

State of the art of what can be done with RNA-seq

Marc DELOGER / Nicolas SERVANT

Gustave Roussy - UMS AMMICA (US 23 INSERM / UMS 3655 CNRS) / Curie
- U900 (INSERM / MINES ParisTech)

DU Bii - March 24th 2021

Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Outline

- **Experimental and sequencing design**
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Experimental and sequencing design

The sequencing of mRNA can be used to address many different biological questions : expression, alternative splicing, RNA editing, fusion gene, repeats profiling, etc...

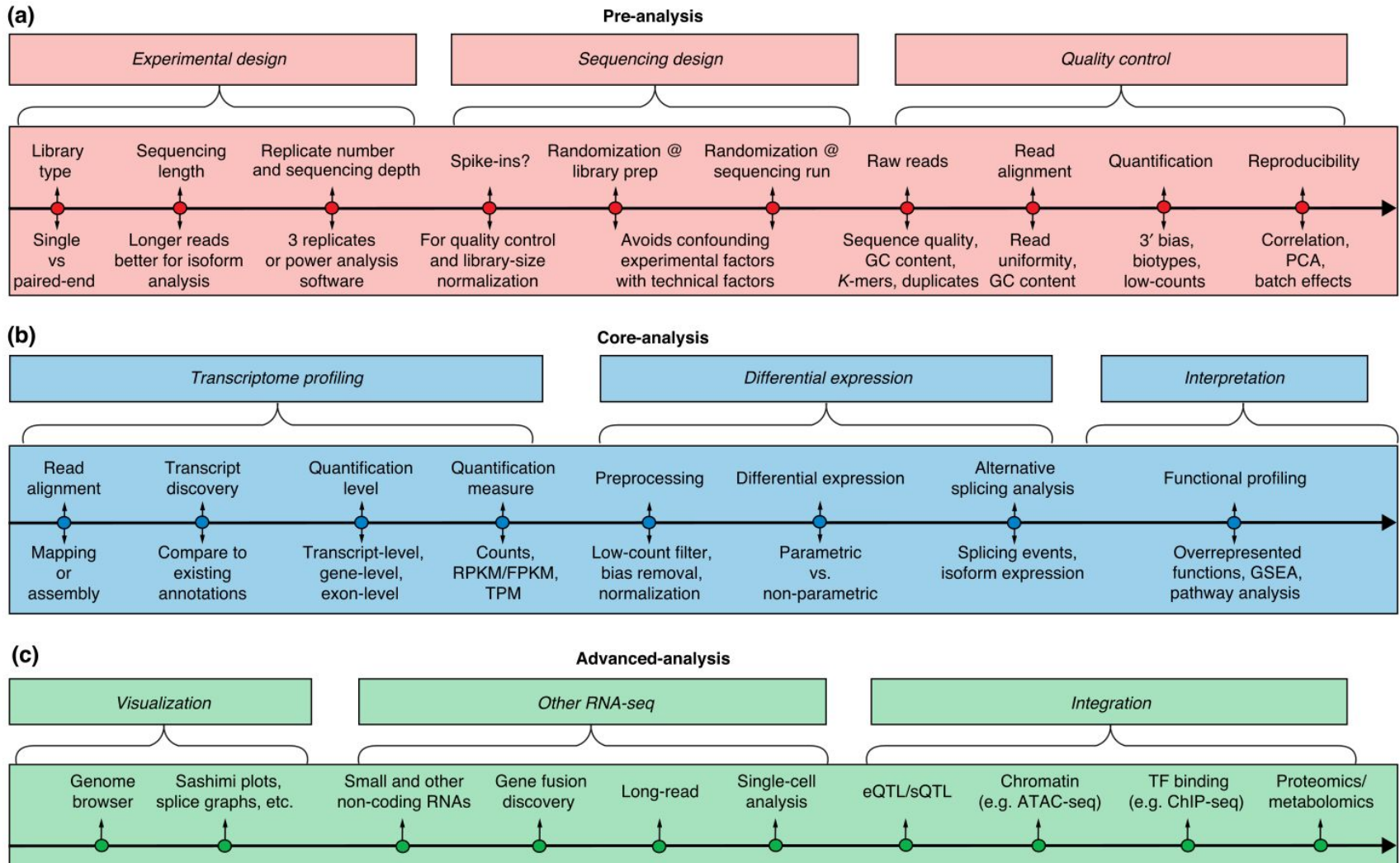
These questions have to be clearly defined **BEFORE** running the sequencing part of the project => Experiment “design”

For each of them, **dedicated** sequencing “design” and bioinformatics pipeline should be used

Experimental and sequencing design

- Citation 1 : "To consult a statistician after an experiment is finished is often merely to ask him to conduct **a post-mortem examination**. He can perhaps say what the experiment died of." (Ronald A. Fisher, Indian Statistical Congress, 1938, vol. 4, p 17)
- Citation 2 : "While **a good design does not guarantee a successful experiment**, a suitably bad design guarantees a failed experiment" (Kathleen Kerr, Atelier Inserm 145, 2003)

Experimental and sequencing design



Experimental and sequencing design

- Samples collection :
 - Number of replicates per condition : achieve enough statistical power to control experiment variability, and so, to be able to answer to your biological question

Replicates : what does that mean ?

- Replicates are mandatory to estimate the biological variability
- The higher the better !
- A biological replicate is not a technical replicate :
 - Technical = Several extractions of the same RNA or Several libraries built from the same RNA extraction or A library sequenced several times

(...) With three biological replicates, nine of the 11 tools evaluated found only 20%–40% of the significantly differentially expressed (SDE) genes identified with the full set of 42 clean replicates (...)

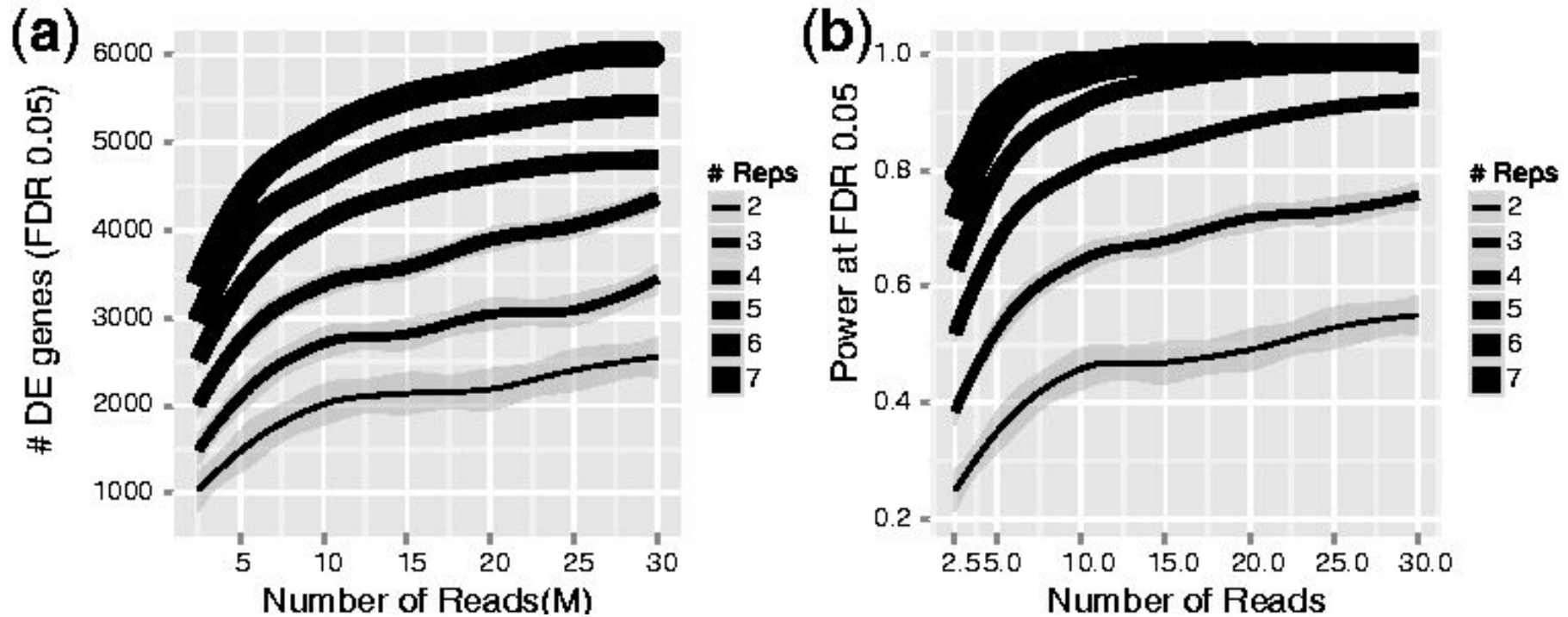
Recommandations for RNA-seq experimental designs :

At least 6 replicates per condition for all experiments.

At least 12 replicates per condition for experiments where identifying the majority of all DE genes is important.

(...)

More reads or more replicates



- (...) Increase in biological replication significantly increases the number of DE genes identified (...)
- (...) Power of detecting DE genes increases with both sequencing depth and biological replication (...)

Experimental and sequencing design

- Samples collection :
 - Number of replicates per condition : achieve enough statistical power to control experiment variability, and so, to be able to answer to your biological question
 - Avoid confounding effects : eg. all the samples for condition 1 processed by technician A (or in year/gender A) and all the samples for condition 2 processed by technician B (or in year/gender B)

Experimental and sequencing design

- Library preparation :
 - Poly-A selection
 - High quantity/quality of input RNA (100-300ng and RIN > 6)
 - Loss of all non-polyadenylated features (majority of non-coding and small RNAs)
 - rRNA depletion
 - Variable depletion procedures efficiency (Petrova OE et al. Sci Rep. 2017)
 - More expensive (+33%)
 - Small RNA / miRNA
 - Unstranded / Stranded protocol => strand-effect control
 - UMI (Unique Molecule Identifier) => PCR-effect control
 - Spike-ins => sensitivity and accuracy of RNA-seq experiments for transcriptome discovery and quantification across different samples (eg. short and/or rare transcripts)

Experimental and sequencing design

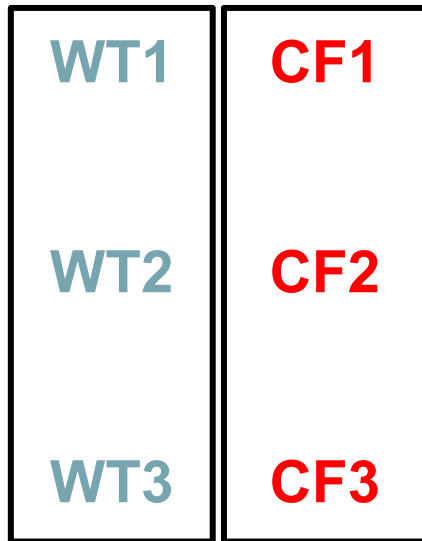
- Sequencing :
 - Paired-End : standard for years now
 - Single-End : no more used except for gene expression only because sequencing price / 2 compared to paired-end and with highly similar gene counts (Pearson $R^2=0.9792$) **so you can multiply by 2 the number of replicates at same price ;-)**
 - Reads length : 50bp for gene expression, >75bp for the rest
 - Insert length : generally ~300bp fragments, the higher the better for splicing/fusion event discovery
 - Depth : single-cell (1M reads or 50K=highly expressed or 20K=cell types), 30M fragments for gene expression, >50/75M fragments for the rest (75M pairs => 150M reads)
 - Run design

Run design

Goal:

Do not add any confounding technical effect (day, lane, run, etc...) to the factor of interest.

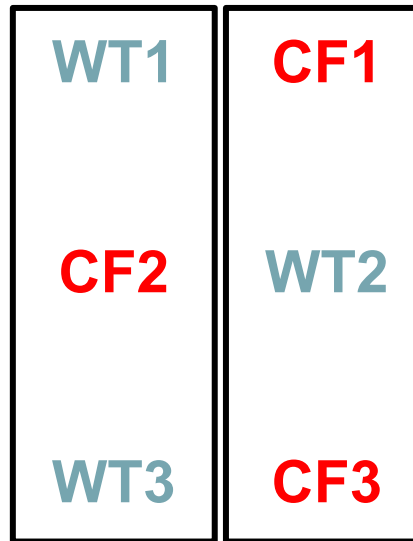
Bad example 



Lane 1

Lane 2

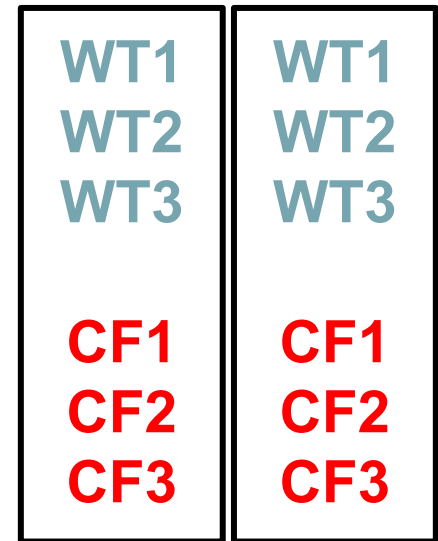
Good example 



Lane 1

Lane 2

Good example 



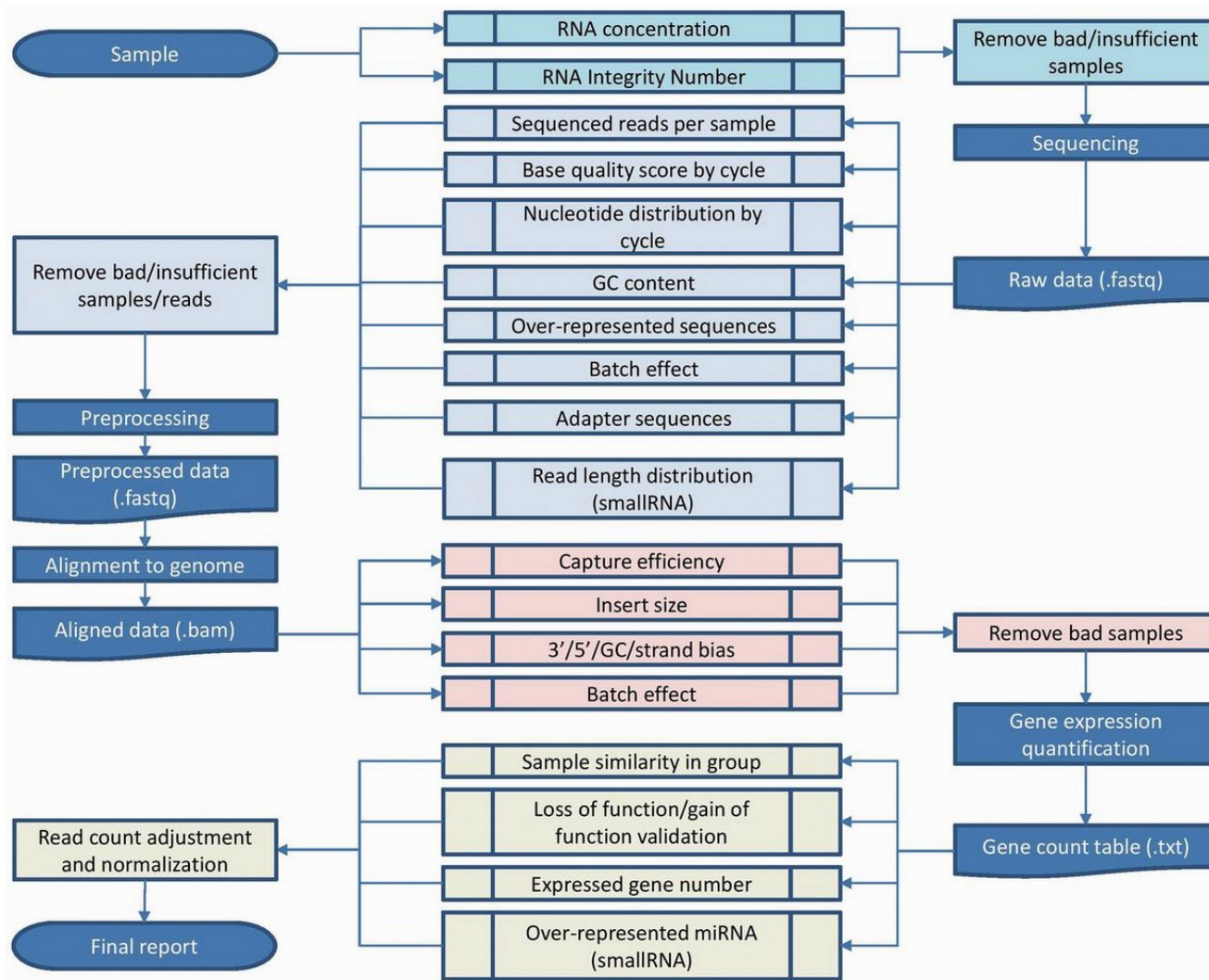
Lane 1

Lane 2

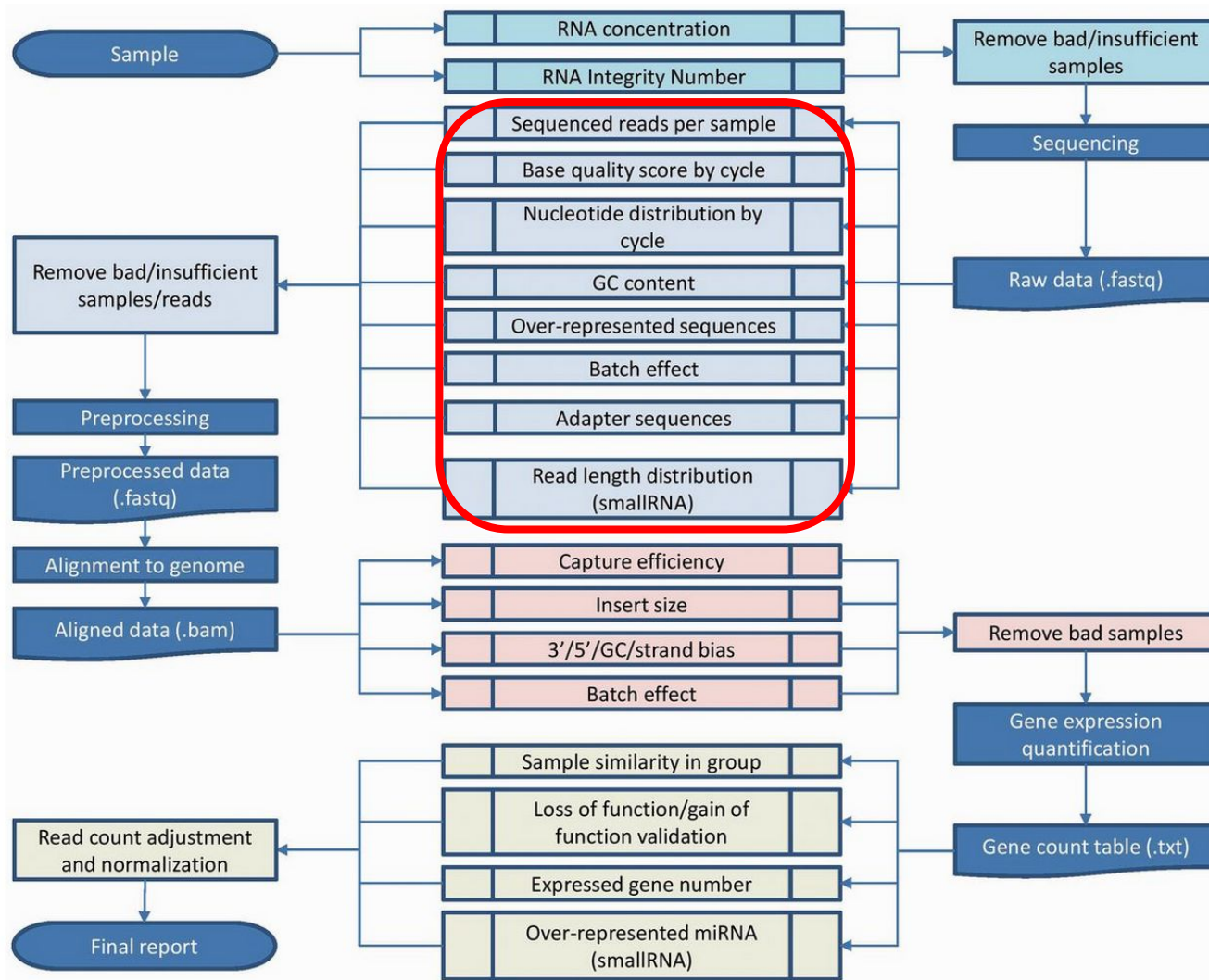
Outline

- Experimental and sequencing design
- **Quality control and mapping**
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Quality control and mapping

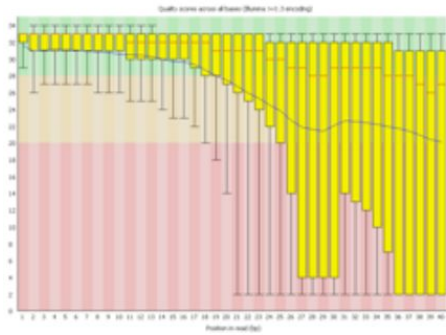


Quality control and mapping

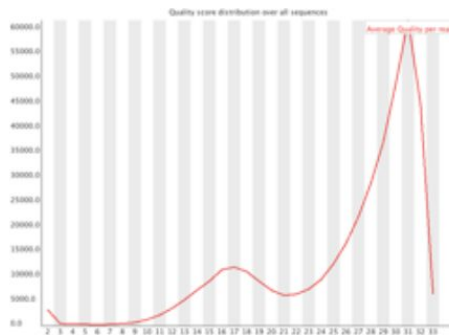


Quality control and mapping

Per Base Sequence Quality



Per Sequence Quality Scores



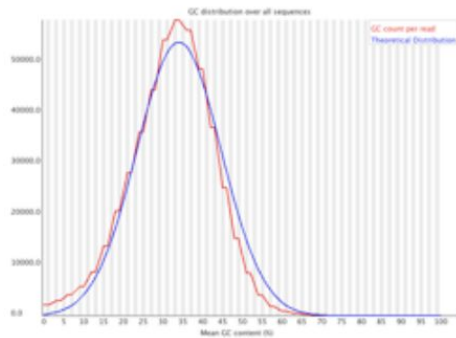
Per Base Sequence Content



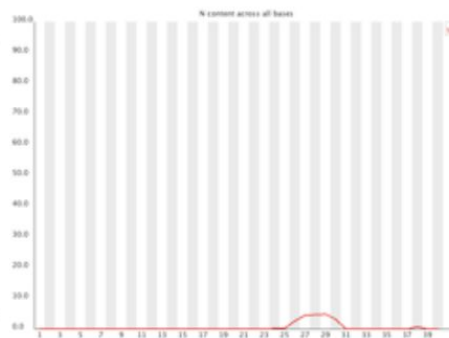
Per Base GC Content



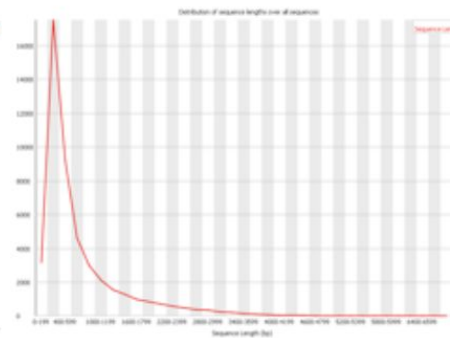
Per Sequence GC Content



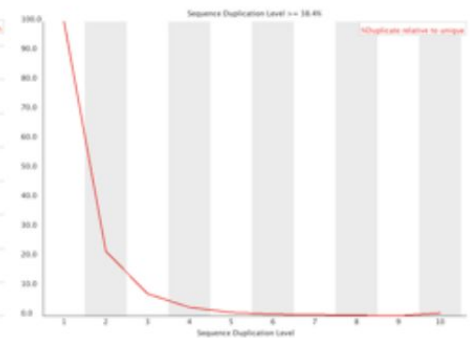
Per Base N Content



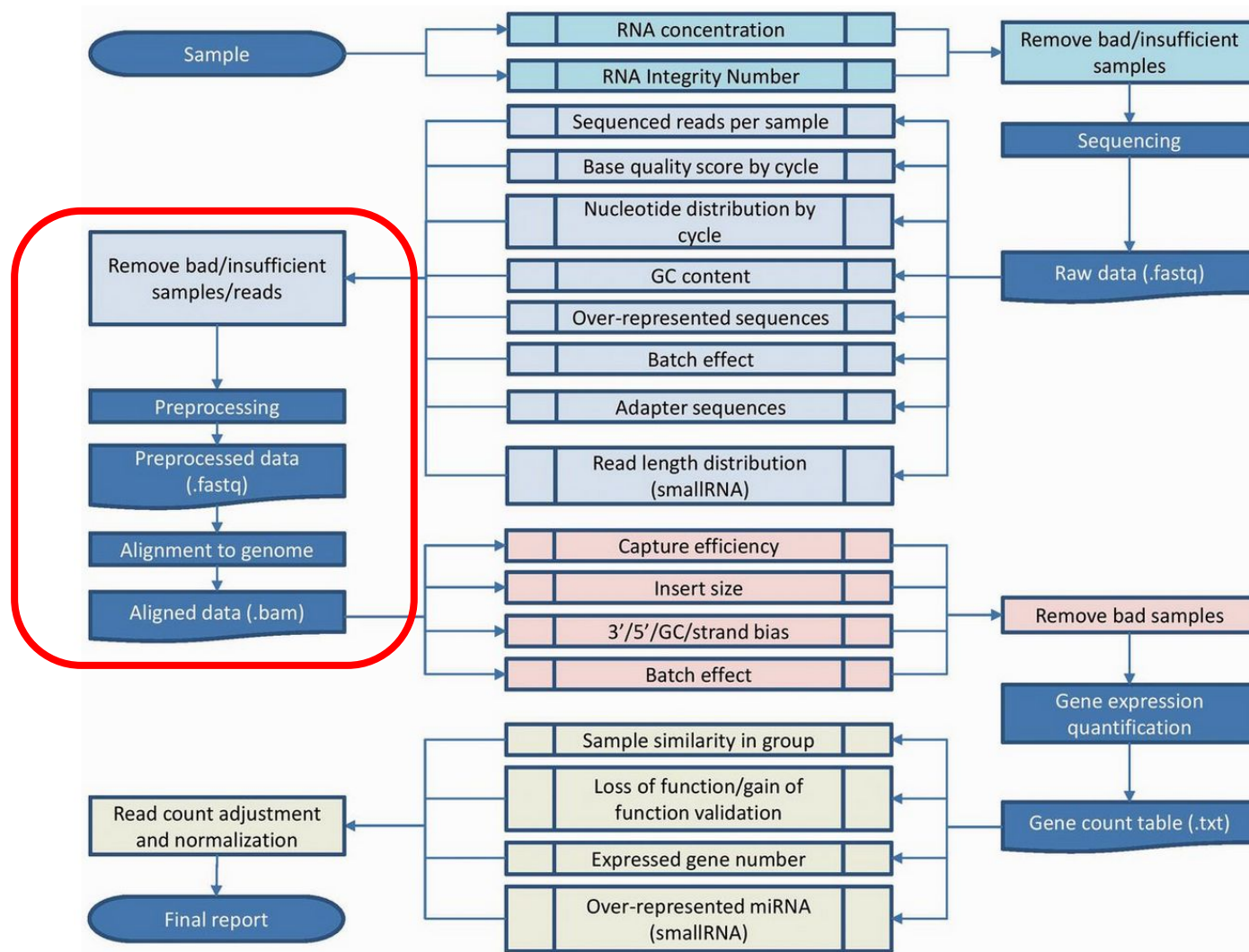
Sequence Length Distribution



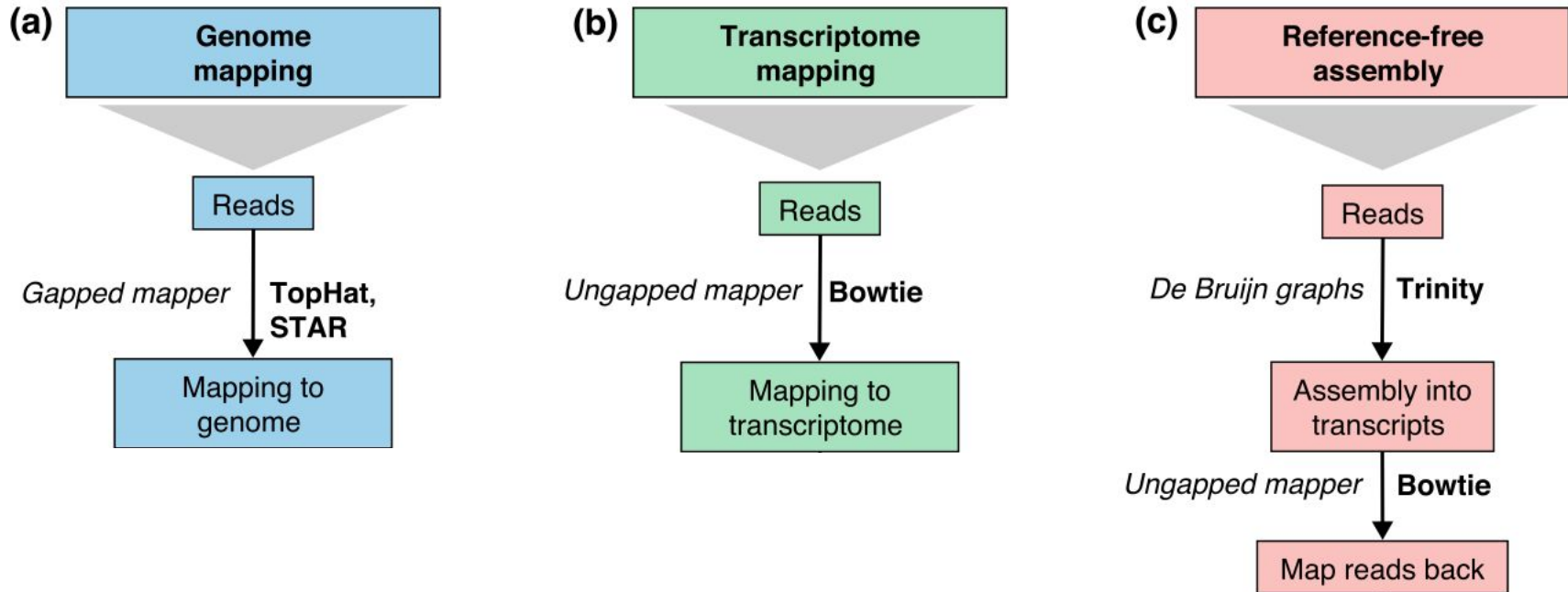
Duplicate Sequences



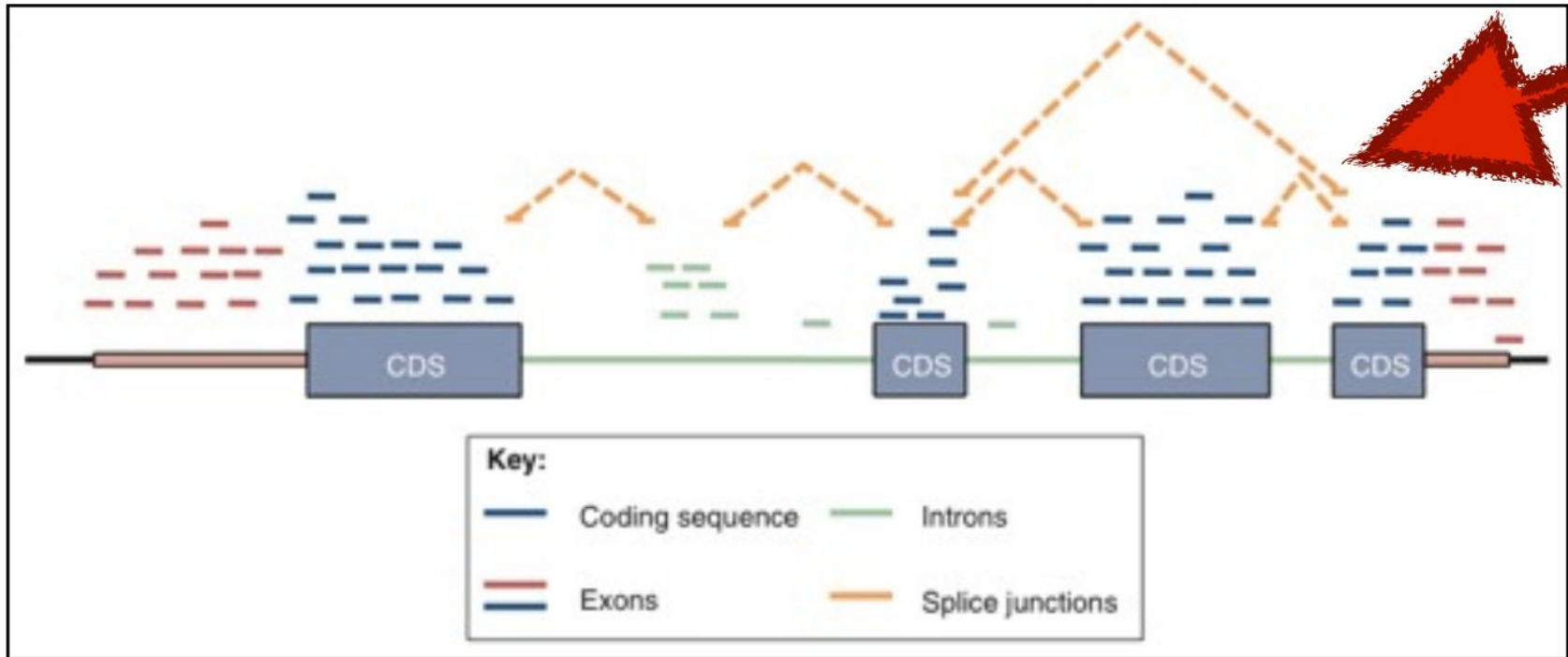
Quality control and mapping



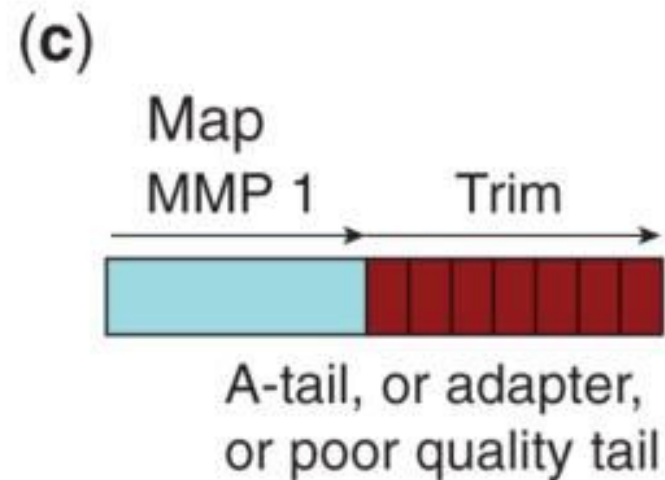
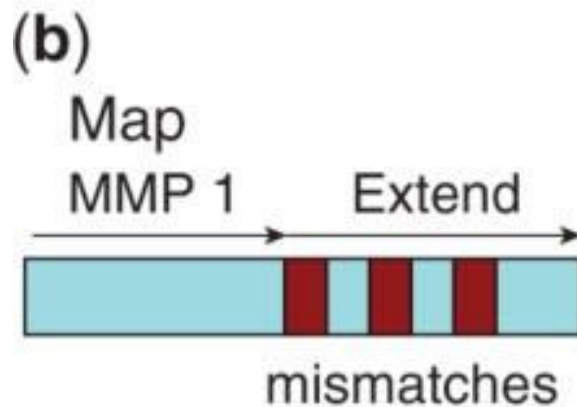
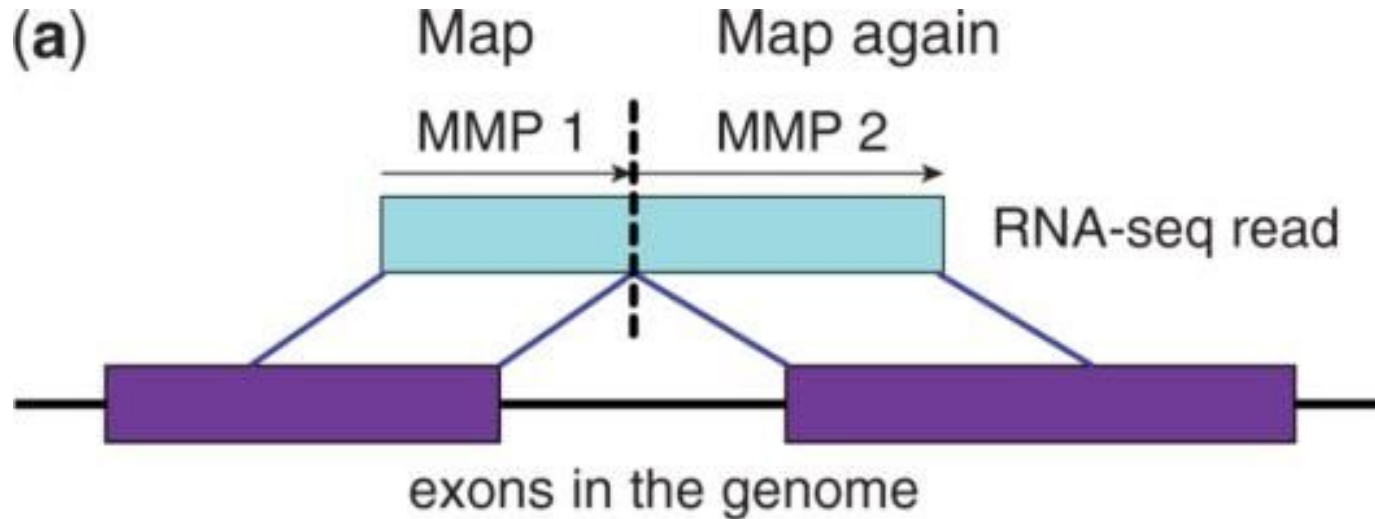
Quality control and mapping



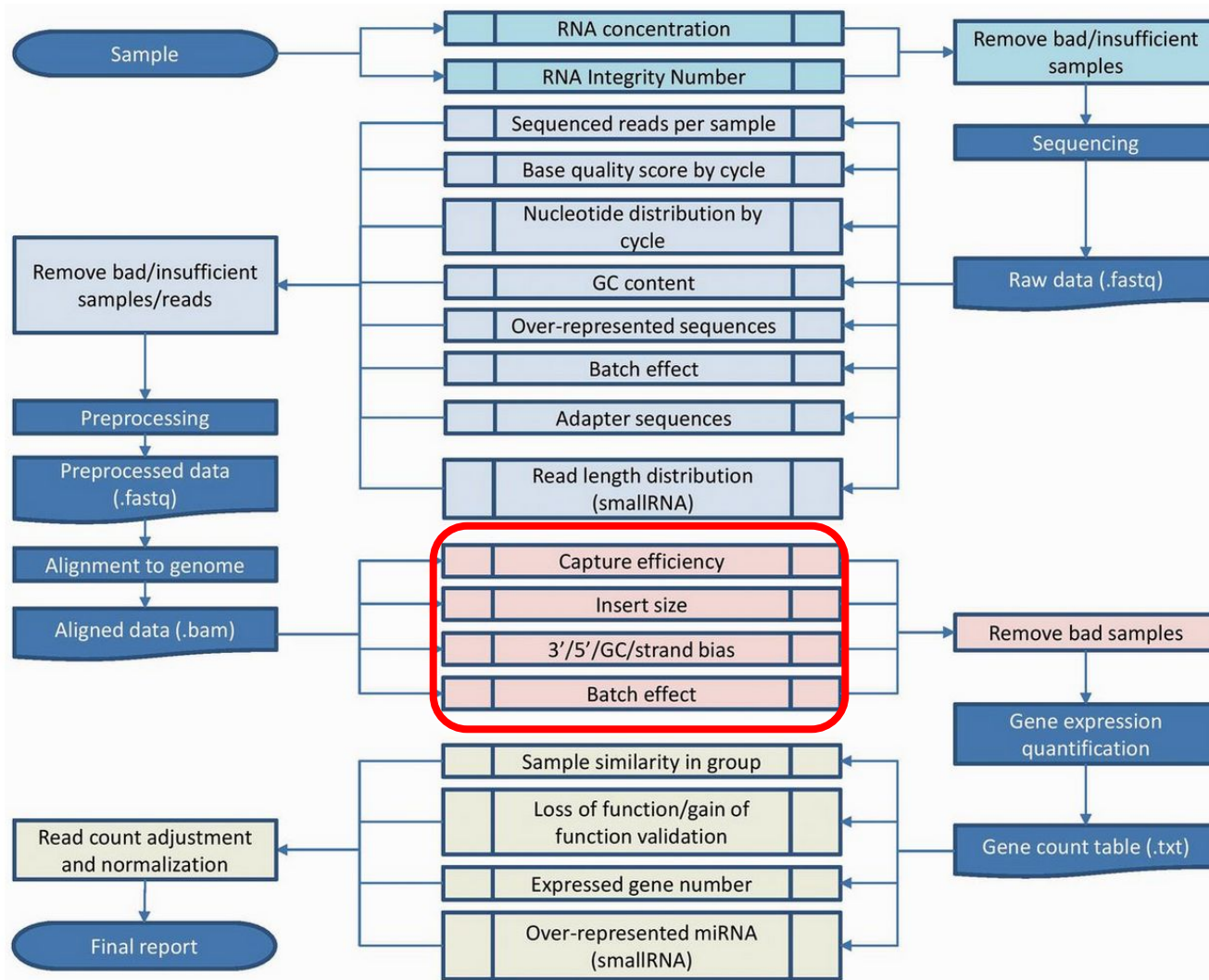
Quality control and mapping



Quality control and mapping



Quality control and mapping



Quality control and mapping

Number of input reads
Average input read length

UNIQUE READS:

Uniquely mapped reads number
Uniquely mapped reads %
Average mapped length
Number of splices: Total
Number of splices: Annotated (sjdb)
Number of splices: GT/AG
Number of splices: GC/AG
Number of splices: AT/AC
Number of splices: Non-canonical
Mismatch rate per base, %
Deletion rate per base
Deletion average length
Insertion rate per base
Insertion average length

MULTI-MAPPING READS:

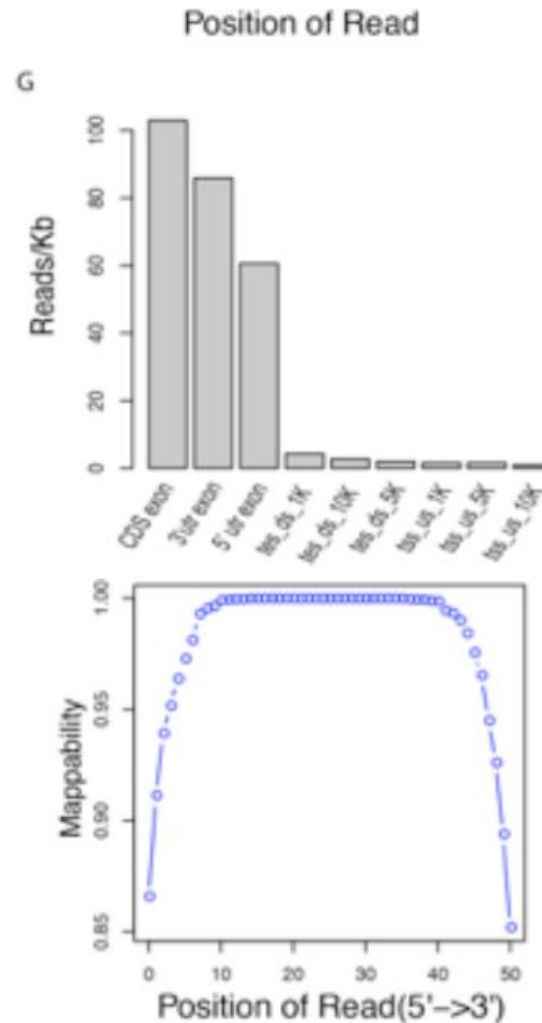
Number of reads mapped to multiple loci
% of reads mapped to multiple loci
Number of reads mapped to too many loci
% of reads mapped to too many loci

UNMAPPED READS:

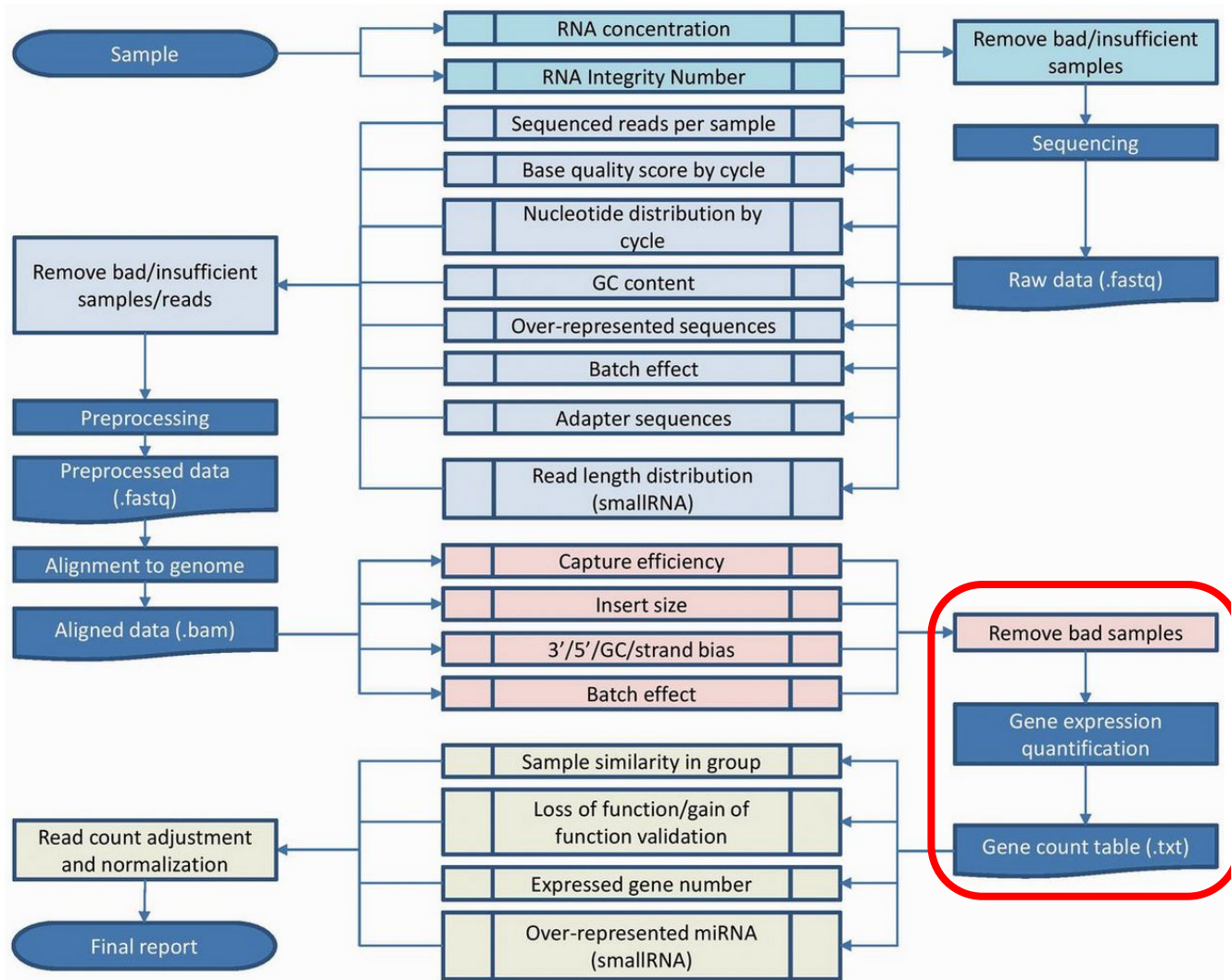
% of reads unmapped: too many mismatches
% of reads unmapped: too short
% of reads unmapped: other

CHIMERIC READS:

Number of chimeric reads
% of chimeric reads



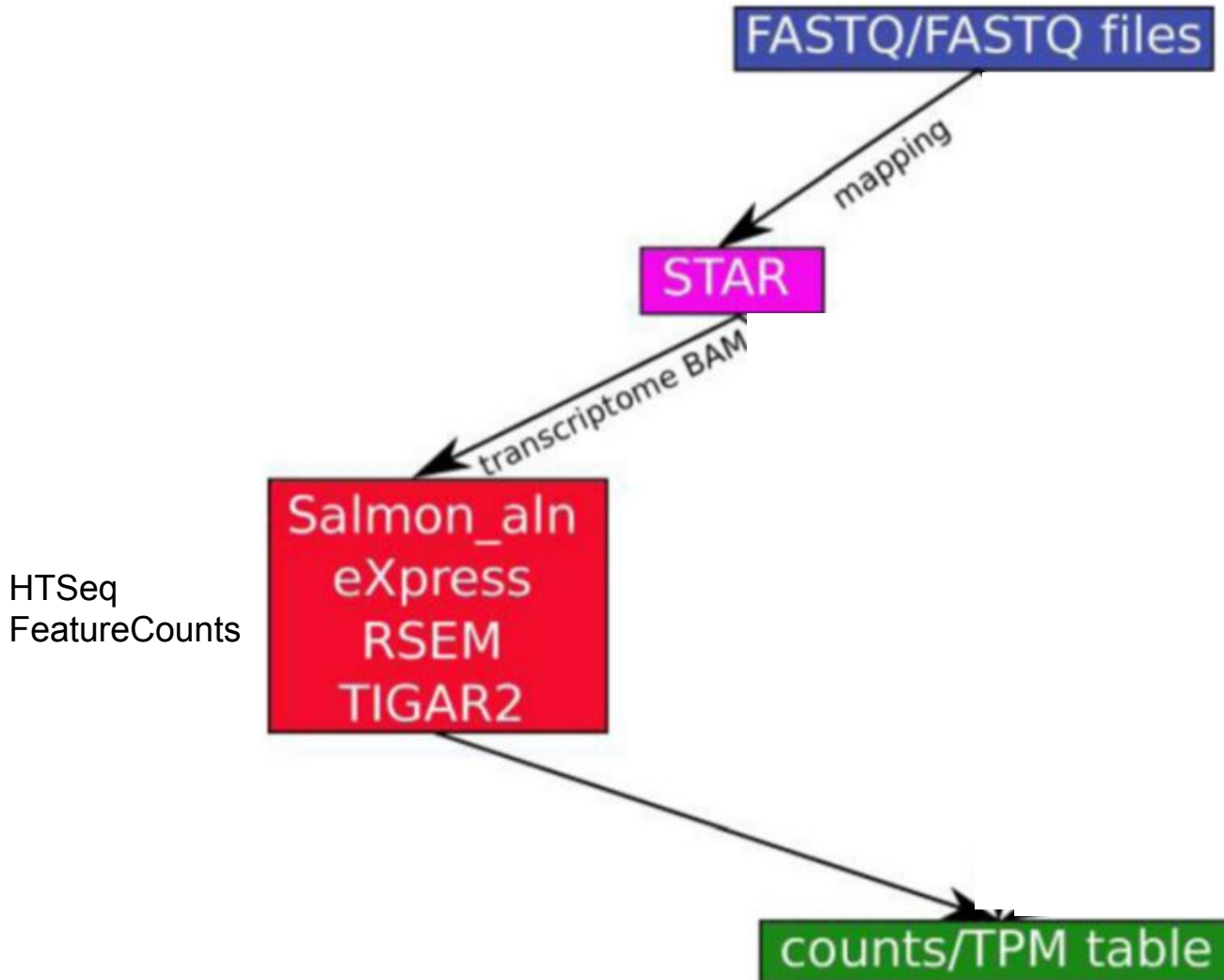
Quality control and mapping



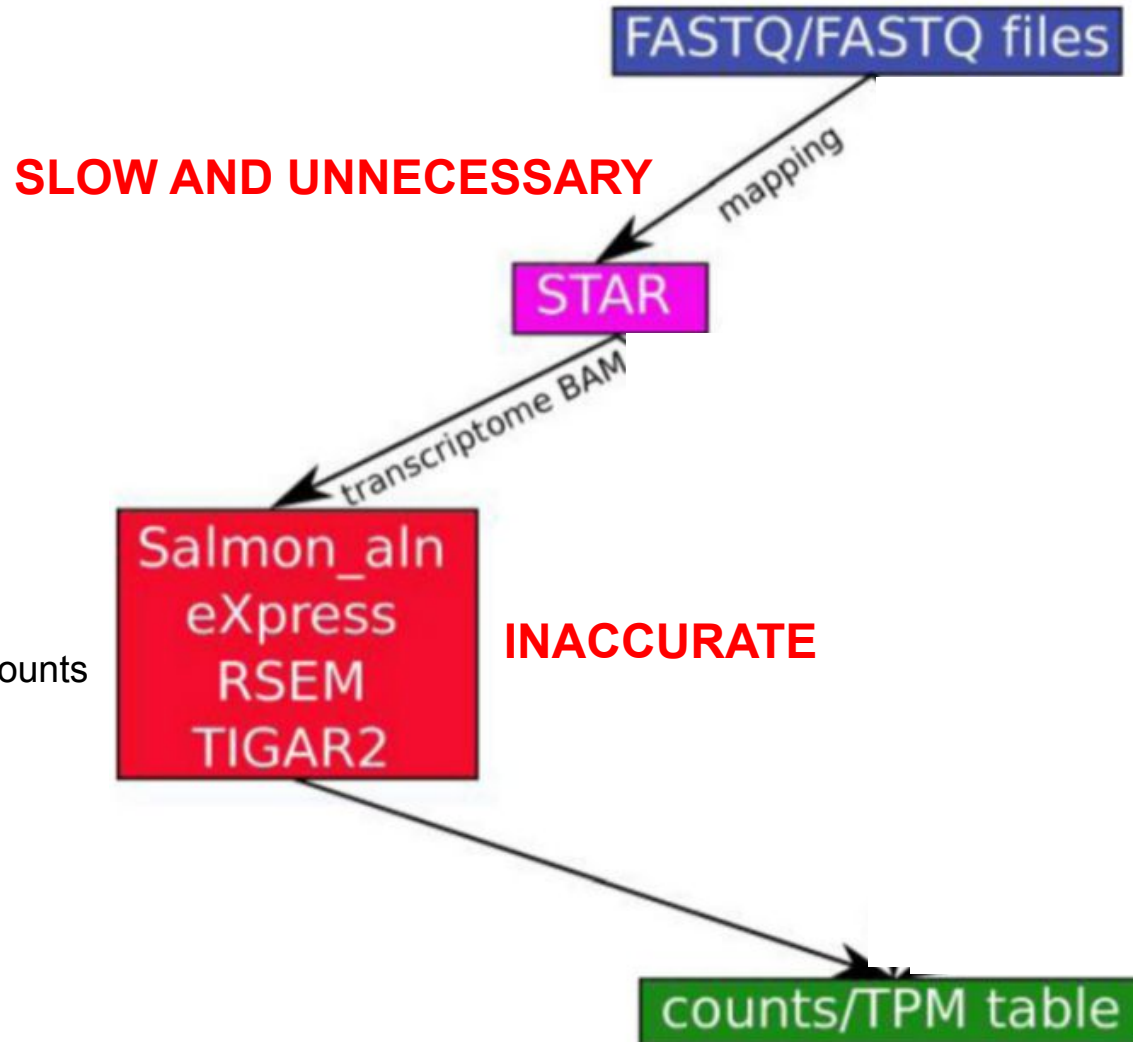
Outline

- Experimental and sequencing design
- Quality control and mapping
- **Genes/transcripts/exons quantification**
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Genes/Transcripts/Exons quantification

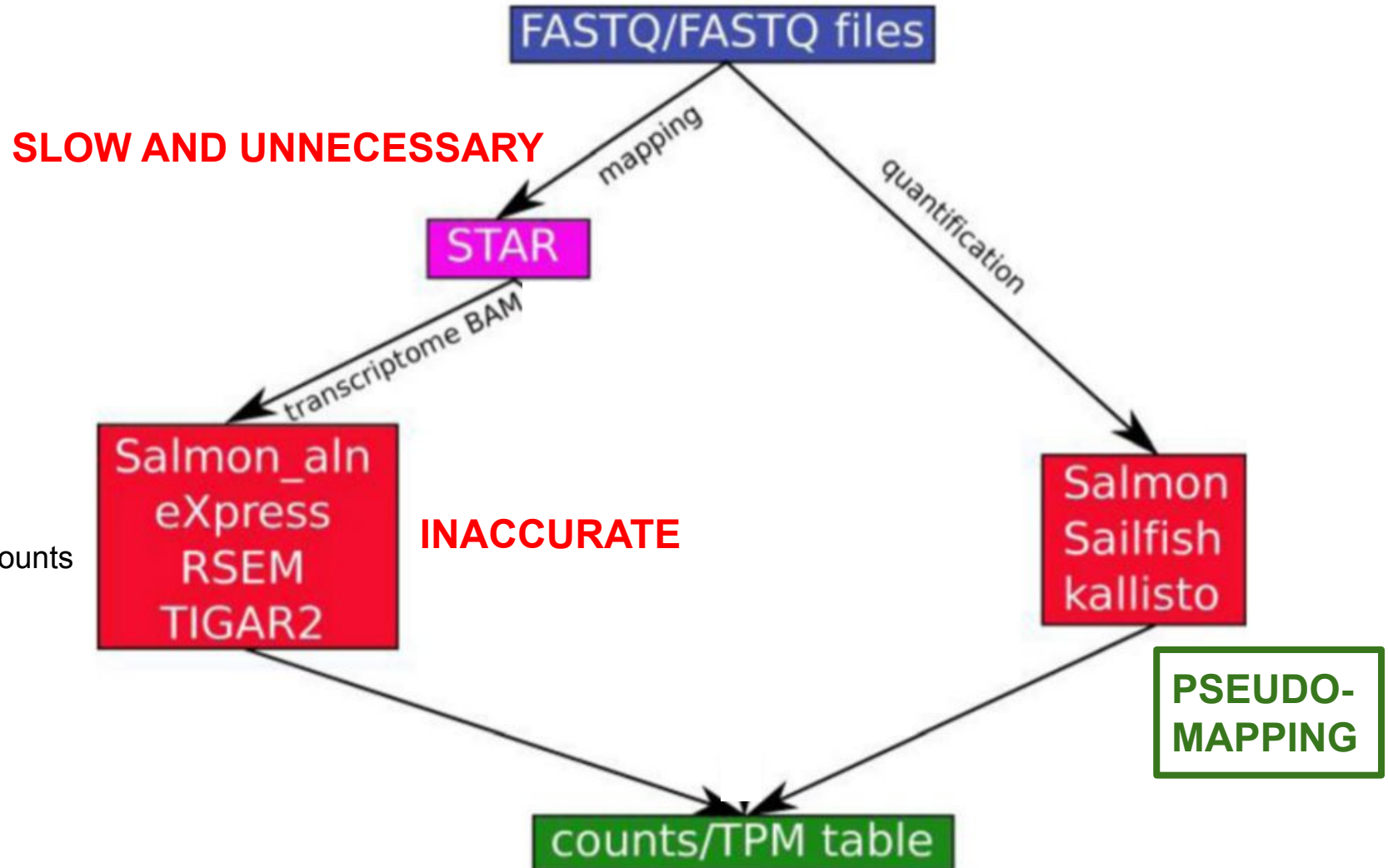


Genes/Transcripts/Exons quantification



HTSeq
FeatureCounts

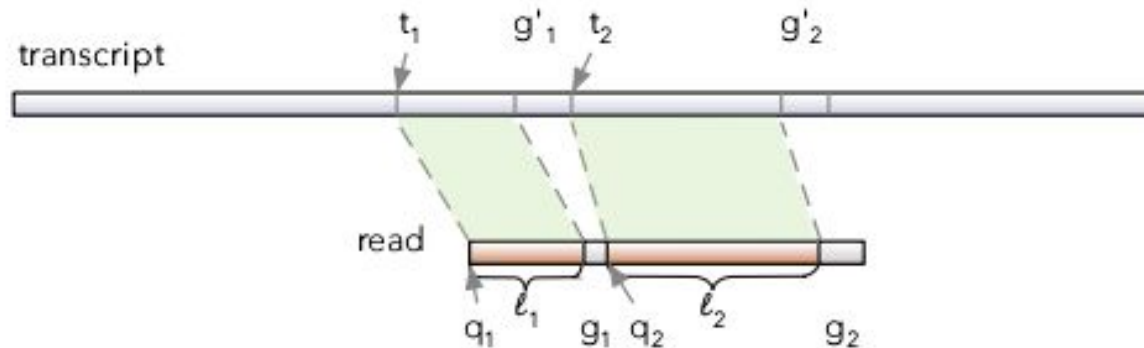
Genes/Transcripts/Exons quantification



Genes/Transcripts/Exons quantification

A simple concept :

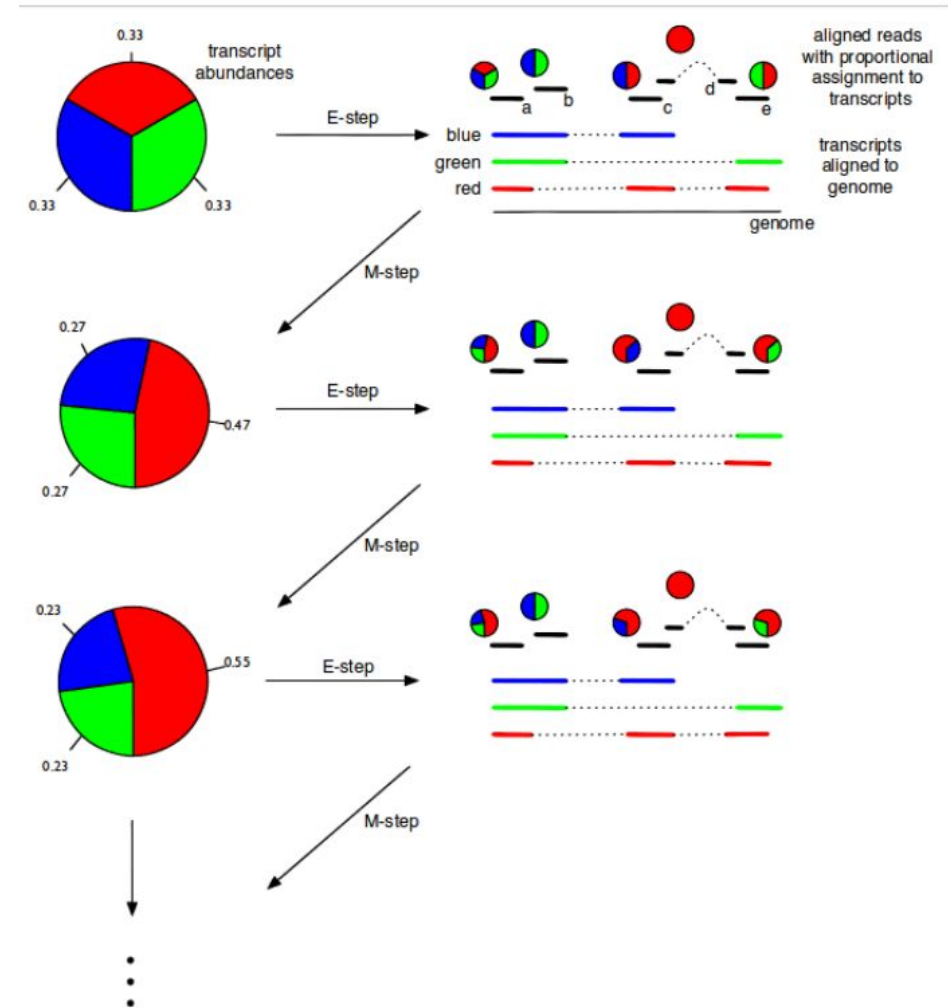
Accurate quantification of transcript abundance from RNA-seq data does not require knowing the optimal alignment for every potential locus of origin. Rather, simply knowing which transcripts (and positions within these transcripts) match the fragments reasonably well is sufficient.



Genes/Transcripts/Exons quantification

First 3 cycles of EM algorithm.
Abundance of **red** isoform
estimated after the 1st M-step:
($\frac{1}{3}$ read a + $\frac{1}{2}$ read c + 1 read
d + $\frac{1}{2}$ read e)/(total read
number), i.e. 0.47
((0.33+0.5+1+0.5)/5)

- proved to converge
- stop criterion: when all probabilities that a fragment is derived from a transcript $\geq 10^{-7}$ have a relative change \leq than 10^{-3}



Genes/Transcripts/Exons quantification

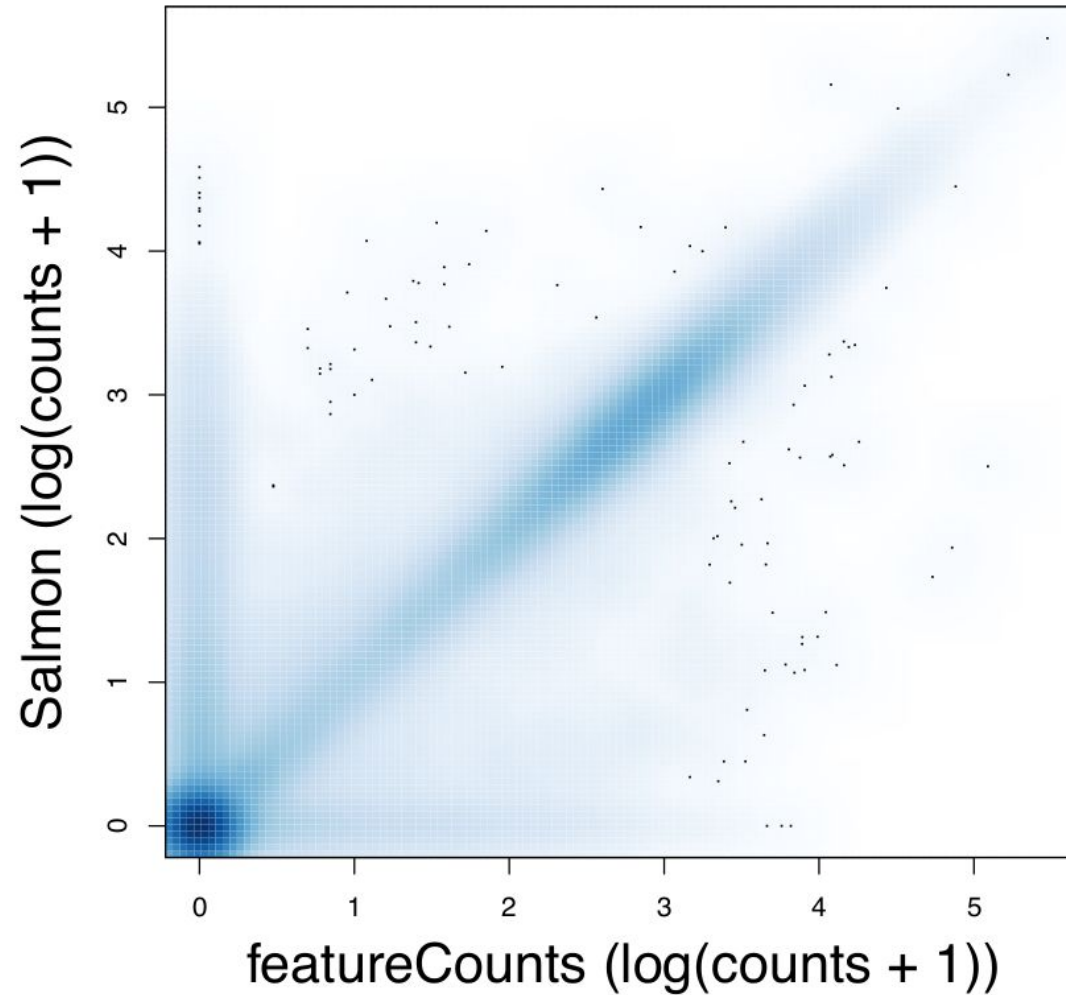
Name	Length	Effective Length	TPM	NumReads
ENST00000424770.1	586	408.443	0	0
ENST00000448070.1	121	12.6016	0	0
ENST00000413156.1	578	400.459	24.2572	5
ENST00000420638.1	685	507.217	151.164	39.465
ENST00000398242.2	1049	869.633	973.655	435.825

OR

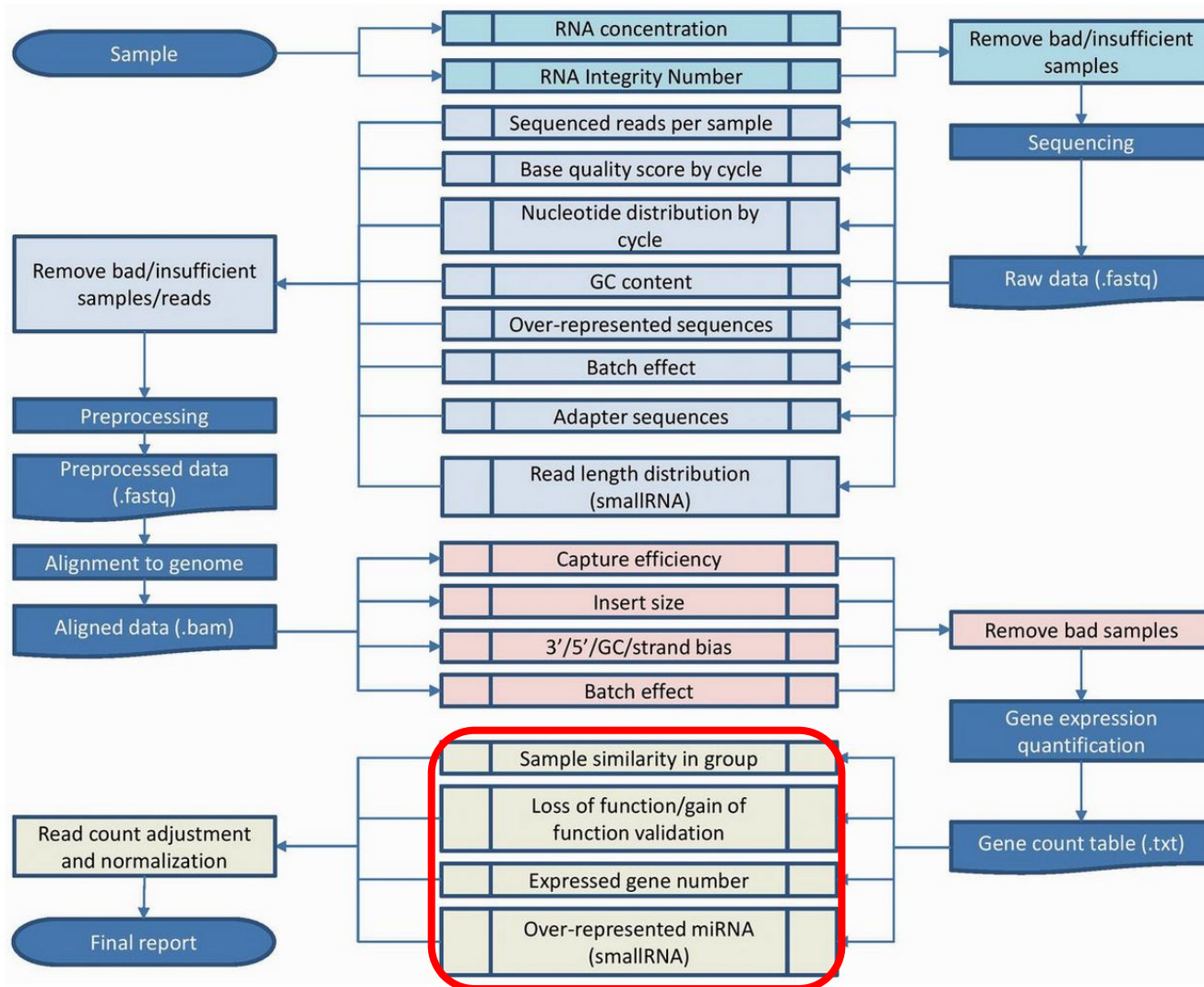
Name	Length	Effective Length	TPM	NumReads
ENSG00000079974.1	176.5	159.427	413.683	356.123
ENSG00000213683.3	308.999	237.082	1808.74	546
ENSG00000254499.1	1606	1432.33	0	0
ENSG00000225929.1	2172.5	1993.13	0	0
ENSG00000212569.1	98	14.2568	0	0

Genes/Transcripts/Exons quantification

SRR1039508



Quality control and mapping



Genes/Transcripts/Exons quantification

RPKM = Reads Per Kilobase per Million

- $(\text{Num_sample_total_reads} / 1,000,000) =$ “per million” scaling factor (SF)
- $(\text{Read_counts} / \text{SF}) =$ This normalizes for sequencing depth, giving you reads per million (RPM)
- $(\text{RPM} / \text{Gene_Length_in_kb}) =$ This gives you RPKM.

TPM = Transcripts per kilobase Per Million

- $(\text{Read_Counts} / \text{Gene_Length_in_kb}) =$ This gives you reads per kilobase (RPK).

$(\text{Sum_all_sample_RPK} / 1,000,000) =$ This is your “per million” scaling factor (SF).

3- $(\text{RPK} / \text{SF}) =$ This gives you TPM.

Genes/Transcripts/Exons quantification

RPKM = Relative Expression Level **WITHIN** Sample

TPM = Normalized Expression Level comparable **BETWEEN**
Samples

Genes/Transcripts/Exons quantification

**EXON-LEVEL
RAW COUNTS**



**DIFFERENTIAL
ANALYSIS
(DEXseq)**

**GENE-LEVEL
RAW COUNTS**



**DIFFERENTIAL
ANALYSIS
(DESeq2)**

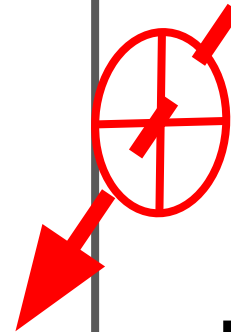
**TRANSCRIPT-LEVEL
RAW COUNTS**



**OPTIONAL :
GENE-LEVEL
AGGREGATION**

**DIFFERENTIAL
ANALYSIS
(Sleuth)**

**FUNCTIONAL ANALYSIS
(clusterProfiler + GSEA + IPA + ROMA)**



Genes/Transcripts/Exons quantification

TABLE 1. RNA-seq differential gene expression tools and statistical tests

Name	Assumed distribution	Normalization	Description	Version	Citations ^d	Reference
<i>t</i> -test	Normal	DEseq ^a	Two-sample <i>t</i> -test for equal variances	–	–	–
log <i>t</i> -test	Log-normal	DEseq ^a	Log-ratio <i>t</i> -test	–	–	–
Mann-Whitney	None	DEseq ^a	Mann-Whitney test	–	–	Mann and Whitney (1947)
Permutation	None	DEseq ^a	Permutation test	–	–	Efron and Tibshirani (1993a)
<i>Bootstrap</i>	Normal	DEseq ^a	Bootstrap test	–	–	Efron and Tibshirani (1993a)
<i>baySeq</i> ^c	Negative binomial	Internal	Empirical Bayesian estimate of posterior likelihood	2.2.0	159	Hardcastle and Kelly (2010)
<i>Cuffdiff</i>	Negative binomial	Internal	Unknown	2.1.1	918	Trapnell et al. (2012)
<i>DEGseq</i> ^c	Binomial	None	Random sampling model using Fisher's exact test and the likelihood ratio test	1.22.0	325	Wang et al. (2010)
<i>DESeq</i> ^c	Negative binomial	DEseq ^a	Shrinkage variance	1.20.0	1889	Anders and Huber (2010)
<i>DESeq2</i> ^c	Negative binomial	DEseq ^a	Shrinkage variance with variance based and Cook's distance pre-filtering	1.8.2	197	Love et al. (2014)
<i>EBSeq</i> ^c	Negative binomial	DEseq ^a (median)	Empirical Bayesian estimate of posterior likelihood	1.8.0	80	Leng et al. (2013)
<i>edgeR</i> ^c	Negative binomial	TMM ^b	Empirical Bayes estimation and either an exact test analogous to Fisher's exact test but adapted to over-dispersed data or a generalized linear model	3.10.5	1483	Robinson et al. (2010)
<i>Limma</i> ^c	Log-normal	TMM ^b	Generalized linear model	3.24.15	97	Law et al. (2014)
<i>NOISeq</i> ^c	None	RPKM	Nonparametric test based on signal-to-noise ratio	2.14.0	177	Tarazona et al. (2011)
<i>PoissonSeq</i> ^c	Poisson log-linear model	Internal	Score statistic	1.1.2	37	Li et al. (2012)
<i>SAMSeq</i> ^c	None	Internal	Mann-Whitney test with Poisson resampling	2.0	54	Li and Tibshirani (2013)

DESeq2 vs edgeR

History

When we started working on DESeq2 in Fall 2012, one of the main differences we were focusing on was a methodological detail* of the dispersion shrinkage steps. Aside from this difference, we wanted to update DESeq to use the GLM framework and to shrink dispersion estimates toward a central value as in edgeR, as opposed to the maximum rule that was previously implemented in DESeq (which tends to overestimate dispersion).

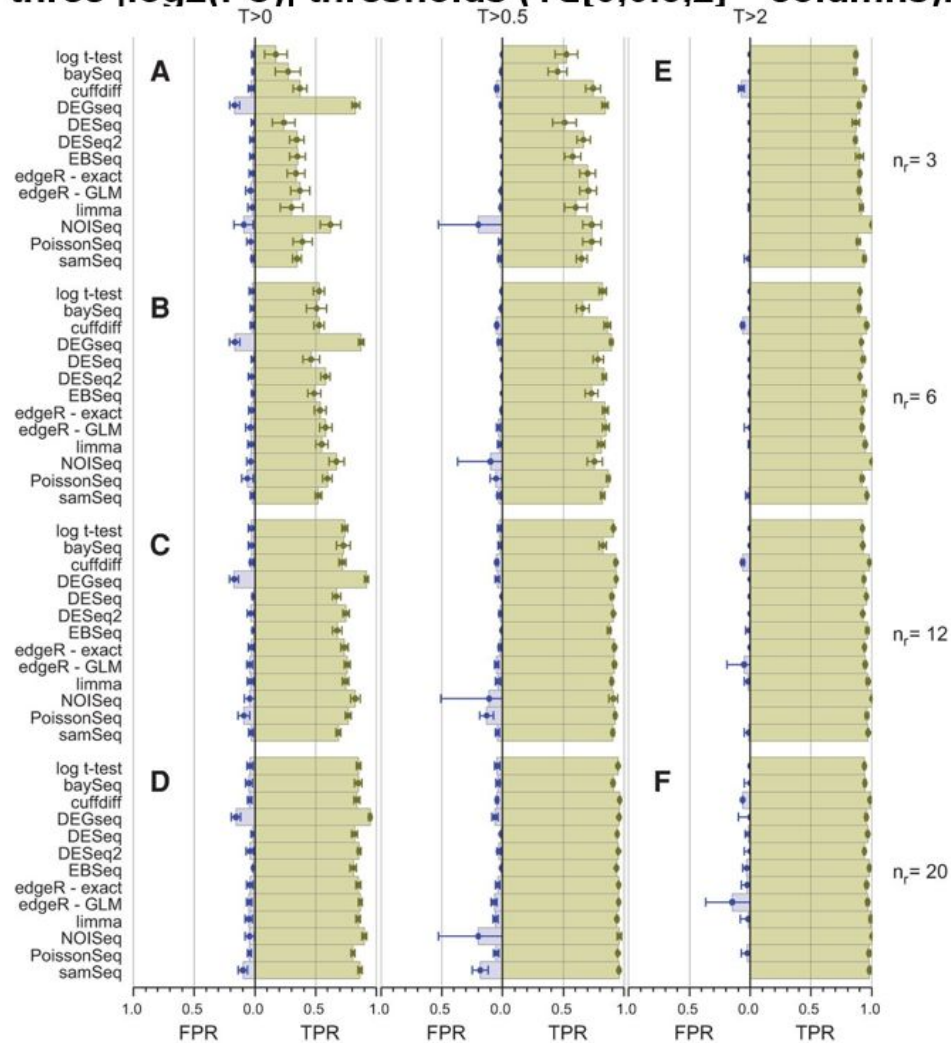
I would say that the difference in dispersion shrinkage didn't make a huge difference in performance compared to edgeR, as can be seen in real and simulated** data benchmarks in our DESeq2 paper published in 2014. From what I've seen in my own testing and numerous third-party benchmarks, the two methods typically report overlapping sets of genes, and have similar performance for gene-level DE testing.

What's different then?

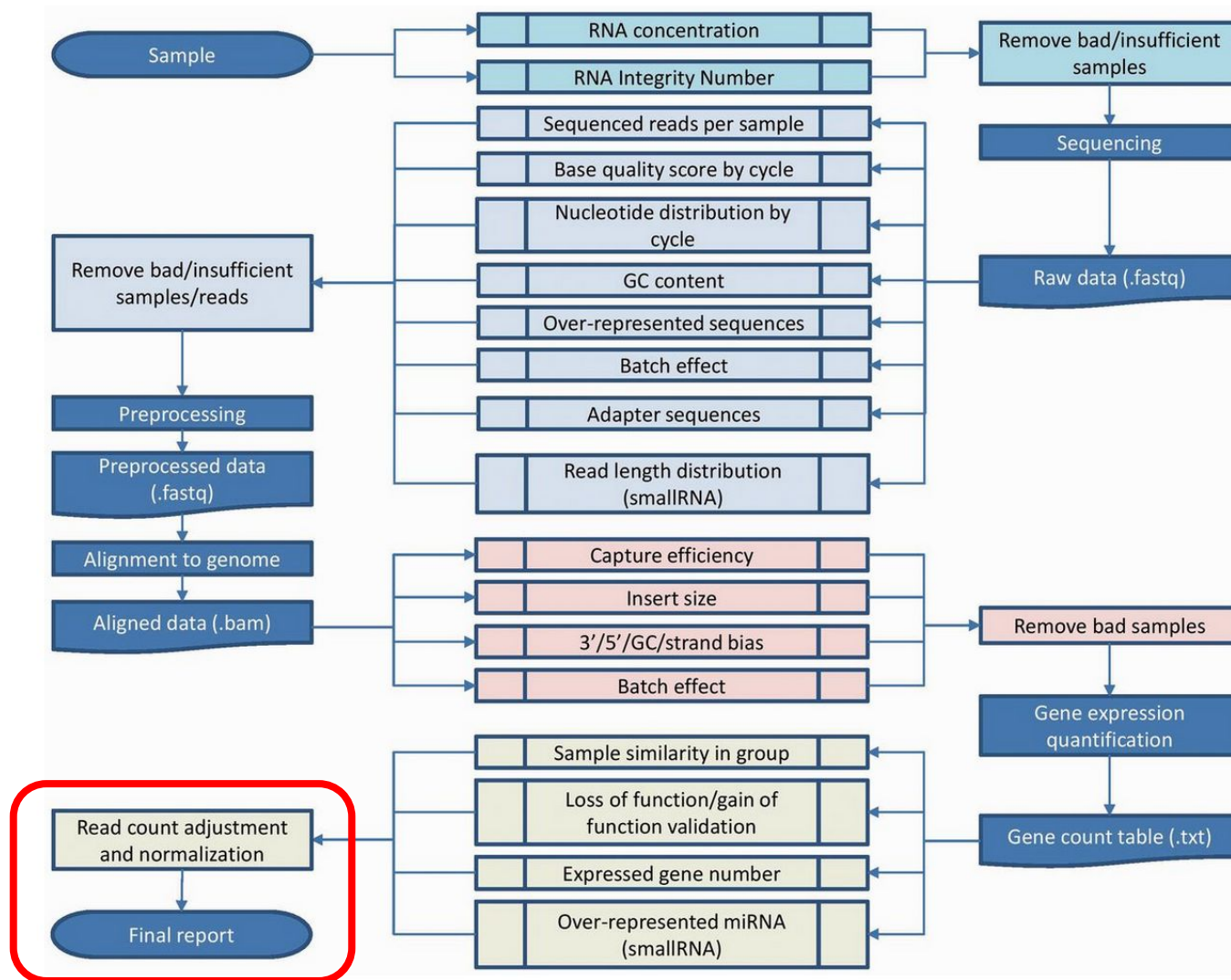
The major differences between the two methods are in some of the defaults. DESeq2 by default does a couple things (which can all optionally be turned off): it finds an optimal value at which to filter low count genes, flags genes with large outlier counts or removes these outlier values when there are sufficient samples per group ($n > 6$), excludes from the estimation of the dispersion prior and dispersion moderation those genes with very high within-group variance, and moderates log fold changes which have small statistical support (e.g. from low count genes). edgeR offers similar functionality, for example, it offers a robust dispersion estimation function, `estimateGLMRobustDisp`, which reduces the effect of individual outlier counts, and a robust argument to `estimateDisp` so that hyperparameters are not overly affected by genes with very high within-group variance. And the default steps in the edgeR User Guide for filtering low counts genes both increases power by reducing multiple testing burden and removes genes with uninformative log fold changes.

DESeq2 vs edgeR

Comparison of the true positive rate (TPR) and false positive rate (FPR) performance for each of the DGE tools on low-, medium-, and highly replicated RNA-seq data ($n_r \in \{3, 6, 12, 20\}$ —rows) for three $|\log_2(\text{FC})|$ thresholds ($T \in \{0, 0.5, 2\}$ —columns).



Quality control and mapping



Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- **Functional Analysis**
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

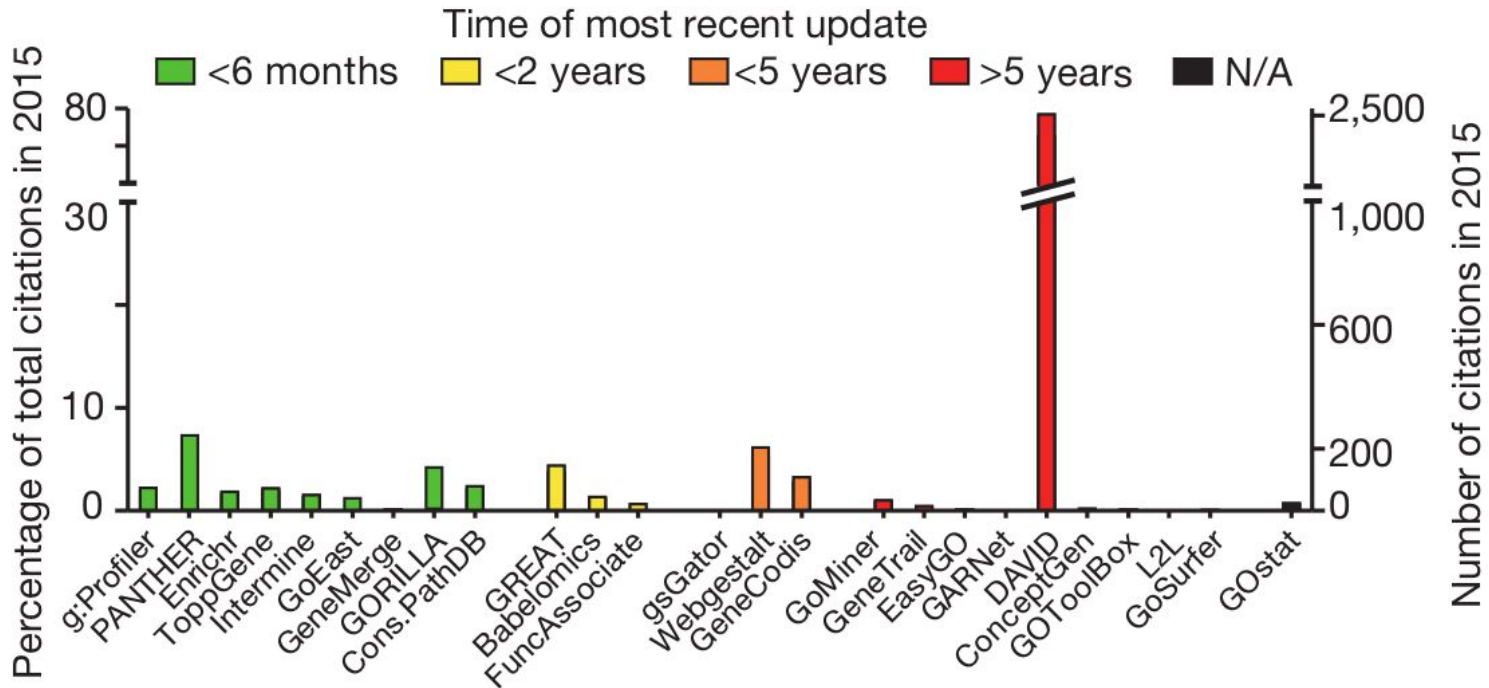
Functional analysis

There are a number of questions that you may answer using your set of DEGs. These include:

- 1) What are the biological process, cellular locations and molecular functions that are particularly over- or under-represented in your set of genes ?**
- 2) What are the pathways that are significantly impacted in your condition ?
- 3) Which are the putative regulatory elements in the promoters of genes that show similar expression patterns (i.e. those which have similar fold-change patterns between different samples) ?

Functional analysis

GO Analysis : There are many tools that allow you to do this...



...but there are also many, many mistakes related to the way the p-values are calculated and corrected for multiple comparisons, the association of genes to multiple GO terms and the implicit redundancy, the choice of the background set of genes, etc... (Rhee et al. 2008. Nat Rev Genet)

Functional analysis

We are in 2021...

DAVID 6.8 Oct. 2016

- The DAVID Knowledgebase completely rebuilt
- Entrez Gene integrated as the central identifier to allow for more timely updates while still incorporating Ensembl and Uniprot as integral data sources
- New GO category (GO Direct) provides GO mappings directly annotated by the source database (no parent terms included)
- New annotation categories
- New list identifier systems added for list uploading and conversion
- A few bugs fixed

DAVID 6.7 Jan. 2010

- The DAVID Knowledgebase completely rebuilt, including the central DAVID id system
- Ensembl Gene included as an integral data source
- DAVID engine completely rebuilt to facilitate future updates and development
- New GO category (GO FAT) filters out very broad GO terms based on a measured specificity of each term (not level-specificity)
- New annotation categories
- New list identifier systems added for list uploading and conversion
- Automatic list naming based on uploaded file name
- Ability to upload expression/other values (some display, but otherwise not used in the analysis at this point)
- A few bugs fixed
- and [more](#)

<https://david.ncifcrf.gov/content.jsp?file=release.html>

Functional analysis

ShinyGO: a graphical gene-set enrichment tool for animals and plants

Steven Xijin Ge ✉, Dongmin Jung, Runan Yao

Bioinformatics, btz931, <https://doi.org/10.1093/bioinformatics/btz931>

Published: 27 December 2019 **Article history** ▼

<http://bioinformatics.sdstate.edu/go/>

Just paste your gene list to get enriched GO terms and other pathways for over 315 plant and animal species, based on annotation from Ensembl (Release 96), Ensembl plants (R. 43) and Ensembl Metazoa (R. 43). An additional 2031 genomes (including bacteria and fungi) are annotated based on STRING-db (v.10). In addition, it also produces KEGG pathway diagrams with your genes highlighted, hierarchical clustering trees and networks summarizing overlapping terms/pathways, protein-protein interaction networks, gene characteristics plots, and enriched promoter motifs.

Functional analysis

There are a number of questions that you may answer using your set of DEGs. These include:

- 1) What are the biological process, cellular locations and molecular functions that are particularly over- or under-represented in your set of genes ?
- 2) What are the pathways that are significantly impacted in your condition ?**
- 3) Which are the putative regulatory elements in the promoters of genes that show similar expression patterns (i.e. those which have similar fold-change patterns between different samples) ?

Functional analysis

- Pathway databases (eg. Reactome, KEGG, etc.) : Here your set of DEG will be mapped on pathways. No real analysis is performed but you see what DEGs are on each pathway. The limitations are obvious: no p-values are calculated, no idea about which pathways are affected beyond random chance, etc.
- Pathway enrichment analysis : Here, pathways are considered as simple sets of genes and an enrichment p-value is calculated for each (e.g. DAVID, GSEA, Ingenuity, etc.). The limitations include the fact that the p-values are calculated based on the assumption that all variables (genes) are independent while the pathways are there precisely to tell you how these genes influence each other. Another limitation is that the pathways are treated as simple bags of genes, disregarding all the phenomena and interactions between genes that they describe. This analysis approach only looks at the number of DE genes and makes no difference between a situation in which a pathway has 3 entry points and all 3 are severely down-regulated thus effectively shutting down the entire pathway and a situation in which 3 other random genes are down regulated on the same pathway.

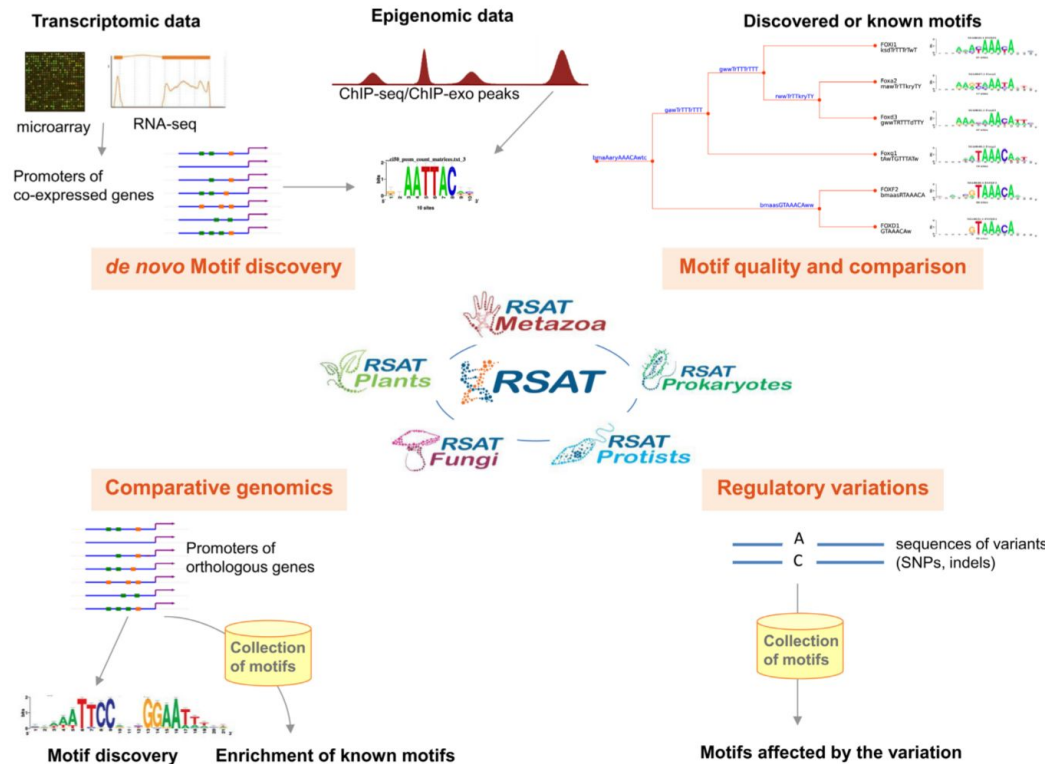
Functional analysis

There are a number of questions that you may answer using your set of DEGs. These include:

- 1) What are the biological process, cellular locations and molecular functions that are particularly over- or under-represented in your set of genes ?
- 2) What are the pathways that are significantly impacted in your condition ?
- 3) Which are the putative regulatory elements in the promoters of genes that show similar expression patterns (i.e. those which have similar fold-change patterns between different samples) ?**

Functional analysis

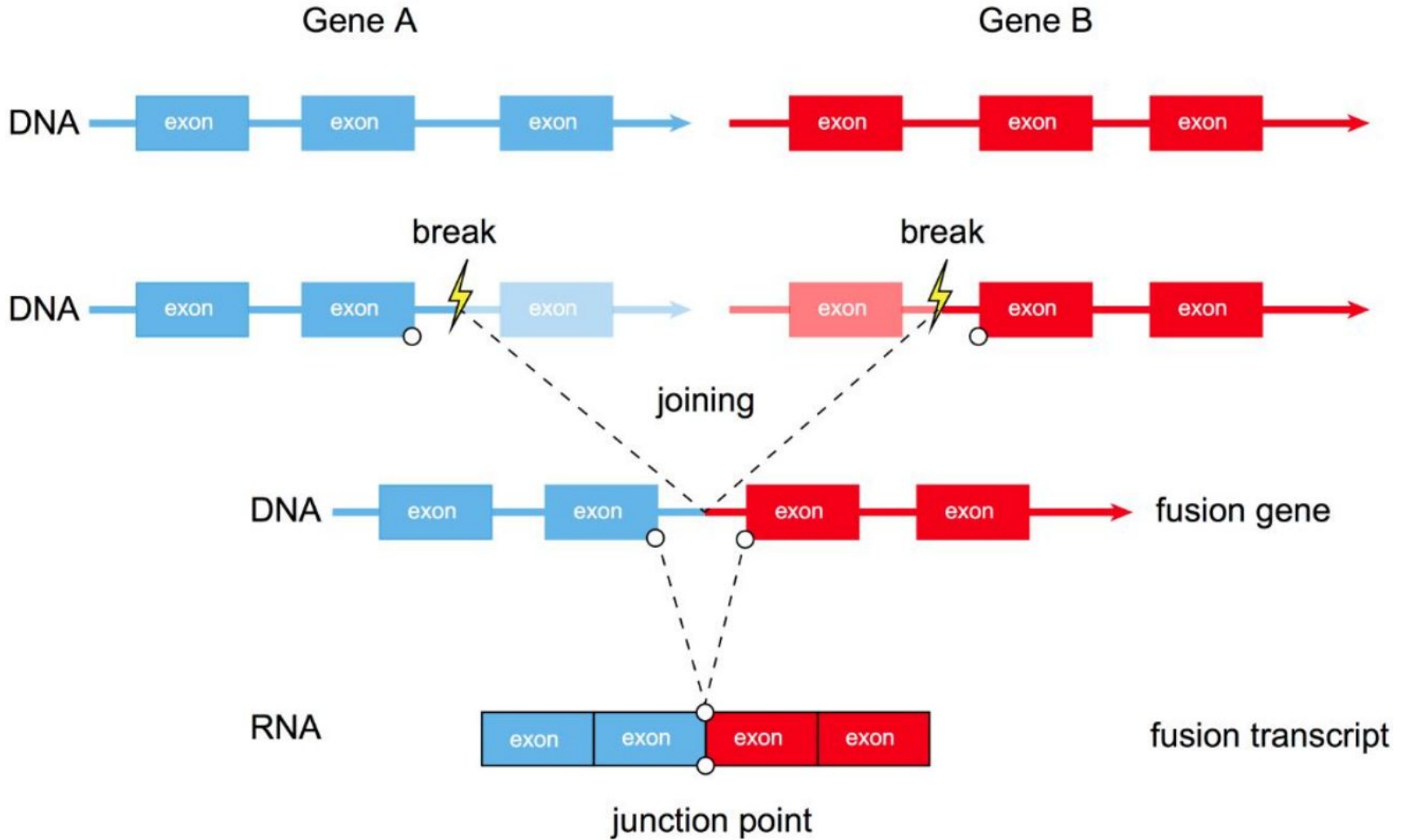
- You can download the putative promoters (usually 1kb upstream of the transcriptional start site is a good starting point) of all the genes you are interested in and search for enriched DNA motifs - you can use the online programs RSAT (<http://www.rsat.eu/>) or MEME (<http://meme.nbcr.net/meme/>). These DNA motifs may be recognised by sequence-specific transcription factors, thereby giving you insight into how this particular transcription profile may be regulated.



Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- **Fusion transcripts / chimera detection**
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

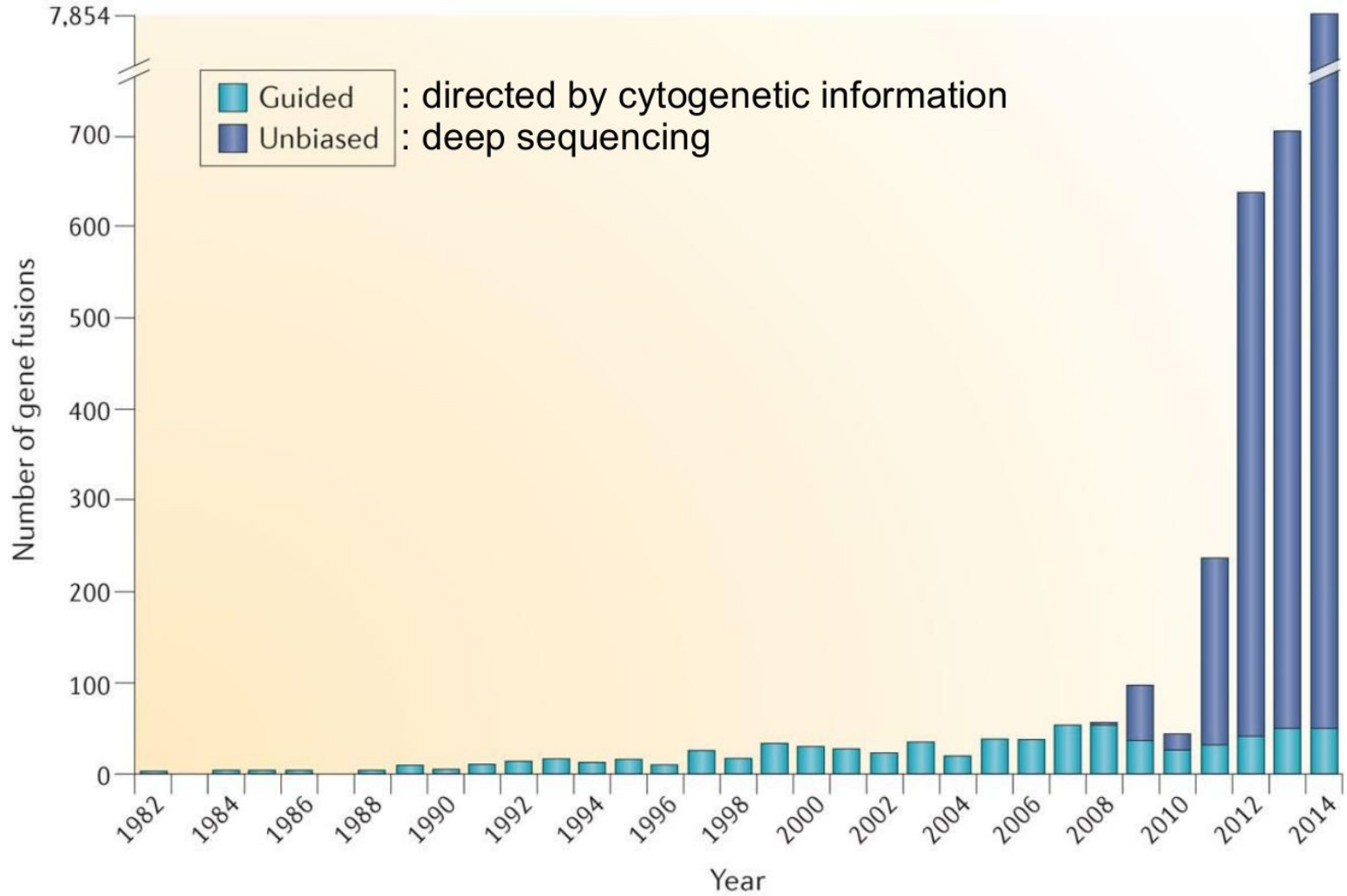
Fusion transcripts



Fusion transcripts

- An important class of cancer-contributing somatic alterations because they can drive the development of cancer
- Attractive as both therapeutic targets and diagnostic tools due to their tumor-specific expression
 - ex. BCR-ABL1 is associated with chronic myeloid leukemia and used as biomarker
- Other classes of chimeric transcripts :
 - Read-through
 - Trans-splicing

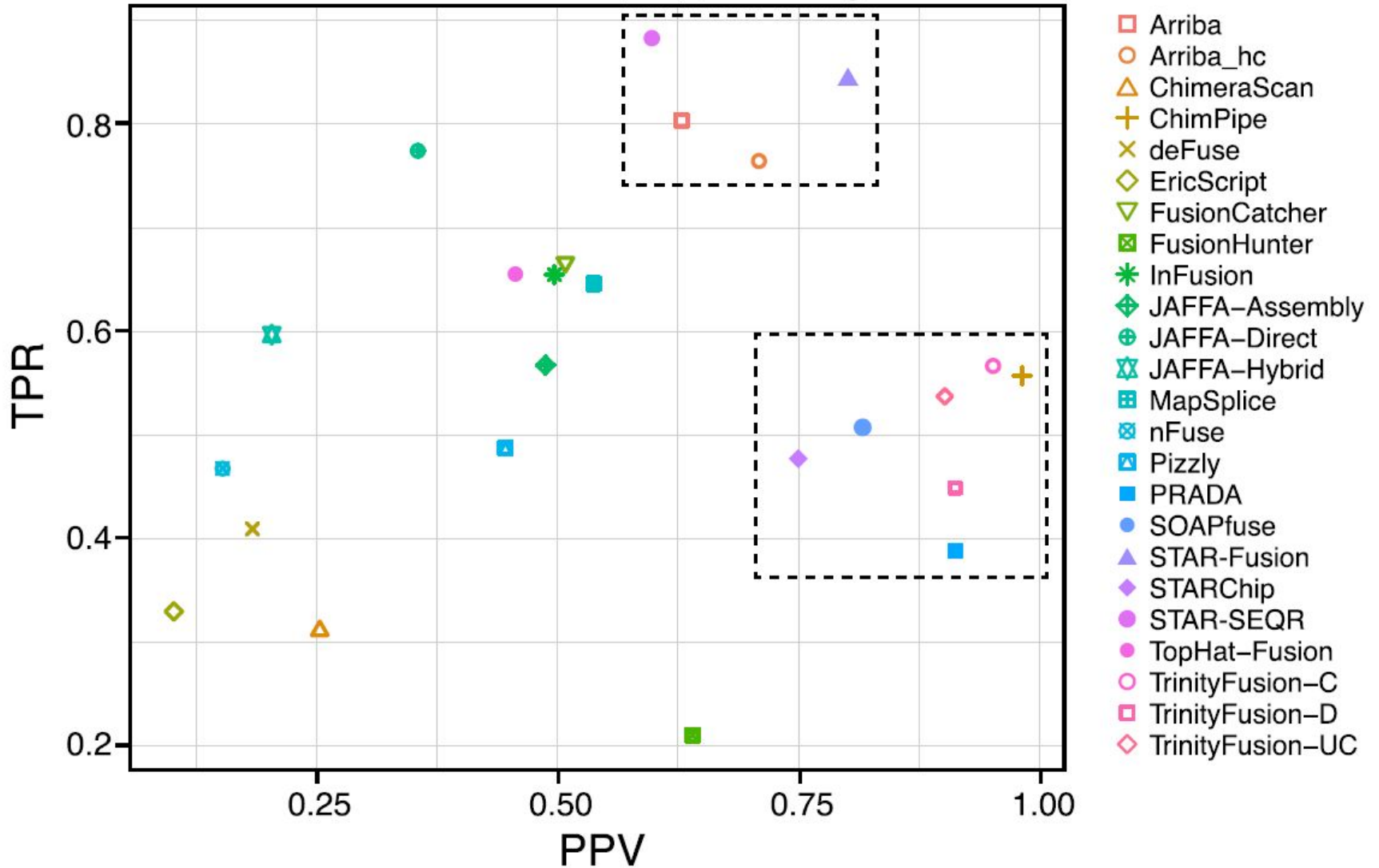
Fusion transcripts



Fusion transcripts

d

TPR versus PPV at max(F1), Min. Agree = 7



Fusion transcripts



Liu et al., 2016 demonstrated the potential of a meta-caller algorithm

\supporting reads\ Tool 1	Tool 2	Tool 3	Tool 4	
Fusion 1	23	100	20	22
Fusion 2	66	130	-	34
Fusion 3	17	-	-	-
Fusion 4	4	-	-	7
...				

Step 1: Keep the fusion transcripts detected by at least 2 tools

\supporting reads\ Tool 1	Tool 2	Tool 3	Tool 4	
Fusion 1	23	100	20	22
Fusion 2	66	130	-	34
Fusion 4	4	-	-	7
...				

Step 2: Rank within each tool (small to large), calculate the sum rank and order it from large to small

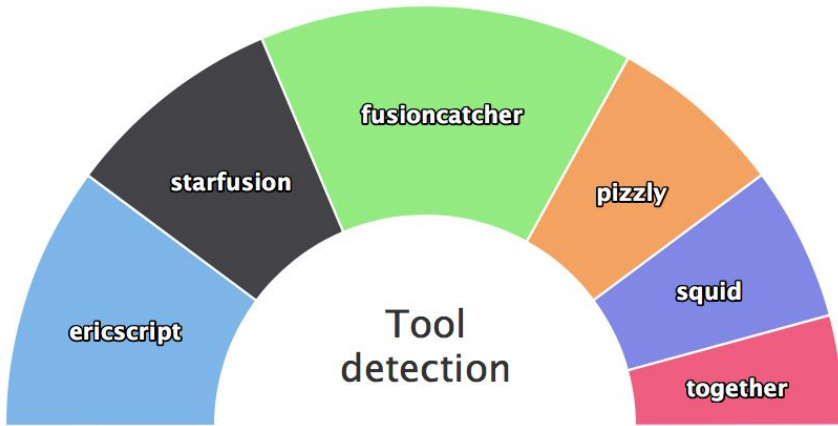
\rank\ Tool 1	Tool 2	Tool 3	Tool 4	Rank sum	Order of rank sum	
Fusion 1	2	1	1	2	6	2
Fusion 2	3	2	0	3	8	1
Fusion 4	1	0	0	1	2	3
...						

🔗 Pipeline overview (v1.2.0 - 2020/07/15)

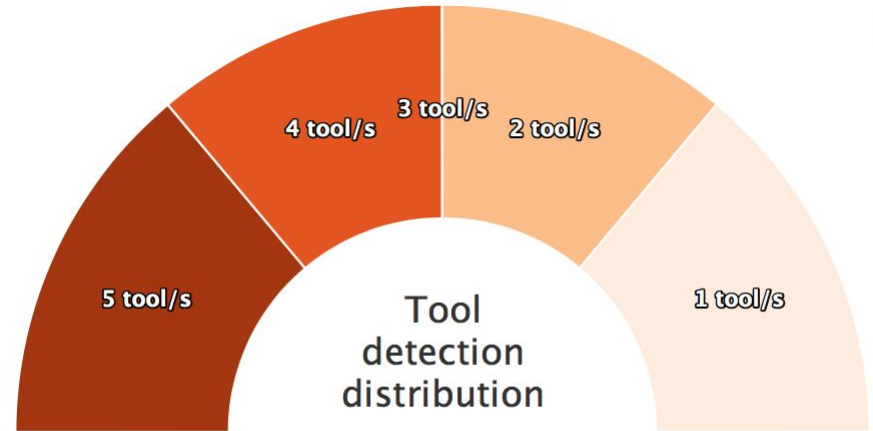
The pipeline is built using [Nextflow](#) and processes data using the following steps:

- [Arriba](#)
- [EricScript](#)
- [FastQC](#) - read quality control
- [FusionCatcher](#)
- [FusionInspector](#)
- [fusion-report](#)
- [MultiQC](#) - aggregate report, describing results of the whole pipeline
- [Pizzly](#)
- [Squid](#)
- [Star-Fusion](#)

Fusion transcripts



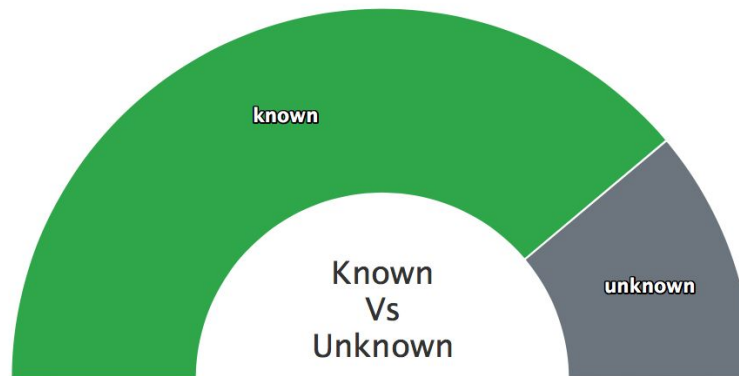
Displays number of found fusions per tool.



Sum of counts detected by different tools per fusion.

- ericscript
- starfusion
- fusioncatcher
- pizzly
- squid
- together

- 5 tool/s
- 4 tool/s
- 3 tool/s
- 2 tool/s
- 1 tool/s



Shows the ration between found and unknown missing fusions in the local database.

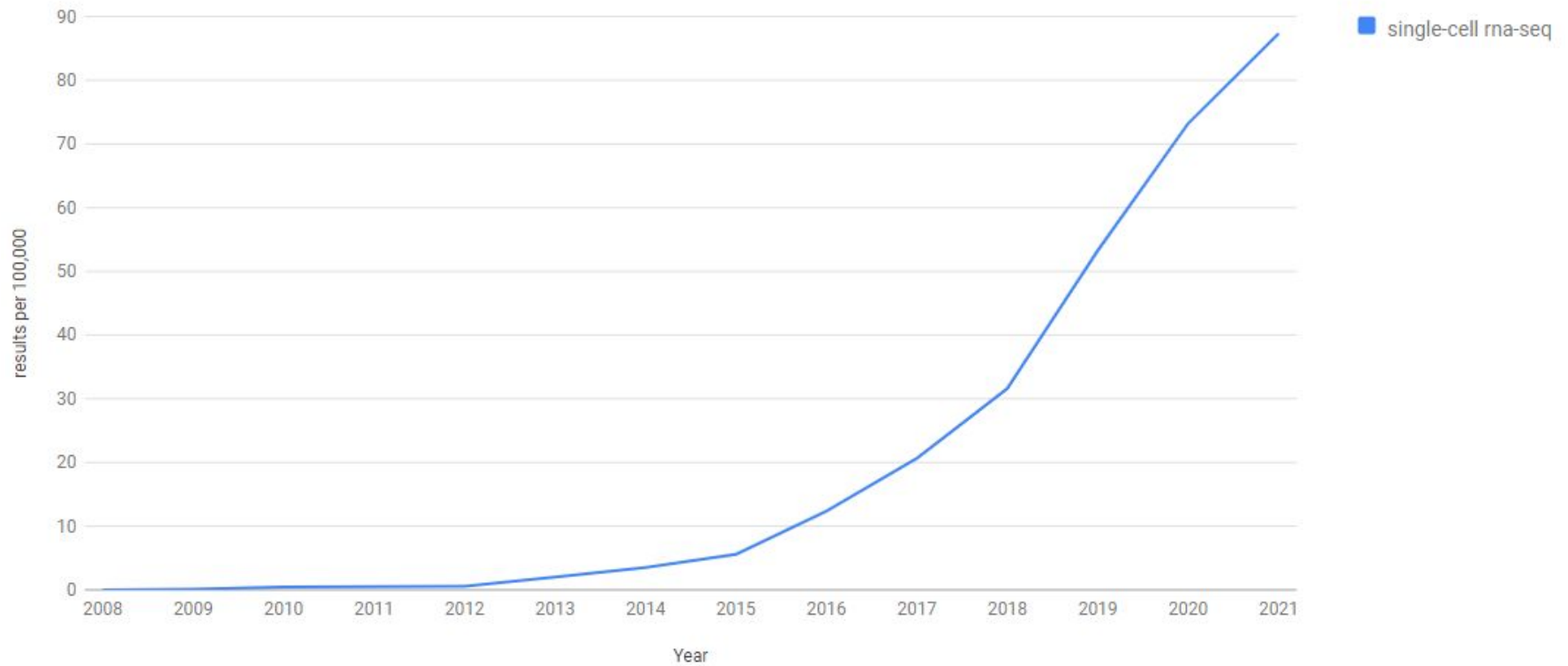
- known
- unknown

Outline

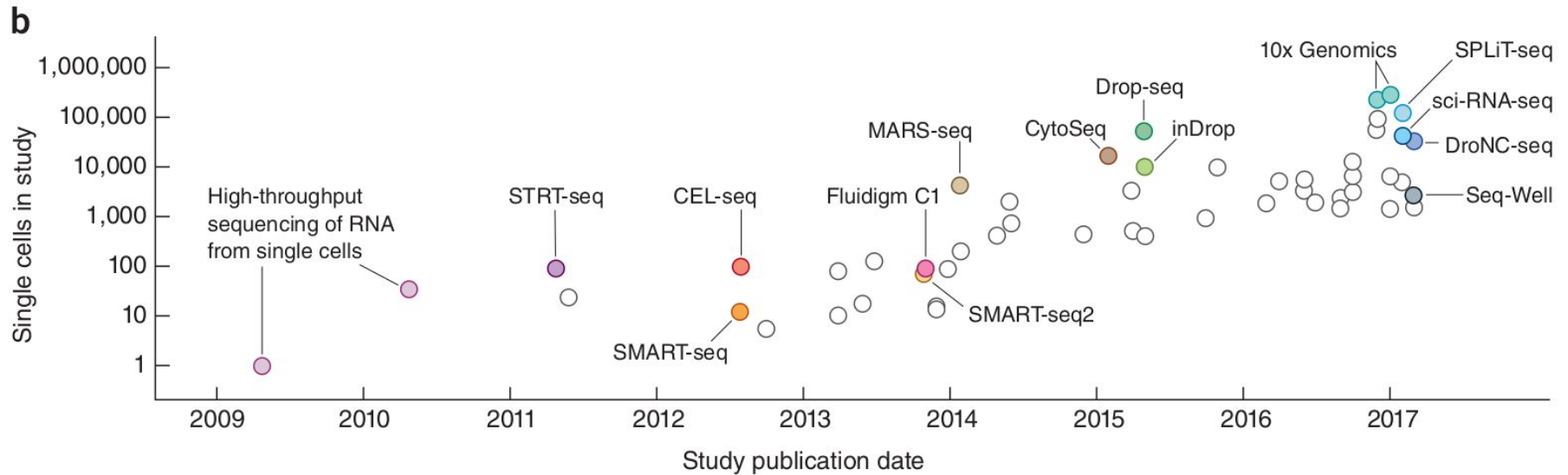
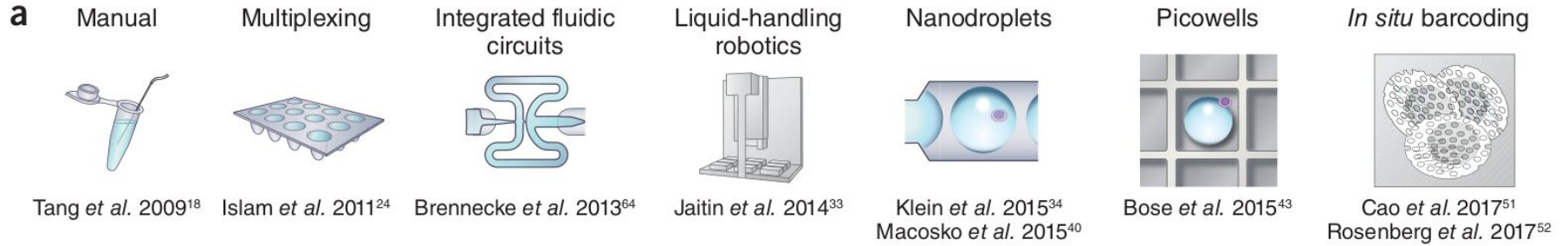
- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- **Single-cell RNA-seq**
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Single-cell RNA-seq

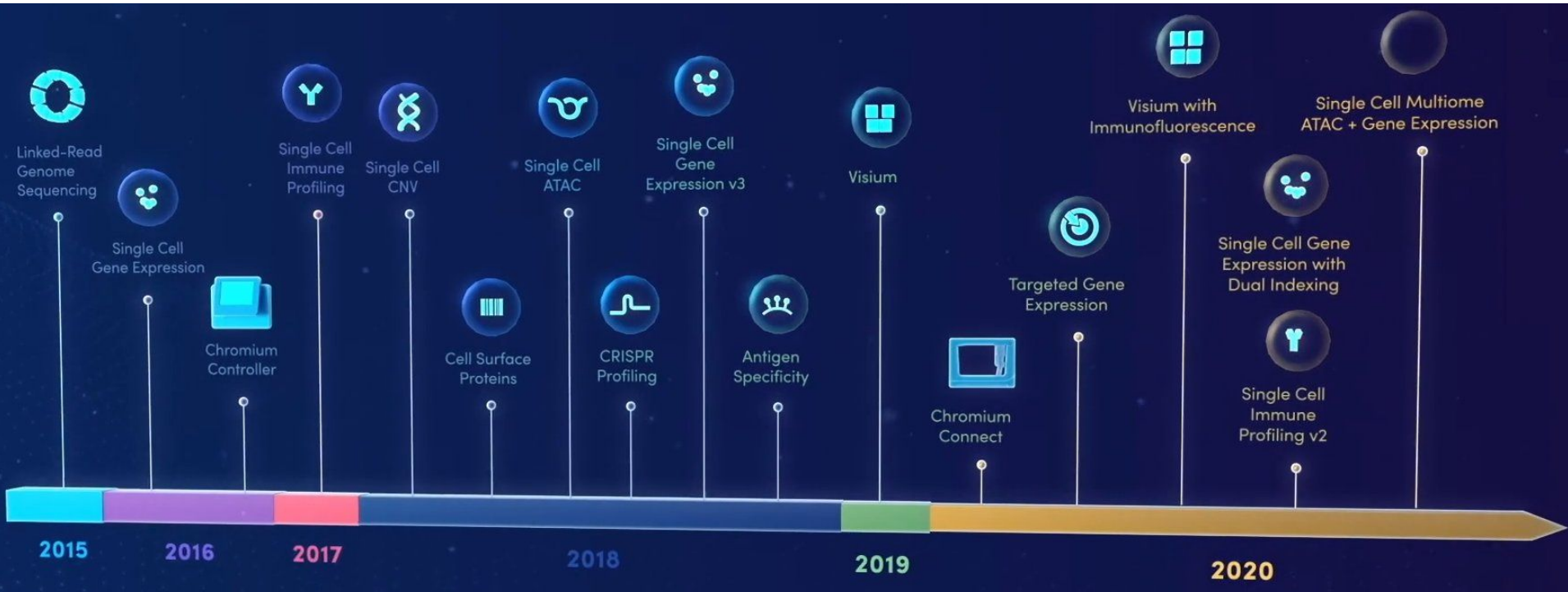
Results per 100,000 citations in PubMed
proportion for each search by year, 2008 to 2021



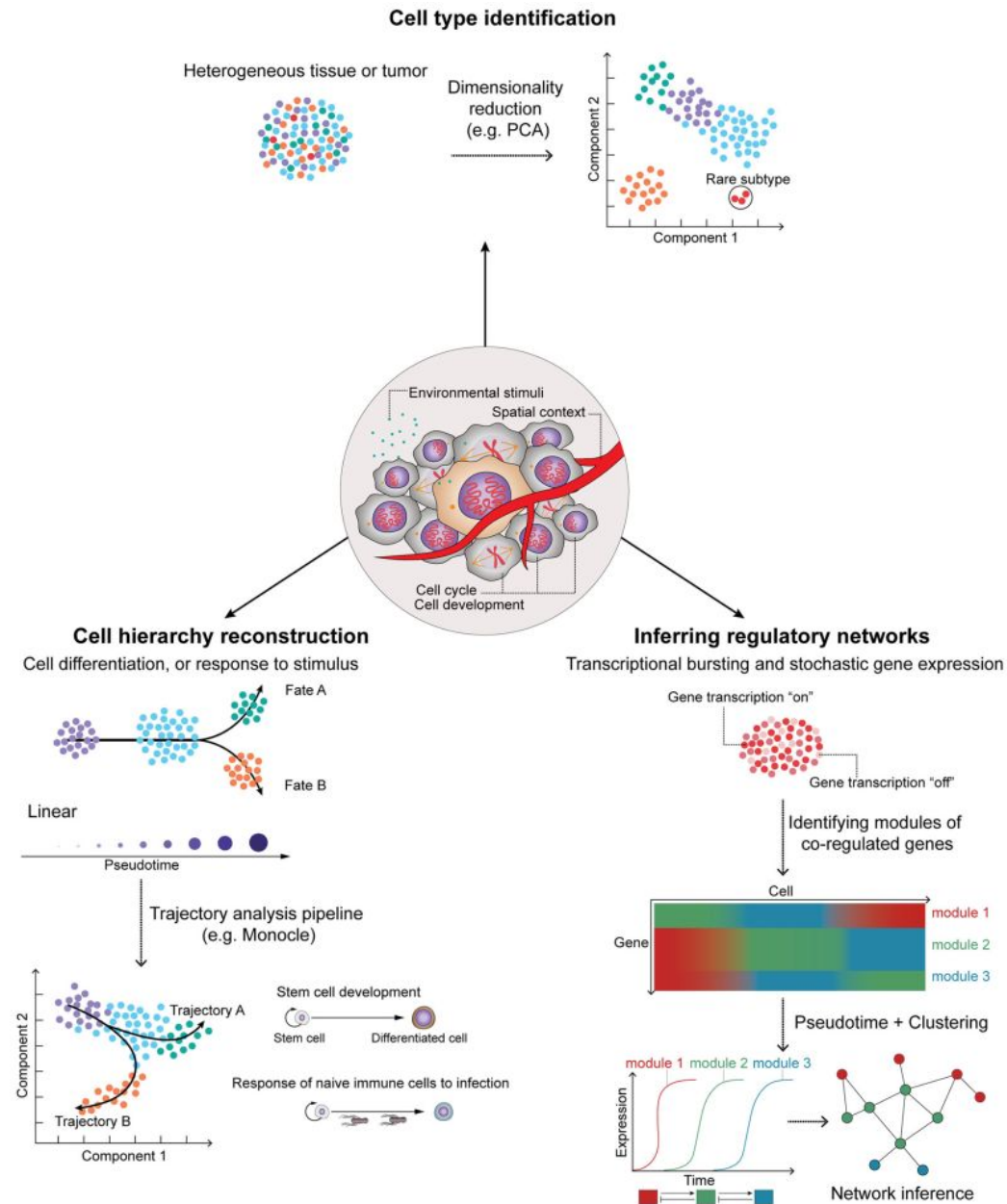
Single-cell RNA-seq



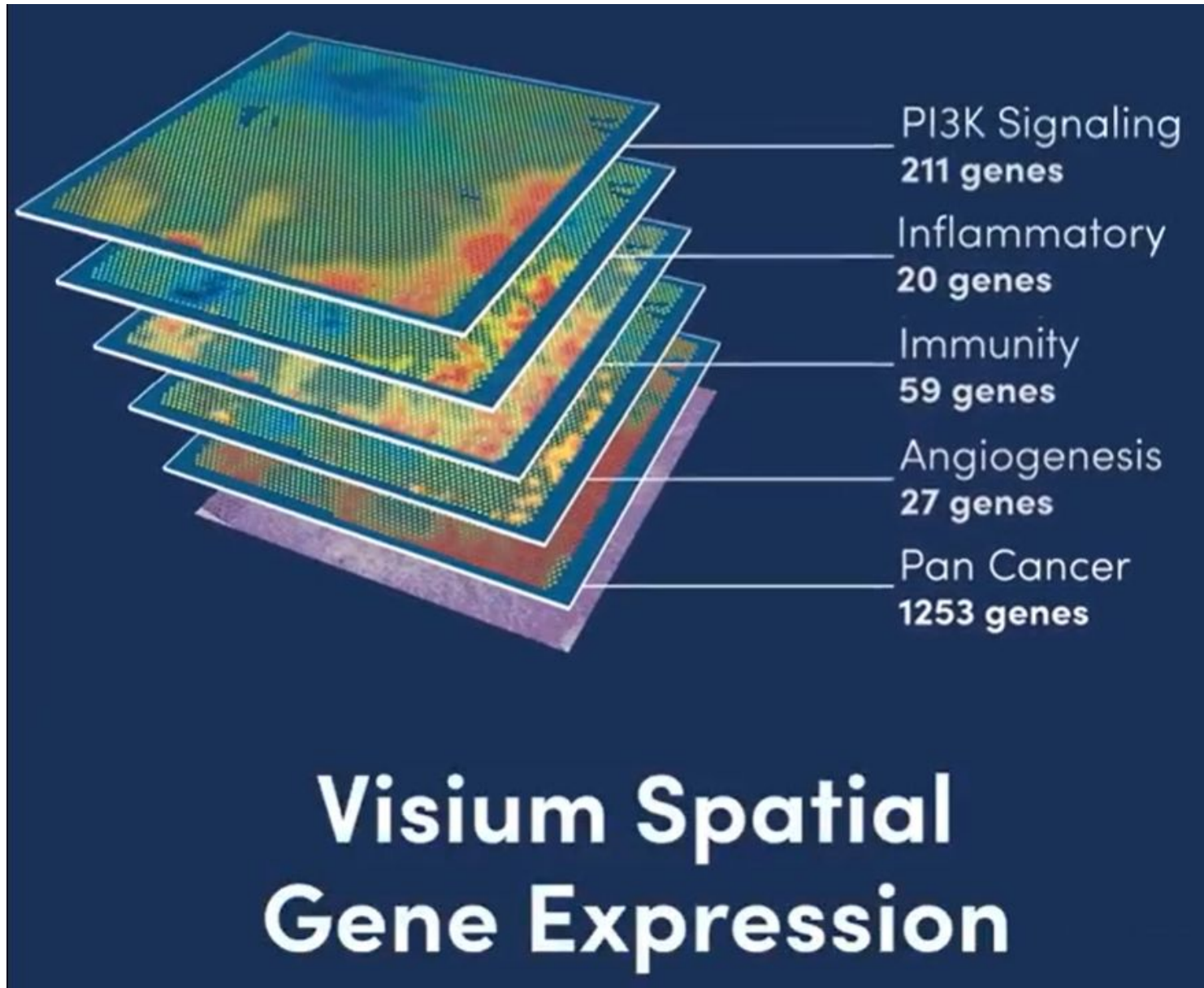
Single-cell RNA-seq



Single-cell RNA-seq



Spatial Single-cell RNA-seq



Single-cell RNA-seq Full-Length

Article | [Open Access](#) | Published: 12 August 2020

High throughput error corrected Nanopore single cell transcriptome sequencing

Kevin Lebrigand , Virginie Magnone, Pascal Barbry  & Rainer Waldmann 

Nature Communications **11**, Article number: 4025 (2020) | [Cite this article](#)

9506 Accesses | **6** Citations | **58** Altmetric | [Metrics](#)



Informations on splicing, fusions, SNPs, editing, imprinting are preserved



PacBio Sequel II



Oxford Nanopore Technology

+++ higher accuracy	+++ high throughput (80M reads / flowcell)
--- low throughput (4M reads / SMRT)	--- lower accuracy
Price (1.400€ / SMRT)	Price (1.400€ / Flowcell)

2 challenges to tackle

- (1) get enough reads to profile molecules (50k reads / cell)
- (2) high accuracy for cell barcode and UMI identification

Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- **Full-length transcripts sequencing using long reads**
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Why ?

SHORT READ TRANSCRIPT ASSEMBLY IS INSUFFICIENT



How insufficient ?

STEIJGER (2013) VS. ANGELINI (2014) VS. CHANG (2014)

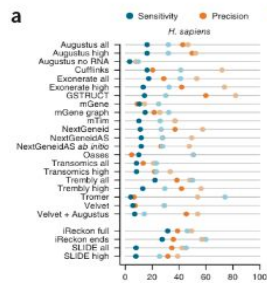
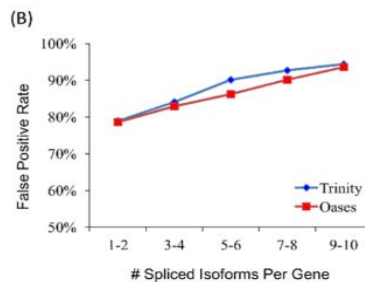
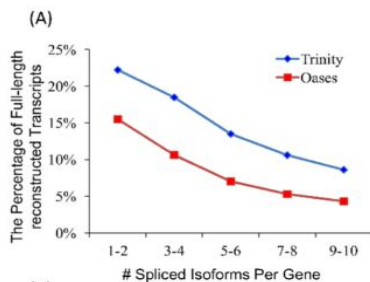
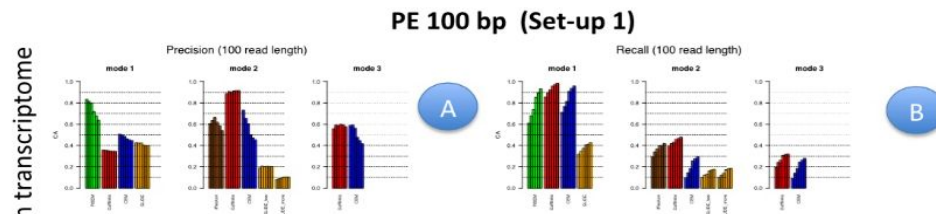


Figure 5 | Transcript assembly performance.
(a) Reference transcripts with a matching submission entry (transcript-level sensitivity, blue) and reported transcripts that match the reference (transcript-level precision, orange).
(b) Transcripts for which various subsets of constituent exons have been reported.

Steijger (2013): recall best at 20%, precision best at 40%
 Performance worse for non-coding transcripts

Angelini (2014): recall best < 30%, precision best ~60%
 Simulated coding transcripts



Chang (2014): recall best < 25%, precision best ~20%
 Simulated coding+non-coding transcripts

Steijger, T. et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Meth* 10, 1177–1184 (2013).
 Angelini, C., Canditii, D. & Feis, I. Computational approaches for isoform detection and estimation: good and bad news. *BMC Bioinformatics* 15, 135–43 (2014).
 Chang, Z., Wang, Z. & Li, G. The Impacts of Read Length and Transcriptome Complexity for De Novo Assembly: A Simulation Study. *PLoS ONE* 9, e94825–8 (2014).

Long-Reads technologies

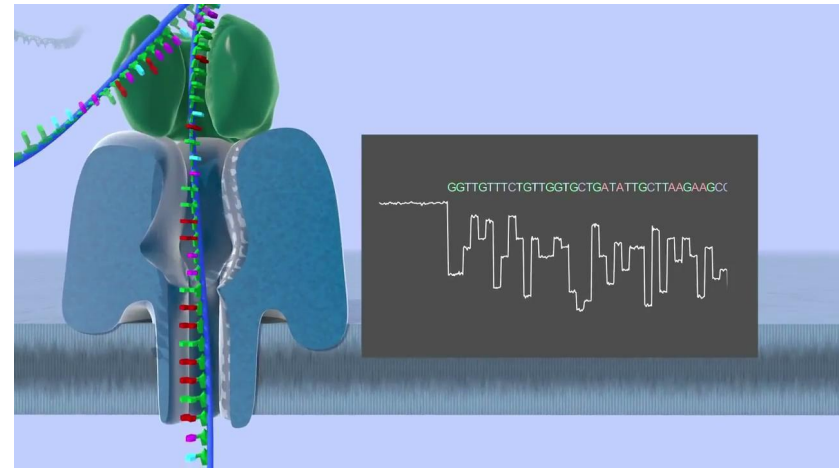
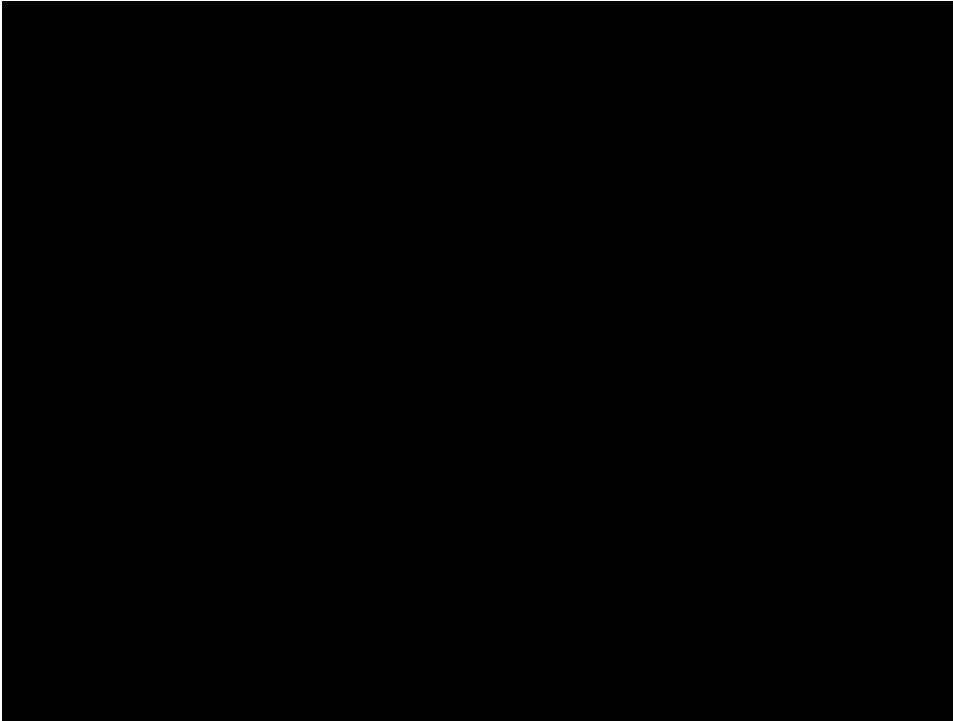


Pacific Bioscience
SEQUEL IIe = 800K€



Oxford Nanopore
PromethION 48 = 600K€

How it works ?

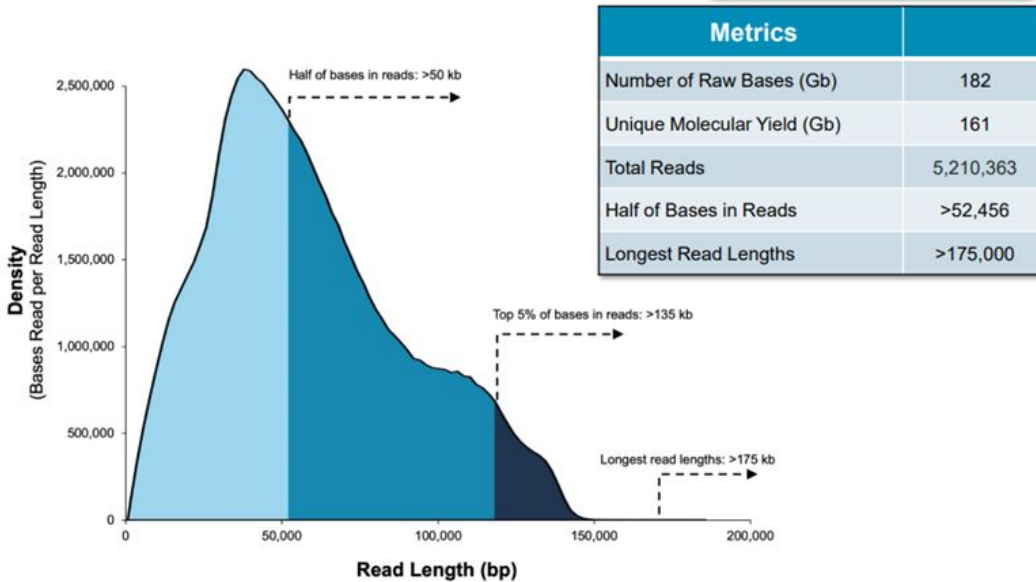


Throughput and lengths

LONG READS

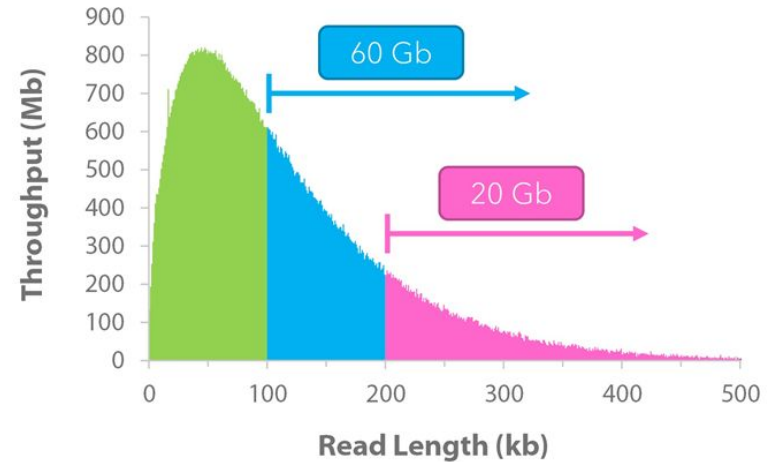
**LONG READ SEQUENCING EXAMPLE:
>35 KB SAMPLE – RAW DATA**

Up to 160 Gb
Average: 100 – 120 Gb



Data shown above from a 35 kb size-selected E. coli library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 15-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

HG02723 on PromethION
130 Gb, 93 kb N50, 1.2 Mb max



Errors are random

i. Generate sequence read:

```
...GAATTCTTAACGTCTGAGACACGACATGCGACCTCTGCACCGGACTCGTCCGCGTCTTTGGCAATCGGGATCAGCTTCGGGAGATGCGGCGCAGCTTGGGGATGGATAGGCGAGCAATGC...
```

ii. Map to reference:

```
...GAATTCTTAACGTC-TGAGACACGACAGCGACCTCTGACCGGACTCG-TCCGCGTCTTTGG-CAATCGGGA-TCAG-CTTCGGGAGATGCGGCGCA-GCTTGGGGATGATAGGCGAGCAATGC...
```

```
...GAATTCTTAACGTCTGAGACACGACAGCGACCTCTGACCGGACTCGTTCGCGTCTTTGGACAATCGGGATTTCAGACTTCGGGGGATGCGGCGCAGGCTTGGGGATGATAGGCGAGCAATGC...
```

iii. Generate consensus (10x coverage):

```
...GAA-TCTTAACGTCTGAGACACG-CAGCGACCTCTGA-CGGACTCGCTCCGCGTCTTTGGACAATC--GATTCAGACTTCGGGAGATGCGGCGCAGGCTTGG-GATGATAGGCGAGCAATGC...
...GAATTCTTAA-GTCTGAGACACGACAGC-ACCTCTGACTTGAC-CGTTCCGCGTCTTTGGACAATCGGGATTTCAGACTTCGGGAGATGCGGCGCAGG-TGGGGATGATAGGCGAGCA-TGC...
...G-ATTC-TAACGTCTGAGACACGACAGCGA-CTCTGACCGGACTCTCTCCGCGGCT-TGGACAA-CGGGATTCAGA-TTCGGGAG-TGGGCGCAGGCTTGGGGATGATAG-CGAGCAATGC...
...GAATTCTTAAC-TCTGAGACACGACAGCGACCTC-GACTGACTCGTTCC-CGTTCTTTGGACAATCGGGA-TCAGACTTCG-GAGATGCGGCGCAGGCTT-GGGATGATAGGCGAGCAATGC...
...GAATTCTTAAACGT-CTGAGACACGACA-CGACCTCTGACTTGA-TCGTTCCGCGTCT-TGGACAATC-GGATTCAGACTTCGGGAGATGCGGCGCAGGC-TGGGGAAGATAGGCGAGCA-TGC...
...GAATTCTT-ACGTCCTGATACACACAGCG-CCTCTG-CGGACTCGTTCGGGCTCTTTGGACAATCGGAGATTCA-ACTTCGGGAGATGCGG-GCAGGCTTGGGGATGATAGGCGAGCAATGC...
...GAATTCTTA-CGTC-CGAGACACGACAGCGACCTCT-ACC-GACTCTCTCCGCGTCTTTGGACAATCGGGATTTCAGACTTCGGGAGATG-GCGCAGGCTTGGGGATGATAGGCGAGCAATGC...
...GAATTCTTAACGTCTGAGACACGACAG-CACCT--GACCTGACTCGT-CCGCGTC-TTGGACAATCGGGATTCAGACTTCGGGAGATGCGGCGCAGGCTTGGGGATGAT-GGCGAGCAATGC...
...GAATTCTTAACGTCTGAGAC-CGACACGACCTCTGACTGACTCGTTCCGCGTCTTTT-ACAATCGG-ATTCAGACT-CGGGAGATGCGGCGCA-GCTTGGGGATG-TAGGCGA-CAATGC...
...GA-TTCTTAACGTCTGAGACAC-ACAGCGACCTCTGACCGGACTCGTTC-GCGTCTTTGG-CAATCGGGATTTCAGACTTCGGGAGATGCGGCGCAGGCTTGGGGATGATAGGCG-GCAATGC...
```

```
...GAATTCTTAACGTCTGAGACACGACAGCGACCTCTGACTGACTCGTTCCGCGTCTTTGGACAATCGGGATTTCAGACTTCGGGAGATGCGGCGCAGGCTTGGGGATGATAGGCGAGCAATGC...
```

```
...GAATTCTTAACGTCTGAGACACGACAGCGACCTCTGACCGGACTCGTTCGCGTCTTTGGACAATCGGGATTTCAGACTTCGGGGGATGCGGCGCAGGCTTGGGGATGATAGGCGAGCAATGC...
```

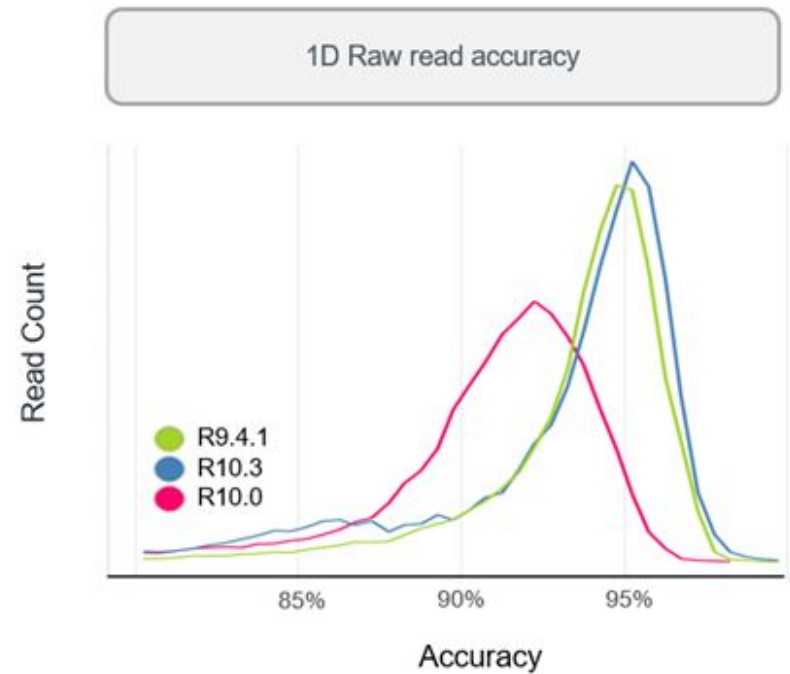
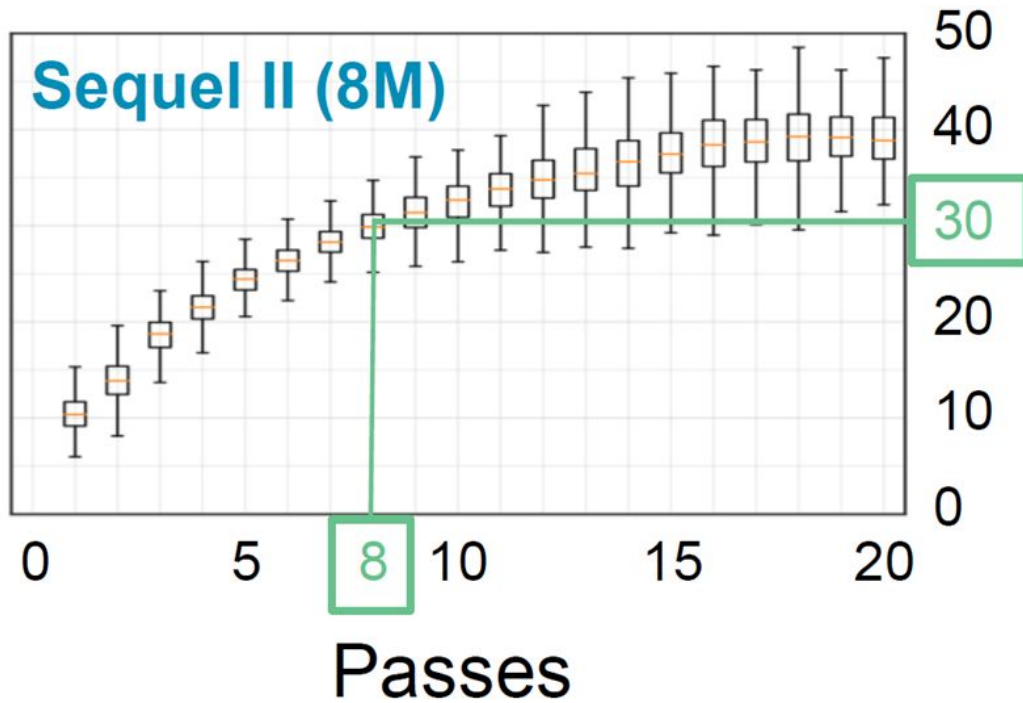
Reference
match

Heterozygous
SNP

Reference

Homozygous
SNP

Errors are random



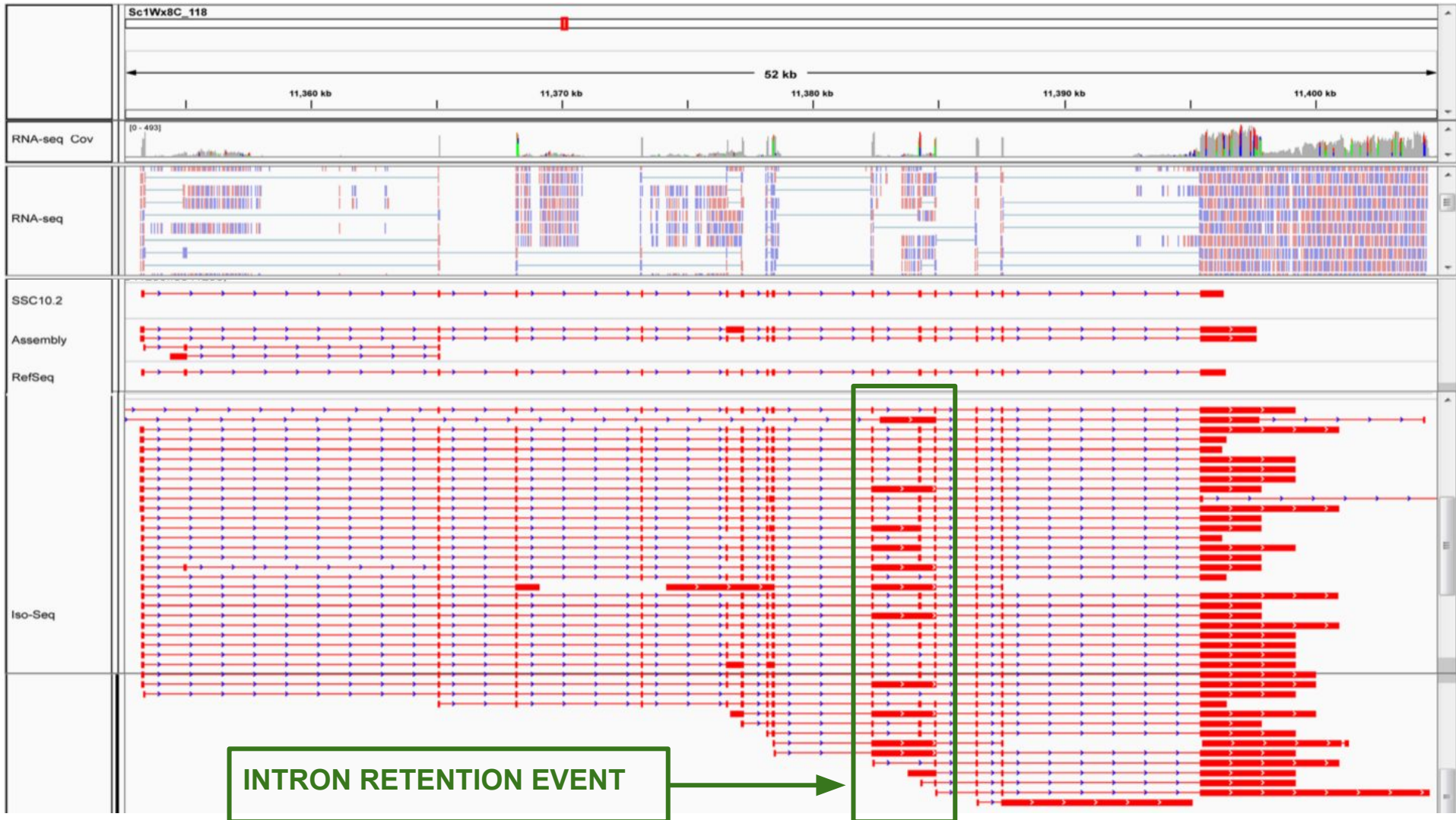
Transcriptome applications

Uncover Novel Insights in Cancer Biology

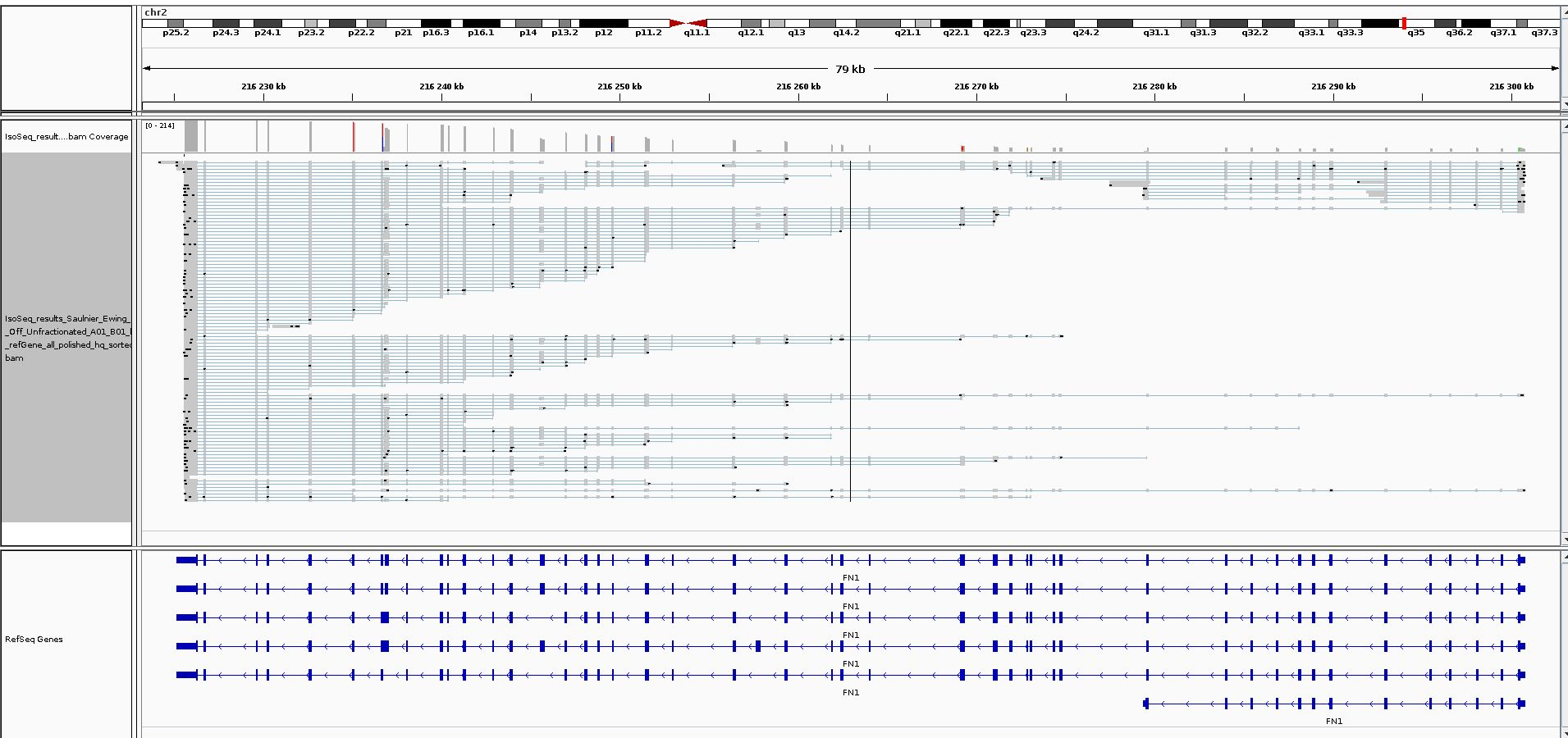
Unlike traditional RNA-Seq techniques, long-read RNA sequencing allows accurate quantification and complete, full-length characterisation of native RNA or cDNA without fragmentation or amplification – streamlining analysis and removing potential sources of bias. Direct RNA sequencing also enables the identification of base modifications alongside nucleotide sequence.

- Full-length transcripts — unambiguous identification of splice variants and gene fusions
- Accurate allele-specific transcript and isoform quantification (thanks to easy phasing)
- Eliminate PCR bias using direct cDNA or direct RNA sequencing
- Detect base modifications alongside nucleotide sequence using direct RNA
- Easy identification of anti-sense transcripts and lncRNA isoforms

Concrete example 1



Concrete example 2



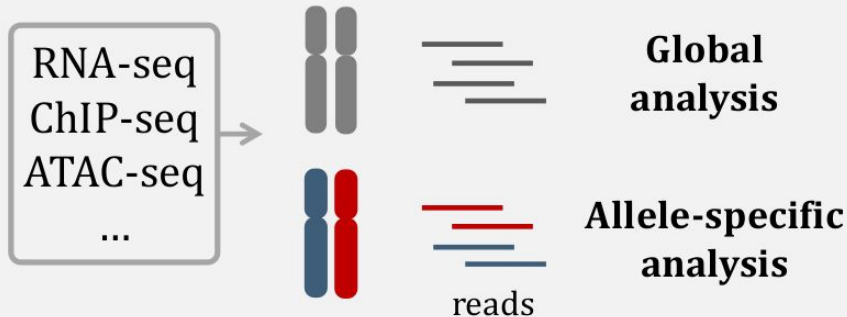
FN1 = Fibronectine1 = 8103 bp

Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- **Allele-specific quantification**
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

Allele-specific analysis

Principle

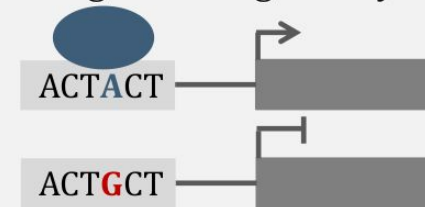


Difficulties

- Need to know haplotypes
- Limited by the number of differential SNPs
- Require dedicated bioinformatics strategy

Domains of application

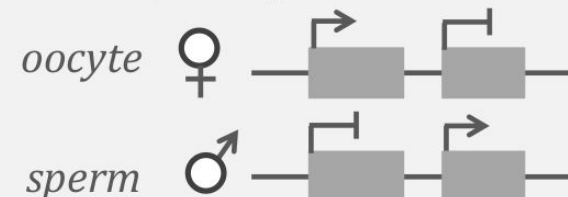
- Effects of genetic regulatory variants



- Non sense mediated decay SNP

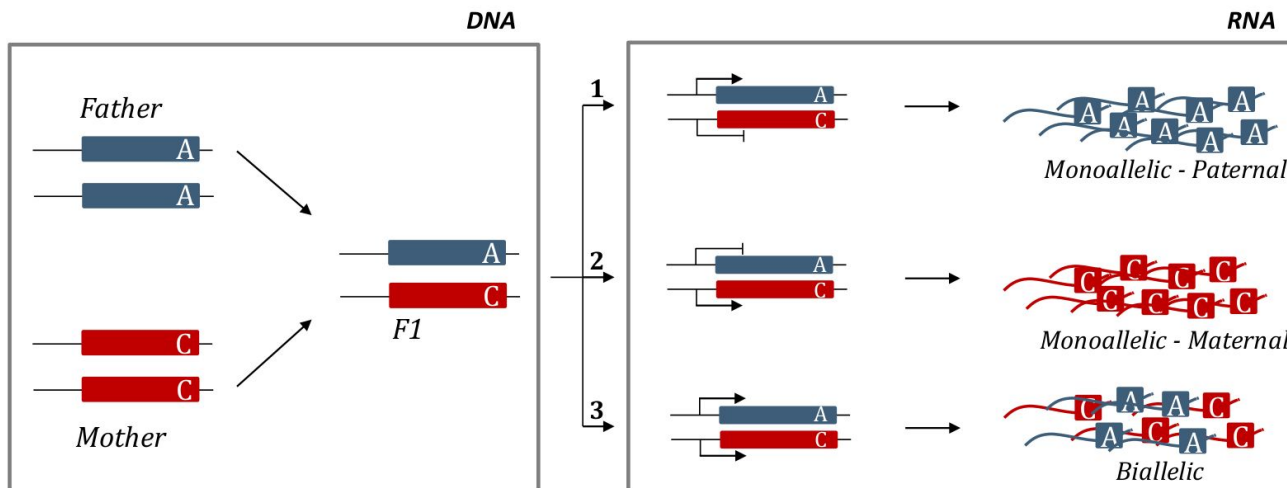
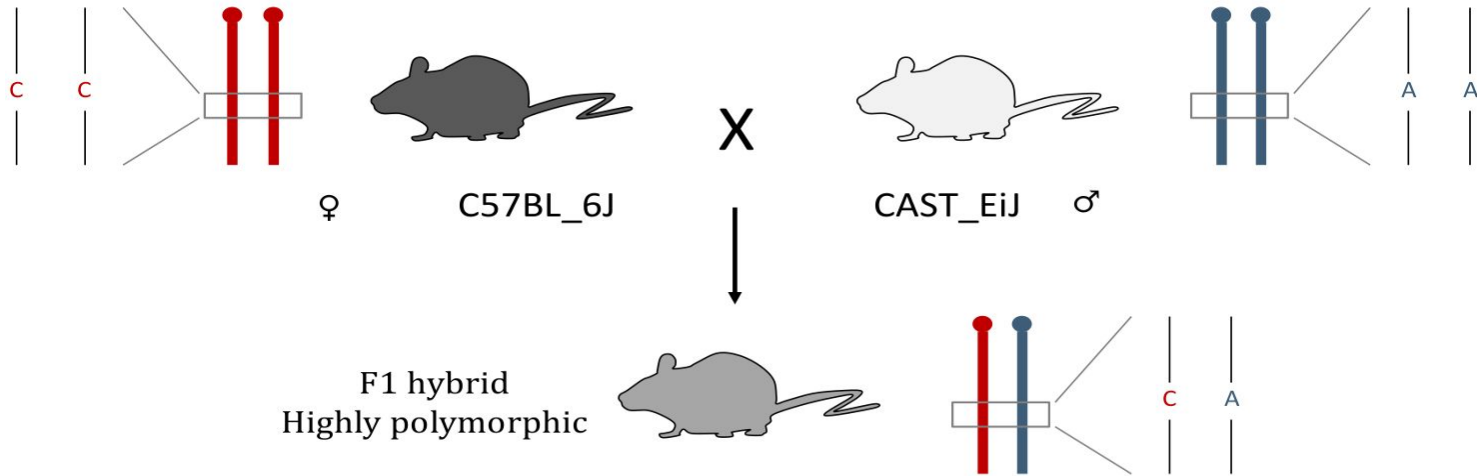


- Parental imprinting



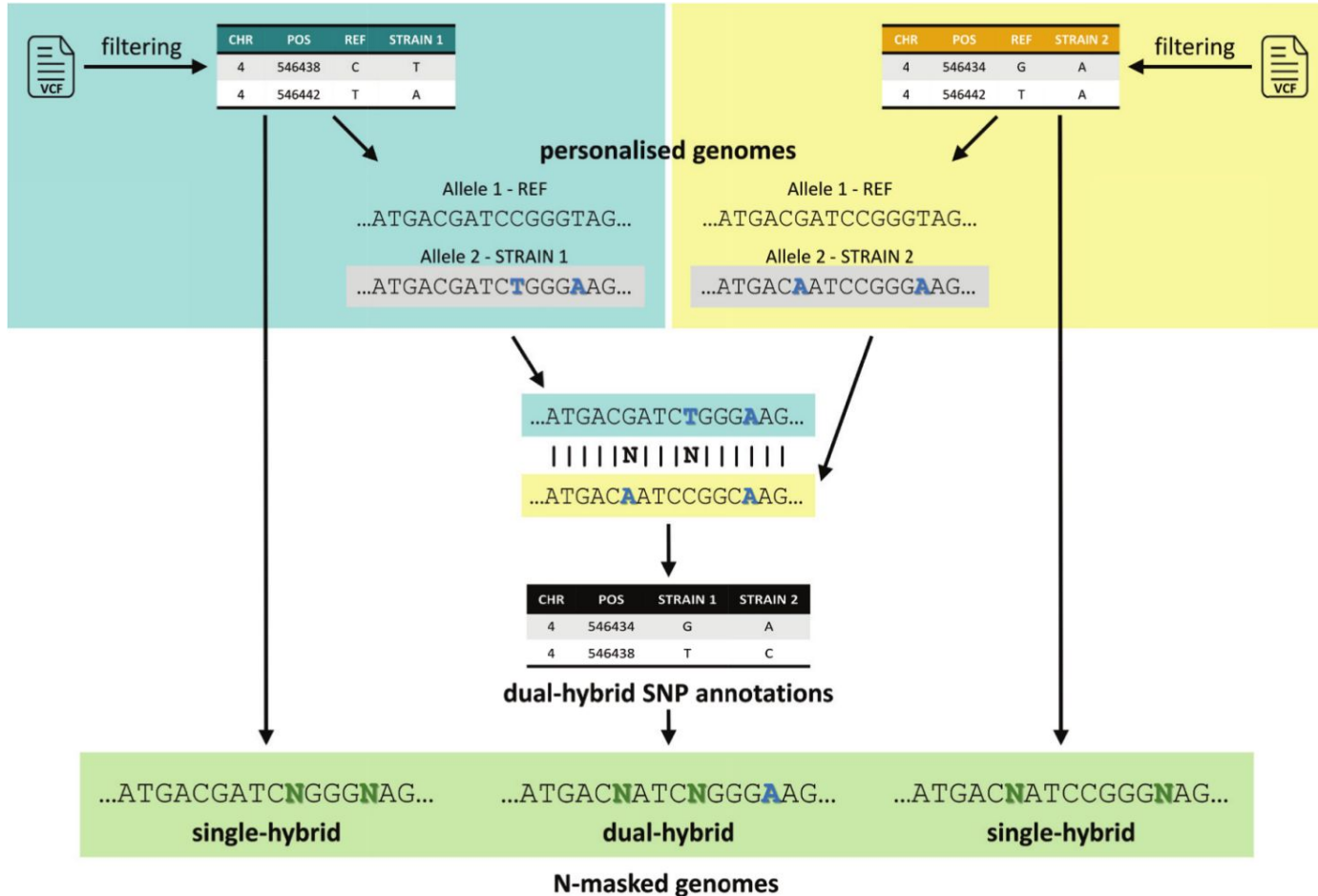
- X chromosome inactivation

Allele-specific analysis



Allele-specific analysis

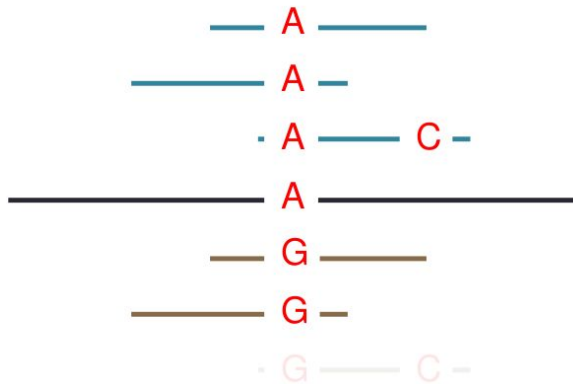
Method 1 : N-masked alignment strategy



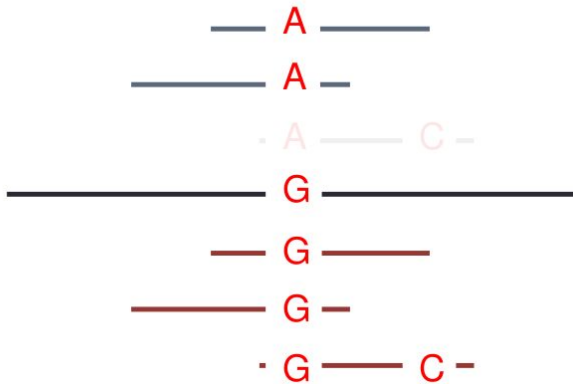
Allele-specific quantification

Method 2 : Alignment to parental genome strategy

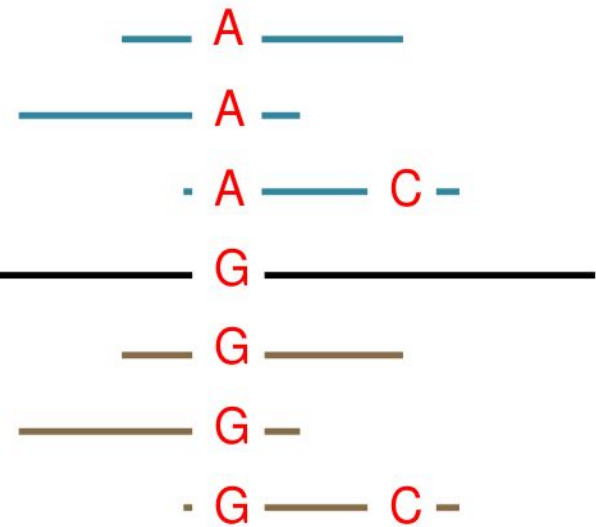
Alignment on genotype 1 genome



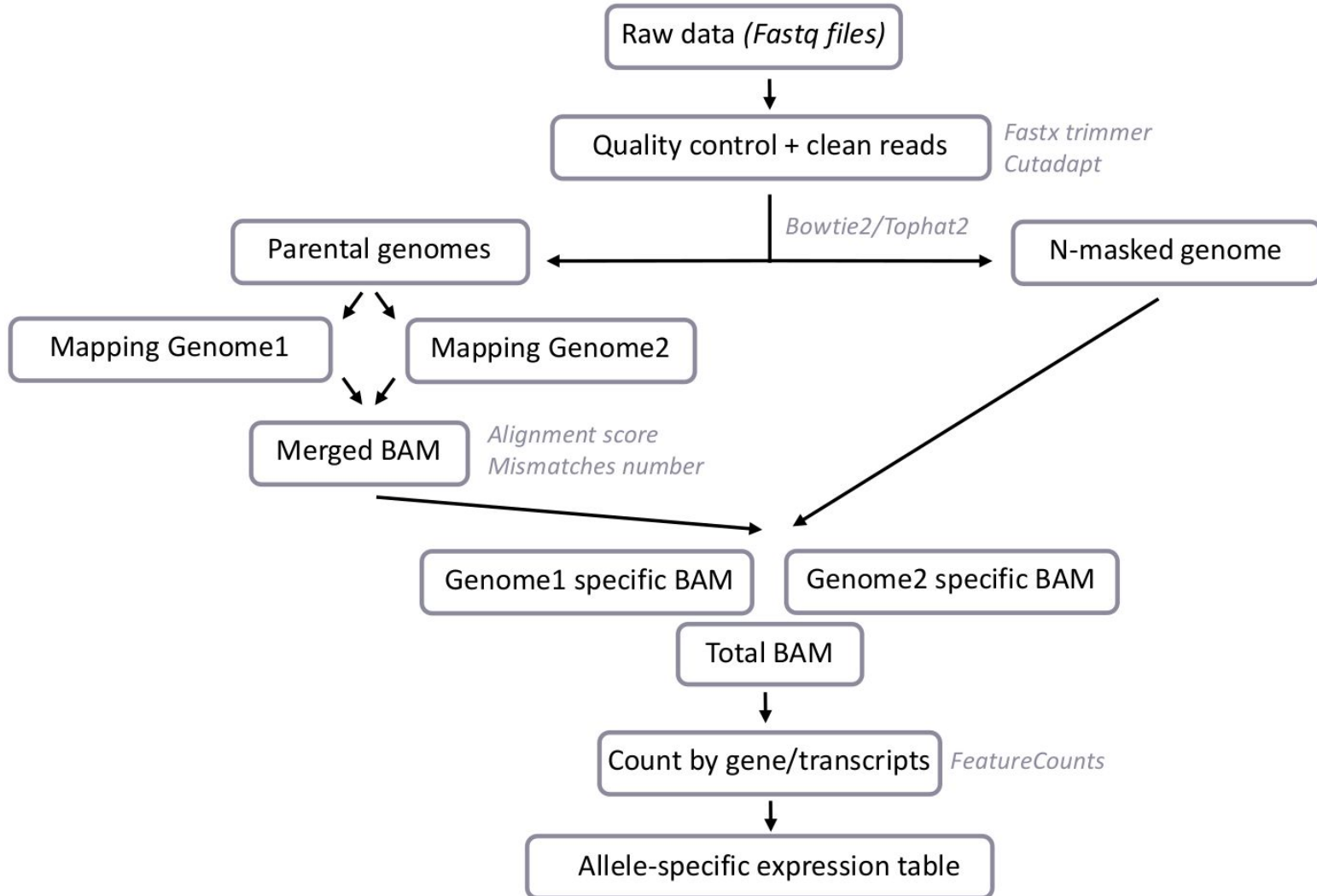
Alignment on genotype 2 genome



Select
Best
Alignments



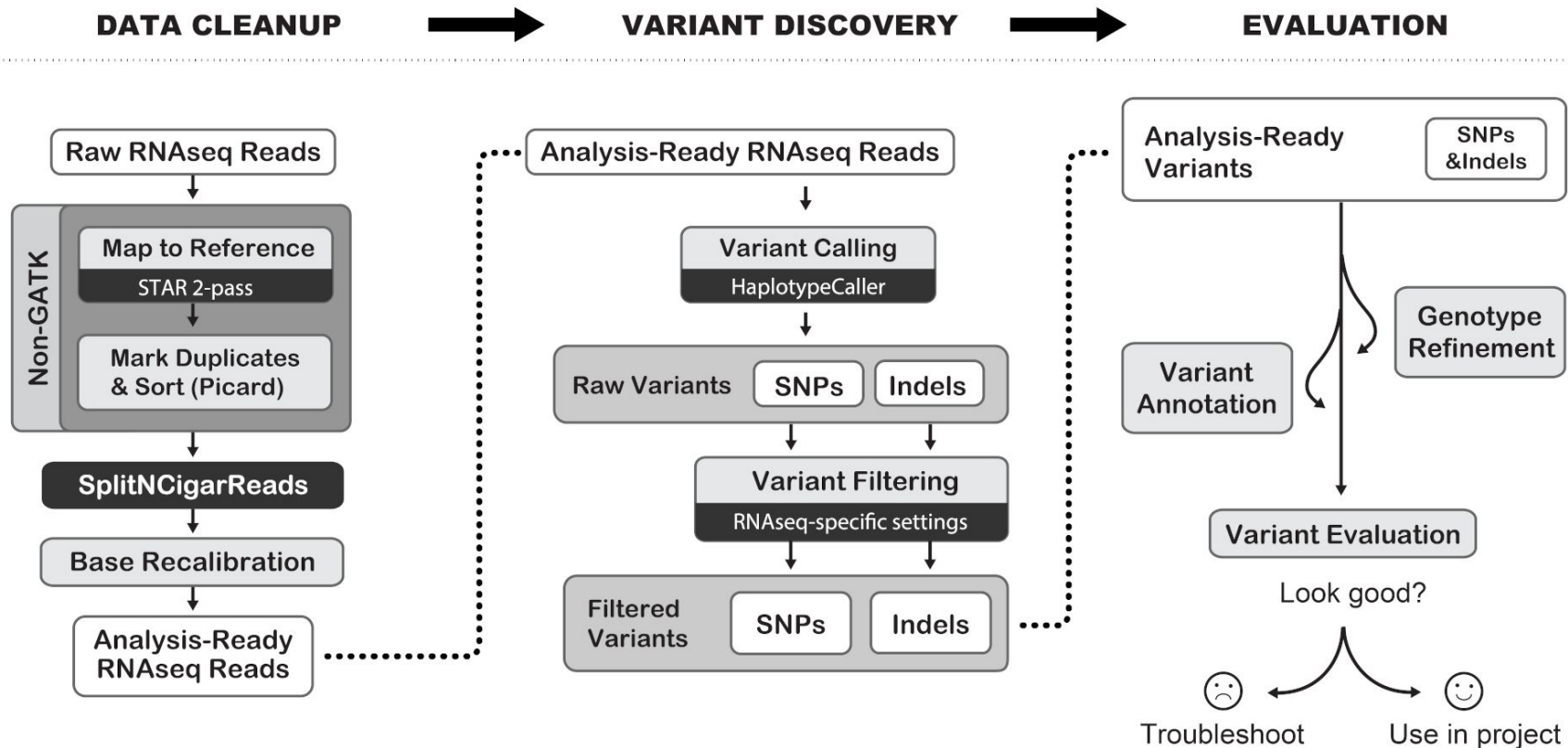
Allele-specific quantification



Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- **Expressed SNVs/indels variants detection**
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

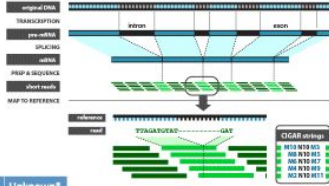
Expressed SNVs/indels



Expressed SNVs/indels

Choosing the best aligner for calling SNPs and indels

Challenge
Aligning RNAseq data to a reference genome is complicated by RNA splicing.



Test results
Of all aligners tested, STAR and TopHat2 performed best, and STAR achieved the highest sensitivity and specificity for both SNPs and indels using a highly curated callset as truth set (CEU sample NA12878).

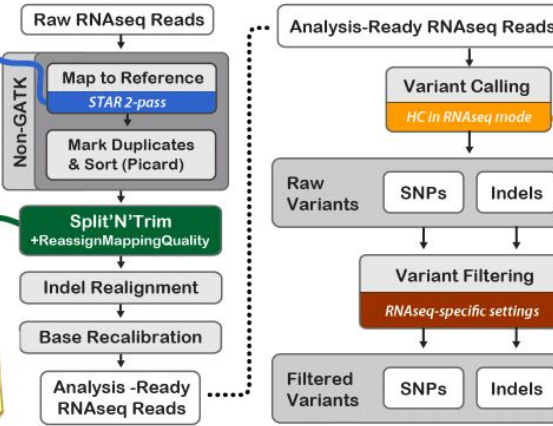
Variant	Aligner	TP	FP	FN	TN	Unknown*
SNP	STAR	8310	65	1851	124422	543
	TopHat2	8148	68	2007	124419	1386
INDEL	STAR	151	21	92	12889	22
	TopHat2	70	20	169	12896	546

* Uncurated variants, mainly false positive calls, some may be RNA editing events.

Selected protocol: 2-pass STAR

- 1) First round of alignment identifies the splice junctions (SJ)
- 2) Use splice junctions to guide the second round of alignment

DATA CLEANUP → VARIANT DISCOVERY



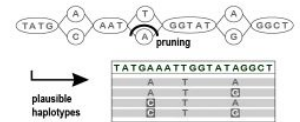
How do we clean up RNAseq data?

What should the variant model be like?

HaplotypeCaller handles graph issues caused by SJs

Challenge

HaplotypeCaller (HC) performs graph-based reassembly in order to compensate for limitations of genome alignment software, but splice junctions cause dangling heads and tails. This in turn causes false negative calls.



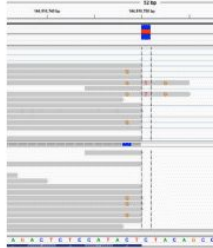
Solution

RNAseq-specific feature of HC allows merging of dangling tails and heads in order to rescue variants located near splice junctions.

Dealing with Ns in read CIGARs

Challenges

Reads mapped across splice junctions have CIGAR strings with Ns and overhangs causing noise.



Solution

New GATK tool called *SplitNCigarReads*

- Splits reads with Ns in the CIGAR string
- Keeps track of grouping information per exon
- Trims overhang reads which cause many false positives (see figure).

CURRENT CHALLENGES & FUTURE WORK

Caveats

This work is in active development. Recommendations are based on limited dataset, single tissue sample with high quality data (low error rates).

Known issues

- Non-optimized filtering process.
- Wrong calls due to allele imbalance & low coverage

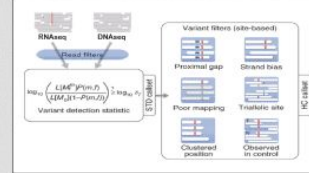
Tweaks for discovery of RNA editing

Primary approach

Contrast calls from both RNAseq and DNAseq

Secondary approach

Use a caller that is less constrained by allele frequency expectations: MuTect, designed to detect somatic (cancer) mutations.



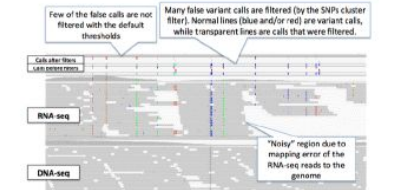
Filtering annotations & thresholds

Challenge

Variant Recalibration (VQSQR) not yet supported for RNAseq.

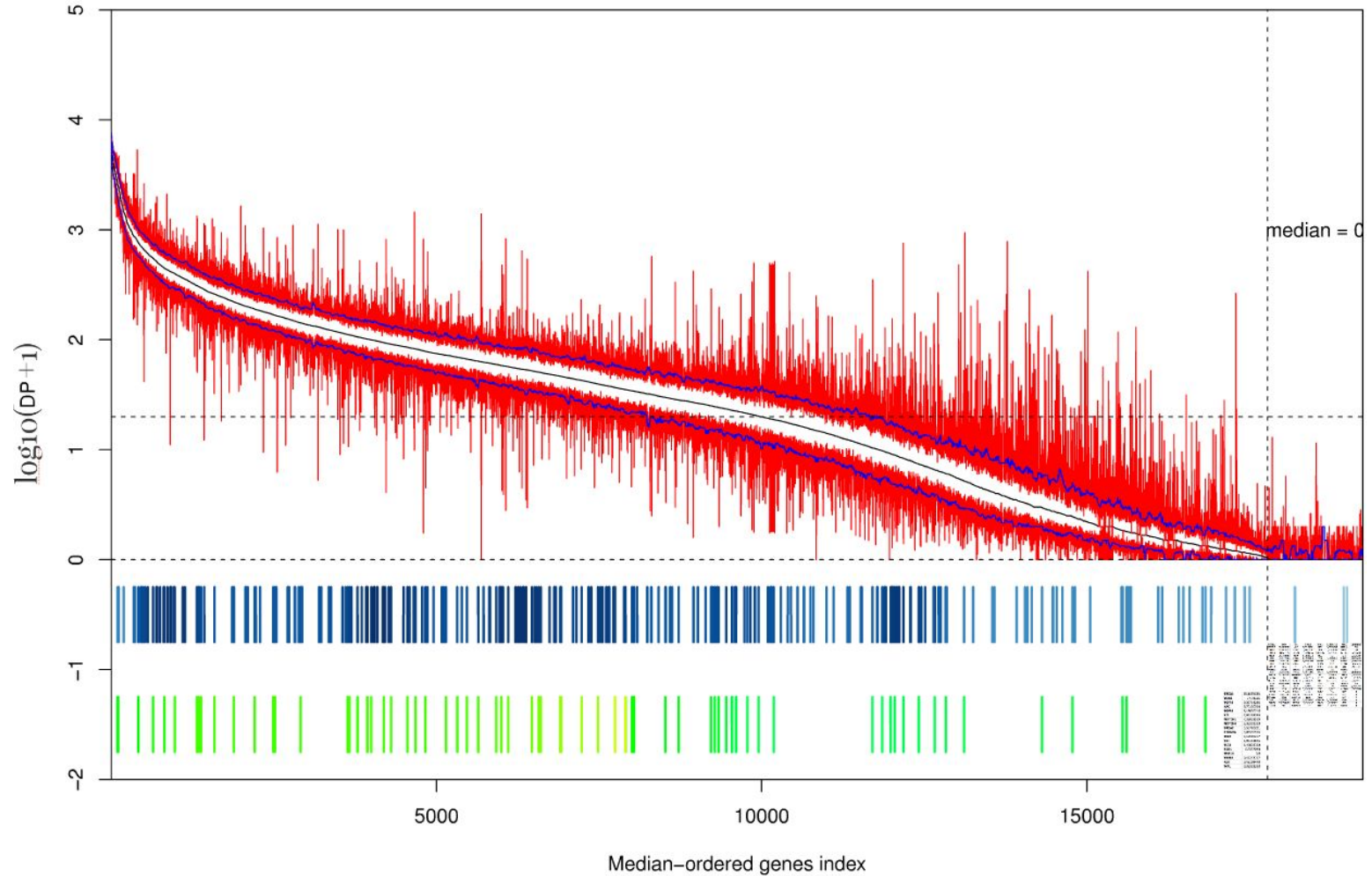
Solution

Manual filtering using hard thresholds on specific annotations (FS for strand bias, QD for quality normalized by depth) and clustering.



Expressed SNVs/indels

Median-ordered median/q25/q75 genes RNA coverage as $\log(DP+1)$
and index of 270 (hotspot) & 73 (PGM) cancer genes



Expressed SNVs/indels

BRCA1	19,1435185
DDR2	17,78125
FGFR3	9,35714286
APC	8,76136364
FGFR4	8,18656716
KIT	7,95238095
NOTCH1	7,33823529
NOTCH4	6,38333333
BRCA2	5,55769231
CDKN2A	5,05555556
JAK3	4,19565217
RET	1,96153846
FLT3	1,45833333
ROS1	0,8372093
HNF1A	0,8
ERBB4	0,47272727
ALK	0,43103448
MPL	0,33333333

BCORL1	19,54	CBL	9,0625	MYCL	5,28571429	FGF10	0,83333333
KMT2A	19,1805556	APC	8,76136364	FLT4	5,06666667	HNF1A	0,8
BRCA1	19,1435185	PTCH1	8,61437909	FLT4	5,06666667	WISP3	0,8
RPTOR	18,328125	EPHA3	8,47916667	CDKN2A	5,05555556	NTRK3	0,77692308
BRD3	18,0454546	FGFR4	8,18656716	BCL2	5	WT1	0,76041667
SUFU	17,9347826	BRIP1	8,10526316	JAK3	4,19565217	GRIN2A	0,56756757
DDR2	17,78125	MITF	8,06521739	GATA3	3,8	NTRK1	0,52941177
BTK	17,2307692	CDK6	8	CEBPA	3,125	ERBB4	0,47272727
CCND2	16,5	KIT	7,95238095	IKZF1	3,06097561	ALK	0,43103448
FCHSD1	16,225	CARD11	7,79166667	ZNF703	2,5	GATA1	0,4
BLM	15,6410256	NTRK2	7,67441861	PAX8	2,27777778	MPL	0,33333333
BARD1	15,0512821	KIAA1549	7,6375	AR	2,21875	FGF19	0,33333333
INHBA	15	HGF	7,5877193	EPHB1	2,09375	NTRK1	0,30106952
RICTOR	14,3962264	NOTCH1	7,33823529	RET	1,96153846	LRP1B	0,24725275
ETV1	14,1643836	CD79A	6,95	IRF4	1,75	NKX2-1	0,2
TET2	13,1	IL7R	6,9375	MYCN	1,75	EPHA5	0,16091954
IRS2	12,5	NOTCH4	6,38333333	ESR1	1,6875	NUTM1	0,13636364
CD79B	11,3235294	PIK3CG	6,25	PLAG1	1,6	FGF23	0
STAT4	11,1521739	GATA2	6,03333333	FLT3	1,45833333	FGF3	0
PRDM1	10,25	SLC34A2	5,625	FGF14	1,4	FGF4	0
DOT1L	10,1785714	BRCA2	5,55769231	CDKN2B	1,16666667	FGF6	0
FGFR3	9,35714286	SS18L1	5,55555556	PAK3	0,85227273	SSX1	0
RNF43	9,1	SLC45A3	5,375	ROS1	0,8372093	SOX10	0

Expressed SNVs/indels

The Author(s) *BMC Genomics* 2017, **18**(Suppl 6):690
DOI 10.1186/s12864-017-4022-x

BMC Genomics

RESEARCH

Open Access

The discrepancy among single nucleotide variants detected by DNA and RNA high throughput sequencing data



Yan Guo^{1*†}, Shilin Zhao^{1†}, Quanhu Sheng¹, David C Samuels² and Yu Shyr^{3*}

From The International Conference on Intelligent Biology and Medicine (ICIBM) 2016
Houston, TX, USA. 08-10 December 2016

Expressed SNVs/indels



Briefings in Bioinformatics, 18(6), 2017, 973–983

doi: [10.1093/bib/bbw069](https://doi.org/10.1093/bib/bbw069)

Advance Access Publication Date: 26 July 2016

Paper

Indel detection from RNA-seq data: tool evaluation and strategies for accurate detection of actionable mutations

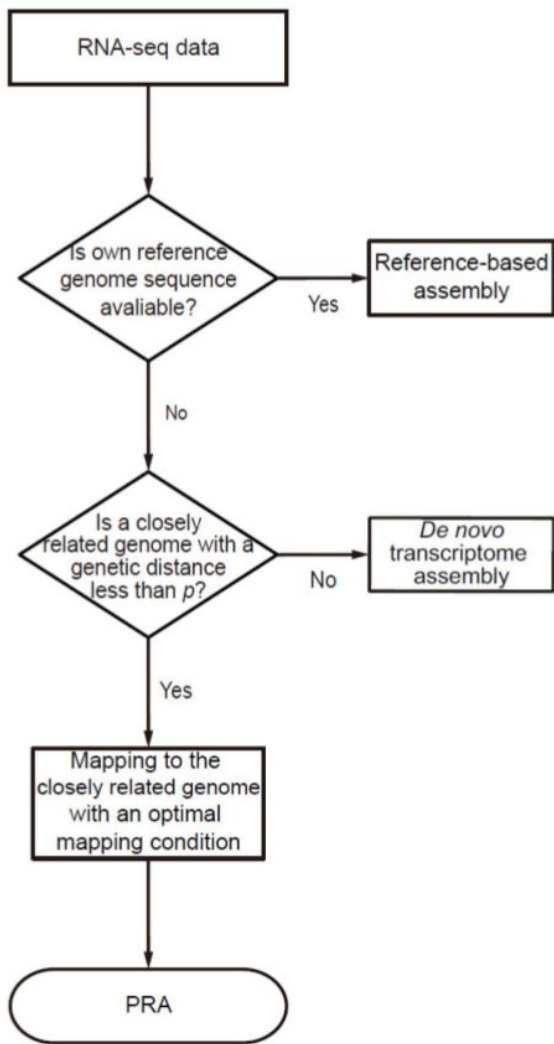
Zhifu Sun, Aditya Bhagwate, Naresh Prodduturi, Ping Yang and Jean-Pierre A. Kocher

Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- **Transcripts reconstruction (assembly)**
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- RNA editing events detection

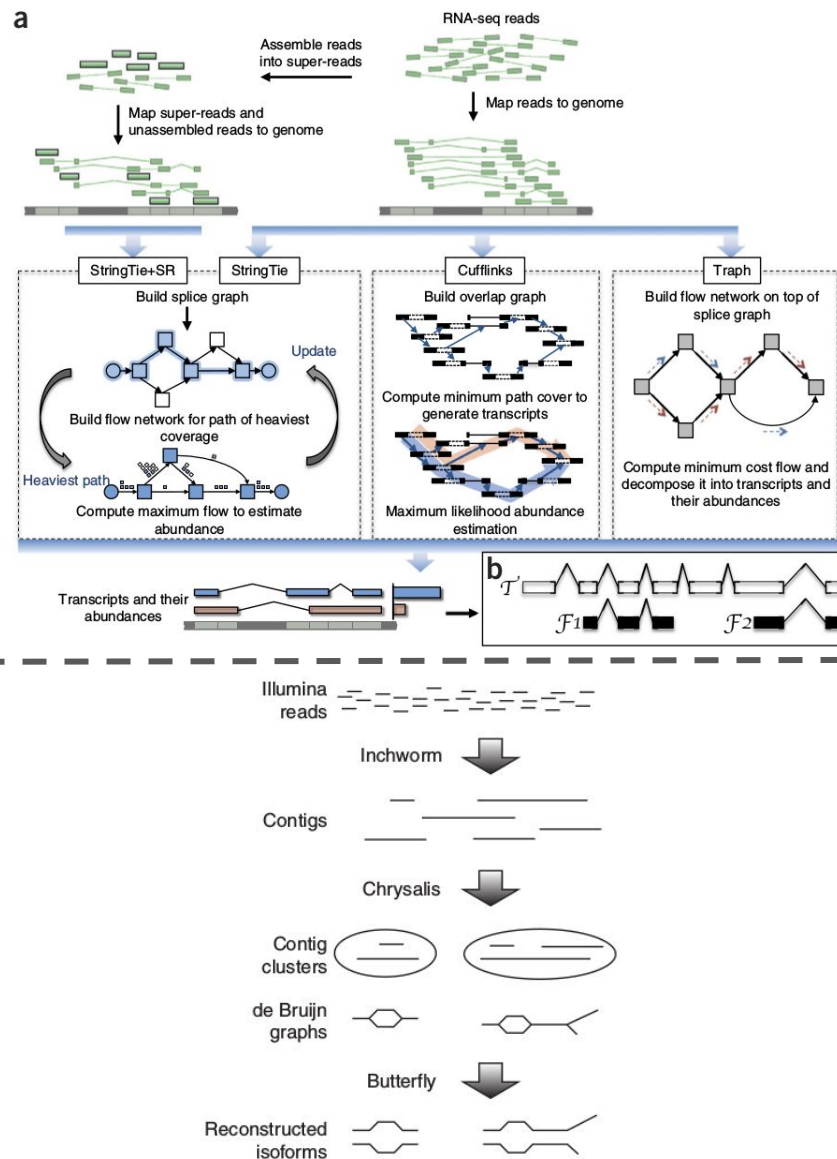
Transcripts reconstruction

Nam K. et al. 2016. Genes



StringTie
(Pertea M. et al. 2015. Nature biotechnology)

Trinity
(Haas BJ. et al. 2013. Nature protocols)



Transcripts reconstruction

BIOINFORMATICS ORIGINAL PAPER

Vol. 30 no. 17 2014, pages 2447–2455
doi:10.1093/bioinformatics/btu317

Gene expression

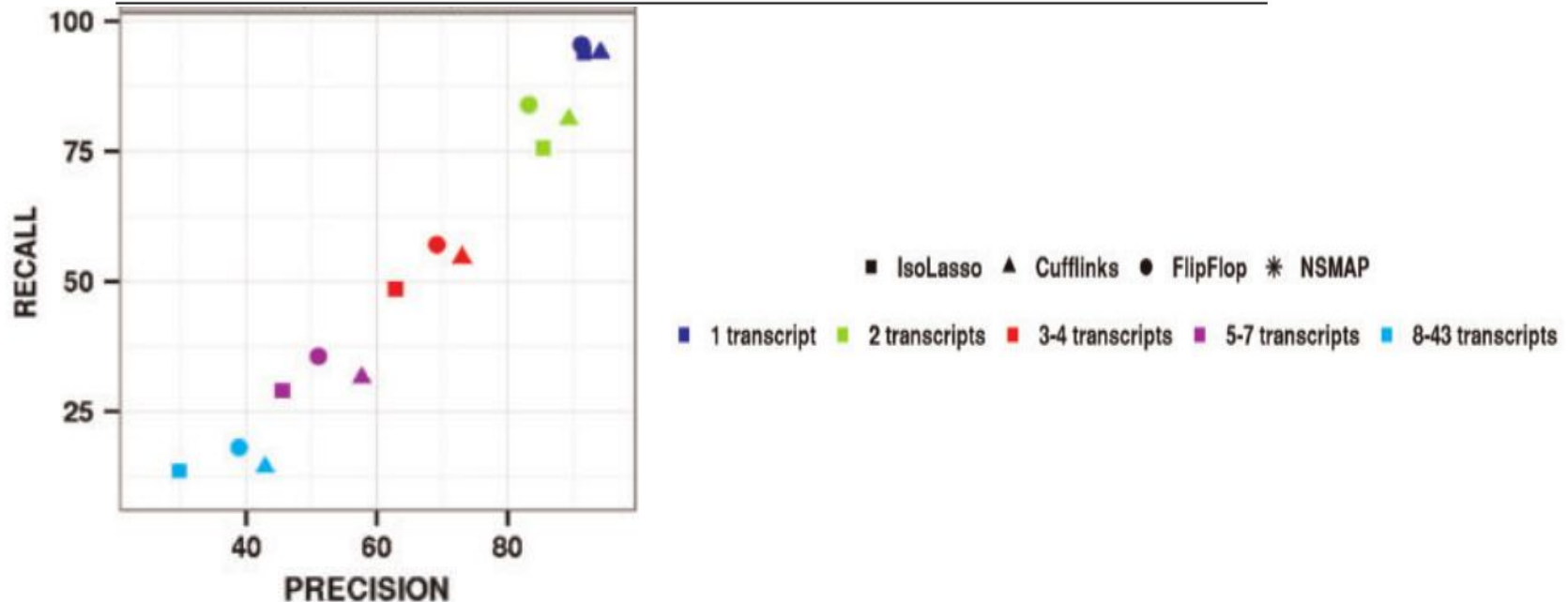
Advance Access publication May 9, 2014

Efficient RNA isoform identification and quantification from RNA-Seq data with network flows

Elsa Bernard^{1,2,3}, Laurent Jacob⁴, Julien Mairal⁵ and Jean-Philippe Vert^{1,2,3,*}

¹Mines ParisTech, Centre for Computational Biology, 77300 Fontainebleau, ²Institut Curie, 26 rue d'Ulm, 75248 Paris Cedex 05, ³INSERM U900, Paris F-75248, France, ⁴Laboratoire Biométrie et Biologie Evolutive, Université de Lyon, Université Lyon 1, CNRS, INRA, UMR5558, Villeurbanne, France and ⁵LEAR Project-Team, INRIA Grenoble Rhône Alpes, 38330 Montbonnot, France

Associate Editor: Janet Kelso



How insufficient

STEIJGER (2013) VS. ANGELINI (2014) VS. CHANG (2014)

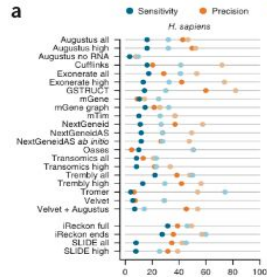
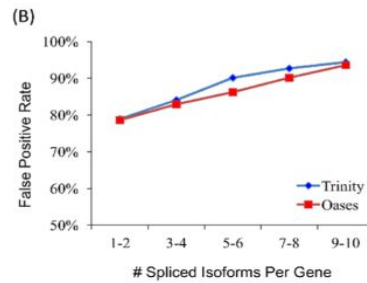
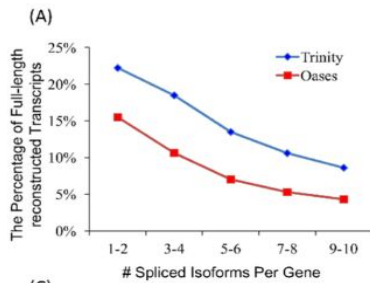
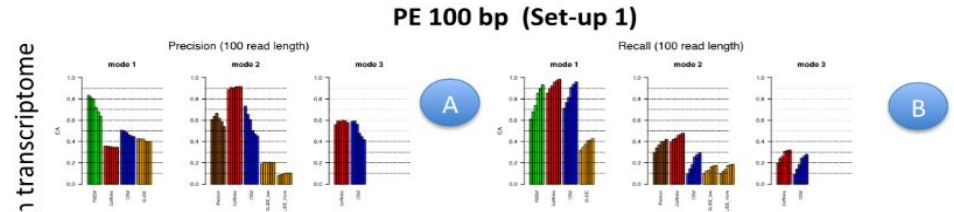


Figure 5 | Transcript assembly performance.
 (a) Reference transcripts with a matching submission entry (transcript-level sensitivity, blue) and reported transcripts that match the reference (transcript-level precision, orange).
 (b) Transcripts for which various subsets of constituent exons have been reported.

Steijger (2013): recall best at 20%, precision best at 40%
 Performance worse for non-coding transcripts

Angelini (2014): recall best < 30%, precision best ~60%
 Simulated coding transcripts



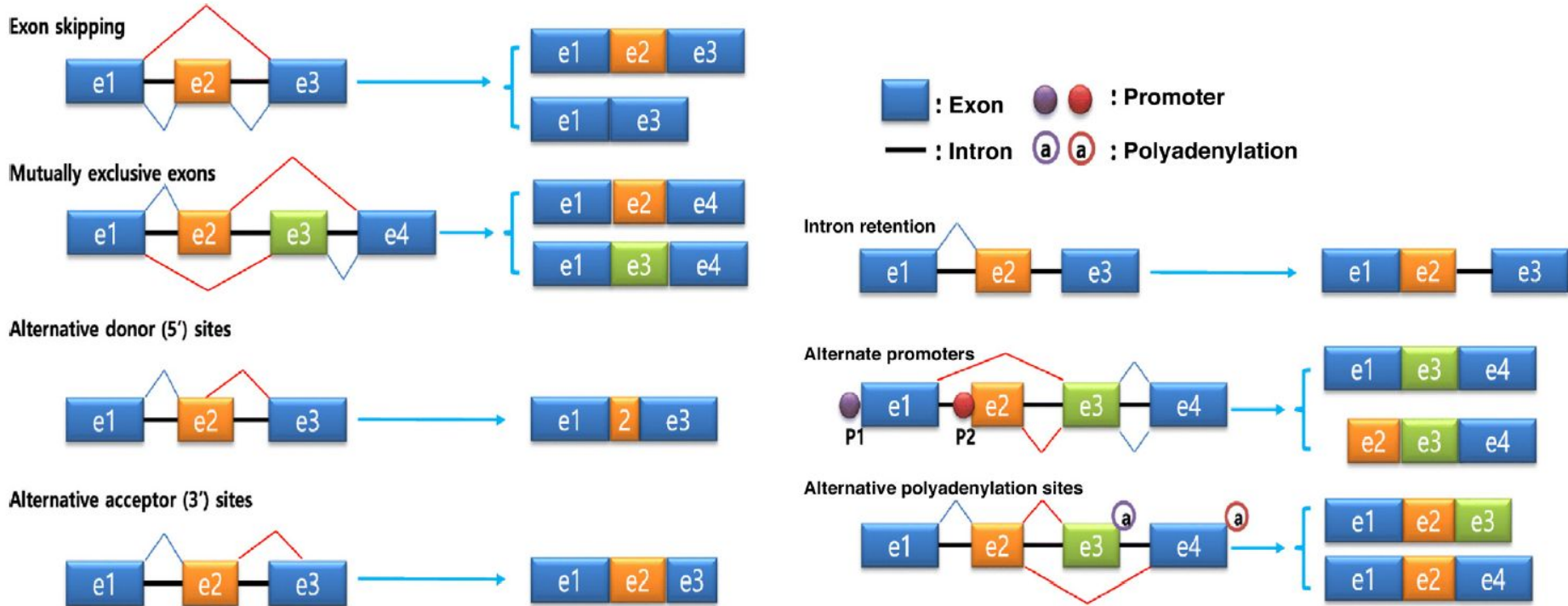
Chang (2014): recall best < 25%, precision best ~20%
 Simulated coding+non-coding transcripts

Steijger, T. et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Meth* 10, 1177–1184 (2013).
 Angelini, C., Canditiis, D. & Feis, I. Computational approaches for isoform detection and estimation: good and bad news. *BMC Bioinformatics* 15, 135–43 (2014).
 Chang, Z., Wang, Z. & Li, G. The Impacts of Read Length and Transcriptome Complexity for De Novo Assembly: A Simulation Study. *PLoS ONE* 9, e94825–8 (2014).

Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- **Alternative splicing / polyadenylation events detection**
- Virus/phages expression detection
- RNA editing events detection

Alternative splicing / polyA



Alternative splicing / polyA

Computational tools to analyze alternatively spliced transcripts can be divided into two categories based on their functionality. Tools in the first category include RSEM [33], Kallisto [34] and Salmon [35], which can be applied to analyze known or annotated transcript isoforms. Tools in the second category include MISO [17], MAJIQ [19], rMATs [36] and LeafCut [37], which can be used to analyze RNA-seq data at the exon level to detect known and novel splicing events. Owing to different definitions on AS events, different software applied to the same datasets often come with different predictions; sometimes the overlap of the software predictions can be very small [38]. A widely used metric for AS is percent spliced [37] (PSI or ψ), which represents the percentage of a gene's mRNA transcripts that include a specific exon or splice site. For a given AS event, the PSI value can be calculated from the number of RNA-seq reads supporting specific exons or splice junctions. However, if the sequencing depth is not deep enough or the expression levels for a particular gene and its isoforms are very low, the calculated ψ is less reliable. Although RNA-seq is indeed now the preferred technology to monitor genome-wide AS, RT-PCR remains the gold standard to validate novel splicing variants and confirm changes in PSI. This is especially true when the depth of reads is not sufficient for high-confidence assessment.

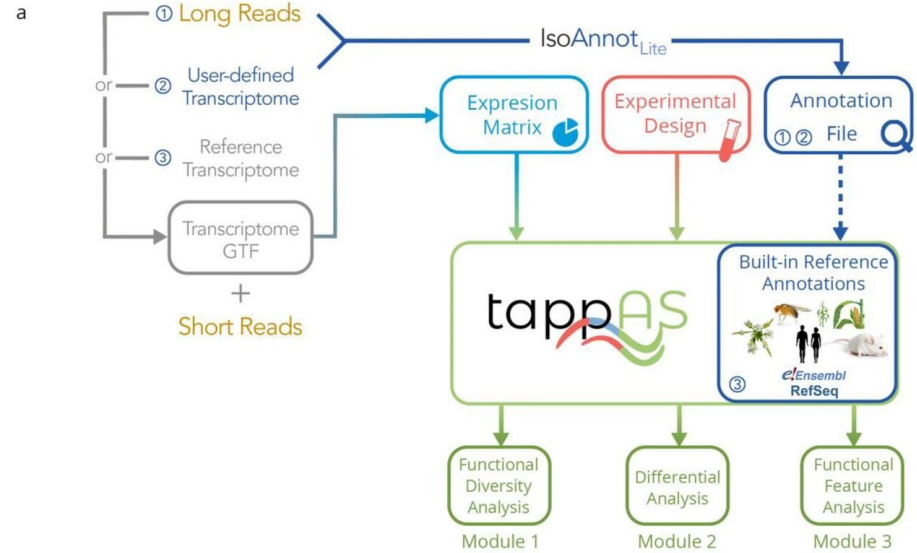
Alternative splicing / polyA

Bioinformatics, 34(15), 2018, 2521–2529
doi: 10.1093/bioinformatics/bty110
Advance Access Publication Date: 23 February 2018
Original Paper

Genome analysis

TAPAS: tool for alternative polyadenylation site analysis

Ashrafal Arefeen¹, Juntao Liu², Xinshu Xiao^{3,*} and Tao Jiang^{1,4,5,*}

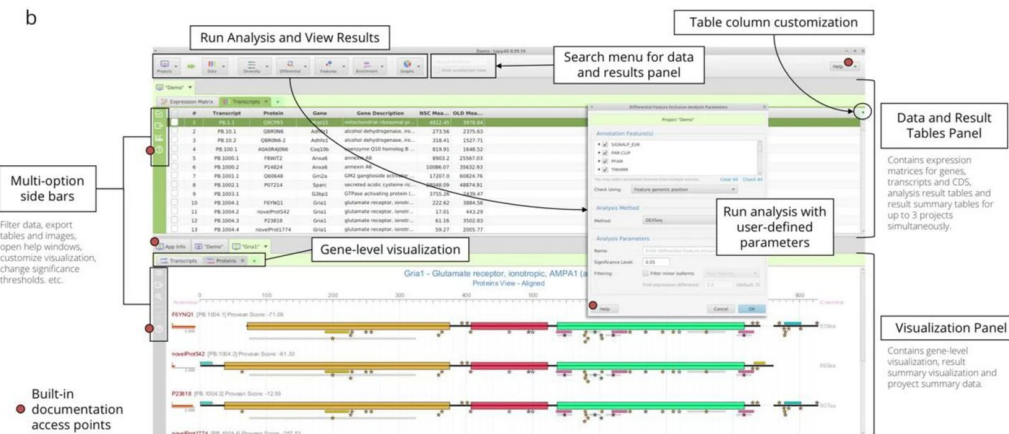


Software | Open Access | Published: 18 May 2020

tappAS: a comprehensive computational framework for the analysis of the functional impact of differential splicing

Lorena de la Fuente, Ángeles Arzalluz-Luque, Manuel Tardaguila, Héctor del Risco, Cristina Martí, Sonia Tarazona, Pedro Salguero, Raymond Scott, Alberto Lerma, Ana Alastrue-Agudo, Pablo Bonilla, Jeremy R. B. Newman, Shunichi Kosugi, Lauren M. McIntyre, Victoria Moreno-Manzano & Ana Conesa

Genome Biology, 21, Article number: 119 (2020) | Cite this article



Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- **Virus/phages expression detection**
- RNA editing events detection

Virus/Phages expression

Table 1. Association of viruses with human cancer.

Cancer	#Samples	HBV	HCV	EBV	CMV	HHV6B	HPV16	HPV38	Other viruses
ACC - adrenocortical carcinoma	79								
BLCA - bladder urothelial carcinoma	119			1	3	1	1		3(HPV6)
BRCA - breast invasive carcinoma	125					1			
CESC - cervical squamous cell carcinoma	91				1	1	59		11(HPV18), 9(HPV45), 3(HPV59), 13(HPV ^{other})
COAD - colon adenocarcinoma	44			2	4	2			2(HPV6)
DLBC - diffuse large B-cell lymphoma	28								
GBM - glioblastoma multiforme	95								
HNSC - head & neck squamous cell carcinoma	123			2	1		14		2(HPV33), 2(HPV35), 1(HPV56), 1(HHV1)
KICH - kidney chromophobe	66								
KIRC - kidney renal clear cell carcinoma	67								1(HPV18), 1(HPV94)
KIRP - kidney renal papillary cell carcinoma	120								
LGG - brain lower grade glioma	100	1							
LIHC - liver hepatocellular carcinoma	115	22	6		1		1		1(HPV18), 1(AAV2)
LUAD - lung adenocarcinoma	125								
LUSC - lung squamous cell carcinoma	125			3			1		1(HPV30)
PRAD - prostate adenocarcinoma	124			1	1				
READ - rectum adenocarcinoma	36			5	4				
SKCM - skin cutaneous melanoma	82				1		1		
THCA - thyroid carcinoma	123								
UCEC - uterine corpus endometrioid carcinoma	168							30*	
UCS - uterine carcinosarcoma	57				1				1(HPV5), 1(HPV38b), 1(HPV133)
Grand Total	2012	23	6	14	17	5	77	30	55

Virus/Phages expression

Article | [Open Access](#) | Published: 16 November 2020

Molecular Diagnostics

Human papilloma virus (HPV) integration signature in Cervical Cancer: identification of *MACROD2* gene as HPV hot spot integration site

Maud Kamal , Sonia Lameiras, Marc Deloger, Adeline Morel, Sophie Vacher, Charlotte Lecerf, Célia Dupain, Emmanuelle Jeannot, Elodie Girard, Sylvain Baulande, Coraline Dubot, Gemma Kenter, Ekaterina S. Jordanova, Els M. J. J. Berns, Guillaume Bataillon, Marina Popovic, Roman Rouzier, Wulfran Cacheux, Christophe Le Tourneau, Alain Nicolas, Nicolas Servant, Suzy M. Scholl, Ivan Bièche & RAIDs Consortium

British Journal of Cancer **124**, 777–785(2021) | [Cite this article](#)

Virus/Phages expression

nf-VIF: A Nextflow-based Virus Insertion Finder

Institut Curie - Bioinformatics Core Facility

nextflow $\geq 0.32.0$ MultiQC 1.6 Install with conda singularity available

Introduction

The pipeline is built using [Nextflow](#), a workflow tool to run tasks across multiple computing infrastructures in a very portable manner. It comes with conda / singularity containers making installation trivial and results highly reproducible, and can be run on a single laptop as well as on a cluster.

The current workflow is based on the nf-core best practice. See the nf-core project from details on [guidelines](#).

Pipeline summary

This pipeline was designed to process Illumina sequencing data from the HPV capture protocol. Briefly, it allows to detect and genotype the HPV strain(s) available in the samples, and to precisely map the insertion sites on the Human genome.

1. Reads cleaning and quality controls ([TrimGalore!](#), [FastQC](#))
2. HPV Genotyping ([Bowtie2](#))
3. Local alignment on detected HPV strain(s) ([Bowtie2](#))
4. Detection of putative HPV breakpoints using soft-clipped reads
5. Soft-clipped reads alignment on Human genome reference ([BLAT](#))
6. Detection of insertion loci and filtering of the results
7. Presentation of results in a dynamic report ([MultiQC](#))

Outline

- Experimental and sequencing design
- Quality control and mapping
- Genes/transcripts/exons quantification
- Functional Analysis
- Fusion transcripts / chimera detection
- Single-cell RNA-seq
- Full-length transcripts sequencing using long reads
- Allele-specific quantification
- Expressed SNVs/indels variants detection
- Transcripts reconstruction (assembly)
- Alternative splicing / polyadenylation events detection
- Virus/phages expression detection
- **RNA editing events detection**

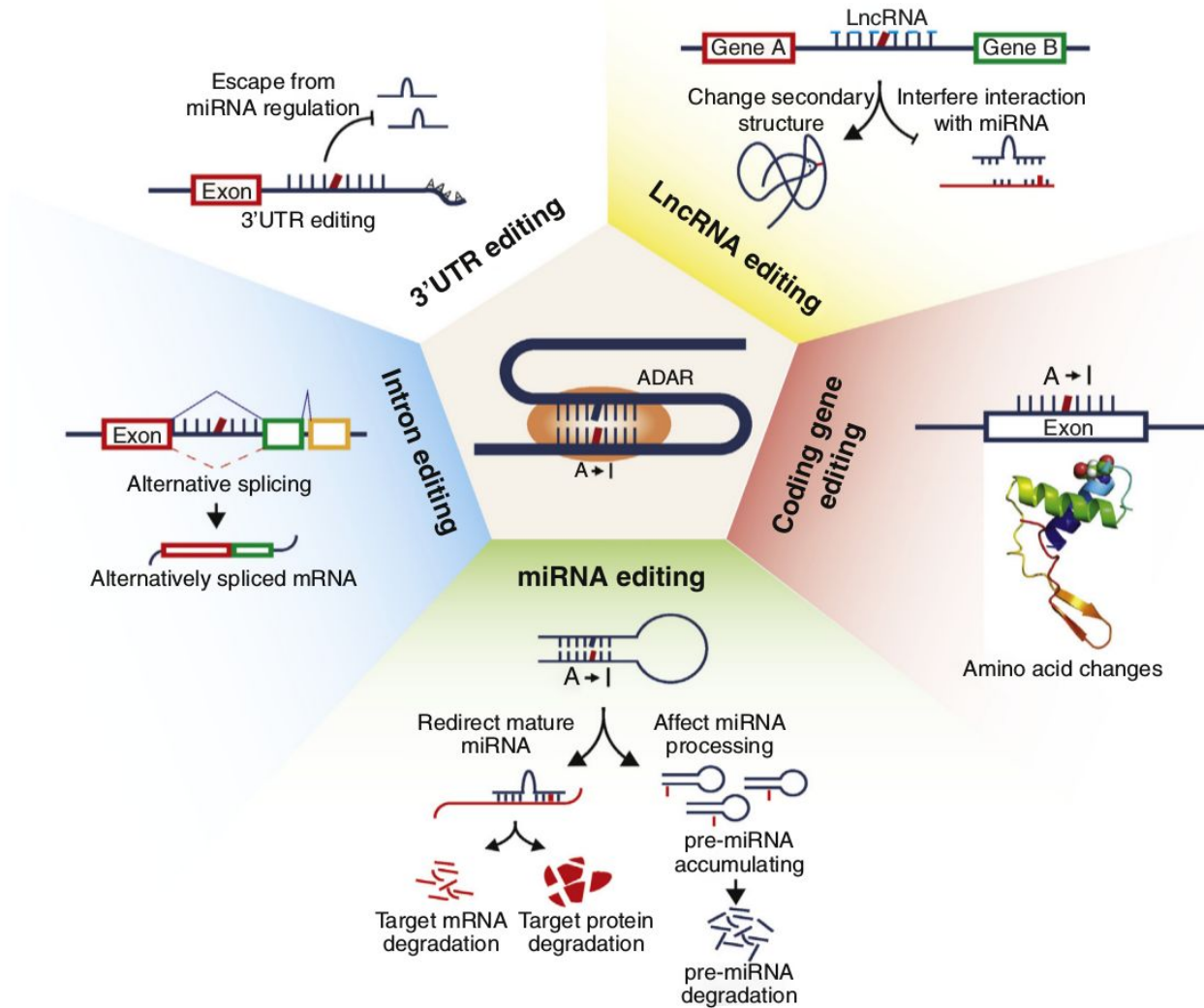
RNA editing

RNA editing (also RNA modification) is a molecular process through which some cells can make discrete changes to specific nucleotide sequences within an RNA molecule after it has been generated by RNA polymerase. It occurs in all living organisms and is one of the most evolutionarily conserved properties of RNAs. RNA editing may include the insertion, deletion, and base substitution of nucleotides within the RNA molecule. RNA editing is relatively rare, with common forms of RNA processing (e.g. splicing, 5'-capping, and 3'-polyadenylation) not usually considered as editing. It can affect the activity, localization as well as stability of RNAs, and has been linked with human diseases.

Type of changes :

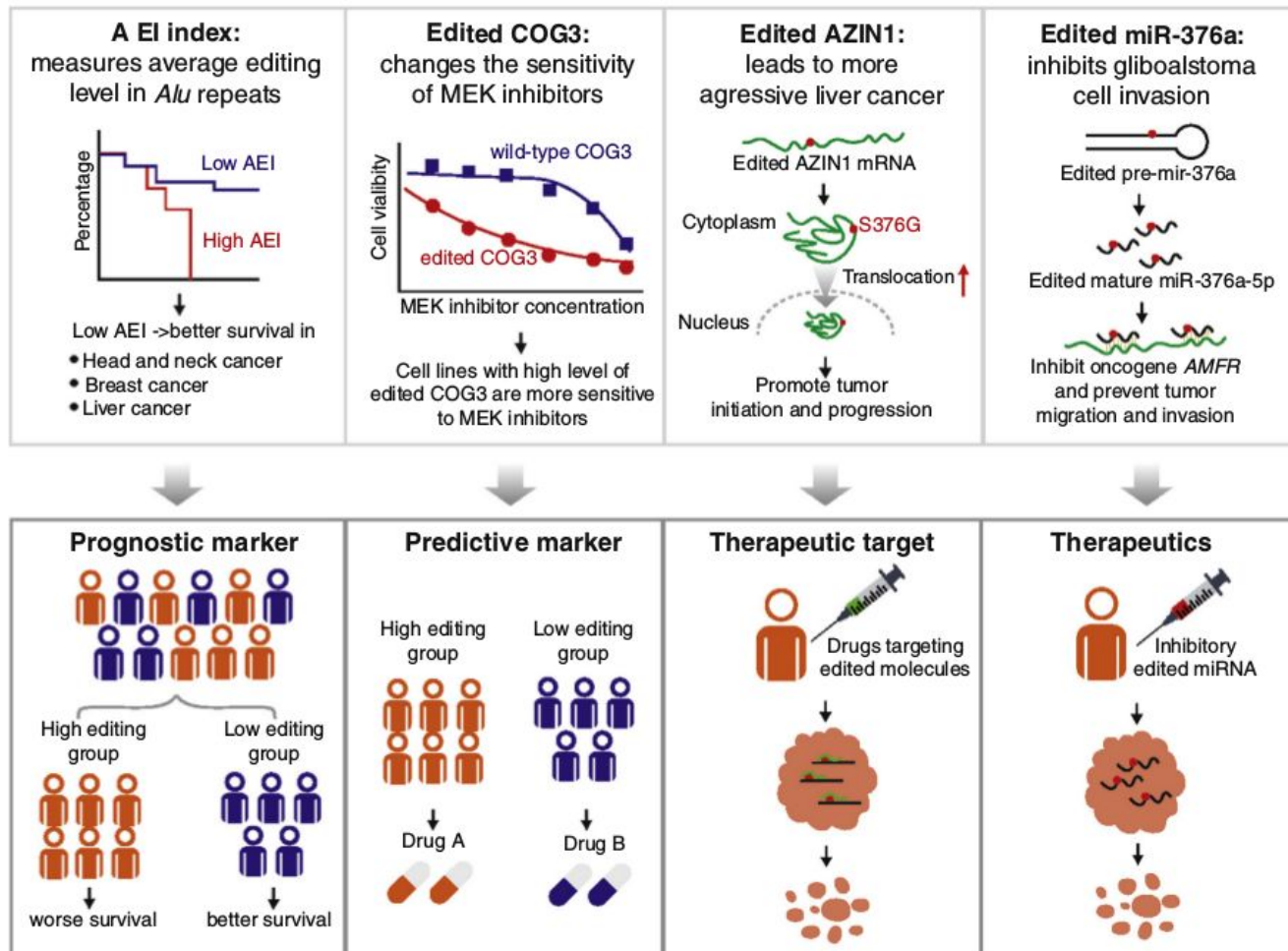
- Editing by insertion or deletion
- Editing by deamination (APOBEC1 and/or ADAR mediated)
 - C-to-U editing
 - A-to-I editing
 - Alternative mRNA editing (U-to-C and G-to-A)

RNA editing



RNA editing

Patients' RNA editing profile



RNA editing

Article | [Open Access](#) | Published: 10 March 2020

ADAR1-mediated RNA editing is a novel oncogenic process in thyroid cancer and regulates miR-200 activity

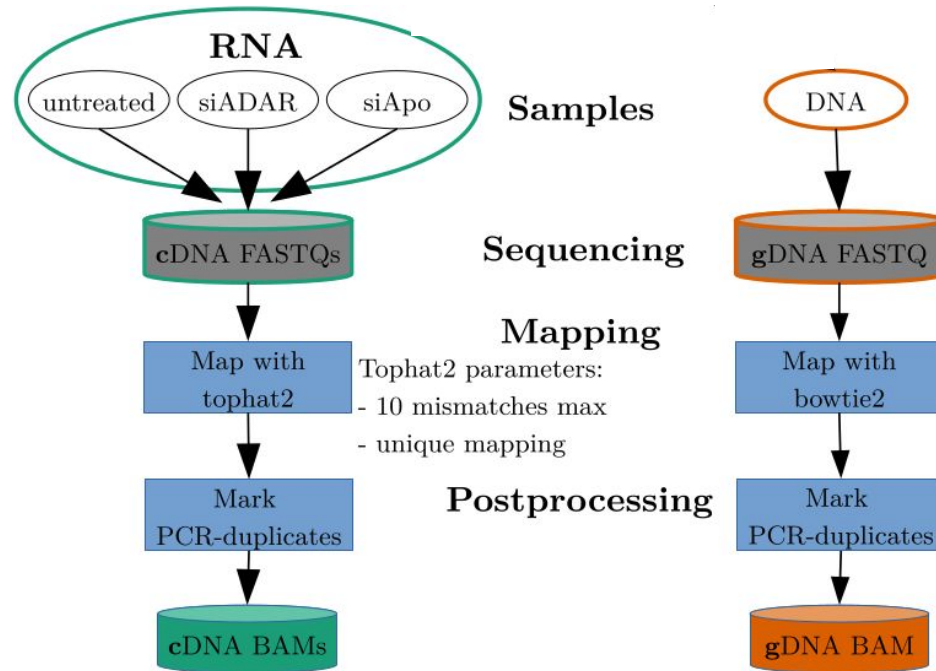
Julia Ramírez-Moya, Allison R. Baker, Frank J. Slack & Pilar Santisteban 

Oncogene **39**, 3738–3753(2020) | [Cite this article](#)

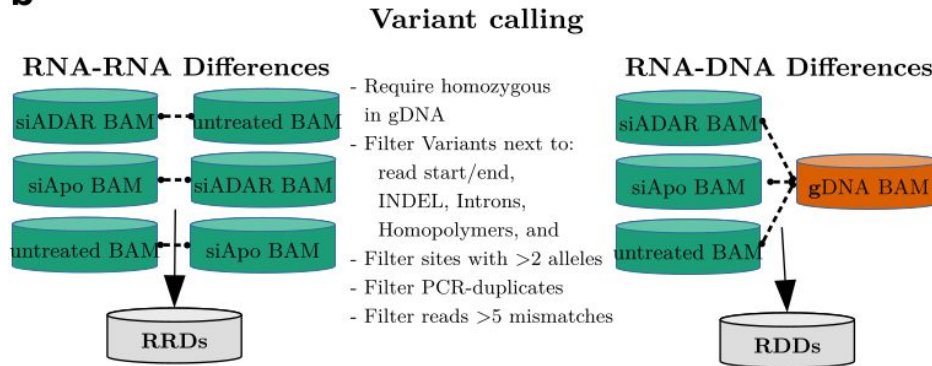
4530 Accesses | **8** Citations | **2** Altmetric | [Metrics](#)

<https://doi.org/10.1038/s41388-020-1248-x>

RNA editing



b



RNA editing

Table 1. Main features of tested tools for RNA editing detection

Name	GIREMI	JACUSA	REDIttools	RES-Scanner	RNAEditor
Required dependencies	HTSlib, SAMtools, R	R for JacusHelper	pysam, BLAT, SAMtools	Perl modules, BWA, SAMtools, BLAT	pysam, pyqt4, matplotlib, numpy, BWA, Picard Tools, GATK, BLAT, BEDtools
Required input files	BAM+ filtered SNVs	BAM	BAM	Fastq or BAM	Fastq
Stranded-oriented samples	Yes	Only single-end reads	Yes	Only FR-stranded (dUTP-protocol)	Yes, but no specified strand
RNA replicates accepted	Yes	Yes	No	No	No
<i>De novo</i> detection of RNA editing sites	Yes	Yes	Yes	Yes	Yes
DNA-RNA comparison	No	Yes	Yes	Yes	No
Works with RNA solo	Yes	Only with replicates	Yes	No	Yes
Multiple types of RNA editing	Yes	Yes	Yes	Yes	Yes
Mapping included	No	No	No	Optional (BWA)	Mandatory (BWA)
SNV calling	No	Yes	Yes (Samtools)	Yes (Samtools)	Yes (GATK)
Annotation of editing sites	No	No	Yes	Yes	Yes
Filtering or P-value provided	P-value	Filtering	Filtering	Filtering+P-value	Filtering

