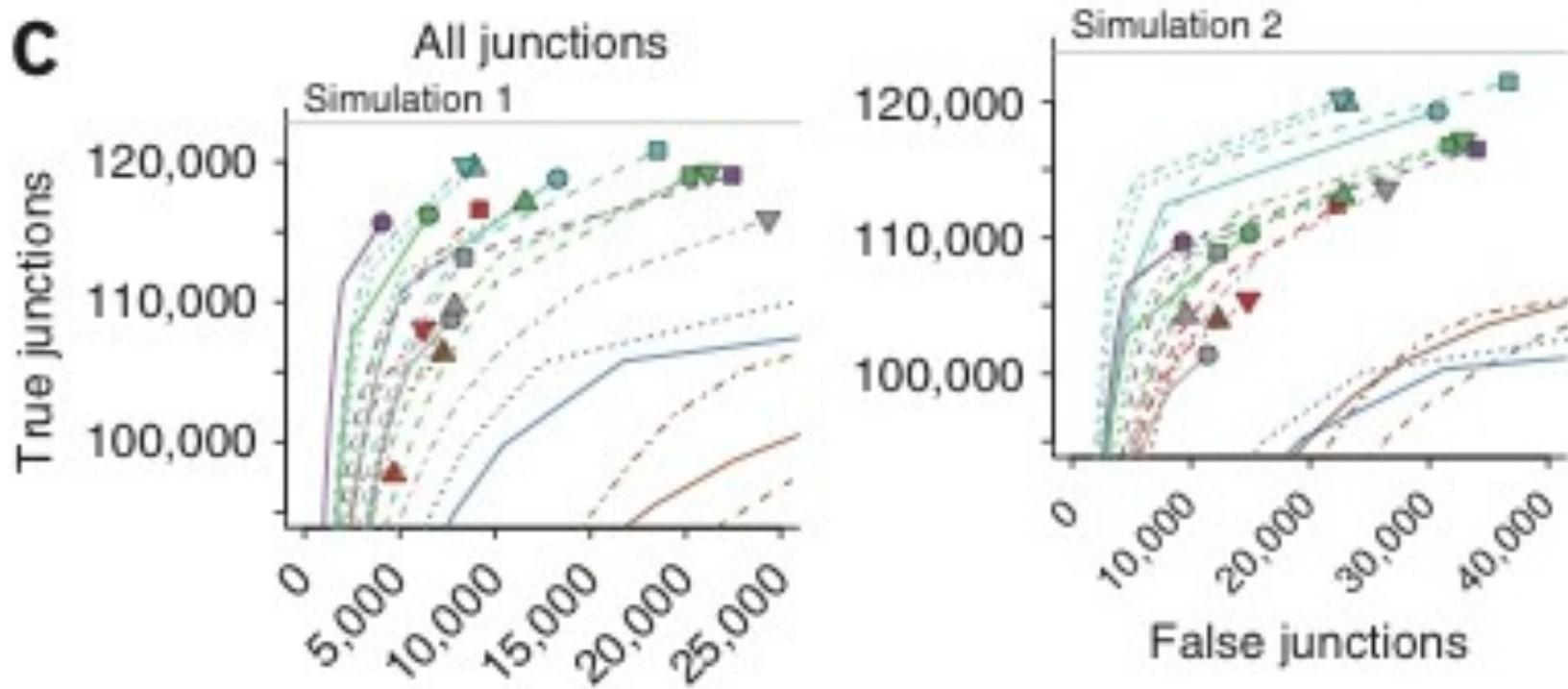
Baget (wrongly placing reads)

SMALT (missing annotated junc., wrongly placing reads)

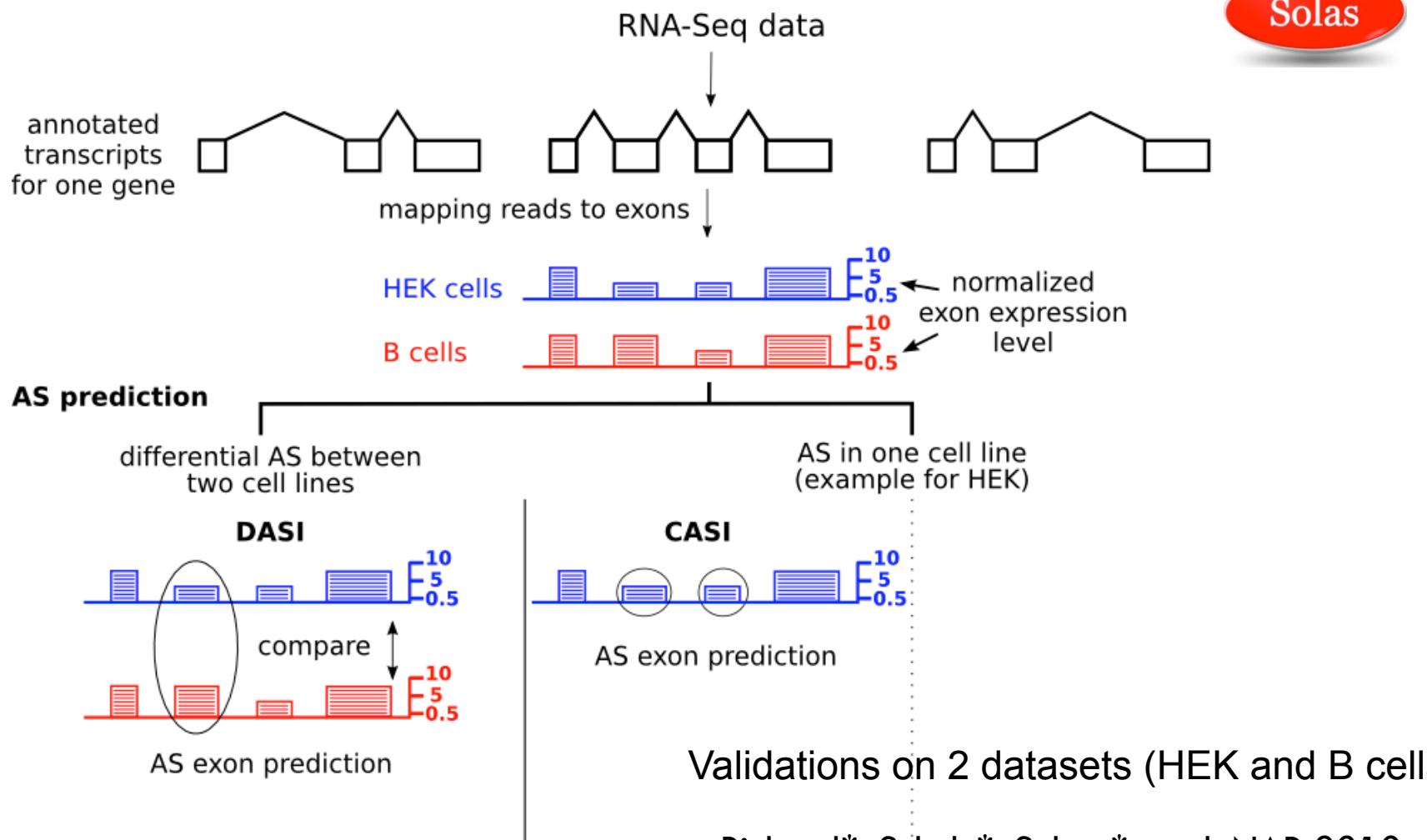
Tophat



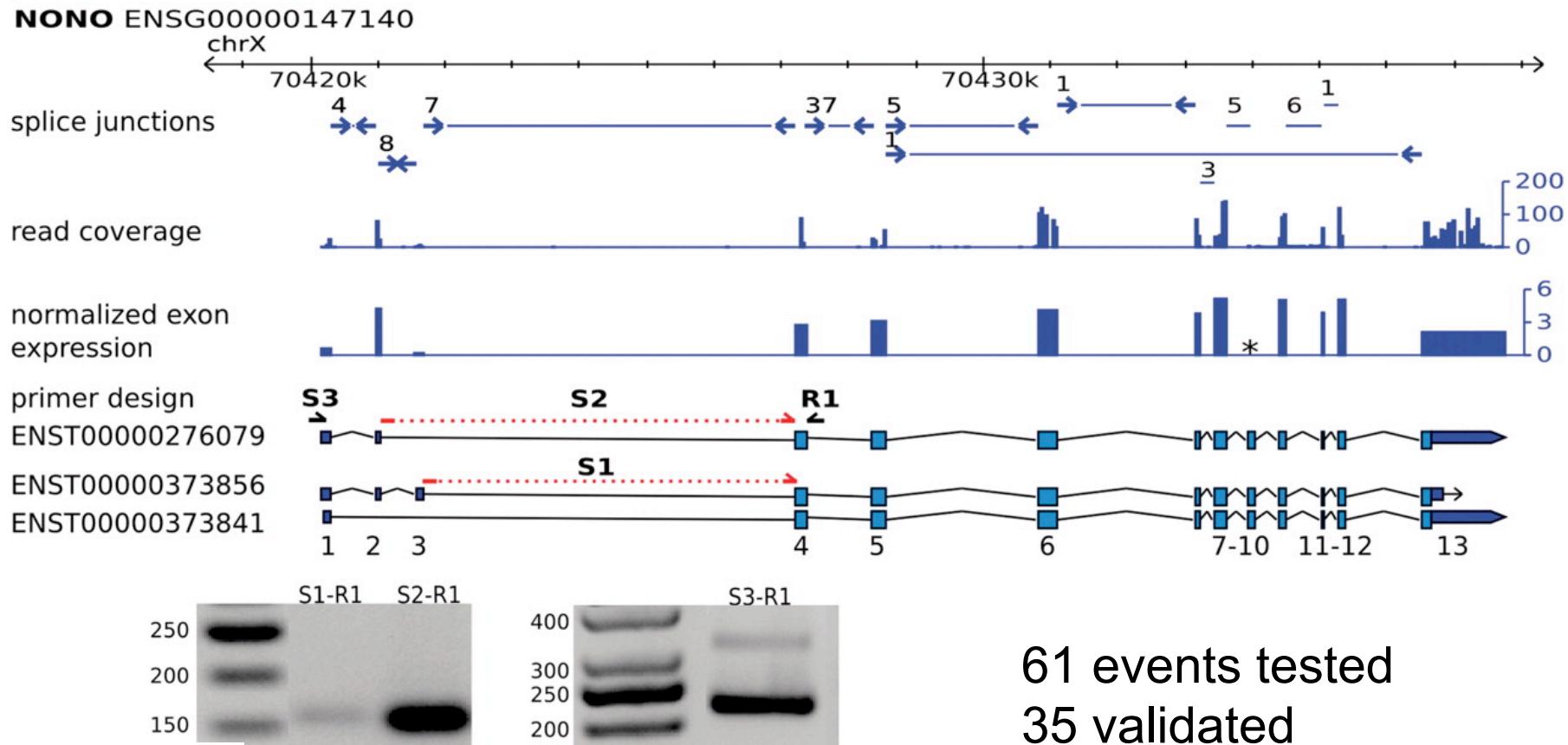
C



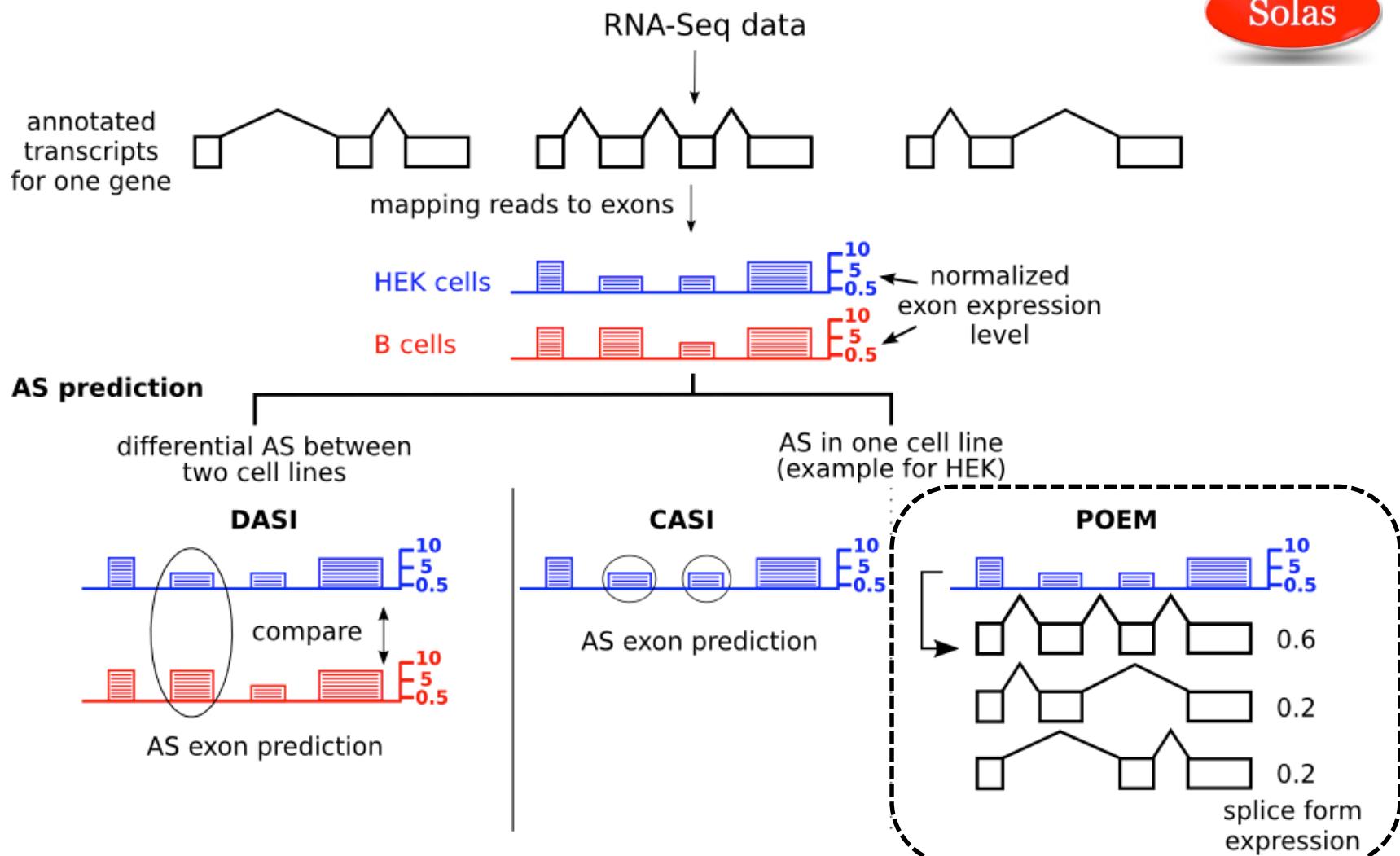
Detecting AEE from Exons Expression



Within cell AS : CASI



Detecting AEE from Exons Expression



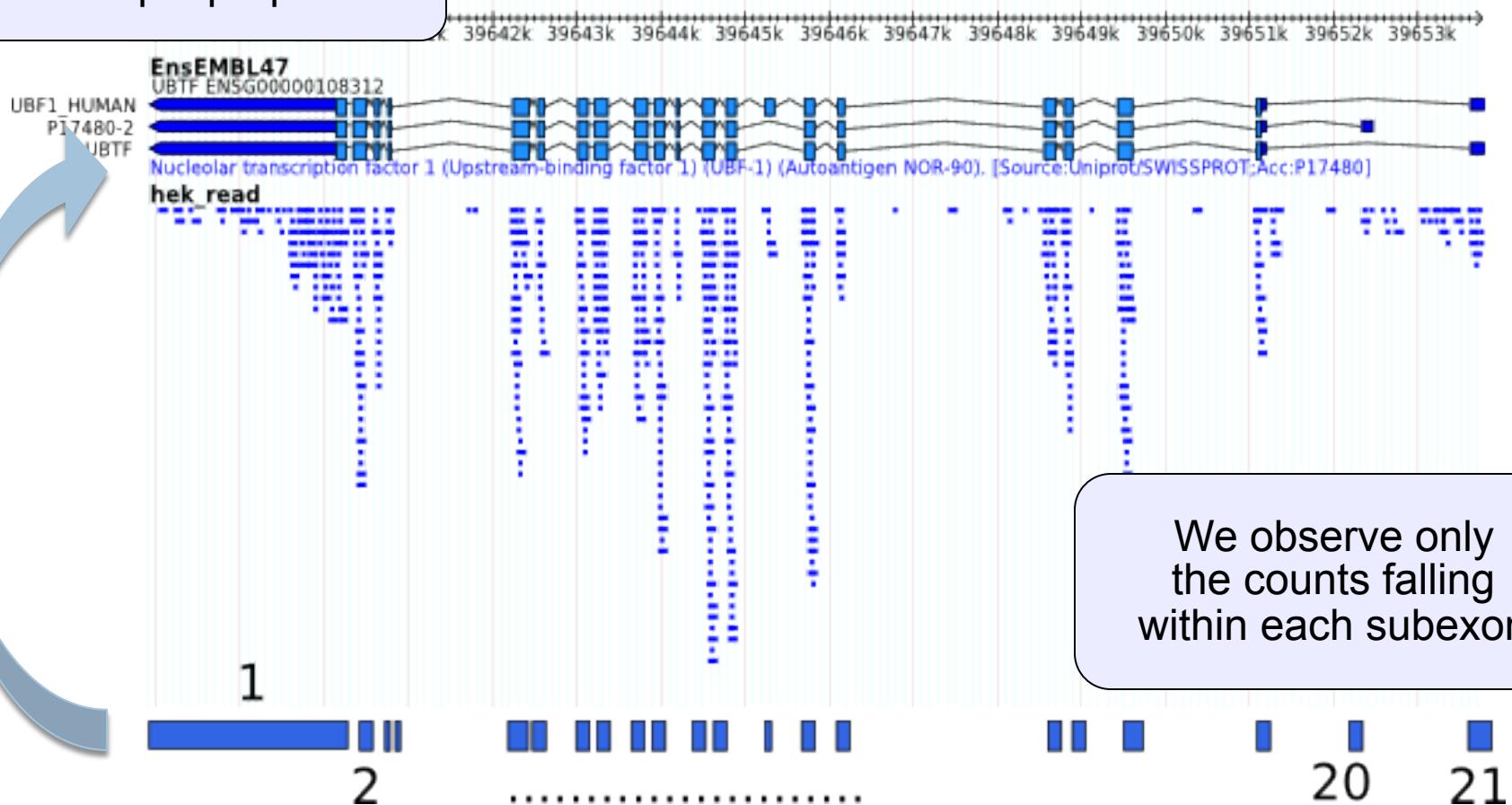
Outline

- Gene/Exon Expression level quantification
 - Sampling model and setup
 - Testing for differential expression
- Analysis of Alternative Splicing/Transcripts isoforms
 - Spliced alignment
 - **Transcript isoforms analysis**
 - Transcripts isoforms quantification
 - Transcriptome reconstruction/*de novo* assembly



Transcripts quantification

We want to infer transcripts proportions

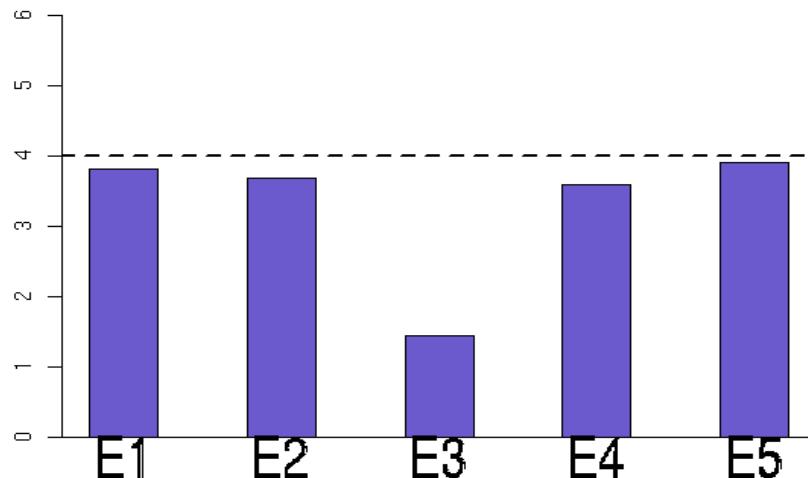
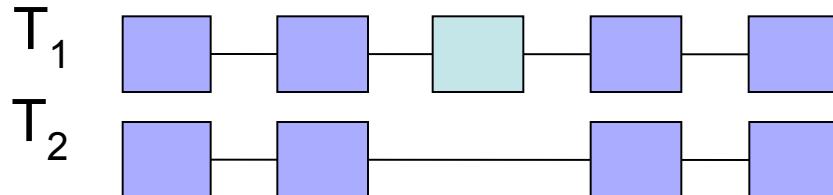


We observe only
the counts falling
within each subexon

Transcripts quantification

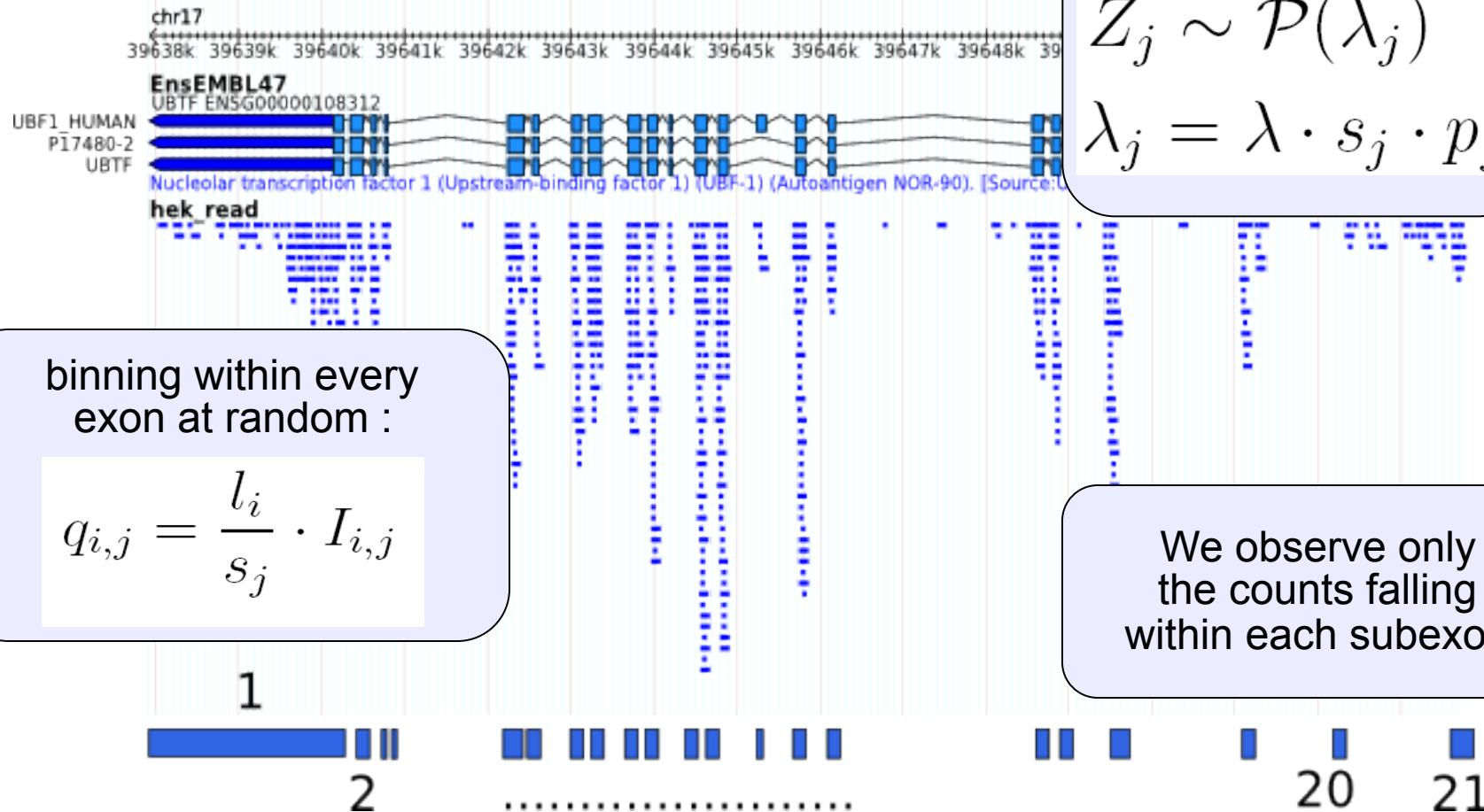
Add the transcript structure:

- Deduce directly transcript proportions
- Like solving linear equations (but with variability)

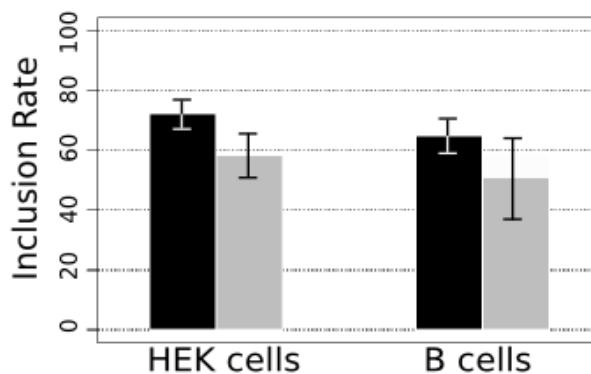
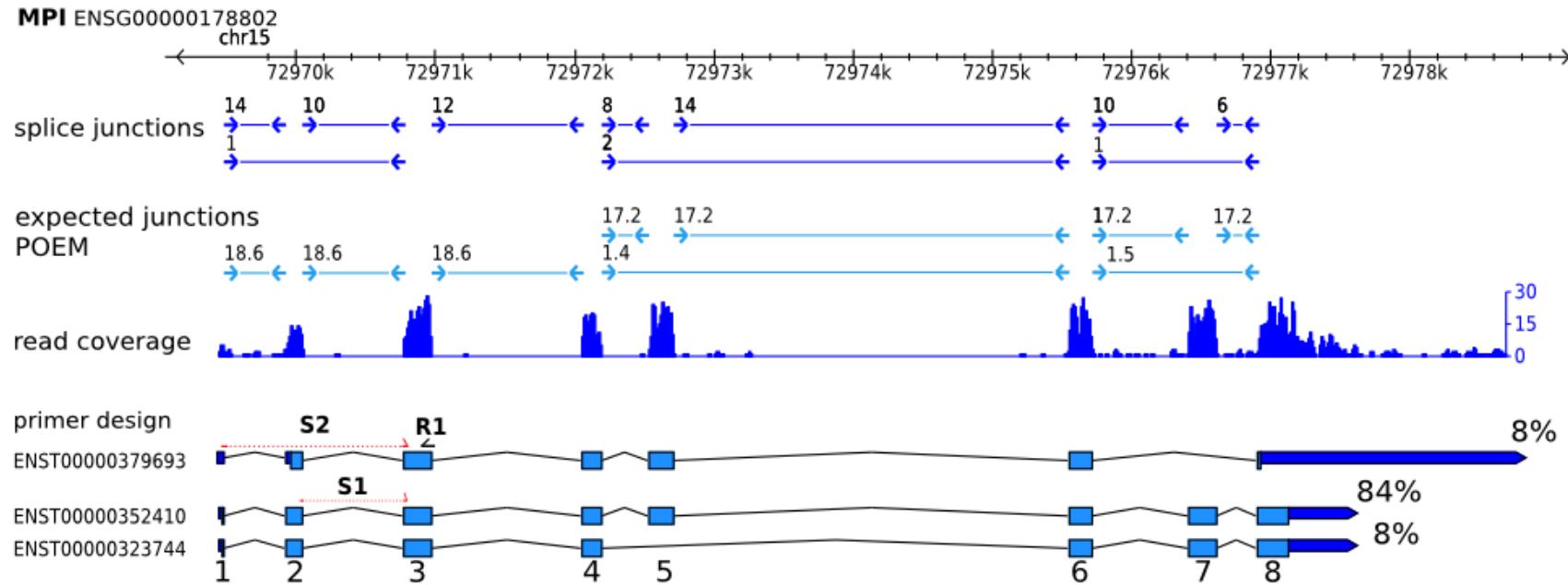


$$T_1 \% \sim \frac{\text{light blue square}}{\text{purple squares}}$$

Transcripts quantification: EM



Isoform quantification: POEM



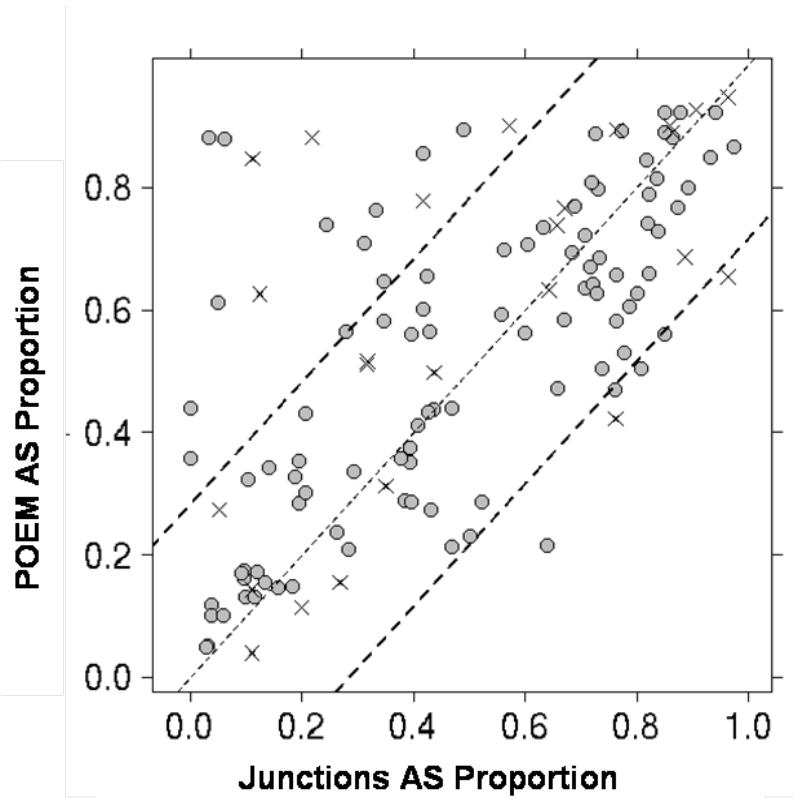
Comparison to:

- qPCR (47 events)
- Estimate from junction counts (267 events)

grey - POEM
black - qPCR

Validation

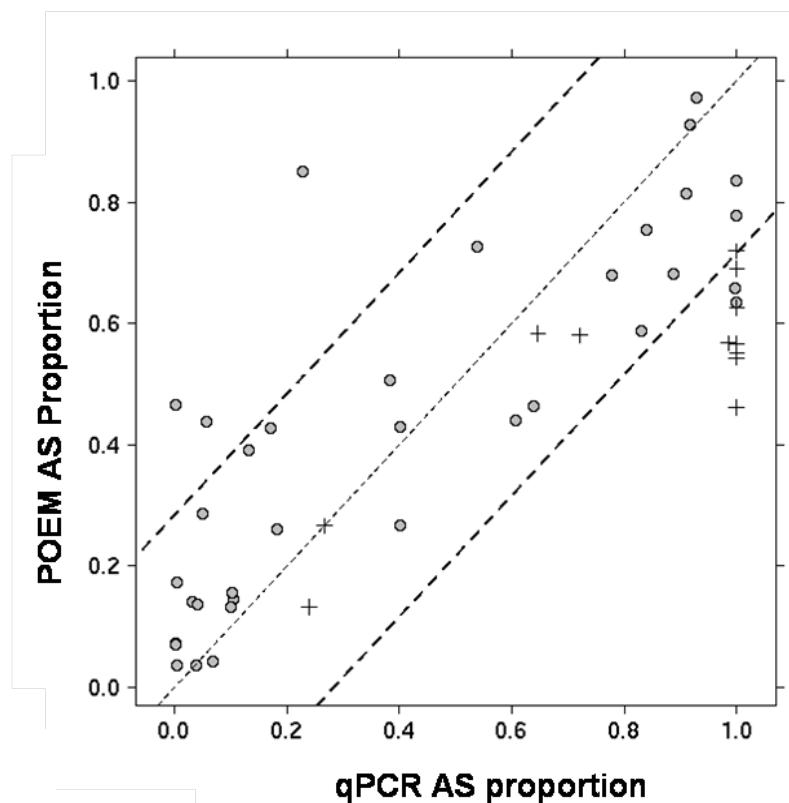
Junction reads



PCC:

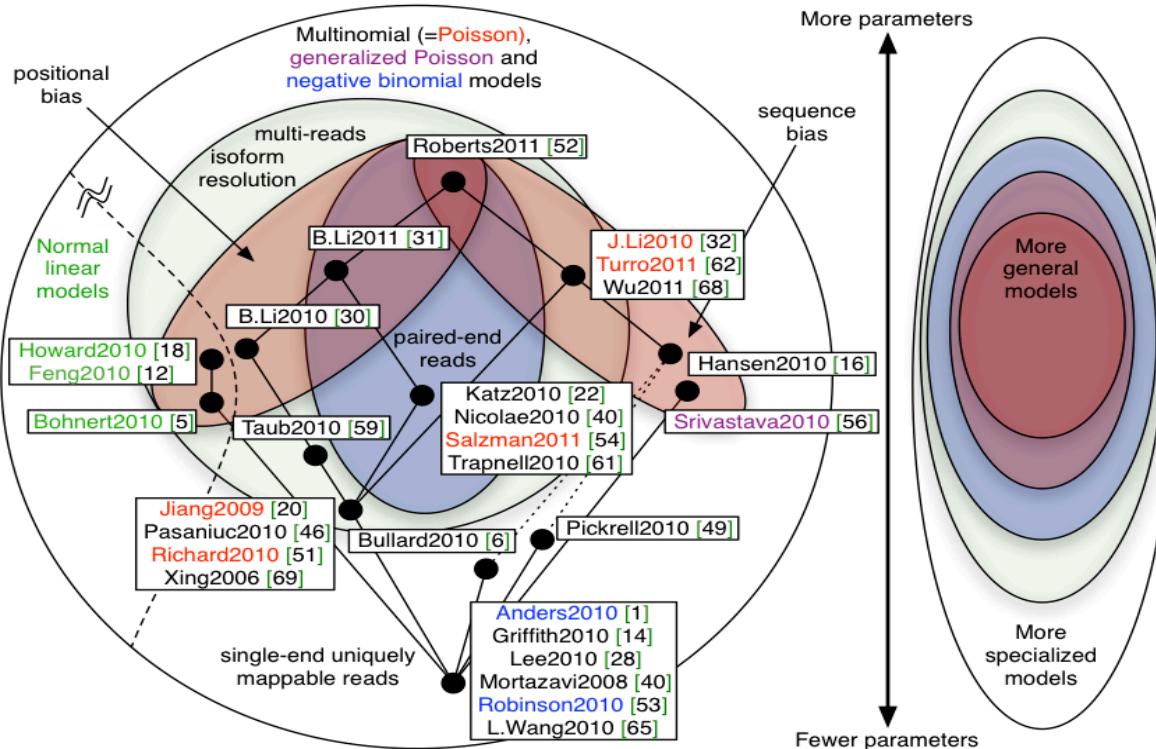
0.65

qRT-PCR



0.81

Transcript Quantification now



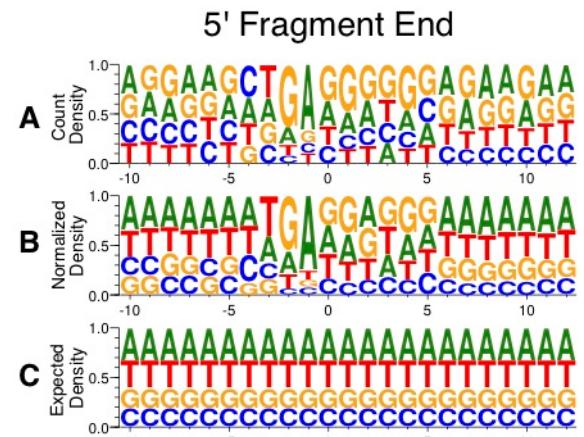
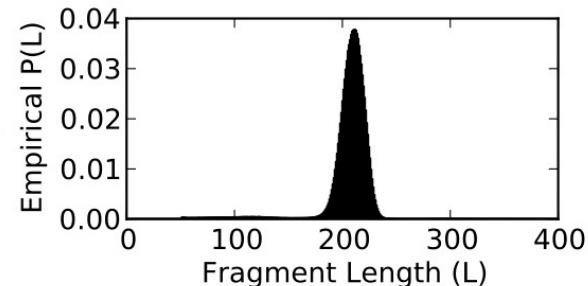
Multiple improvements lately :

- Multi-reads
- Sequence/longitudinal artifacts
- Paired-end reads

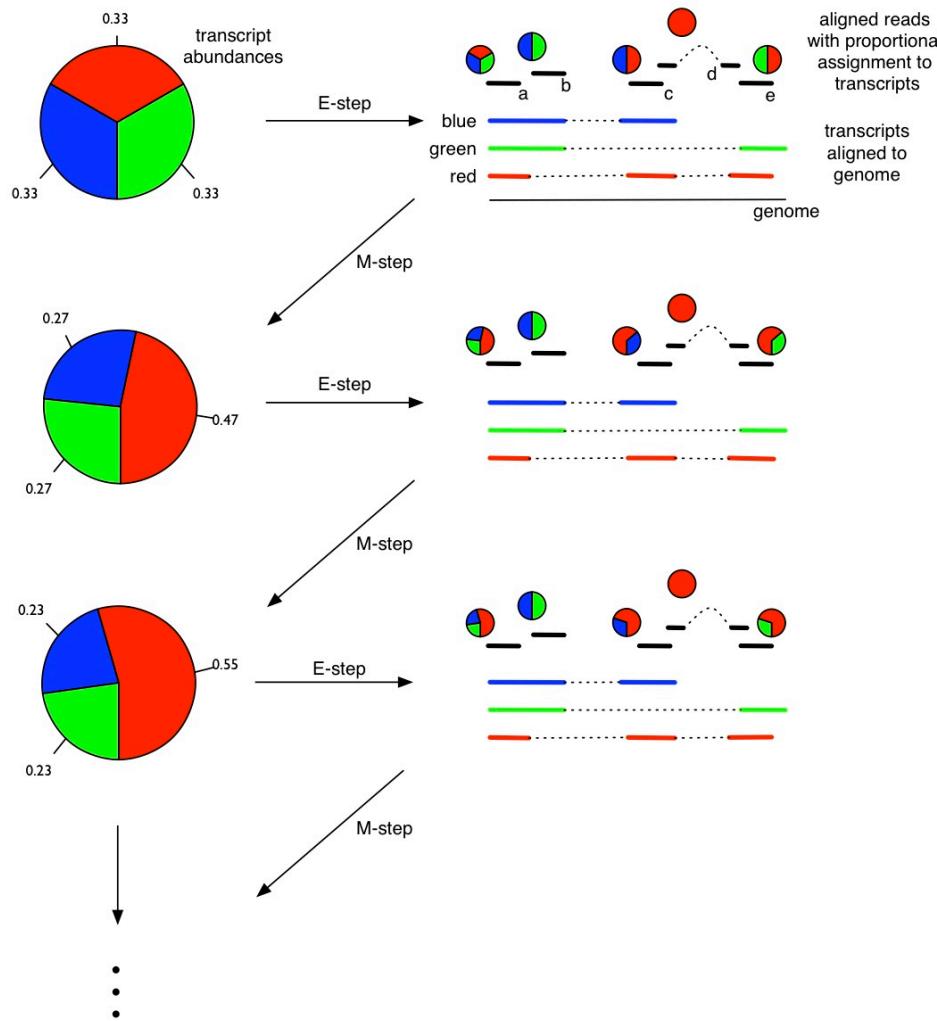
FIGURE 1. Models for RNA-Seq. The figure shows a Venn diagram (and the partially ordered set it induces) representing relationships among models. More general models are nested inside simpler models.

RSEM and other

- Incorporate library artifacts
 - fragment length distribution
 - positional biases
 - context information
- Model fragments to transcripts assignations
 - RSEM, Cufflinks, Express, IsoEM, Ireckon...
 - EM to estimate parameters and transcripts abundances



Batch EM for transcript proportions

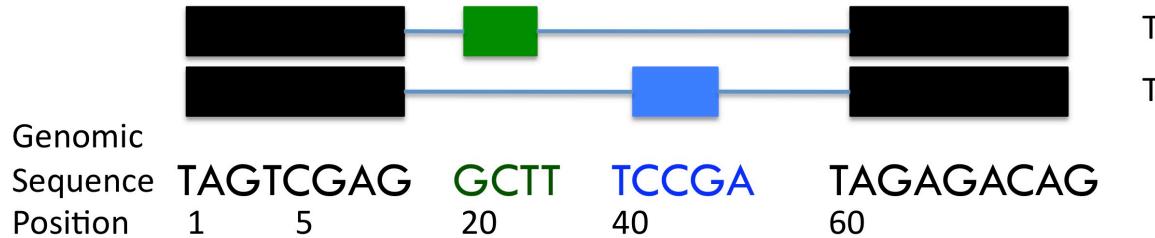


Transcripts reconstruction

- The reconstruction problem is non trivial
 - Combination of AS events can be **non identifiable** given the limited read length.
 - Paired end information is critical
 - Enumerating all possible isoforms is exponential in the number of AS events
 - Need a regularization strategy (Parcimony, Lasso)
- Can be formulated:
 - *ab initio* with a genome (Cufflinks, Ireckon, Slider, IsoLasso...)
 - *de novo* assembly (Trinity, Oases, TransAbyss...)



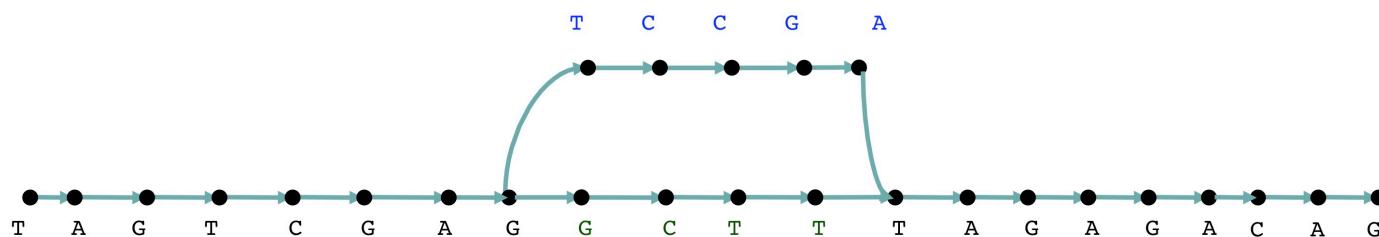
Splicing graph



Definition: Splicing graph $G=(V,E)$ (G is a Directed Acyclic Graph)

$$V = T_1 \cup T_2$$

edge (v,w) in E , if v and w are consecutive positions in T_1 or T_2



Splicing graph over genomic positions of T_1 and T_2 with representative letter at that position as node label

Heber et al. Bioinformatics, 2001

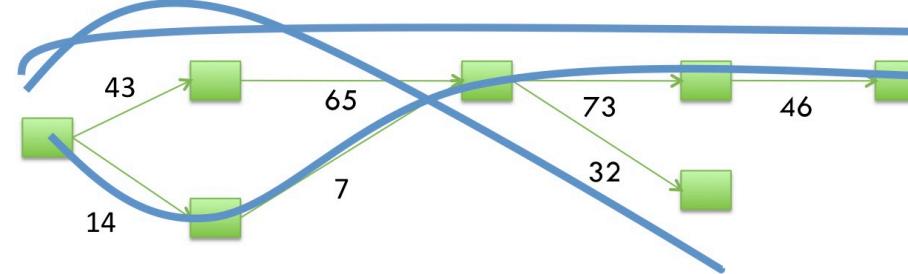
Courtesy M.Schulz



Reconstruct transcripts

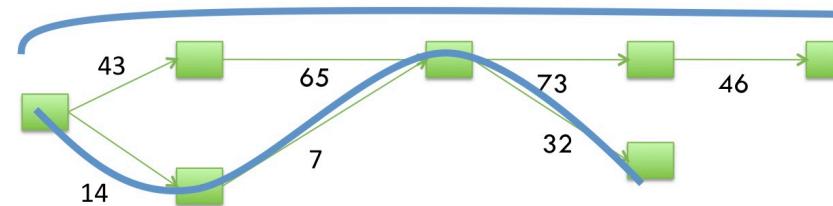
□ Maximum Likelihood approach:

assume independent AS event and a majority of constitutive exons



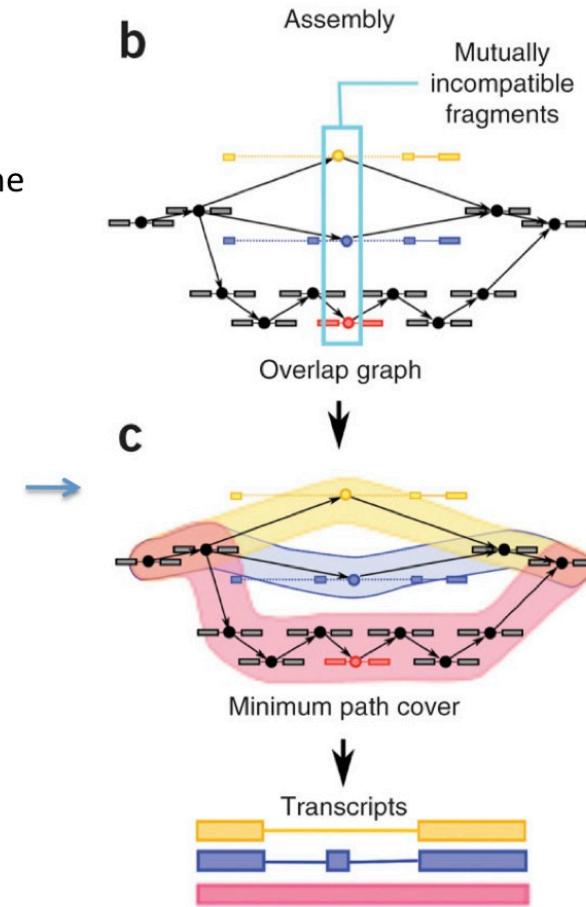
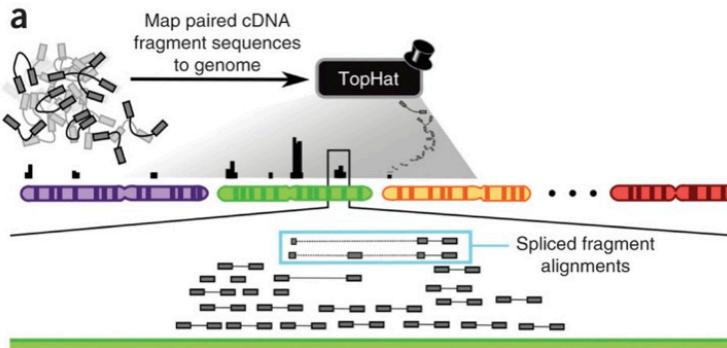
□ Maximum parsimony

constrain the number of transcripts



Cufflinks

spliced alignment of paired-end reads to genome



modified Trapnell et al. 2010



Comparing each approach

Approach	Advantages	Disadvantages
<i>ab initio</i> Reference based	<ul style="list-style-type: none">- alignment tolerate seq. errors- repeat detected through alignment- grouping by genomic proximity	<ul style="list-style-type: none">- need the reference seq.- assumes collinearity of the transcripts with the genome
<i>de novo</i>	<ul style="list-style-type: none">- no reference needed- detection of non-collinear transcripts- handles micro-exons	<ul style="list-style-type: none">- lowly expressed genes are indistinguishable from seq. errors- missassemblies due to repeats



Comparing approaches on mouse RNA-Seq

Method	transfrag > 100 bps	N50	total in mb	Nucleotide Sensitivity			Nucleotide Specificity
				all genes	lowly exp. genes RPKM < 1	highly exp. genes RPKM > 20	
Cufflinks*	72,745	2,613	~ 73	45.3	24.1	71.4	67
Oases\$	73,357	1,287	~ 64	28.3	0.9	64.8	85.3

*Trapnell et al. *Nat. Biotech.* 2010

\$Schulz, Zerbino et. al, submitted

Results for Ensembl 57 annotation

Courtesy M.Schulz

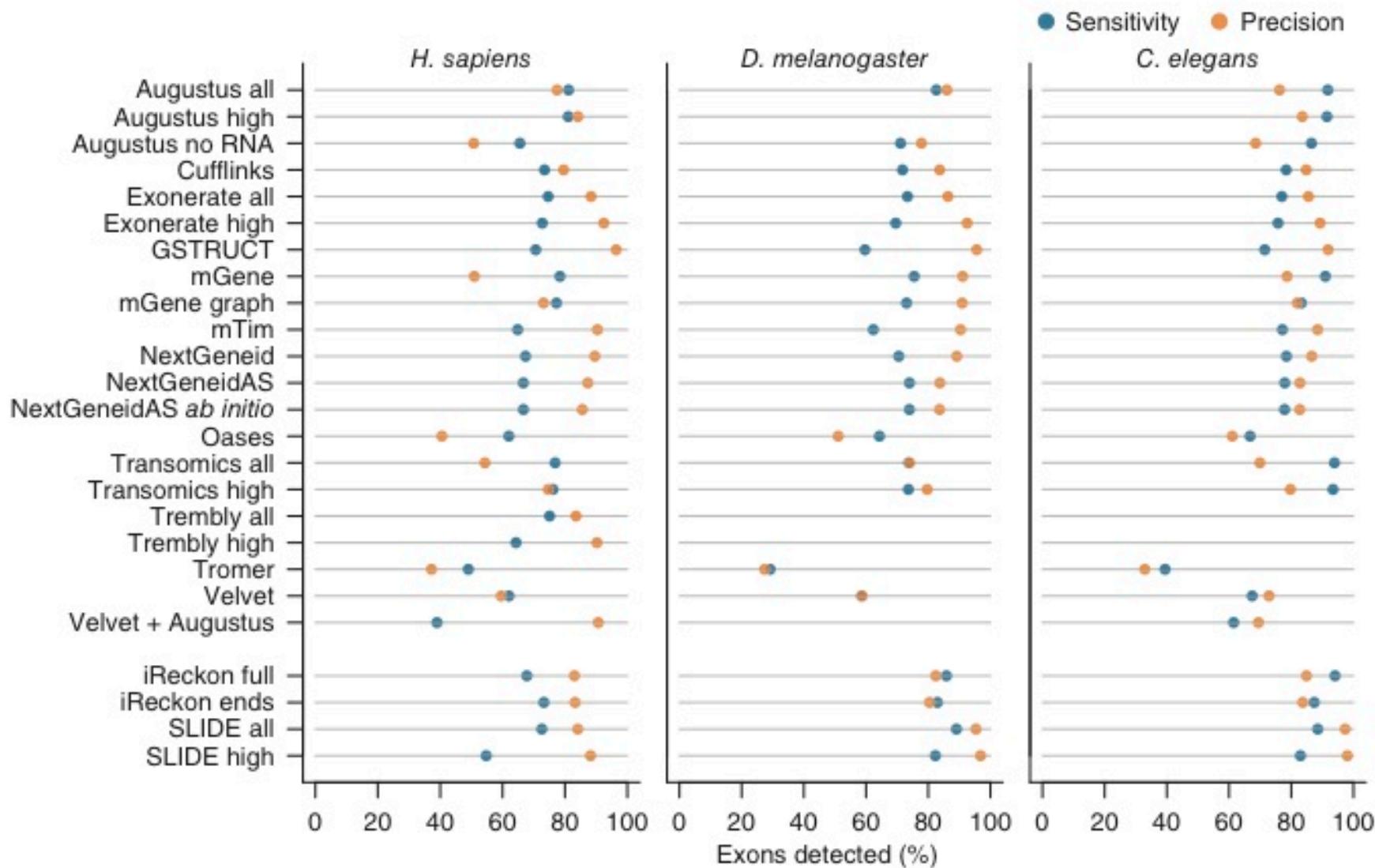


RGASP – transcript reconstruction

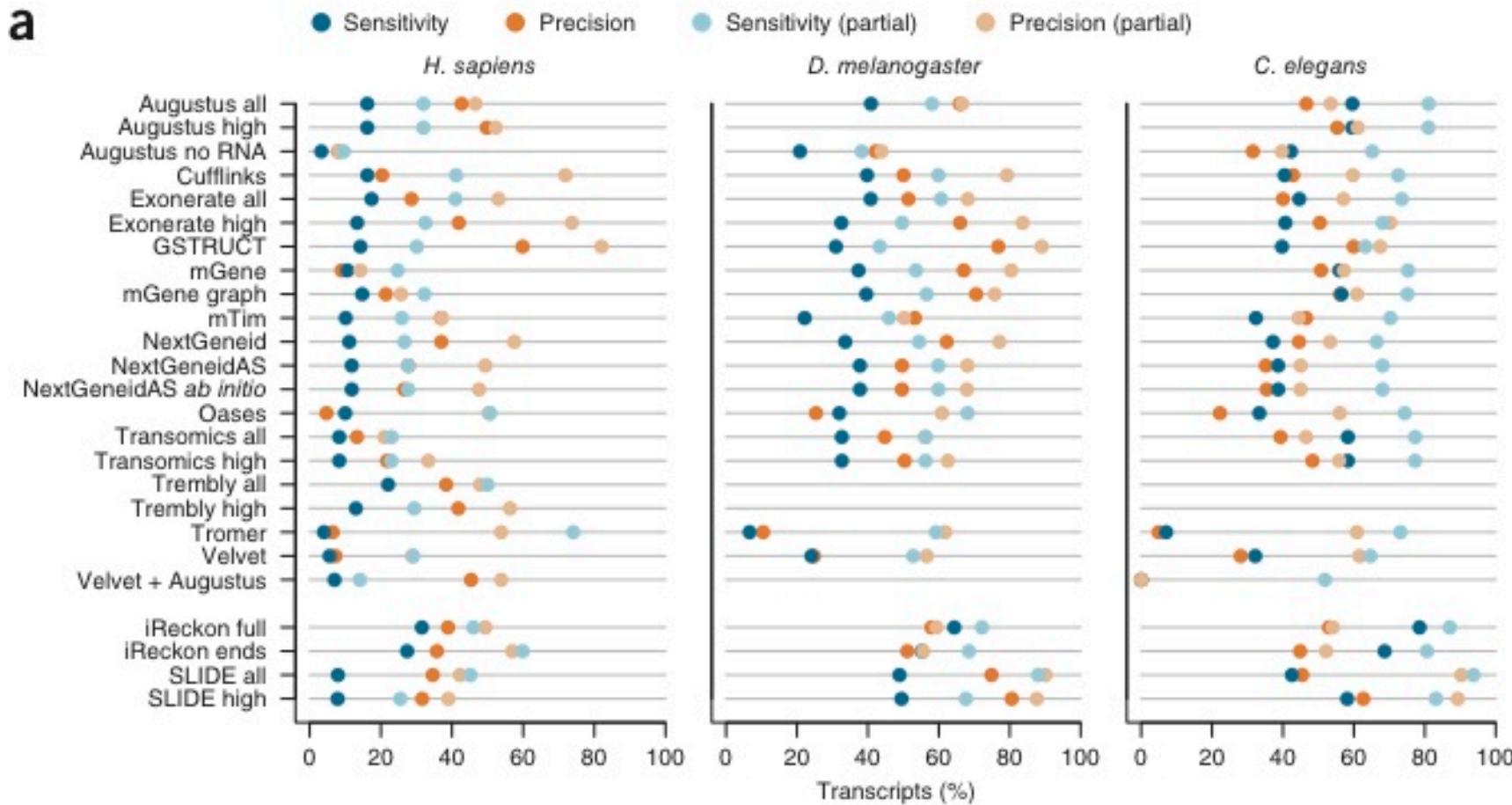
- Evaluate 25 transcripts reconstruction methods:
 - Ab initio methods:
Cufflinks, iReckon, SLIDE, Augustus, mGene, Transomics
 - De novo:
Oases, Velvet + Augustus.



results- exon level



results- transcript level



Summary transcriptome

- Transcriptome is more complicated than the genome
 - Often still use a genome proxy
 - Quantification from digital data (counts)
 - Identification of exons boundary is a special alignment problem with large gaps
- Transcripts quantification and reconstruction:
 - Inferring transcript proportions: linear system or models with Hidden variables
 - When transcripts are unknown, reconstructing transcript structure is an additional challenge.



