

# Mapping

Isheng Jason Tsai

Introduction to NGS Data and Analysis  
Lecture 3 v2020

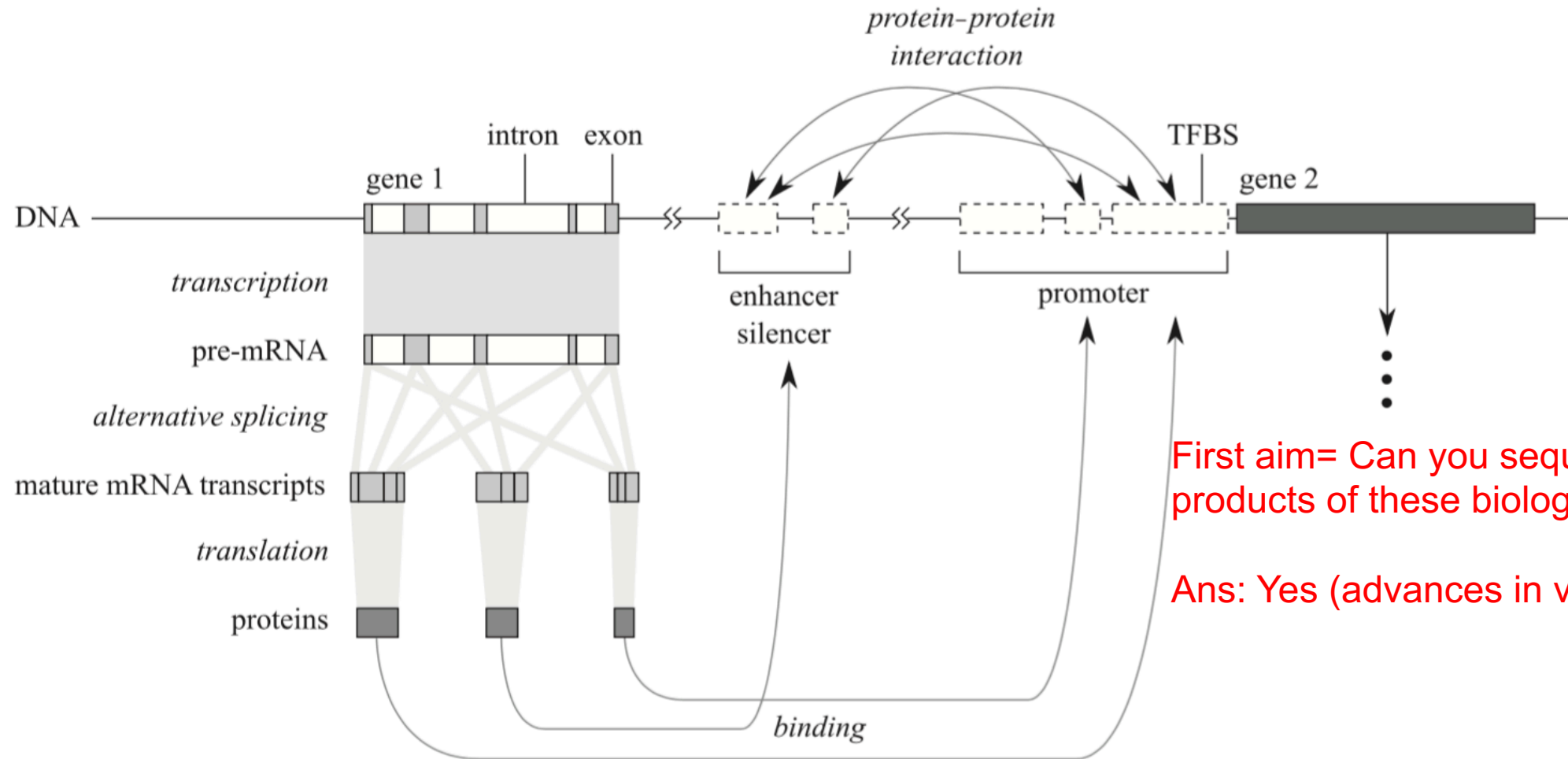


# Lecture outline

1. Background
2. Mapping algorithms
3. Mapping processes
4. Variant calling

Background

# High throughput sequencing applications

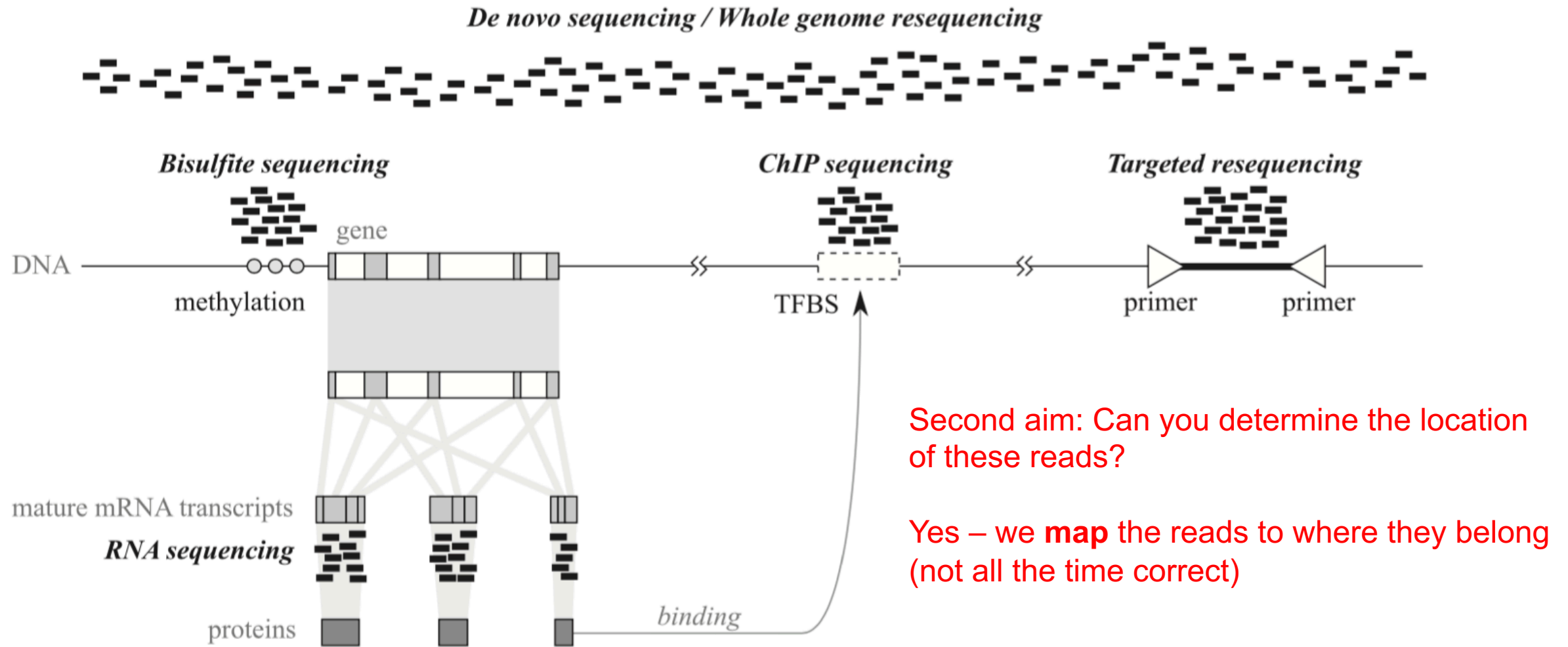


First aim= Can you sequence the products of these biological processes?

Ans: Yes (advances in various experiments)

**Figure 1.1** A schematic illustration of the central dogma. Gene 1 has three alternatively spliced transcripts. The relative expression of such transcripts affects the regulatory modules of gene 2, and eventually its expression. Definitions are given in Section 1.1.

# High throughput sequencing applications



**Figure 1.2** A schematic summary of high-throughput sequencing applications. Details are described in Section 1.3.

# *De novo* vs mapping approach

**Mapping** is less complicated and more intuitive

Can gather lots of information from many individuals given a good reference

But, information on repeats/ gene families / *de novo* genes / large structural variants are more difficult to detect

**Assembly** is powerful but also computationally demanding

And is your question worth the trouble to assemble 100 strains?

In practice, people do a combination of both approaches

In humans, *de novo* genomes of references and cancer cells are being generated. In butterflies, many assemblies to reveal super gene

# Recommended paper

## Technology dictates algorithms: Recent developments in read alignment

Mohammed Alser, Jeremy Rotman, Kodi Taraszka, Huwenbo Shi, Pelin Icer Baykal, Harry Taegyun Yang, Victor Xue, Sergey Knyazev, Benjamin D. Singer, Brunilda Balliu, David Koslicki, Pavel Skums, Alex Zelikovsky, Can Alkan, Onur Mutlu, Serghei Mangul

*(Submitted on 28 Feb 2020)*

”...Our review provides a survey of **algorithmic foundations and methodologies across alignment methods for both short and long reads**. We provide rigorous experimental evaluation of 11 read aligners to demonstrate the effect of these underlying algorithms on speed and efficiency of read aligners. We separately **discuss how longer read lengths produce unique advantages and limitations to read alignment techniques**. We also **discuss how general alignment algorithms have been tailored to the specific needs of various domains in biology, including whole transcriptome, adaptive immune repertoire, and human microbiome studies**. “

<https://arxiv.org/ftp/arxiv/papers/2003/2003.00110.pdf>

Note: a preprint

# Preface



**Nick Loman** @pathogenomenick · Mar 11

Got a talk at ECCMID entitled: "So you've sequenced your (bug) genome ... what now?" Crowdsourcing best answers please, will acknowledge!



RETWEETS

7

FAVORITES

2



11:56 AM - 11 Mar 2015 · Details

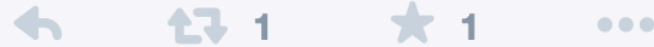


# We all know...



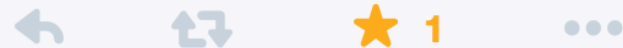
**Alan McNally** @alanmcn1 · Mar 11

@pathogenomenick @biomickwatson in that case "give it to someone who knows what they are doing!"



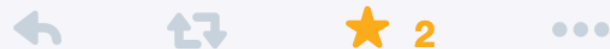
**Nicki Fawcett** @DrNJFawcett · Mar 11

@alanmcn1 @pathogenomenick Clinician thirthing/fourthing 'Give it to someone who knows what they're doing'. #ooohYersiniaInEverything



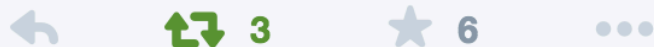
**Mick Watson** @BioMickWatson · Mar 11

@pathogenomenick ah. Clinician clinicians? Give the data to someone who knows what to do with it, then ;-)

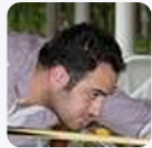


**azizipeasie** @AzizAboobaker · Mar 11

@pathogenomenick send it to your bioinformation friend and give them a week to send back a paper with themselves as a middle author.

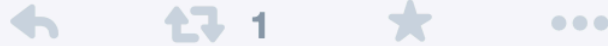


# Logical answer



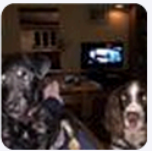
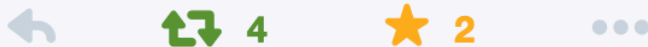
**azizipeasie** @AzizAboobaker · Mar 11

@pathogenomenick sequence some more while your thinking.



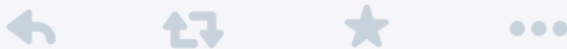
**Esther Robinson** @ilovechocagar · Mar 11

@pathogenomenick first law of doing a lab test: don't unless you know what your question is



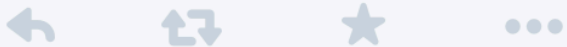
**ruth massey** @bowsermassey · Mar 11

@WvSchaik @pathogenomenick determine ID, resistance profile and dare I say it....virulence potential!

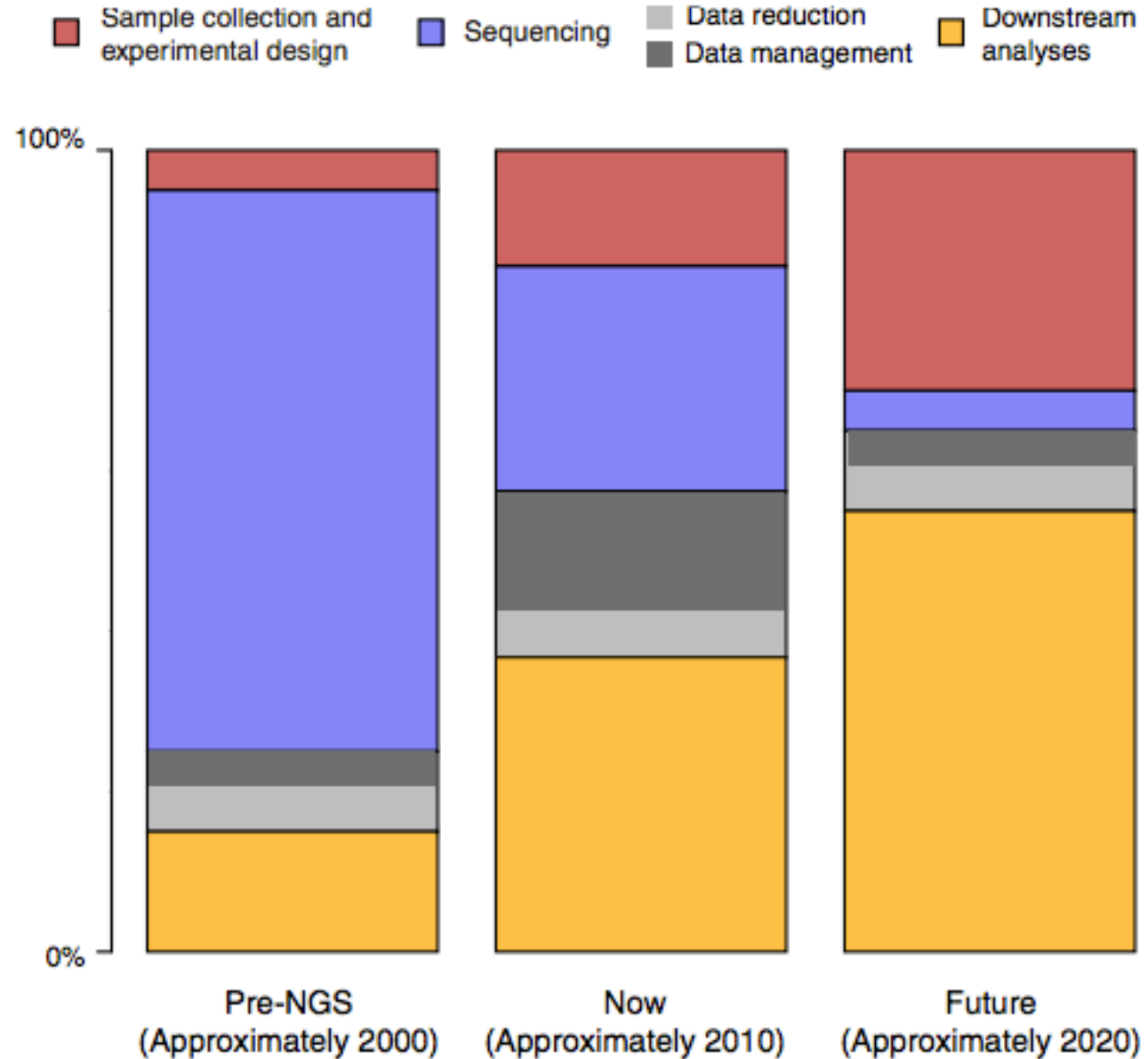
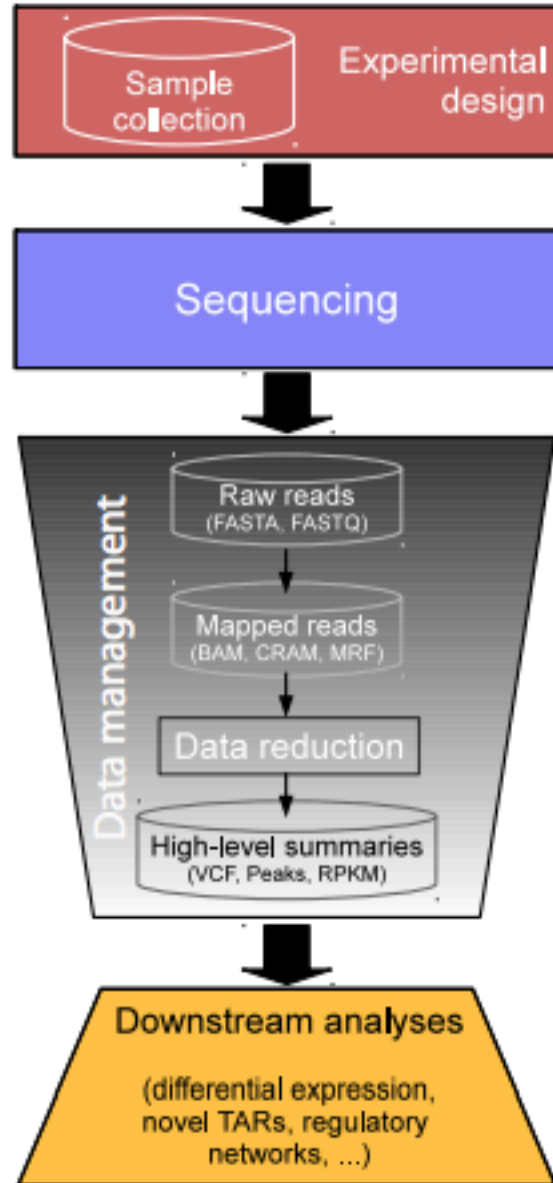


**Bill Hanage** @BillHanage · Mar 11

@pathogenomenick you've had many good suggestions but it completely depends on what you are interested in. Resistance? Epi? Something else?

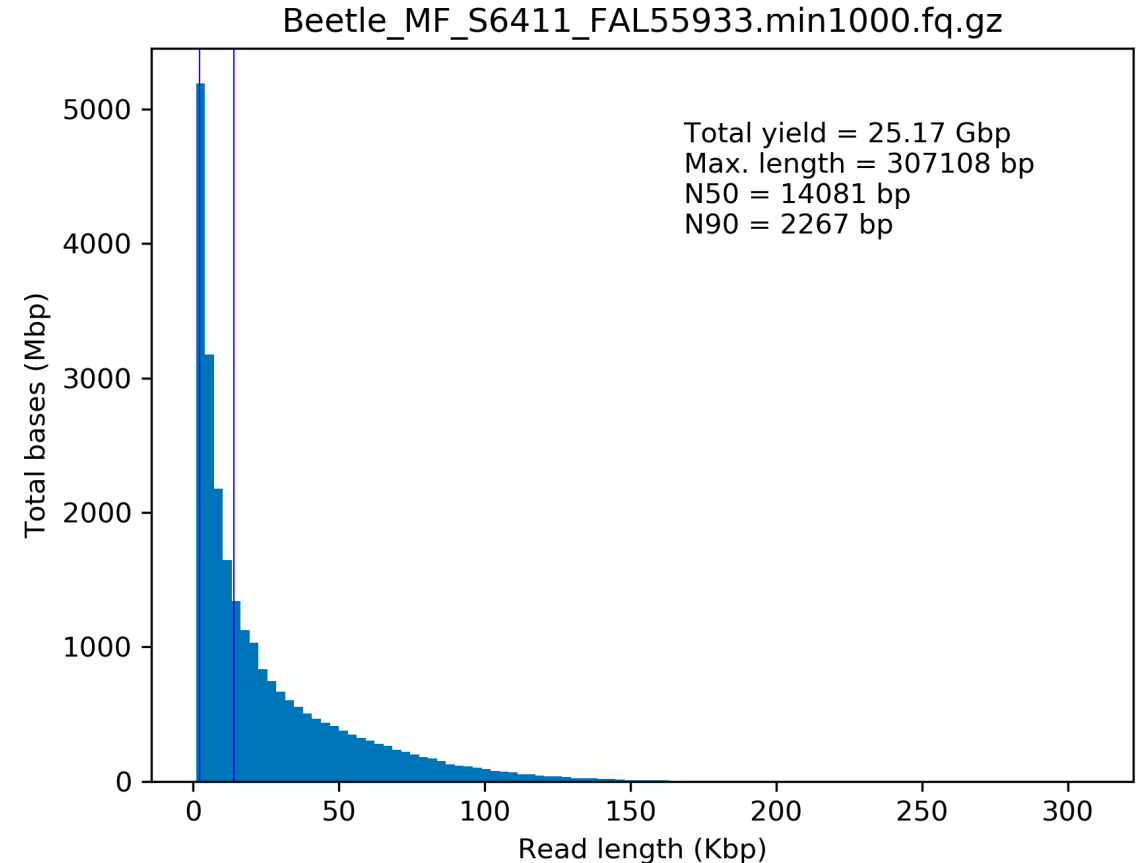


# The real cost of sequencing



# Long reads are now common

- Most users should and will be mainly analyzing Paired-end Illumina reads (typically 150 bases)
- Pacific Biosciences or Oxford Nanopore (long reads) are increasingly very common



# Mapping

Mapping is **aligning** the read to where **the most likely origin** within the reference/assembly

Sequence alignment has not changed and will remain a classic problem  
Tradeoffs of speed, accuracy and sensitivity

## **Sequence data we want to map:**

- Mostly nucleotide

**Very short evolutionary distances** (human to reference, isolate/strain to reference, 'slightly diverged' strain will result in less mapped reads)

**A lot of reads**— needs faster processing per read (BLAST is too slow!)

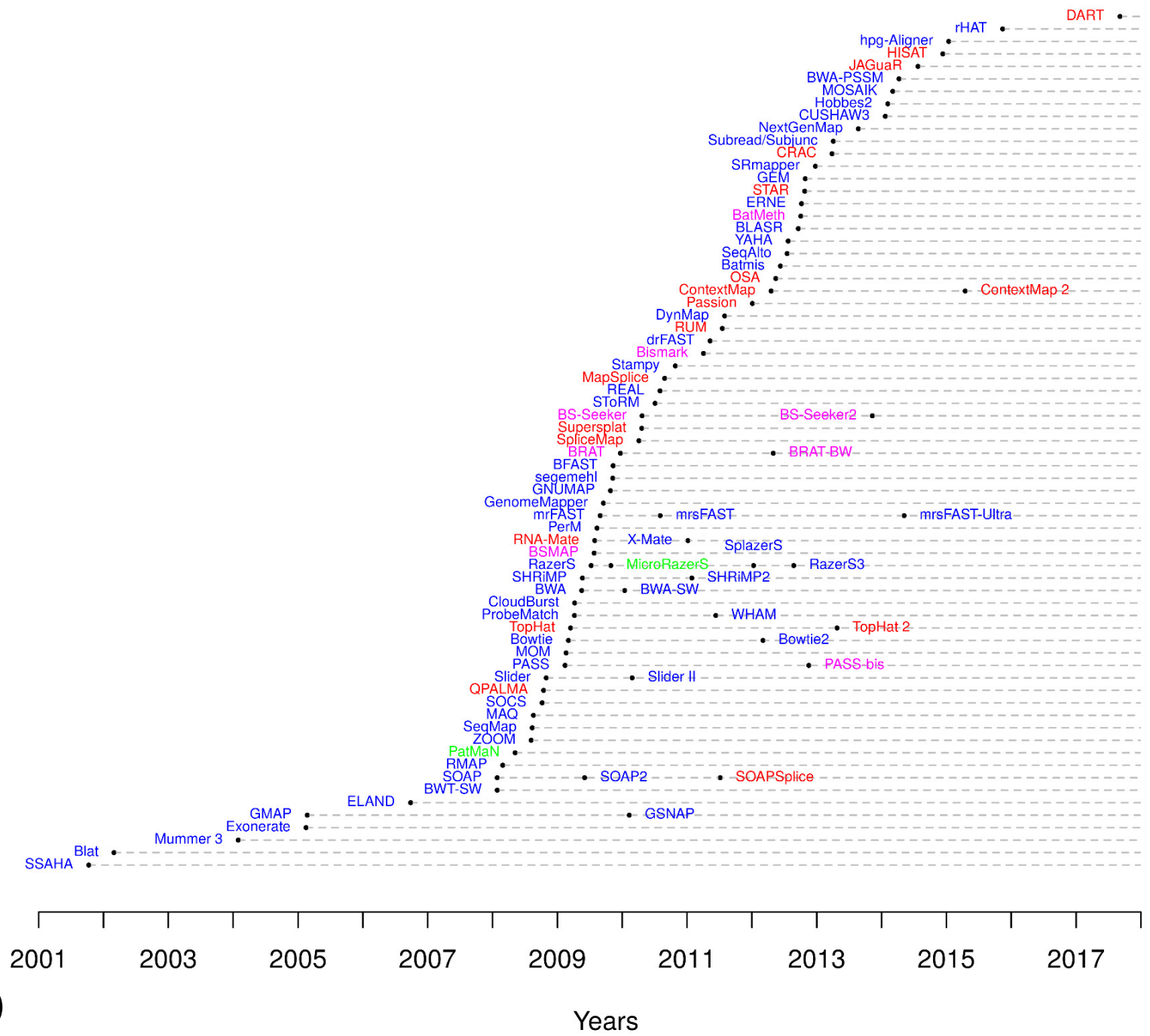
There are some assumptions to make alignment process faster  
(like allows most 2 mismatches)

# How to map billions of short reads onto genomes

Cole Trapnell & Steven L Salzberg

**Mapping the vast quantities of short sequence fragments produced by next-generation sequencing platforms is a challenge. What programs are available and how do they work?**

# All the mappers!



[http://www.ebi.ac.uk/~nf/hts\\_mappers/](http://www.ebi.ac.uk/~nf/hts_mappers/)  
(link no longer working; last updated 2018)

Mapper	Data	Availability	Version	O.S.	Number Citations	Seq.Plat.	Input	Output
<b>BWA</b>	DNA	OS	0.6.2	Linux,Mac,Windows	<b>13341</b>	I,So,4,Sa,P	FASTA/Q	SAM
<b>Bowtie</b>	DNA	OS	0.12.7	Linux,Mac,Windows	<b>11207</b>	I,So,4,Sa,P	(C)FAST(A/Q)	SAM TSV
<b>Bowtie2</b>	DNA	OS	2.0beta5	Linux,Mac,Windows	<b>8586</b>	I,4,Ion	FASTA/Q	SAM TSV
<b>Blat</b>	DNA	OS	34	Linux,Mac	<b>6252</b>	N	FASTA	TSV BLAST
<b>TopHat</b>	RNA	OS	1.4.1	Linux,Mac	<b>3764</b>	I	FASTA/Q GFF	BAM
<b>BWA-SW</b>	DNA	OS	0.6.2	Linux,Mac,Windows	<b>3494</b>	I,4,Sa,Hel,Ion,P	FASTA/Q	SAM
<b>MAQ</b>	DNA	OS	0.7.1	Linux,Mac	<b>2592</b>	I,So	(C)FAST(A/Q)	TSV
<b>Mummer 3</b>	DNA	OS	3.23	Linux,Mac	<b>2446</b>	N	FASTA	TSV
<b>SOAP2</b>	DNA	OS	2.21	Linux	<b>1655</b>	I	FASTA/Q	SAM TSV
<b>SOAP</b>	DNA	OS	1.11	Linux,Mac	<b>1284</b>	I	FASTA/Q	TSV
<b>GSNAP</b>	DNA	OS	2012-04-27	Linux,Unix,Mac,Windows	<b>1156</b>	I,4,Sa,Hel,Ion,P	FASTA/Q	SAM Native
<b>TopHat 2</b>	RNA	OS	2.0.8	Linux,Mac	<b>1102</b>	I	FASTA/Q	BAM
<b>Exonerate</b>	DNA	OS	2.2	Linux,Mac	<b>918</b>	N	FASTA	TSV
<b>Bismark</b>	Bisulfite	OS	0.7.3	Linux,Mac	<b>887</b>	I	FASTA/Q	SAM
<b>SSAHA2</b>	DNA	Bin	2.5.5	Linux,Mac	<b>874</b>	I,4,Sa	FASTA/Q	SAM
<b>SSAHA</b>	DNA	OS	3.1	Linux,Mac	<b>874</b>	N	FASTA/Q	TSV
<b>GMAP</b>	DNA	OS	2012-04-27	Linux,Unix,Mac,Windows	<b>868</b>	I,4,Sa,Hel,Ion,P	FASTA/Q	SAM GFF Native
<b>CloudBurst</b>	DNA	OS	1.1	Linux,Mac,Windows	<b>650</b>	N	FASTA	TSV
<b>MapSplice</b>	RNA	OS	1.15.2	Linux	<b>610</b>	I	FASTA/Q	SAM BED
<b>STAR</b>	RNA	OS	2.3.0	Linux,Unix,Mac	<b>602</b>	I,4,Sa,Ion,P	FASTA/Q	SAM
<b>mrFAST</b>	DNA	OS	2.5.0.1	Linux,Unix	<b>602</b>	I	FASTA/Q	SAM DIVET
<b>SHRIMP</b>	DNA	OS	1.3.2	Linux,Mac	<b>573</b>	I,So,4,Hel	(C)FAST(A/Q)	TSV
<b>BFAST</b>	DNA	OS	0.7.0	Linux,Mac	<b>553</b>	I,So,4, Hel	(C)FAST(A/Q)	SAM TSV
<b>HISAT</b>	RNA	OS	1	Windows, Linux, Unix, Mac	<b>480</b>	I	FASTA/Q	SAM



# Heng Li has contributed a dozen of those mappers



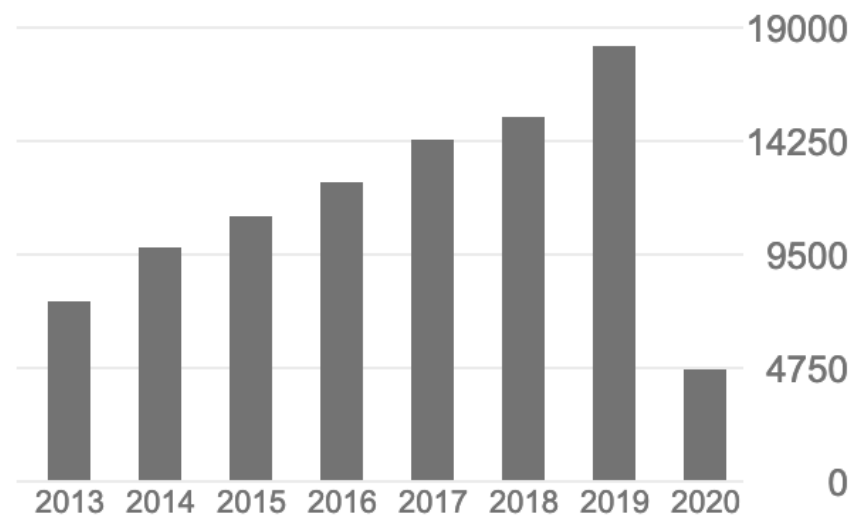
Heng Li

Dana-Farber Cancer Institute & [Harvard University](#)  
在 [jimmy.harvard.edu](mailto:jimmy.harvard.edu) 的電子郵件地址已通過驗證 - [首頁](#)  
[Computational Biology](#) [Bioinformatics](#) [Genomics](#)

Author of maq, bwa, bwa-mem, wtdbg2  
and minimap2...

<http://www.liheng.org/>

	全部	自 2015 年
引文	108019	76196
H 指數	53	50
i10 指數	68	62



# How?

Brute force comparison

Smith-Waterman

Suffix Tree

Burrows-Wheeler Transform

# Brute force

TCGATCC  
    ↘  
GACCTCATCGATCCACTG

1.

TCGATCC  
X  
GACCTCATCGATCCACTG

2.

TCGATCC  
X  
GACCTCATCGATCCACTG

3.

TCGATCC  
| | X  
GACCTCATCGATCCACTG

4.

TCGATCC  
| | | | | | |  
GACCTCATCGATCCACTG

# Exact matching

What's a simple algorithm for exact matching?

*P*: word

*T*: There would have been a time for such a word

word word word word word word word word word word


word word word word word word word word

word word word word word word word word

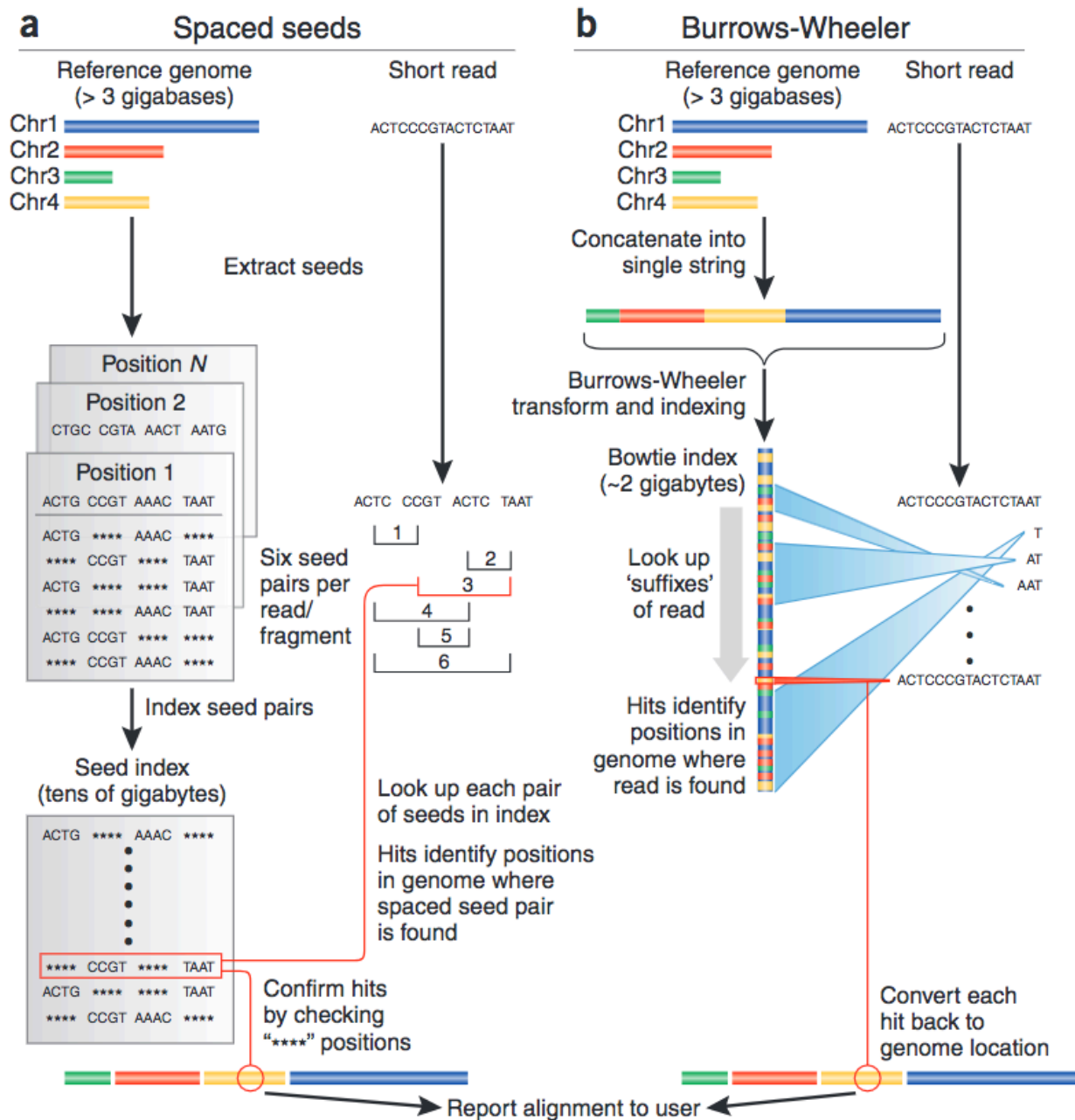
word word word word word word word word

word word word word word word word word

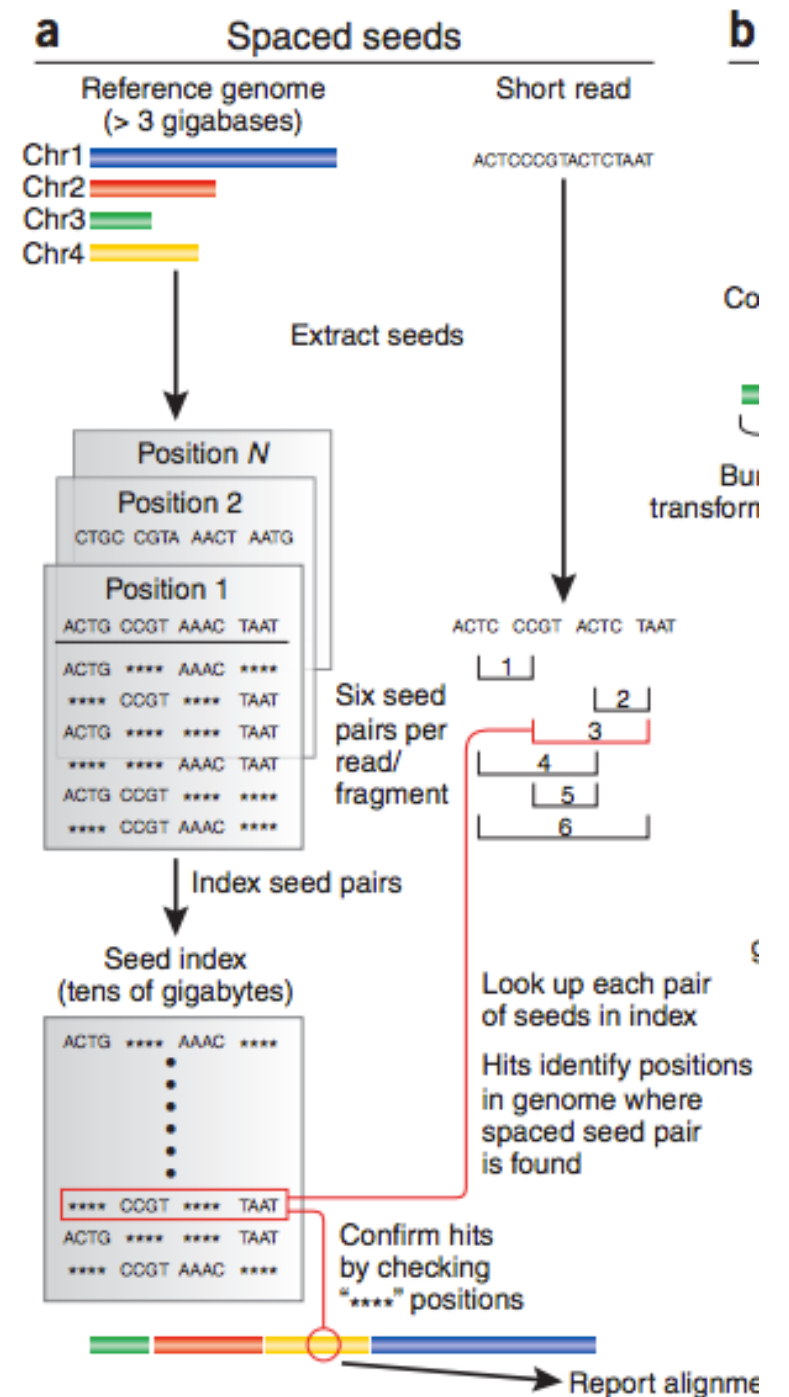
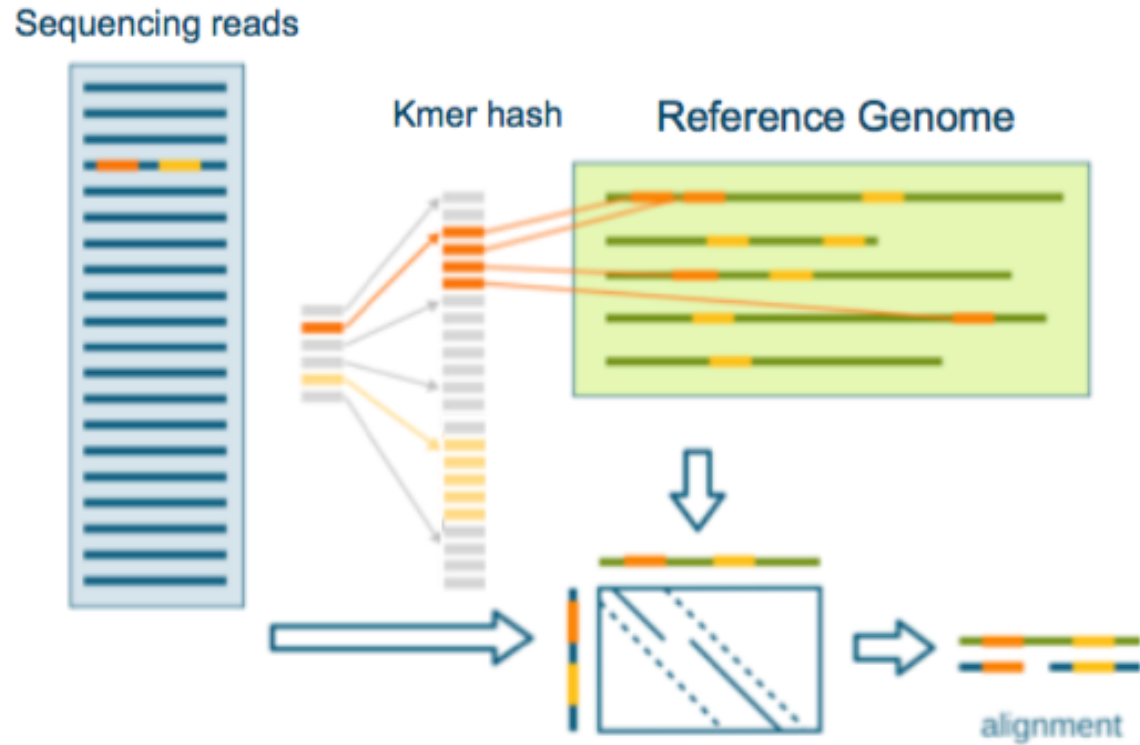
One occurrence



Try all possible alignments. For each, check whether it's an occurrence. "Naïve algorithm."



# Mapping (hash table)



- Identify all the seeds in the index
- Determine the most likely location
- Perform Smith-Waterman alignment to fully align
- Output (important)

Example: BLAST, MAQ (Heng Li 2008)

# Suffix tree

GACCTCATCGATCCCACTG

A				C							G		T					
C	T			A	C	G	T		A	\$		C		G				
C	T	C		C	T	A	C	T	A	C	G	C	T		A	C	G	\$
T	G	C	G	T	C	C	A	C	T	A	\$	C	C		T	C	A	
C	\$	C	A	G	G	T	C	A	C	T		T	C		C	A	T	
A		A	T	\$	A	G	T	T	C	C		C	C		G	C	C	
T		C	C		T	\$	G	C	C	G		A	A		A	T	C	
C		T	C		C		\$	G	A	A		T	C		T	G	C	
G		G	C		C			A	C	T		C	T		C	\$	A	
A		\$	A		C			T	T	C		G	G		C		C	
T			C		A			C	G	C		A	\$		C		T	
C			T		C			C	\$	C		T			A		G	
C			G		T			C		A		C		C		\$		
C			\$		G			A		C		C		T				
A					\$			C		T		C		G				
C								T		G		A		\$				
T								G		\$		C						
G								\$				T						
\$												G						
												\$						

But suffix can be very very big if data structure not considered carefully!

# Burrows-Wheeler Transform

A transformation that will result in many repeated characters

This means it's easy to compress

And an elegant way to search!

Transformation				
Input	All Rotations	Sorting All Rows into Lex Order	Taking Last Column	Output Last Column
<code>^BANANA  </code>	<code>^BANANA  </code> <code>  ^BANANA</code> <code>A   ^BANAN</code> <code>NA   ^BANA</code> <code>ANA   ^BAN</code> <code>NANA   ^BA</code> <code>ANANA   ^B</code> <code>BANANA   ^</code>	<code>ANANA   ^B</code> <code>ANA   ^BAN</code> <code>A   ^BANAN</code> <code>BANANA   ^</code> <code>NANA   ^BA</code> <code>NA   ^BANA</code> <code>^BANANA  </code> <code>  ^BANANA</code>	<code>ANANA   ^B</code> <code>ANA   ^BAN</code> <code>A   ^BANAN</code> <code>BANANA   ^</code> <code>NANA   ^BA</code> <code>NA   ^BANA</code> <code>^BANANA  </code> <code>  ^BANANA</code>	<code>BNN^AA   A</code>



Original sequence	All permutations	Alphabetical ordering of rows	Output of last column
>BONOBO*	>BONOBO* *>BONOBO O*>BONOB BO*>BONO OBO*>BON NOBO*>BO ONOBO*>B BONOBO*>	BONOBO*> BO*>BONO NOBO*>BO OBO*>BON ONOBO*>B O*>BONOB >BONOBO* *>BONOBO	> O O N B B * O

Inverse transformation using Burrows-Wheeler transform

Add cycle 1	Sort cycle 1	Add cycle 2	Sort cycle 2
>	B	>B	BO
O	B	OB	BO
O	N	ON	NO
N	O	NO	OB
B	O	BO	ON
B	O	BO	O*
*	>	*>	>B
O	*	O*	*>
Add cycle 3	Sort cycle 3	Add cycle 4	Sort cycle 4
>BO	BON	>BON	BONO
OBO	BO*	OBO*	BO*>
ONO	NOB	ONOB	NOBO
NOB	OBO	NOBO	OBO*
BON	ONO	BONO	ONOB
BO*	O*>	BO*>	O*>B
*>B	>BO	*>BO	>BON
O*>	*>B	O*>B	*>BO

GACCTCATCGATCCCACTG\$  
ACCTCATCGATCCCACTG\$G  
CCTCATCGATCCCACTG\$GA  
CTCATCGATCCCACTG\$GAC  
TCATCGATCCCACTG\$GACC  
CATCGATCCCACTG\$GACCT  
ATCGATCCCACTG\$GACCTC  
TCGATCCCACTG\$GACCTCA  
CGATCCCACTG\$GACCTCAT  
GATCCCACTG\$GACCTCATC  
ATCCCACTG\$GACCTCATCG  
TCCCACTG\$GACCTCATCGA  
CCCACTG\$GACCTCATCGAT  
CCACTG\$GACCTCATCGATC  
CACTG\$GACCTCATCGATCC  
ACTG\$GACCTCATCGATCCC  
CTG\$GACCTCATCGATCCCA  
TG\$GACCTCATCGATCCCAC  
G\$GACCTCATCGATCCCACT  
\$GACCTCATCGATCCCACTG

Sort →

ACCTCATCGATCCCACTG\$G  
ACTG\$GACCTCATCGATCCC  
ATCCCACTG\$GACCTCATCG  
ATCGATCCCACTG\$GACCTC  
CACTG\$GACCTCATCGATCC  
CATCGATCCCACTG\$GACCT  
CCACTG\$GACCTCATCGATC  
CCCACTG\$GACCTCATCGAT  
CCTCATCGATCCCACTG\$GA  
CGATCCCACTG\$GACCTCAT  
CTCATCGATCCCACTG\$GAC  
CTG\$GACCTCATCGATCCCA  
GACCTCATCGATCCCACTG\$  
GATCCCACTG\$GACCTCATC  
G\$GACCTCATCGATCCCACT  
TCATCGATCCCACTG\$GACC  
TCCCACTG\$GACCTCATCGA  
TCGATCCCACTG\$GACCTCA  
TG\$GACCTCATCGATCCCAC  
\$GACCTCATCGATCCCACTG



TCGATCC  
↓ ?  
GACCTCATCGATCCCACTG

GAC  
CAC  
GAT  
CAT  
CCA  
→ TCA  
CCC  
→ TCC  
ACC  
→ TCG  
CCT  
ACT  
\$GA  
CGA  
→ TGS  
CTC  
ATC  
ATC  
CTG  
G\$G

- Start with the transform column
- My read starts with a T, so I want rows with Ts in them
- This column gives me all the single nucleotide counts
- Sort the single nucleotide counts to get the alphabetically first column
- Now these two columns give me all the dinucleotide counts
- Sort those to get the alphabetically first two columns
- Now there is only one place my read can match

# BWT – a summary

Stores all possible suffixes to enable fast string matching

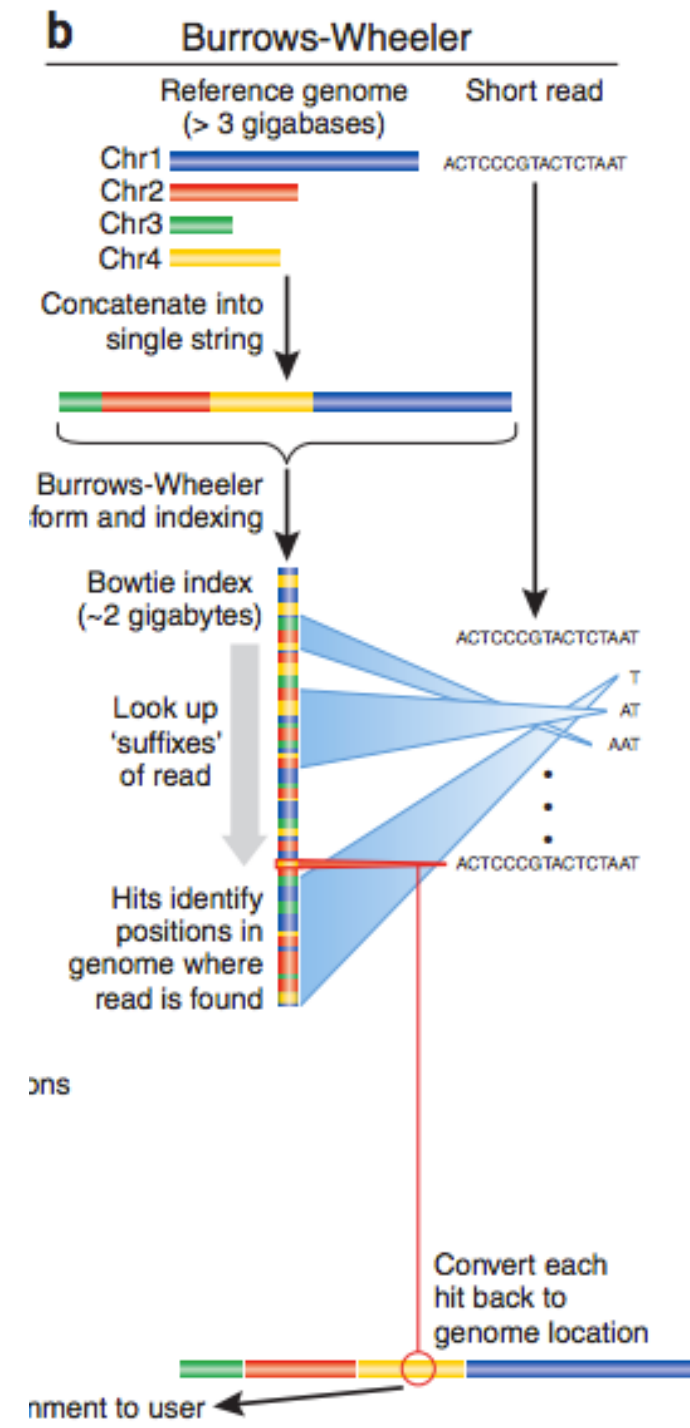
Much smaller memory footprint than hash table  
(hash table need to store all different kmers)

Examples:

MUMMER, bwa, bowtie2

Still need local alignment in final step

Trapnell *et al* (2009)



# Hash table vs. BWT

**Table 1 A selection of short-read analysis software**

Program	Website	Open source?	Handles ABI color space?	Maximum read length
Bowtie	<a href="http://bowtie.cbcb.umd.edu">http://bowtie.cbcb.umd.edu</a>	Yes	No	None
BWA	<a href="http://maq.sourceforge.net/bwa-man.shtml">http://maq.sourceforge.net/bwa-man.shtml</a>	Yes	Yes	None
Maq	<a href="http://maq.sourceforge.net">http://maq.sourceforge.net</a>	Yes	Yes	127
<del>Mosaik</del>	<del><a href="http://bioinformatics.bc.edu/marthlab/Mosaik">http://bioinformatics.bc.edu/marthlab/Mosaik</a></del>	<del>No</del>	<del>Yes</del>	<del>None</del>
<del>Novoalign</del>	<del><a href="http://www.novocraft.com">http://www.novocraft.com</a></del>	<del>No</del>	<del>No</del>	<del>None</del>
<del>SOAP2</del>	<del><a href="http://soap.genomics.org.cn">http://soap.genomics.org.cn</a></del>	<del>No</del>	<del>No</del>	<del>60</del>
<del>ZOOM</del>	<del><a href="http://www.bioinfor.com">http://www.bioinfor.com</a></del>	<del>No</del>	<del>Yes</del>	<del>240</del>

BWT

Hash table

# Hash table vs. BWT strengths and weaknesses

## **Burrows-Wheeler**, e.g. bwa, bowtie

- Fast, esp. (multiple) exact matches
- High sensitivity at repetitive regions
- less robust at high genomic variation (because you need to retry with a substitution)

## **Hashing** (overlapping k-mer words, e.g. SMALT, Stampy)

- Slower (more memory hungry)
- Less sensitivity at repetitive regions
- tolerate high genomic variation
- partial alignments (junction reads) easier
- Flexible (multiple sequencing platforms)

# Choose an mapper/ aligner

Hash based approaches are more suitable for divergent alignments

General rule:

<2% divergence -> BWT

E.g. human samples

>2% divergence -> hash based approach

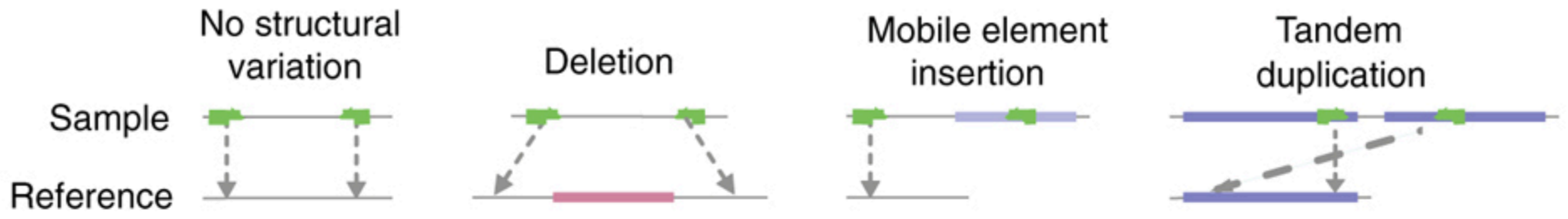
E.g. wild sample alignments ;

Watch out for latest advancement ; and don't stay at one for too long

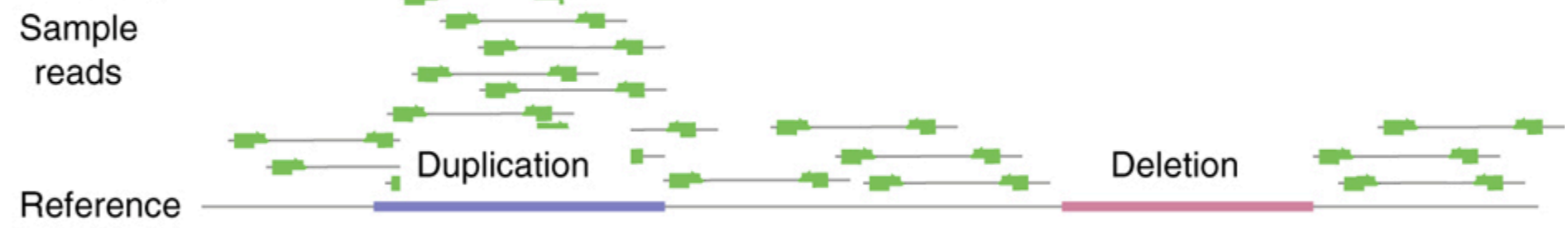


# Detecting structural variations (ideally assembly is probably better)

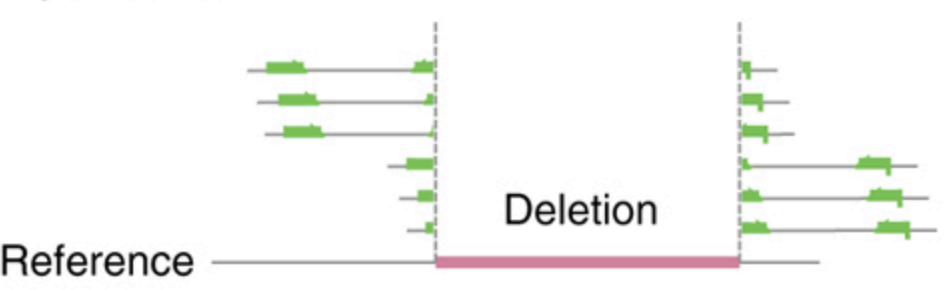
Read pairs



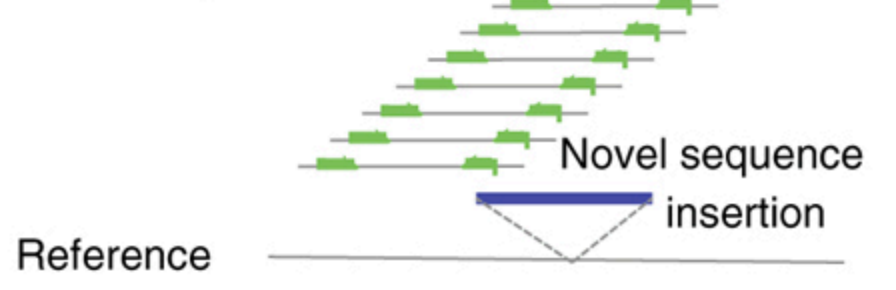
Read depth



Split reads



Assembly



# What to do with repetitive (multi) reads?

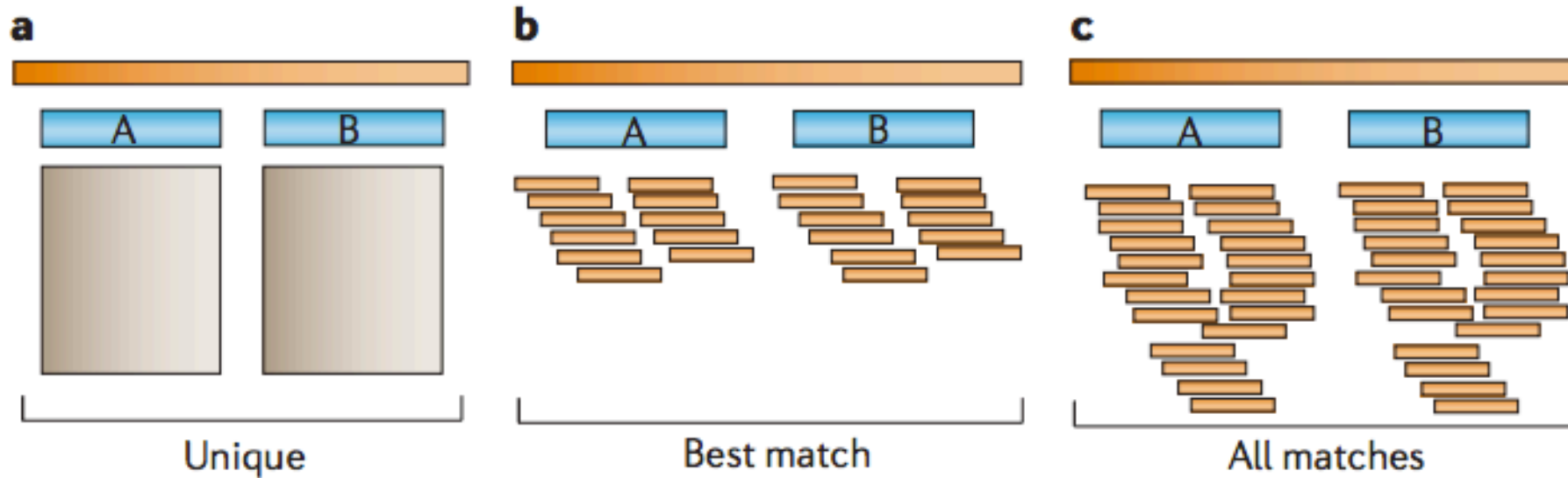
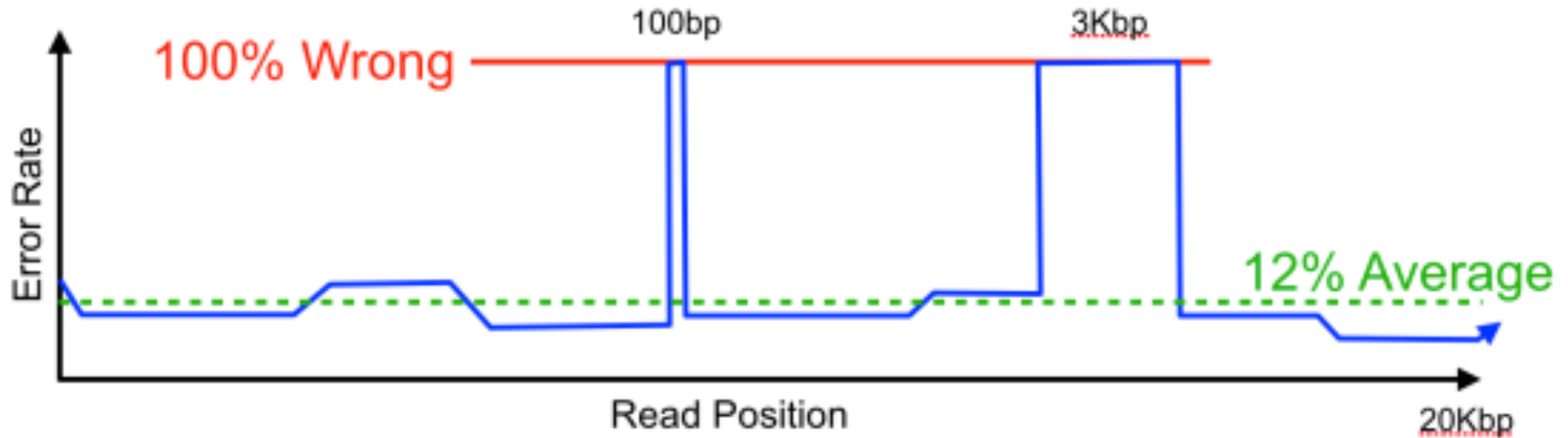


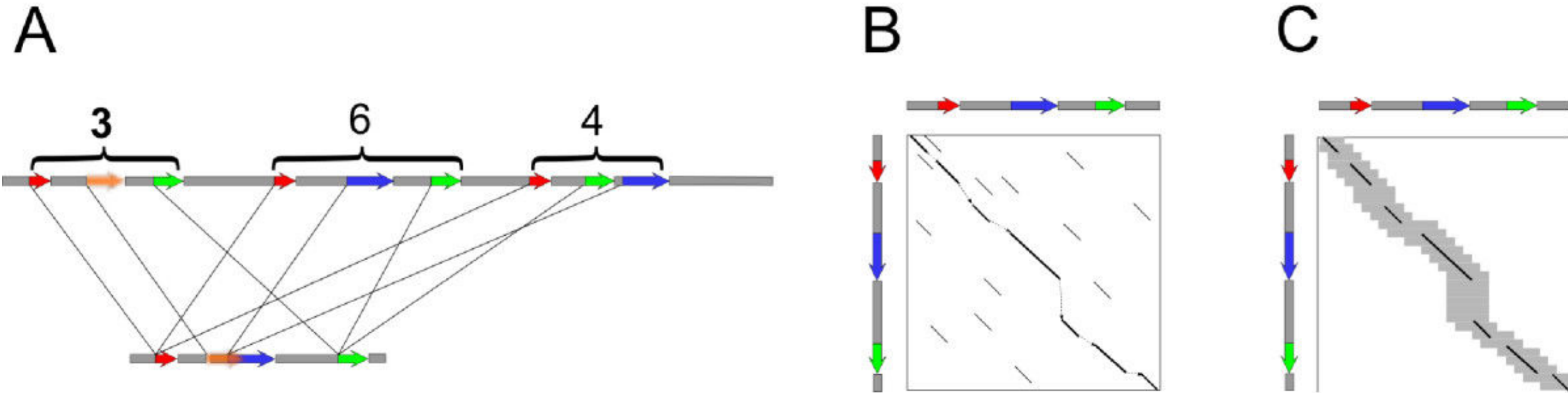
Figure 2 | **Three strategies for mapping multi-reads.** The shaded rectangles at the top represent intervals along a chromosome. The two blue rectangles below each region represent an identical two-copy repeat containing the paralogous genes A and B. The small orange bars represent reads aligned to specific positions. **a** | The ‘unique’ strategy reports only those reads that are uniquely mappable. Because A and B are identical, no alignments are reported. **b** | The ‘best match’ alignment strategy reports the best possible alignment for each read, which is determined by the scoring function of the alignment algorithm. In the case of ties, this strategy randomly distributes reads across equally good loci, as shown here. **c** | The ‘all matches’ strategy simply reports all alignments for each multi-read, including lower-scoring alignments.

# What about long read mapping?



- **BLASR** and **Daligner** designed for long error-prone (but random) reads (PacBio)
- Now there's alternative such as GMAP and minimap2 which are much faster

# What about long read mapping?



- **BLASR**

- One of the first tools in long read alignments (meaning it's easier to understand)
- Combines multiple methods
- Starts by finding short exact matches using suffix or B-W
- Next locally identifies a linear chain of shorter exact matches
- Performs banded Smith-Waterman constrained by the shorter exact matches

# Minimap2 – the most popular tool to use in long read alignment

1. Read `-I [=4G]` reference bases, extract `(-k,-w)`-minimizers and index them in a hash table.
2. Read `-K [=200M]` query bases. For each query sequence, do step 3 through 7:
3. For each `(-k,-w)`-minimizer on the query, check against the reference index. If a reference minimizer is not among the top `-f [=2e-4]` most frequent, collect its the occurrences in the reference, which are called *seeds*.
4. Sort seeds by position in the reference. Chain them with dynamic programming. Each chain represents a potential mapping. For read overlapping, report all chains and then go to step 8. For reference mapping, do step 5 through 7:
5. Let  $P$  be the set of primary mappings, which is an empty set initially. For each chain from the best to the worst according to their chaining scores: if on the query, the chain overlaps with a chain in  $P$  by `--mask-level [=0.5]` or higher fraction of the shorter chain, mark the chain as *secondary* to the chain in  $P$ ; otherwise, add the chain to  $P$ .
6. Retain all primary mappings. Also retain up to `-N [=5]` top secondary mappings if their chaining scores are higher than `-p [=0.8]` of their corresponding primary mappings.
7. If alignment is requested, filter out an internal seed if it potentially leads to both a long insertion and a long deletion. Extend from the left-most seed. Perform global alignments between internal seeds. Split the chain if the accumulative score along the global alignment drops by `-z [=400]`, disregarding long gaps. Extend from the right-most seed. Output chains and their alignments.
8. If there are more query sequences in the input, go to step 2 until no more queries are left.
9. If there are more reference sequences, reopen the query file from the start and go to step 1; otherwise stop.

## 1 Indexing

## 2-4 Collect and sort seeds

## 5-7 Reference mappings

Note: It is the backbone of many assemblers and applications

# What about even longer mappings (genome vs genome)

RESEARCH ARTICLE

## MUMmer4: A fast and versatile genome alignment system

Guillaume Marçais<sup>1,2\*</sup>, Arthur L. Delcher<sup>3</sup>, Adam M. Phillippy<sup>4</sup>, Rachel Coston<sup>3</sup>, Steven L. Salzberg<sup>3,5</sup>, Aleksey Zimin<sup>1,3\*</sup>

Aligner	Graphical User Interface	Multi-platform Windows/Linux	Multi-threaded	Callable from C++, scripting languages	Whole genome aln.	Short read aln.	Long read aln.	SAM format output	P-value output
MUMmer4			✓	✓	✓	✓	✓	✓	
MUMmer3					✓				
Blast	✓	✓	✓		✓				✓
Blat					✓				✓
Mauve	✓	✓			✓				
LASTZ					✓			✓	✓
bwa-mem			✓		-	✓	✓	✓	
Bowtie2			✓		-	✓	-	✓	
BLASR			✓		-	-	✓	✓	✓

# What about even longer mappings (genome)

RESEARCH ARTICLE

## MUMmer4: A fast and versatile genome alignment system

Guillaume Marçais<sup>1,2\*</sup>, Arthur L. Delcher<sup>3</sup>, Adam M. Phillippy<sup>4</sup>, Rachel Coston<sup>3</sup>, Steven L. Salzberg<sup>3,5</sup>, Aleksey Zimin<sup>1,3\*</sup>

		Arabidopsis	Tardigrade	Human/Chimp
nucmer4	Wall time (min)	3.7	4.0	207
	CPU time (min)	22	26	2897
	Memory (GB)	4.6	4.9	66
Mauve	Wall time (min)	41	273	> 2 days
	CPU time (min)	38.6	268	> 2 days
	Memory (GB)	3.3	4.0	> 2 days
LASTZ default	Wall time (min)	1122	> 2 days	> 2 days
	CPU time (min)	1113	> 2 days	> 2 days
	Memory (GB)	1.3		
LASTZ match	Wall time (min)	66	77	> 2 days
	CPU time (min)	66	76	> 2 days
	Memory (GB)	0.6	0.4	

# Mapping algorithm – a summary

Feature	Hash table index tools	BWT tools
<b>Speed</b>	Slower	Faster
<b>Memory</b>	Higher	Lower
<b>Sensitivity</b>	Higher	Lower

Build an index of your reference

Align your reads to your index

Choose an aligner!

Bowtie2, BWA-MEM, SMALT

Minimap2, GMAP, MUMMER4 (Pacbio or Nanopore)

As reads get longer, there seems to be a new generation of mappers arriving

Use the output to do subsequent analysis

What's the output?

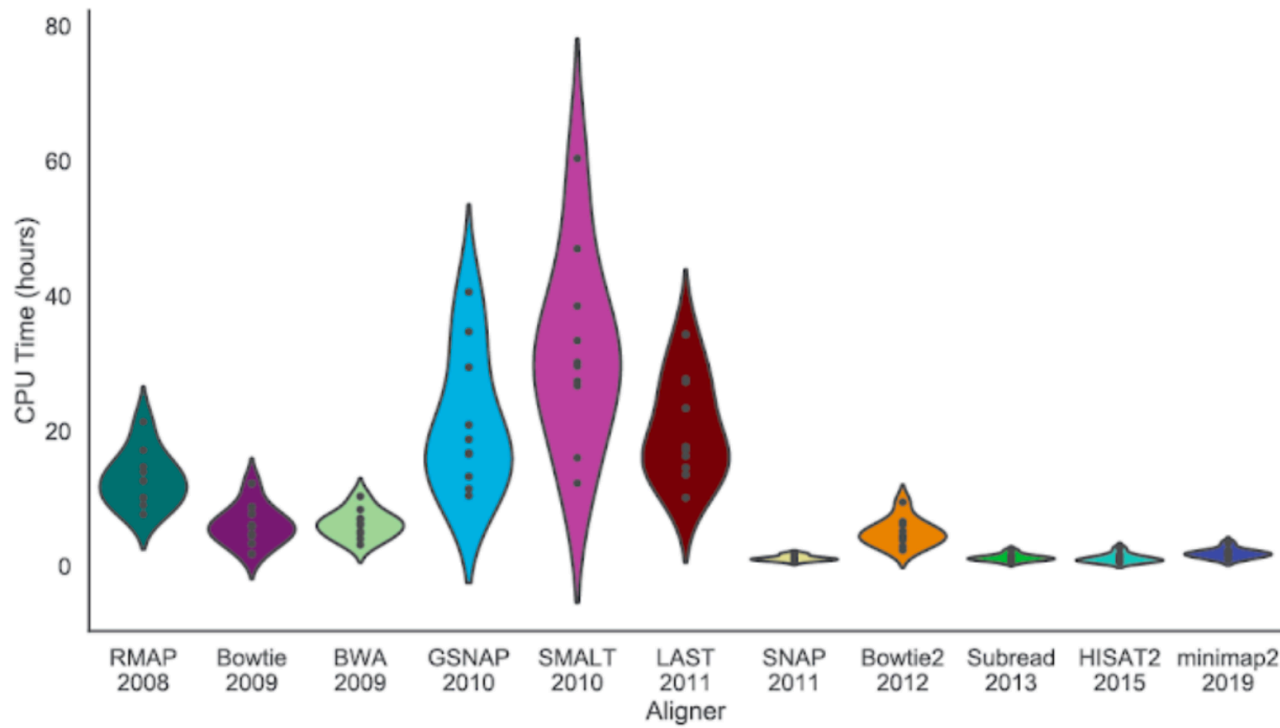
How to use this output?



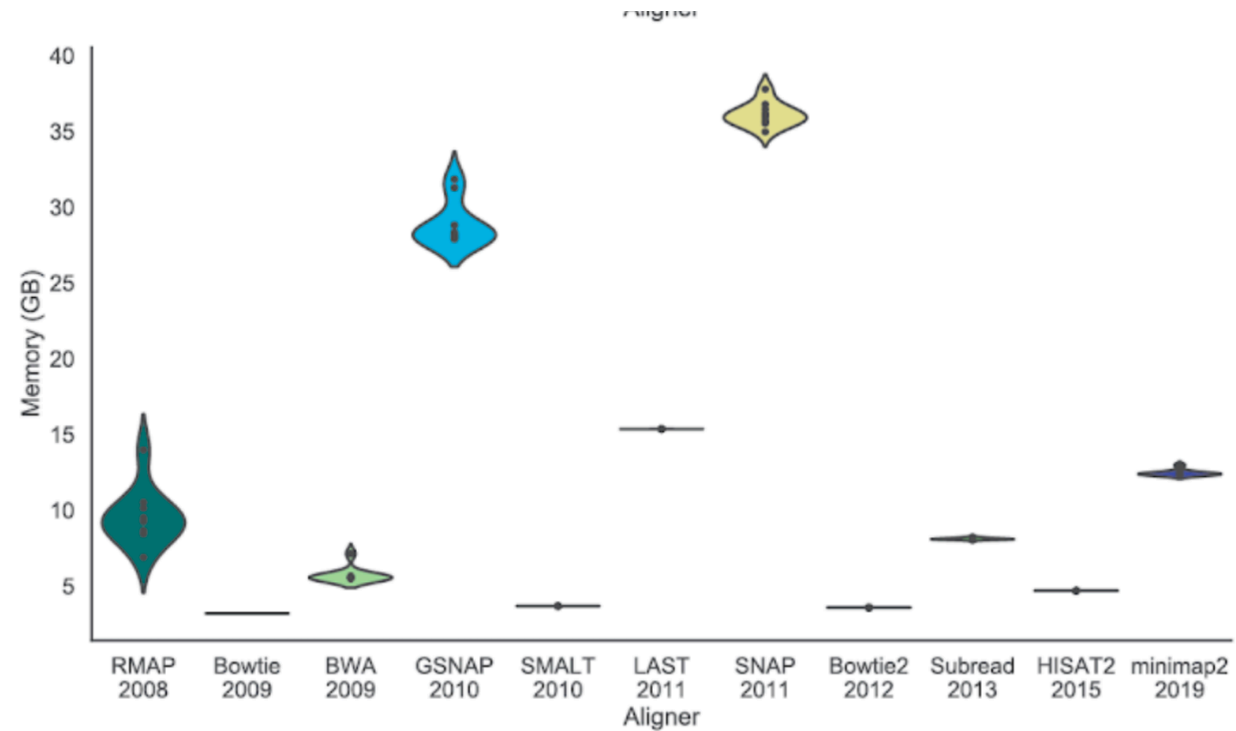
# Mapping algorithm – a summary

- Improvement in both indexing and seed searching
- Allows more applications to be developed

## CPU time



## Memory (GB)



<https://arxiv.org/ftp/arxiv/papers/2003/2003.00110.pdf>

Note: a preprint

# For your reference – speedups in various alignment algorithms including BLAST

## Fast and sensitive protein alignment using DIAMOND

Benjamin Buchfink<sup>1</sup>, Chao Xie<sup>2,3</sup> & Daniel H Huson<sup>1,2</sup>

**The alignment of sequencing reads against a protein reference database is a major computational bottleneck in metagenomics and data-intensive evolutionary projects. Although recent tools offer improved performance over the gold standard BLASTX, they exhibit only a modest speedup or low sensitivity.**

**We introduce DIAMOND, an open-source algorithm based on double indexing that is 20,000 times faster than BLASTX on short reads and has a similar degree of sensitivity.**

- Improvement in both indexing and seed searching
- Allows more applications to be developed

Most sequence comparison programs, including BLASTX and RAPSearch2, use single consecutive seeds, which need to be short (length 3–6 amino acids) to ensure sensitivity. To increase speed without losing sensitivity, DIAMOND uses spaced seeds—that is, longer seeds in which only a subset of positions are used<sup>9,10</sup>. The number and exact layout of those positions are called the weight and shape of the spaced seed, respectively. To achieve high sensitivity, DIAMOND uses a set of four carefully chosen shapes<sup>11</sup> of length 15–24 and weight 12 by default. The most sensitive version of DIAMOND uses 16 shapes of weight 9. In addition, DIAMOND uses a reduced amino acid alphabet of size 11 to enhance sensitivity<sup>12</sup>. A simple exact match criterion determines which seeds are passed on to the extension phase, in which a Smith-Waterman alignment<sup>13</sup> is computed.

In a recent metagenomic study of 12 permafrost samples<sup>14</sup>, a BLASTX comparison of 176 million high-quality DNA reads against the KEGG reference database<sup>3</sup> was reported to require 800,000 CPU hours at a supercomputing center<sup>15</sup>. When we used DIAMOND with its default settings, the analysis of all 246 million reads took 2.3 h on a single workstation, producing a total of 568.9 million alignments on 43 million reads.

Mapping process

# Back to the beginning: FASTQ

```
@HISEQ:409:HA7CJADXX:1:1101:1202:2113 1:N:0:GCNAAT
AAAAAAGTTTCCATACAATTACAAGCATCACACTGTGGGCATGCACTTGGGAAAGAAG
+
==?DBD@<AA<ADAFHGGE<ECHHCG+:1::?D;G4::?BBGCFHI<BCCC;FCGC96
```

Read ID

Sequence

Quality score

Quality Score	Error Probability
Q40	0.0001 (1 in 10,000)
Q30	0.001 (1 in 1,000)
Q20	0.01 (1 in 100)
Q10	0.1 (1 in 10)

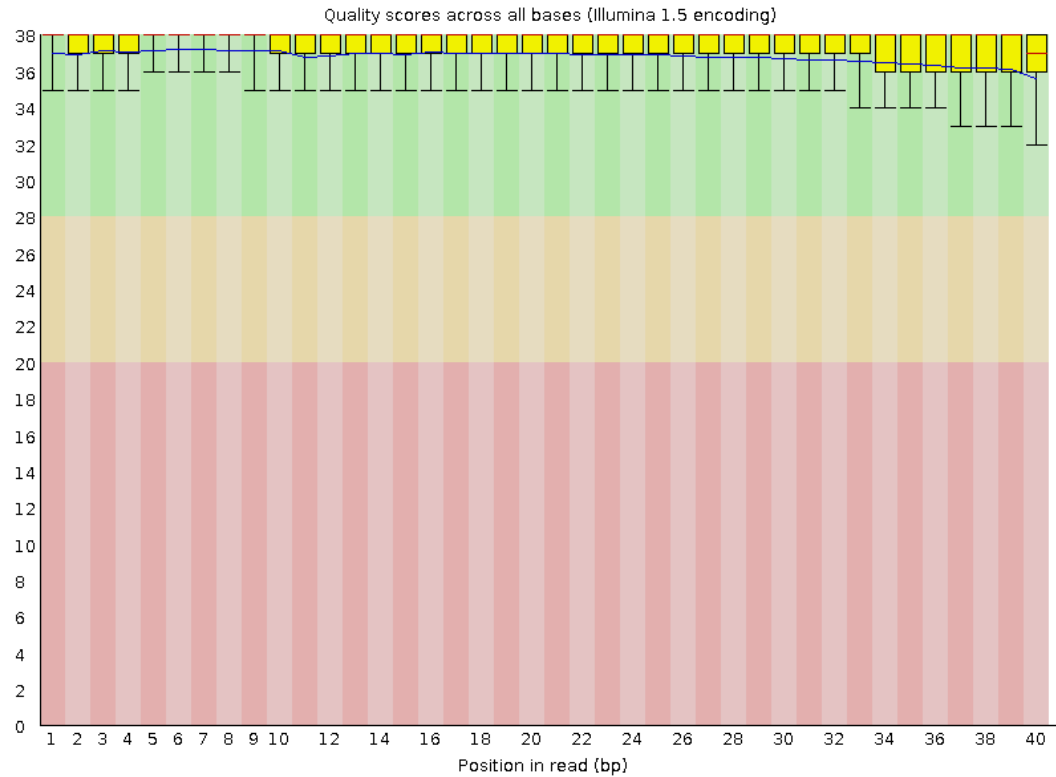
Q-Score Bins	Example of Empirically Mapped Q-Scores*
N (no call)	N (no call)
2-9	6
10-19	15
20-24	22
25-29	27
30-34	33
35-39	37
≥ 40	40

[http://www.illumina.com/content/dam/illumina-marketing/documents/products/technotes/technote\\_understanding\\_quality\\_scores.pdf](http://www.illumina.com/content/dam/illumina-marketing/documents/products/technotes/technote_understanding_quality_scores.pdf)

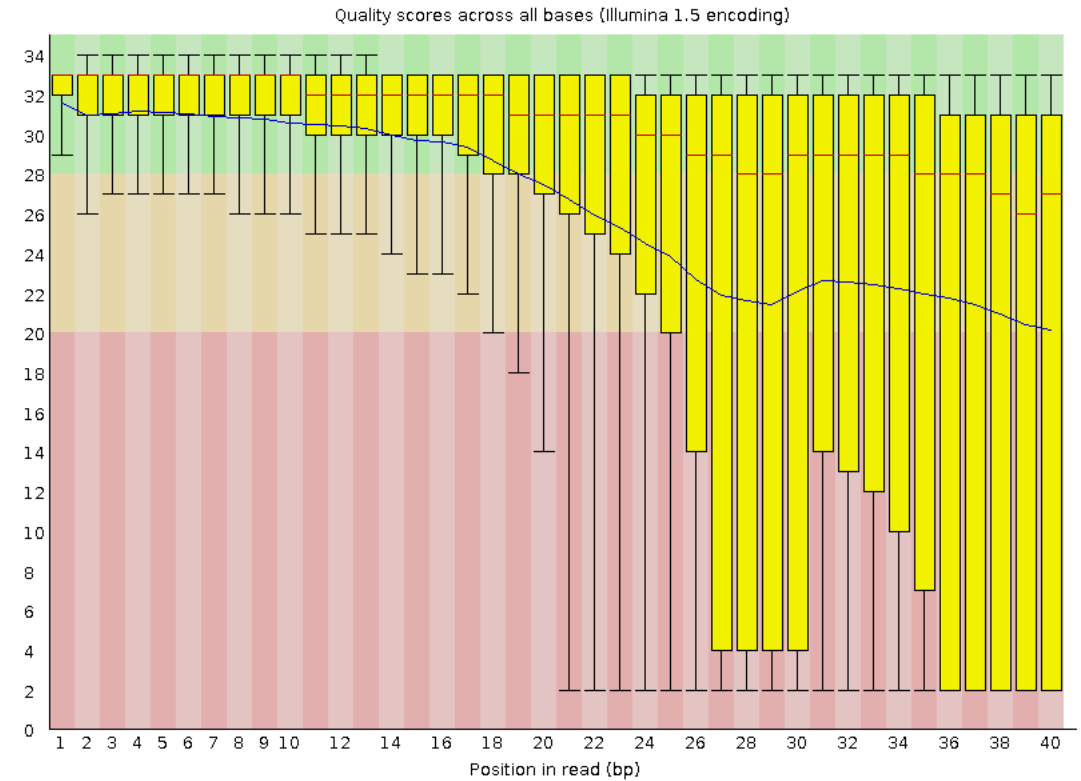
# QC first - always always the first step

- Contamination! \*
- Is it of good quality?
  - Read quality
  - Adaptor contamination
  - Insert size distribution
  - PCR duplicate rate
- Is it your species or someone else's (sample swap)?

# Sequence quality - FastQC



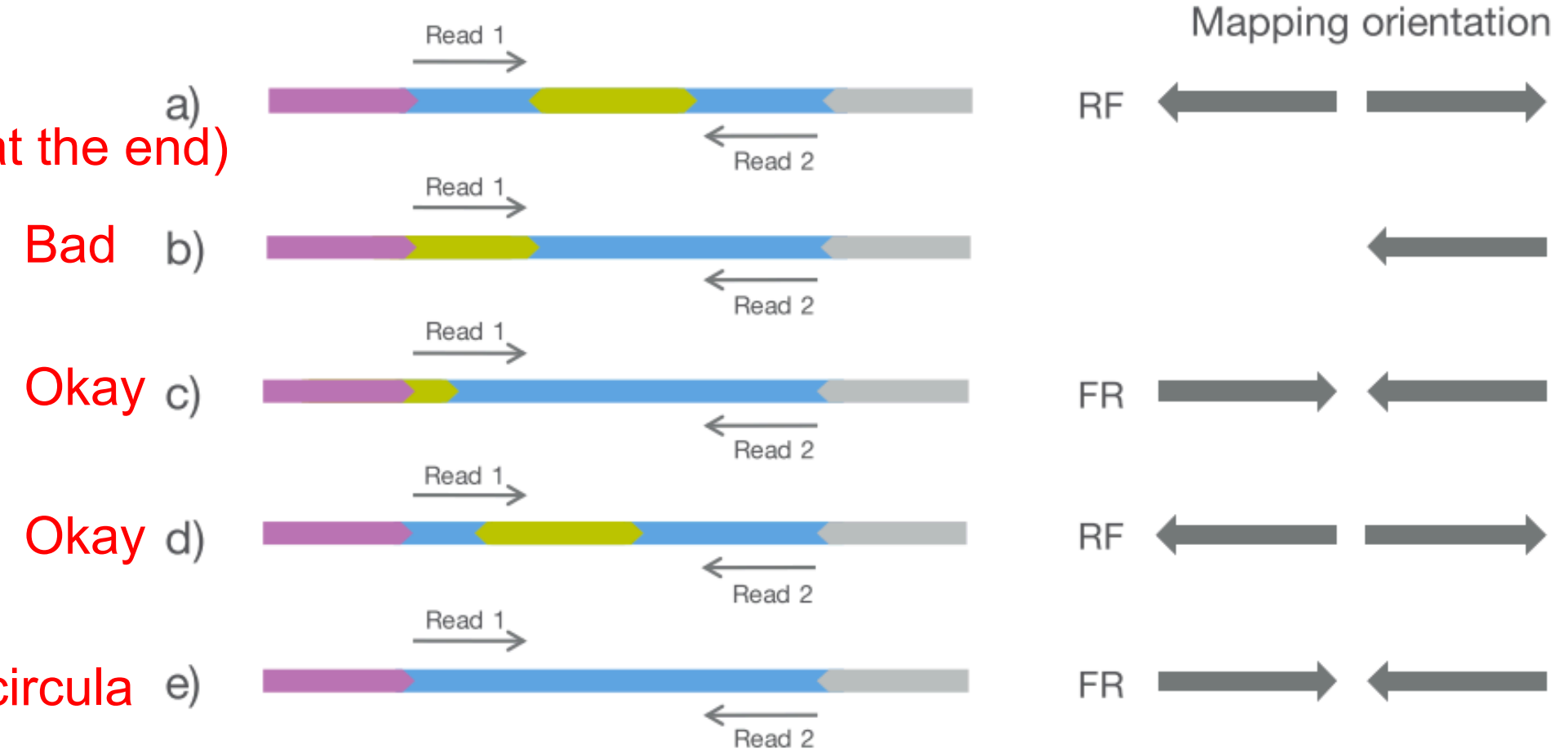
Good (unlikely)



Bad

# Basically the adaptor sequence can appear everywhere (but in a logic way)

Best case  
(a bit of adaptor at the end)



Bad

Okay

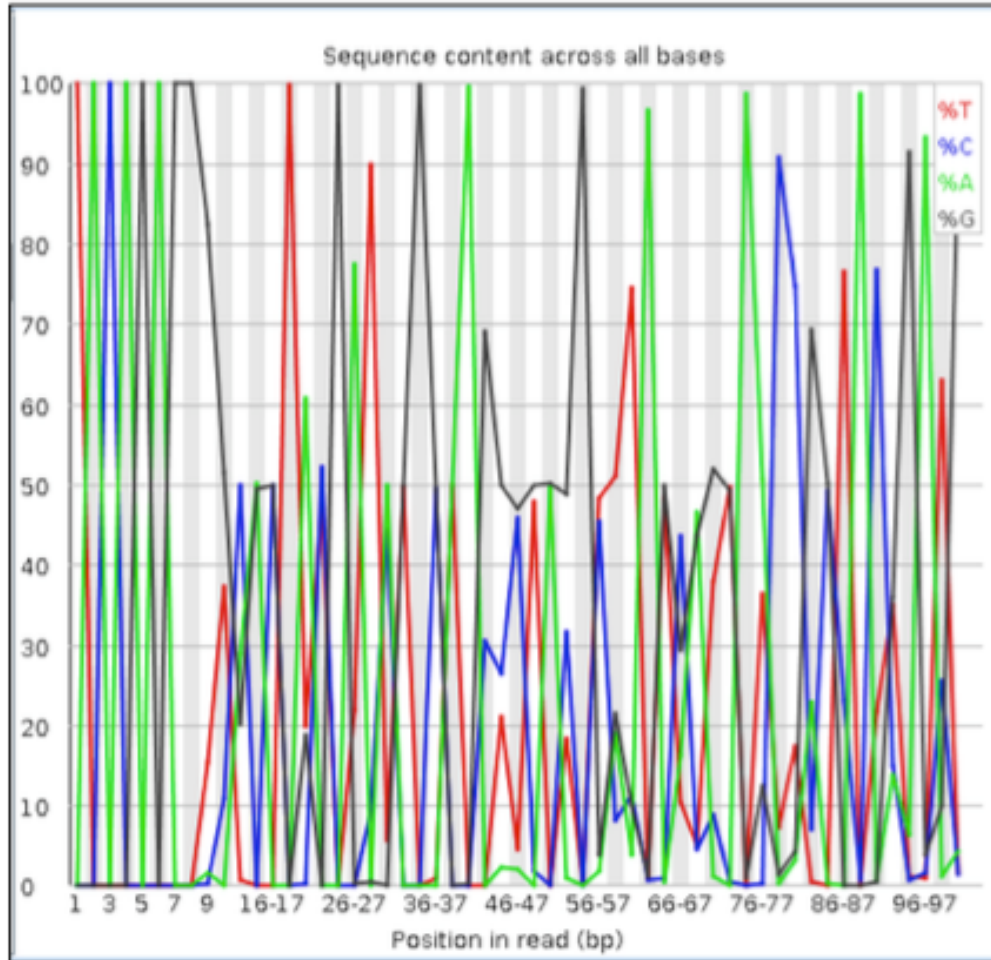
Okay

Totally fail to circula

Composition of example library templates from a mate pair experiment. For each example (a–e), the position of the junction adapter sequence is shown in green and the mapping orientation (either FR or RF, 'forward-reverse' and 'reverse-forward', respectively) of the resulting read pairs is shown to the right. Sections of genomic DNA sequence are shown in blue and the TruSeq adapter sequences are shown in purple and grey. Amplification/sequencing primer adapters are shown in grey and purple.

# FastQC will offer some insights in adaptor

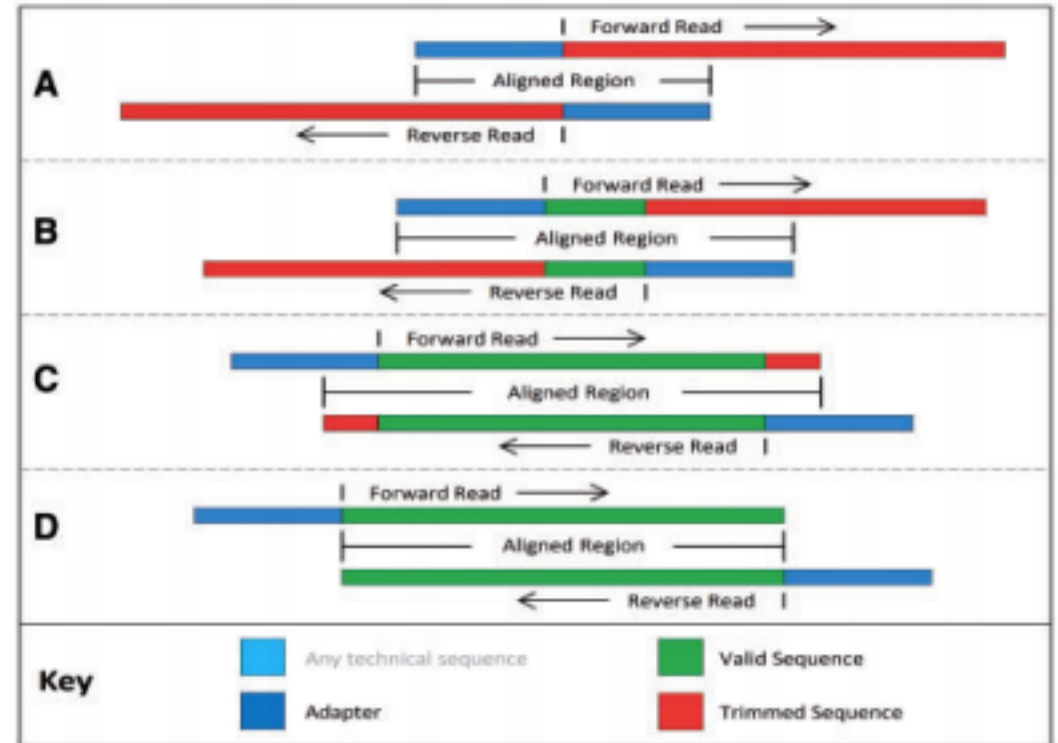
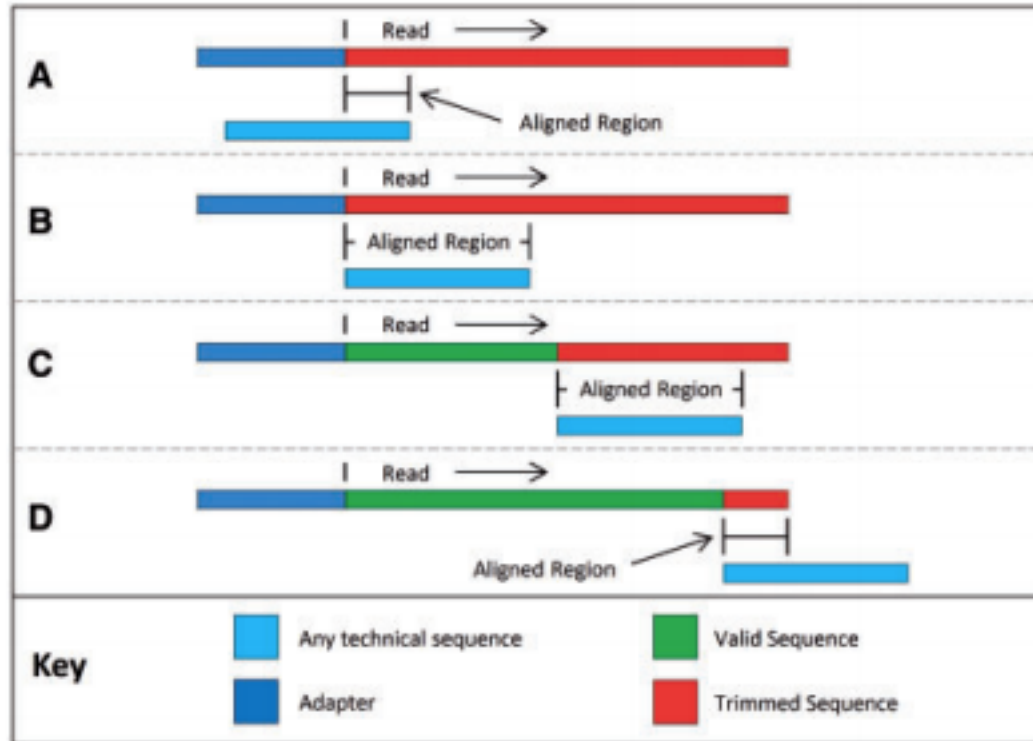
TACAGAGG overrepresented – what is it?



Overrepresented sequences			
Sequence	Count	Percentage	Possible Source
AGCAGCATTGTACA...	3398	3.398	No Hit
TACAGTCCGACGAT...	1814	1.814	Illumina PCR Prime...
TCTACAGTCCGACG...	1570	1.57	RNA PCR Primer, In...
TATTGCACTTGTCCC...	1421	1.421	No Hit
TTCTACAGTCCGAC...	1181	1.181	RNA PCR Primer, In...
CTACAGTCCGACGA...	1168	1.168	Illumina PCR Prime...
CATTGCACTTGTCTC...	839	0.839	No Hit
ACAGTCCGACGATC...	835	0.835	RNA PCR Primer, In...
AGTTCTACAGTCCG...	648	0.648	Illumina PCR Prime...
AAAGTGCTGCGACA...	491	0.491	No Hit
TCGTATGCCGTCTT...	465	0.465	Illumina Single En...
CAGTCCGACGATCT...	436	0.436	Illumina PCR Prime...
TNNNNNNNNNNNN...	392	0.392	No Hit
TAGCTTATCAGACT...	388	0.388	No Hit
TATTGCACTCGTCC...	366	0.366	TruSeq Adapter, I..
ACCGGGCGGAAAC...	357	0.357	No Hit
ANNNNNNNNNNN...	355	0.355	No Hit
GTTCTACAGTCCGA...	353	0.353	Illumina PCR Prime...
AAGTGCTGCGACAT	341	0.341	No Hit



# Trimmomatic for quality and adaptor trimming (many other tools also exist)



Trimmomatic: a flexible trimmer for Illumina sequence data - NCBI - NIH

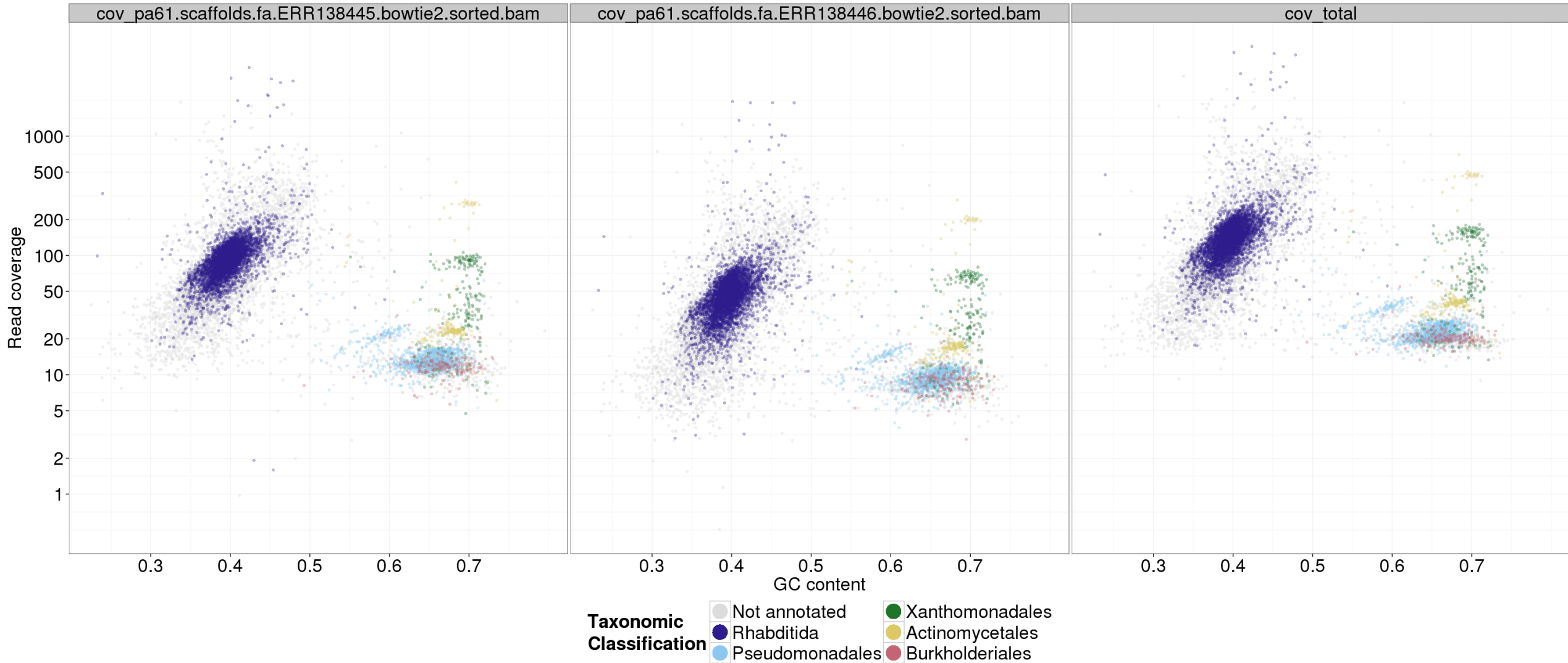
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4103590/> ▾ 翻譯這個網頁

由 AM Bolger 著作 - 2014 - 被引用 5183 次 - 相關文章

2014年4月1日 - Motivation: Although many next-generation sequencing (NGS) read preprocessing tools already existed, we could not find any tool or combination of tools that met our requirements in terms of flexibility, correct handling of paired-end data and high performance. We have developed **Trimmomatic** as a more ...

Bolger *et al.*, (2014)

# Check what your samples contain - Blobology



# Source of contamination

- Difficult to remove (gut from microorganisms)
- Fail to remove
- Not careful
- Bad company
- Sequencer carry over (from previous run)
- Sample (barcode) mix up

- Or simply bad day (not your fault)

Salter *et al.* *BMC Biology* 2014, **12**:87  
<http://www.biomedcentral.com/1741-7007/12/87>



RESEARCH ARTICLE

Open Access

Reagent and laboratory contamination can critically impact sequence-based microbiome analyses

# Source of contamination

Mukherjee et al. *Standards in Genomic Sciences* 2015, **10**:18  
<http://www.standardsingenomics.com/content/10/1/18>



Standards in  
Genomic Sciences

COMMENTARY

Open Access

## Large-scale contamination of microbial isolate genomes by Illumina PhiX control

Supratim Mukherjee<sup>1\*</sup>, Marcel Huntemann<sup>1</sup>, Natalia Ivanova<sup>1</sup>, Nikos C Kyrpides<sup>1,2</sup> and Amrita Pati<sup>1</sup>

....In this study we screened over **18,000 publicly available microbial isolate genome sequences in the Integrated Microbial Genomes database and identified more than 1000 genomes that are contaminated with PhiX, a control frequently used during Illumina sequencing runs.** .....The presence of PhiX contamination in several publicly available isolate genomes can result in additional errors when such data are used in comparative genomics analyses. **Such contamination of public databases have far-reaching consequences in the form of erroneous data interpretation and analyses, and necessitates better measures to proofread raw sequences before releasing them to the broader scientific community.**

# Sample storage matters (case of humans)

## 3 months storage resulted in less efficient DNA extraction

High fragmentation: loss of material

Decrease in library complexity

High increase in PCR duplicates, 60-85% for FFPE vs. 30% for FF

## C > U deamination is a common cause of artifacts

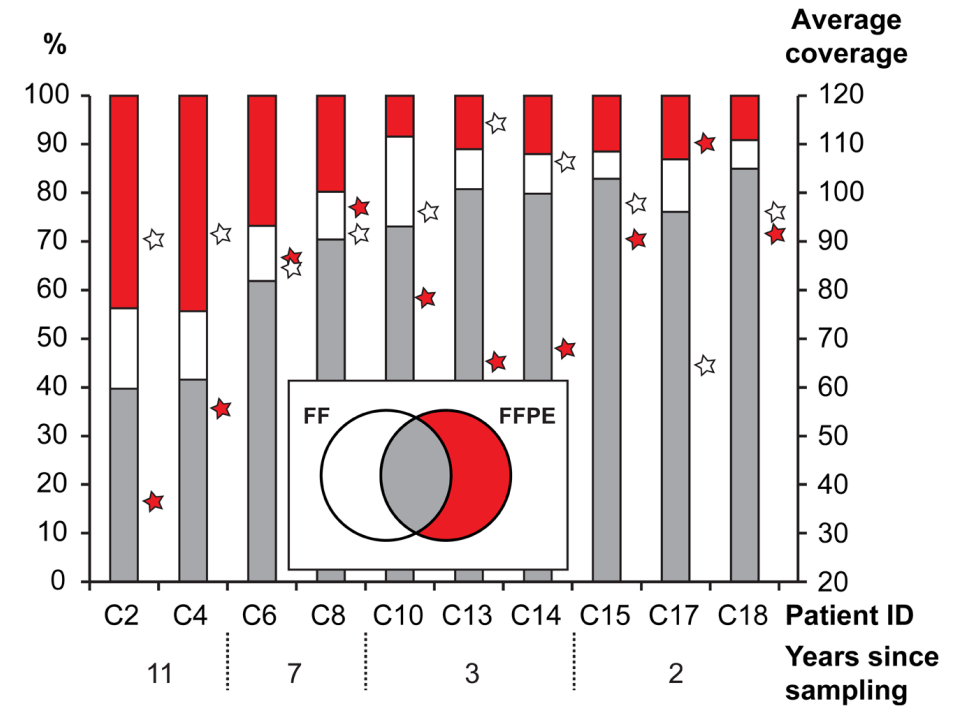
U-tolerant polymerase didn't help

Pattern, T <> C, A <> G transition

## The fraction of mapped reads decreases with storage time

Increase in partial mappings

Increase in gapped mappings



Hedegaard et al. 2014

# Mapping output format: SAM/BAM

Spec defined by maq/bwa/samtools author Heng Li

SAM: text version

tab-delimited

Exome (GBs) ; Whole genome (TBs)

BAM: binary/compressed version

indexed so it's faster to look up using samtools

Exome (1-2GBs) ; Whole genome (GBs)

# SAM file header

```
@HD VN:1.4 S0:coordinate
@SQ SN:PNOK.scaff0001.C LN:7761079
@SQ SN:PNOK.scaff0002.C LN:4533150
@SQ SN:PNOK.scaff0003.C LN:3409659
@SQ SN:PNOK.scaff0004.0 LN:3380754
@SQ SN:PNOK.scaff0005.0 LN:2749859
@SQ SN:PNOK.scaff0006.0 LN:2613677
@SQ SN:PNOK.scaff0007.0 LN:1690816
@SQ SN:PNOK.scaff0008 LN:1673160
@SQ SN:PNOK.scaff0009.0 LN:1538597
@SQ SN:PNOK.scaff0010 LN:1377172
@SQ SN:PNOK.scaff0011 LN:633856
@SQ SN:PNOK.scaff0012 LN:52253
@SQ SN:PNOK.mito LN:163443
@PG ID:smalt VN:0.7.4 CL:/h
```

Always start with @

Contains “background” information

@HD = Header

@SQ = Sequence dictionary

# SAM file header

Very detailed in how one should specify the headers

Subsequent programs (like variant calling) will use these info

<http://samtools.github.io/hts-specs/SAMv1.pdf>

Tag	Description
<b>@HD</b>	The header line. The first line if present.
<b>VN*</b>	Format version. <i>Accepted format:</i> /^[0-9]+\.[0-9]+\$/.
<b>SO</b>	Sorting order of alignments. <i>Valid values:</i> <b>unknown</b> (default), <b>unsorted</b> , <b>queryname</b> and <b>coordinate</b> . For coordinate sort, the major sort key is the RNAME field, with order defined by the order of @SQ lines in the header. The minor sort key is the POS field. For alignments with equal RNAME and POS, order is arbitrary. All alignments with '*' in RNAME field follow alignments with some other value but otherwise are in arbitrary order.
<b>GO</b>	Grouping of alignments, indicating that similar alignment records are grouped together but the file is not necessarily sorted overall. <i>Valid values:</i> <b>none</b> (default), <b>query</b> (alignments are grouped by QNAME), and <b>reference</b> (alignments are grouped by RNAME/POS).
<b>@SQ</b>	Reference sequence dictionary. The order of @SQ lines defines the alignment sorting order.
<b>SN*</b>	Reference sequence name. Each @SQ line must have a unique SN tag. The value of this field is used in the alignment records in RNAME and RNEXT fields. Regular expression: [!-]+-<->-[!-]*
<b>LN*</b>	Reference sequence length. <i>Range:</i> [1,2 <sup>31</sup> -1]
<b>AS</b>	Genome assembly identifier.
<b>M5</b>	MD5 checksum of the sequence in the uppercase, excluding spaces but including pads (as '*'s).
<b>SP</b>	Species.
<b>UR</b>	URI of the sequence. This value may start with one of the standard protocols, e.g http: or ftp:. If it does not start with one of these protocols, it is assumed to be a file-system path.
<b>@RG</b>	Read group. Unordered multiple @RG lines are allowed.
<b>ID*</b>	Read group identifier. Each @RG line must have a unique ID. The value of ID is used in the RG tags of alignment records. Must be unique among all read groups in header section. Read group IDs may be modified when merging SAM files in order to handle collisions.
<b>CN</b>	Name of sequencing center producing the read.
<b>DS</b>	Description.
<b>DT</b>	Date the run was produced (ISO8601 date or date/time).
<b>FO</b>	Flow order. The array of nucleotide bases that correspond to the nucleotides used for each flow of each read. Multi-base flows are encoded in IUPAC format, and non-nucleotide flows by various other characters. <i>Format:</i> /\* [ACMGRSVTWYHKDBN]+/
<b>KS</b>	The array of nucleotide bases that correspond to the key sequence of each read.
<b>LB</b>	Library.
<b>PG</b>	Programs used for processing the read group.
<b>PI</b>	Predicted median insert size.
<b>PL</b>	Platform/technology used to produce the reads. <i>Valid values:</i> <b>CAPILLARY</b> , <b>LS454</b> , <b>ILLUMINA</b> , <b>SOLID</b> , <b>HELICOS</b> , <b>IONTORRENT</b> , <b>ONT</b> , and <b>PACBIO</b> .
<b>PM</b>	Platform model. Free-form text providing further details of the platform/technology used.
<b>PU</b>	Platform unit (e.g. flowcell-barcode.lane for Illumina or slide for SOLiD). Unique identifier.
<b>SM</b>	Sample. Use pool name where a pool is being sequenced.
<b>@PG</b>	Program.
<b>ID*</b>	Program record identifier. Each @PG line must have a unique ID. The value of ID is used in the alignment PG tag and PP tags of other @PG lines. PG IDs may be modified when merging SAM files in order to handle collisions.
<b>PN</b>	Program name
<b>CL</b>	Command line







# SAM file spec

Col	Field	Type	Regex/Range	Brief description
1	QNAME	String	[!-?A-~]{1,254}	Query template NAME
2	FLAG	Int	[0,2 <sup>16</sup> -1]	bitwise FLAG
3	RNAME	String	\*  [!-( )+-<>-~] [!-~]*	Reference sequence NAME
4	POS	Int	[0,2 <sup>31</sup> -1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0,2 <sup>8</sup> -1]	MAPping Quality
6	CIGAR	String	\*  ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	\* =  [!-( )+-<>-~] [!-~]*	Ref. name of the mate/next read
8	PNEXT	Int	[0,2 <sup>31</sup> -1]	Position of the mate/next read
9	TLEN	Int	[-2 <sup>31</sup> +1,2 <sup>31</sup> -1]	observed Template LENgth
10	SEQ	String	\*  [A-Za-z=.]+	segment SEQUENCE
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33



# SAM flags

hexadecimal	decimal	binary bit; 0=no, 1=yes	position of bit	description
0x1	1	"0000 0000 0001"	1	paired-end (or multiple-segment) sequencing technology
0x2	2	"0000 0000 0010"	2	each segment properly aligned according to the aligner
0x4	4	"0000 0000 0100"	3	segment unmapped
0x8	8	"0000 0000 1000"	4	next segment in the template unmapped
0x10	16	"0000 0001 0000"	5	SEQ is reverse complemented
0x20	32	"0000 0010 0000"	6	SEQ of the next segment in the template is reverse complemented
0x40	64	"0000 0100 0000"	7	the first segment in the template
0x80	128	"0000 1000 0000"	8	the last segment in the template
0x100	256	"0001 0000 0000"	9	secondary alignment
0x200	512	"0010 0000 0000"	10	not passing quality controls
0x400	1024	"0100 0000 0000"	11	PCR or optical duplicate
0x800	2048	"1000 0000 0000"	12	supplementary alignment

<http://gatkforums.broadinstitute.org/gatk/discussion/7019/sam-flags-down-a-boat>

# CIGAR String a few examples

ATCGATCGATCGATCG      Reference

ATGGACGATTTCGTGAA      Read mapping = 5M1D3M1I3M4S

Soft clip usually the result of lower mapping quality

- 200M
- 203M1D4M1I48M1I13M
- 164M
- 232M
- 159M
- 162M1D4M1I43M
- 101M7S
- 227M
- 155M
- 105M

op	Description
M	Alignment match (can be a sequence match or mismatch)
I	Insertion to the reference
D	Deletion from the reference
N	Skipped region from the reference
S	Soft clip on the read (clipped sequence present in <seq>)
H	Hard clip on the read (clipped sequence NOT present in <seq>)
P	Padding (silent deletion from the padded reference sequence)

# CIGAR string of long reads

```
9368d6b3-3242-4583-a37c-8ecbfa44ccee 0 ref2.scaff0001 1 60 1574S76M1D12M6D17M1D10M1D3M1D6M4D2M1D6M1D16M2D20M1I10M1I5M1I5M4D4M1D3M1I12M1D42M2D7M1I13M1I15M2I5M1I4M1I8M3D10M3D9M1D4M1I6M1D4M1D5M1I11M1D13M1D10M2D18M4I11M2I13M1D6M1D2M3D13M1D14M2D2M5D26M2D4M2D11M1I8M1I4M1I6M1D8M1D3M2D3M1I41M1I14M1D10M1D36M1D54M2D2M1D7M1I6M5D12M2I2M1I12M2I2M1D16M1I1M2I5M1I19M1D9M1D2M1D7M1I15M1D9M2D12M1I7M2I2M1D3M1D8M1I13M1D5M1D6M1D21M3I4M1I8M1I15M3D1M2D18M3D10M1D5M1I9M1D21M1I22M4D2M1D3M2D14M1I1M1I15M2I5M1D2M1D10M1D11M2D6M1I11M5I11M1I6M1D19M8I14M1I10M1I27M3I14M2D5M1D6M1I4M3D2M1I28M1D17M2I13M3D23M2D5M4D1M2D9M3I16M1D17M2D11M1I14M1D16M1D2M1D7M2D10M1I1M1I11M1I22M1I54M1D9M1I9M2I9M1I3M1I14M2D7M1D4M1D7M5D1M1D14M1I15M1D3M1I9M2D13M1I4M1D5M2D10M2D1M2D13M3D16M2I4M1I5M1D3M1I17M2D6M5D46M4D14M5D4M1D21M3D45M1D17M2D5M2D1M1D2M4D30M1D29M3I9M1D7M1D15M1I42M2I7M1D10M1I9M1D3M1D3M1D5M4D16M4I13M1D29M1D4M3D6M1I35M6D20M6I4M1I6M3D16M1D7M1I34M1D11M4D21M1D23M2D1M1D35M4D12M1D6M1I2M1I25M1D13M1I35M2D11M3D20M5D2M1D21M1I12M2D3M2D21M3I31M1D37M3D6M1D13M1D35M1I17M1D5M3I16M1I71M1I1M1I2M1I4M1I5M1I19M4D35M2D32M1D28M1I7M2I8M1I6M1I14M1I16M1D58M2D12M1I6M1D11M1D17M3I28M2I2M1D18M2I6M2D3M1D9M1D15M2D1M1D16M1I5M1I30M4I6M2I9M1I5M1D26M2D2M3D1M1D4M1I13M2D21M5D5M2D10M4D17M1D24M1D11M1D6M4D15M1I20M1D32M1D16M5D22M2I11M1I35M2D7M1D43M2D10M1D6M2I8M1D5M1D9M5D2M1D9M1I9M1D35M3I1M1I9M1D9M3D5M3D5M3D19M1I23M2D17M1I4M1I8M1D13M2I11M3I1M1I20M1I1M1I22M2D36M3D19M2I19M2D9M2D6M2I1M1I9M1I30M1I10M1D3M1D23M1I31M1I20M2I14M2D1M1D15M1D22M2D3M1D9M1D10M6D3M1D2M1D21M1I22M1I16M1I13M2I2M1I7M1D38M5D6M3D7M1D8M1D32M1I3M3D39M5D25M1D1M1D9M1I8M1D2M2D14M1D47M1D1M1D7M2I20M1D3M1D8M2D39M2D4M1D51M1D5M4D8M2I5M1D12M1D6M72S *
```

```
606ec45a-9559-4dcc-8901-f39b24c67e21 2048 ref2.scaff0001 1 60 15175H9M2D7M1D5M1D2M1D25M1D22M1D8M1D30M1D9M2D14M2D12M1D29M1D9M1I4M2D1M1D2M5D12M1D7M1I7M2D12M1I28M6D14M1D27M2D15M2D13M1I18M2D1M3D1M1D14M2D28M1D8M1I7M2D13M1D30M1D46M1I31M1I4M3D8M1D14M1I11M2I7M1I22M1D6M1D7M1D15M3D46M1D5M1I5M1D2M1D8M2D16M8D39M6D24M2D6M1I6M1I12M1D15M11D13M1D20M1D11M1D3M1D2M1D4M1I23M1I22M1D7M2D12M1I5M1D51M1I11M4D7M1D12M1I2M1D4M1D9M1D43M1D50M1I12M1D1M1D22M1D11M1I4M3I3M1D13M1I7M1I2M1D20M1D13M1D7M1D7M5D30M1D5M1D7M3D1M1D14M2I4M1I21M1I72M1I12M1I28M1I6M1D28M1D8M2D5M1I8M3D2M5D6M1D58M1I29M1D17M1D3M1D10M1D11M2D2M1D11M3D9M1D2M4D8M2D11M2I12M5D47M8D9M3D24M1I16M1I22M1D55M1D24M1I20M2D13M2D2M1D61M1I11M1D13M1D7M1I4M5D6M2D42M1I17M1I3M1D20M1I10M1I14M3D25M2I5M1D22M1D11M8D11M1D1M1D10M5D3M1D12M1D30M2I5M1D3M1D10M1I9M2I13M1D6M1I31M2I8M1D4M1D1M3D29M2D20M1I6M3D7M5D13M1I24M1D12M3D1M3D22M1D9M1D25M1D13M1D1M1D4M1D28M1I40M1D5M2D9M3D13M1D5M1I19M2I17M2D16M1D2M1D7M2D11M1I35M1D7M1D29M1D5M1D2M1I22M1D5M2D19M1D5M3D19M1D5M6D13M3D8M3D9M1D25M1D1M1D18M1D3M1D9M1D25M2D8M1D27M2D7M2D3M1D6M1D10M1D1M3D3M2D14M2D16M1I19M1D4M2I3M1D4M1D25M1D8M3D8M4D1M5D8M1I15M2D16M2D21M1D5M3D51M1I5M1I10M1I14M1D2M1I27M1D47M1D35M2D17M3D1M2D4M1D5M2D18M4D6M3D2M3D7M4D10M1I2M1D19M2D21M1I8M3I16M1I12M1D6M2D23M1D12M1D3M3I32M2I8M2I11M4I2M1D7M1D33M2D17M2D2M2D9M1I12M5D40M2D35M1D11M1I11M3D1M2D41M1D4M1D30M1D7M5D2M1D22M1D1M2D12M5D14M2D6M1I35M1I25M1D9M2D4M2D31M1D2M1D11M2D7M1I8M1D5M2D15M1D14M4D4M1I24M1I55M1D9M1D19M3D32M1D1M2D26M2D10M1D5M1D5M1D26M1I18M1I9M1D10M6D2M3D64M2D12M1D8M1D12M1I12M3I24M1I20M1D16M1I18M2D19M1I31M2D4M1D11M2D15M1D8M4D11M1I13M1D22M1I28M1D7M1I46M1D8M2D26M1D1M2D5M2I9M1I2M1D40M2D20M4D15M1I12M1D5M2D13M1D10M2D5M1I4M1I22M2D40M1I6M3D19M1D17M1I20M5D27M2D7M1D5M1D3M3D5M2I16M1D12M1D4M5D12M1I4M1D6M1D5M1D3M2D17M3I4M1I27M1I22M1D11M1D14M1I35M1D23M3D13M1D10M14D16M1I10M1D14M1I6M1D3M1D28M1I7M1I11M1I7M2D12M4D27M1D7M1D1M3D5M1D9M3I24M2D7M3D21M1D8M1D30M1I14M7D9M2D39M1I27M3I4M1D15M1I2M2I25M1D8M2D13M1D29M4D3M6D43M3D1M3D36M1D2M1D16M1D10M1D5M1D8M1D5M1I20M2D6M1D7M4D3M1D10M1I13M1D8M1D2M1I1M1I4M1D24M2D12M1D10M1D2M3D8M1I10M2D5M1D10M1D9M1I13M1I4M1I3M3D7M2I13M1I6M1I7M2D3M1D24M4I5M1I9M1D8M1I13M1I8M1D35M1D22M6D7M3D3M8D6M3D22M20148
```

# Mapping quality

Probability that a read is mapped incorrectly

Useful for calling SNP later on

Function of

- Uniqueness

- Number of mismatches

- Number of indels

- Quality of bases in read

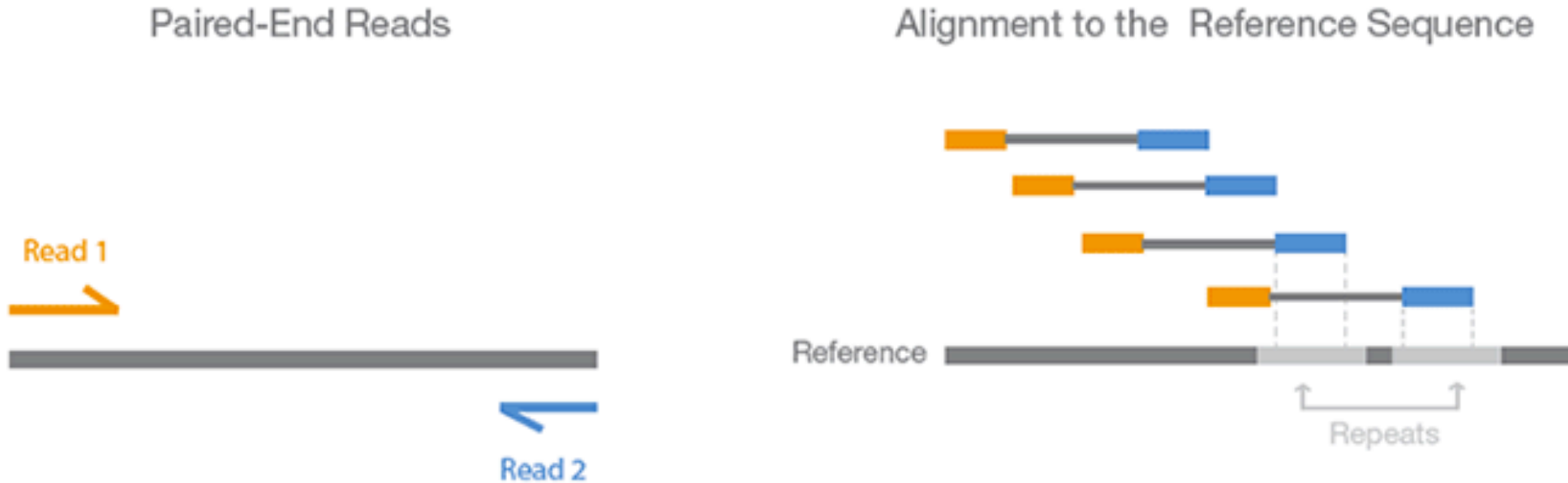
MQ30 = 1 in 1000 alignment is wrong

MQ40 = 1 in 10000 alignment is wrong



# Post mapping QC: insert size in PE mapping?

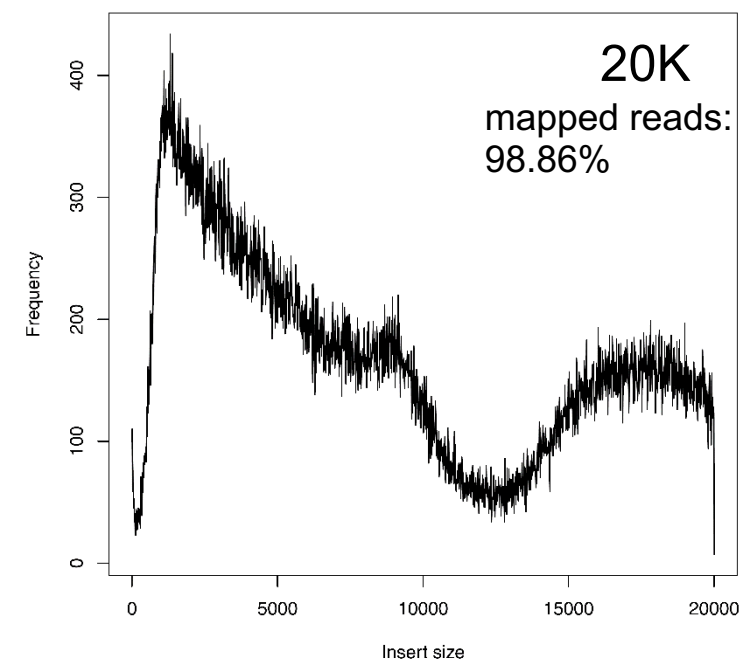
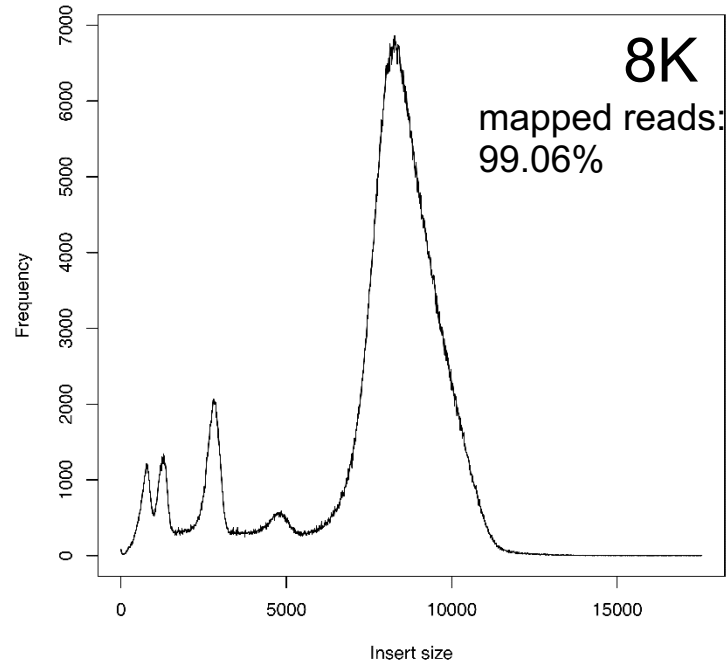
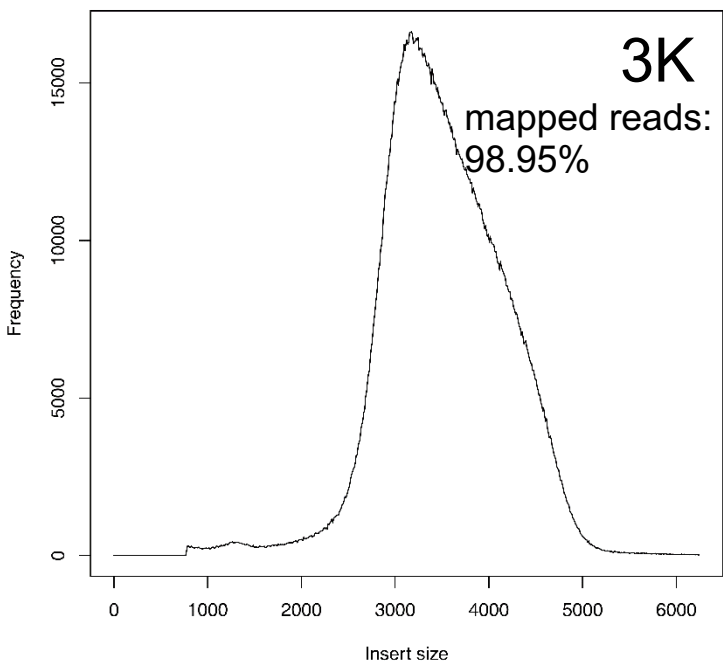
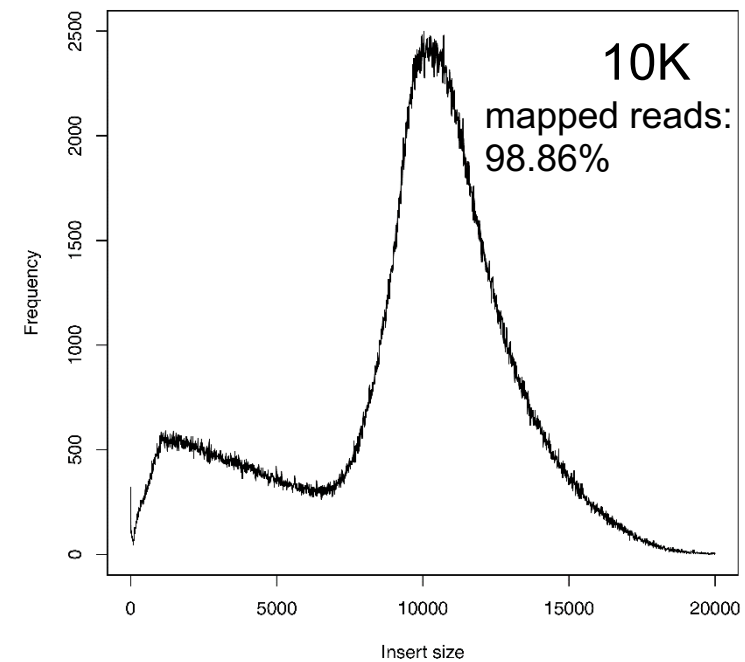
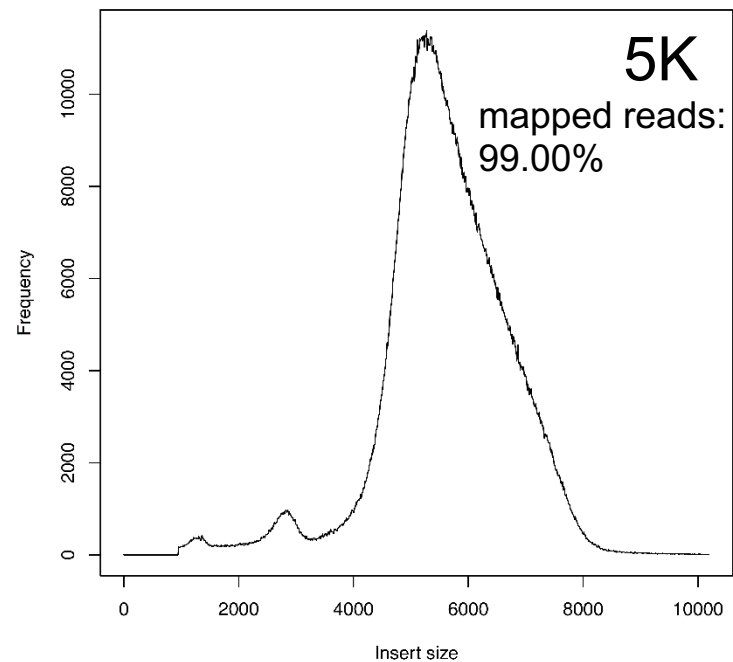
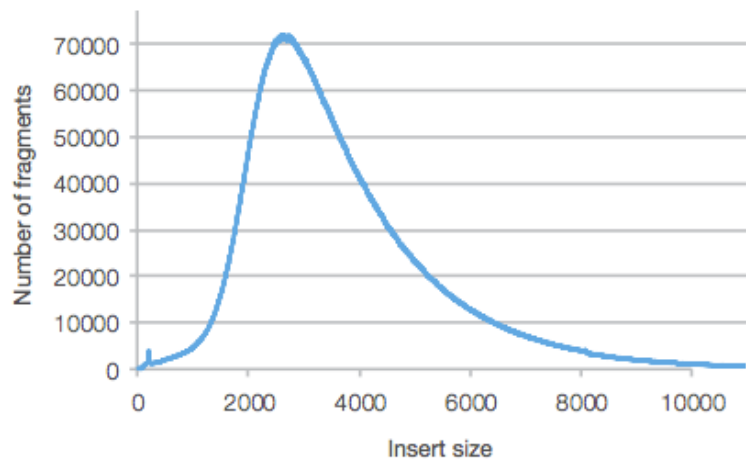
Figure 4. Paired-End Sequencing and Alignment



Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

DNA fragment length should be longer than most repeat size in your genome  
No point to boost up coverage if your fragment len < repeat length

# Insert size



# Post mapping QC – how much coverage?

```
*****  
Stats for BAM file(s):  
*****  
  
Total reads:      2963812  
Mapped reads:    2926492      (98.7408%)  
Forward strand:  1463708      (49.386%)  
Reverse strand:  1462784      (49.3548%)  
Failed QC:       0      (0%)  
Duplicates:      8469 (0.285747%)  
Paired-end reads: 2963812      (100%)  
'Proper-pairs':  2808018      (94.7435%)  
Both pairs mapped: 2901342      (97.8922%)  
Read 1:          1481906  
Read 2:          1481906  
Singletons:      25150      (0.848569%)  
Average insert size (absolute value): 808.327  
Median insert size (absolute value): 466
```

2963812 reads  
x 300 bp per read  
/ 32000000bp genome  
= **27.8X**

This number is overestimated because

1. ~1.3% not mapped
2. Trimmed reads (not all reads have now 300bp)

# 1 million dollar question: how much coverage is better

In mapping:

- ~15X for SNP calling in bacteria
- ~30X for SNP calling in diploid (to delineate heterozygous bases)
- >50X for exome (because you need to be sure)
- No point with >100X in the Illumina world

# PCR duplicates

PCR duplicates during sample prep

= the same fragment is sequenced again and again and again

Some worse than others (because starting material is not good)

< 5% is good

High duplication rate will lead to problems in downstream analysis

Example: 30X ; 1 out of ~30 fragment get duplicated 15 times

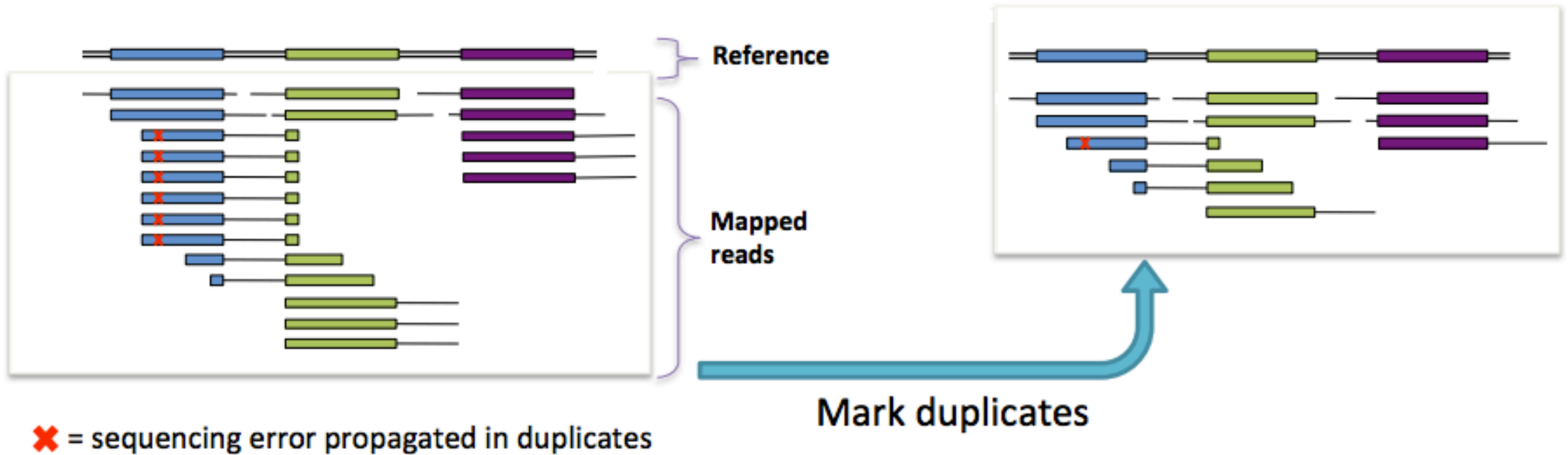
= skew allele frequency

= false SNP discovery

Can be detected (and removed) by read pairs map at the complete position. **We usually keep one copy only**

# PCR duplicates

Can be detected (and removed) by read pairs map at the complete position.  
**We usually keep one copy only**



# De- duplication



**Showing duplicate reads**



**Hiding duplicate reads**

# Case study: Check lane quality and assembly

	Total reads	Mapped reads (%)	Duplicates	Proper-pairs	Both pairs mapped	Median
BRC PE	217,190,726	95.80%	1.47%	53.97%	92.85%	968
Old PE	249,742,439	4.40%	1.51%	2.81%	3.08%	59
Old PE	1,167,521,211	98.21%	11.81%	68.97%	97.18%	465
Old PE	917,638,787	97.97%	5.12%	75.54%	96.99%	261
Company hmm	38,508,236	94.15%	7.51%	48.22%	90.53%	1681
Company hmm	76,992,221	95.09%	10.75%	48.57%	92.43%	1675
Company hmm	26,348,302	93.54%	6.23%	47.58%	89.29%	1681
Company hmm	398,746,361	98.42%	79.36%	57.28%	97.23%	1500
Company hmm	396,241,991	98.42%	79.03%	57.31%	97.24%	1500
Company hmm	39,879,176	92.45%	29.40%	40.66%	88.55%	4623
Company hmm	43,010,934	92.10%	31.27%	40.40%	87.99%	4623
Company hmm	316,963,201	97.71%	84.14%	57.79%	96.11%	410
Company hmm	71,118,483	96.00%	70.88%	42.97%	93.67%	283
Company hmm	61,803,780	94.60%	73.18%	45.36%	91.55%	285



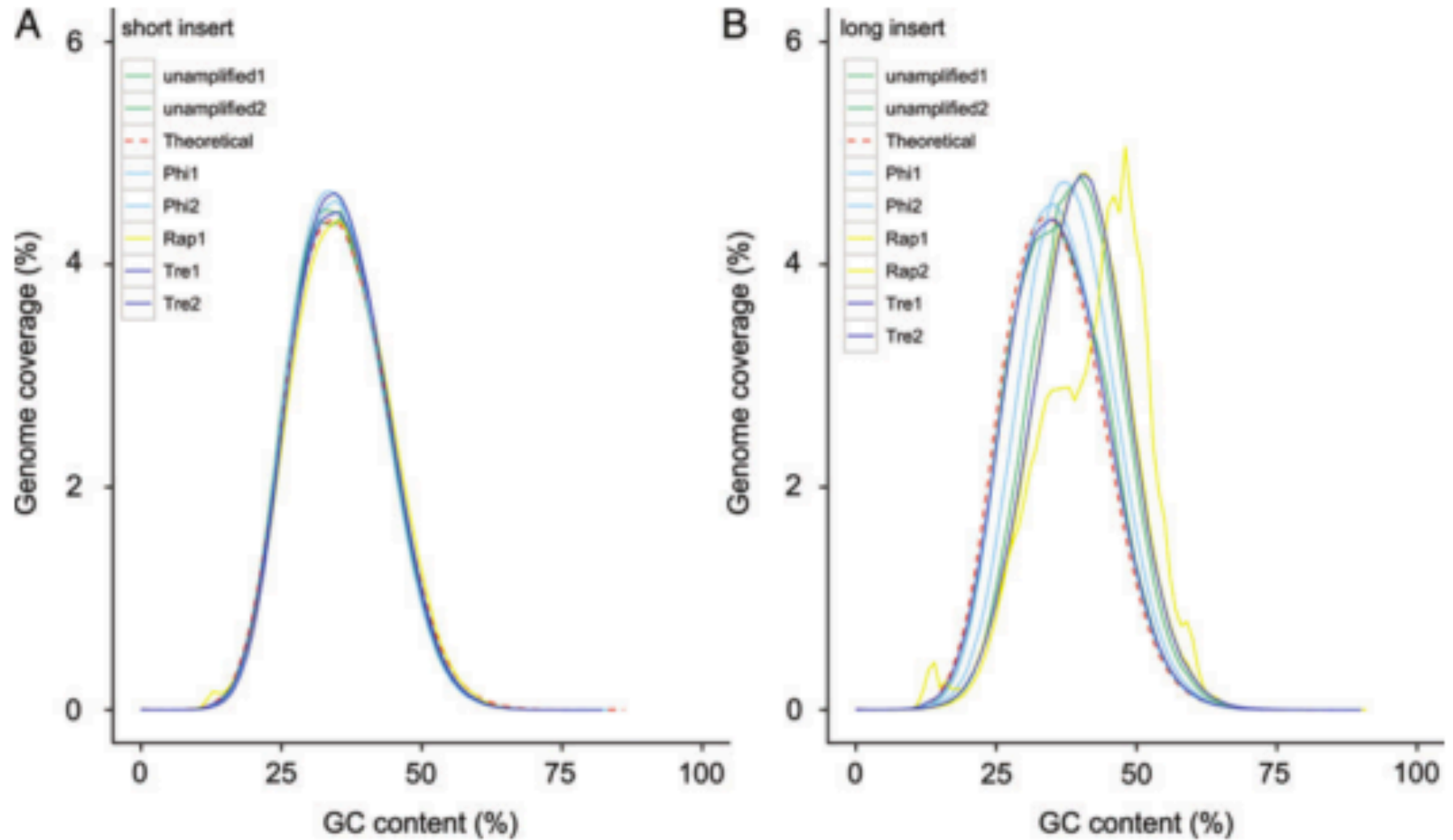
# Case study: Check lane quality and assembly

PE from one of my projects											
Directory	TotalReads	MappedReads	MappedReads(%)	ForwardStrand	ForwardStrand(%)	ReverseStrand	ReverseStrand(%)	FailedQC	FailedQC(%)	Duplicates	Duplicates(%)
map.718-S1	3171334	2920504	92.09%	1461181	46.07%	1459323	46.02%	0	0%	28099	0.89%
map.KPN91	2963812	2926683	98.75%	1463794	49.39%	1462889	49.36%	0	0%	8464	0.29%
map.KPN92	38811800	37864479	97.56%	18931099	48.78%	18933380	48.78%	0	0%	614696	1.58%
map.NTU	151505774	140827017	92.95%	70410162	46.47%	70416855	46.48%	0	0%	72381797	47.77%
MP from a company in Japan (Nextera)											
3kb.map	29432914	28598764	97.17%	14280601	48.52%	14318163	48.65%	0	0%	2369419	8.05%
5kb.map	8887196	8436683	94.93%	4208676	47.36%	4228007	47.57%	0	0%	1455786	16.38%
MP from another institute											
Directory	TotalReads	MappedReads	MappedReads(%)	ForwardStrand	ForwardStrand(%)	ReverseStrand	ReverseStrand(%)	FailedQC	FailedQC(%)	Duplicates	Duplicates(%)
MP.2kb	87149392	61884157	71.01%	31033842	35.61%	30850315	35.40%	0	0%	6893810	7.91%
MP.4kb	92488172	60082343	64.96%	30124954	32.57%	29957389	32.39%	0	0%	6688542	7.23%
MP.6kb	79969510	50558184	63.22%	25273991	31.60%	25284193	31.62%	0	0%	3919754	4.90%
MP.9kb	63262972	44161175	69.81%	22132740	34.99%	22028435	34.82%	0	0%	6938278	10.97%
a project from BRC											
Directory	TotalReads	MappedReads	MappedReads(%)	ForwardStrand	ForwardStrand(%)	ReverseStrand	ReverseStrand(%)	FailedQC	FailedQC(%)	Duplicates	Duplicates(%)
MP10kb.map	4809196	4757184	98.92%	2380384	49.50%	2376800	49.42%	0	0%	31589	0.66%
MP15kb.map	4557418	4492023	98.57%	2247623	49.32%	2244400	49.25%	0	0%	101159	2.22%
MP4kb.map	5349212	5266083	98.45%	2633803	49.24%	2632280	49.21%	0	0%	26721	0.50%
MP6kb.map	5185824	5129611	98.92%	2566177	49.48%	2563434	49.43%	0	0%	30809	0.59%





# Experiment biases



**Figure 5.** Distribution of GC content in sequenced reads of (A) short- and (B) long-insert libraries.

# Mapping output: A summary

There is a lot you can do from the initial mapping output

- Post mapping QC

- Assembly QC

At this point you should decide whether

- it's a good run and you can go ahead to the next stage

- you need additional run

- you need to abandon the whole run

Variant calling

# Variant calling

You have just:

- Mapped the reads to where they belong (supposedly)
- Provided accurate mapping quality scores

Next:

- Give the correct file (**BAM**) to variant callers

How to determine the above are correct?

# SNP discovery

## Heterozygous and homozygous SNP

10X

```
ATCGATGACTGACTGAATGGTTGAC
ATCGATGACTGACTGAATGGTTGAC
ATCCATGACTGACTGAATGGTTGAC
ATCGATGACTGACTGAATGGTTGAC
ATCGATAACTGACTGAATGGTTGAC
ATCGATGACTGAGTGAATGGTTGAC
ATCGATGACTGAGTGAATGGTTGAC
ATCGATGACTGAGTGAATGGTTGAC
ATCGATGACTGAGTGAATGGTTGAC
ATCGATGACTGAGTGAATGGTTGAC
```

---

...ATCGATGACTGACTGACTGGTTGAC...

reference



# INDELS (insertion deletions) and Structural variations

## Indel examples

wild-type sequence

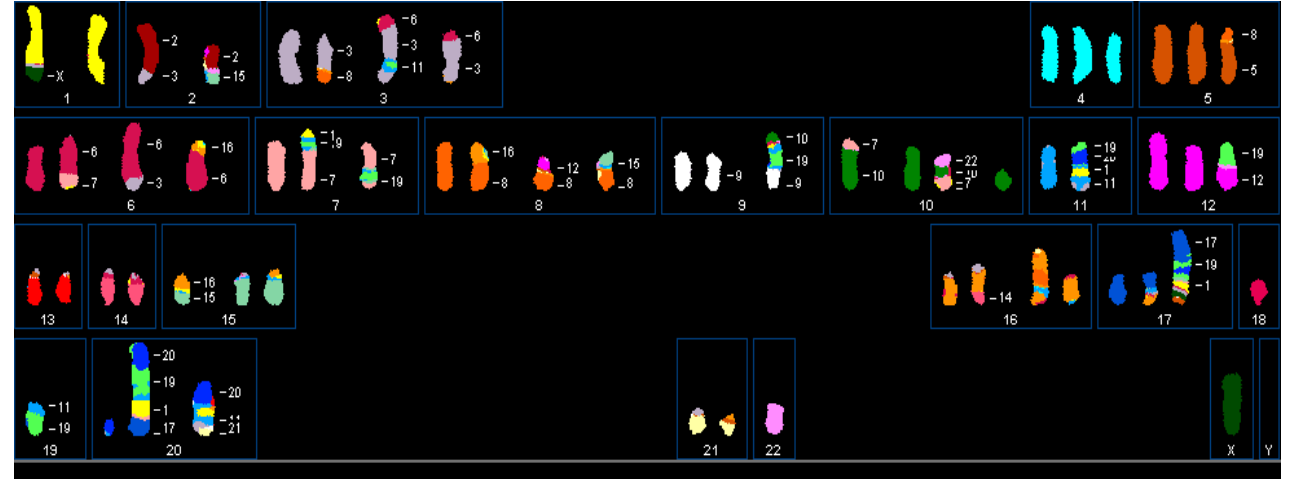
```
ATCTTCAGCCATAAAAAGATGAAGTT
```

3 bp deletion

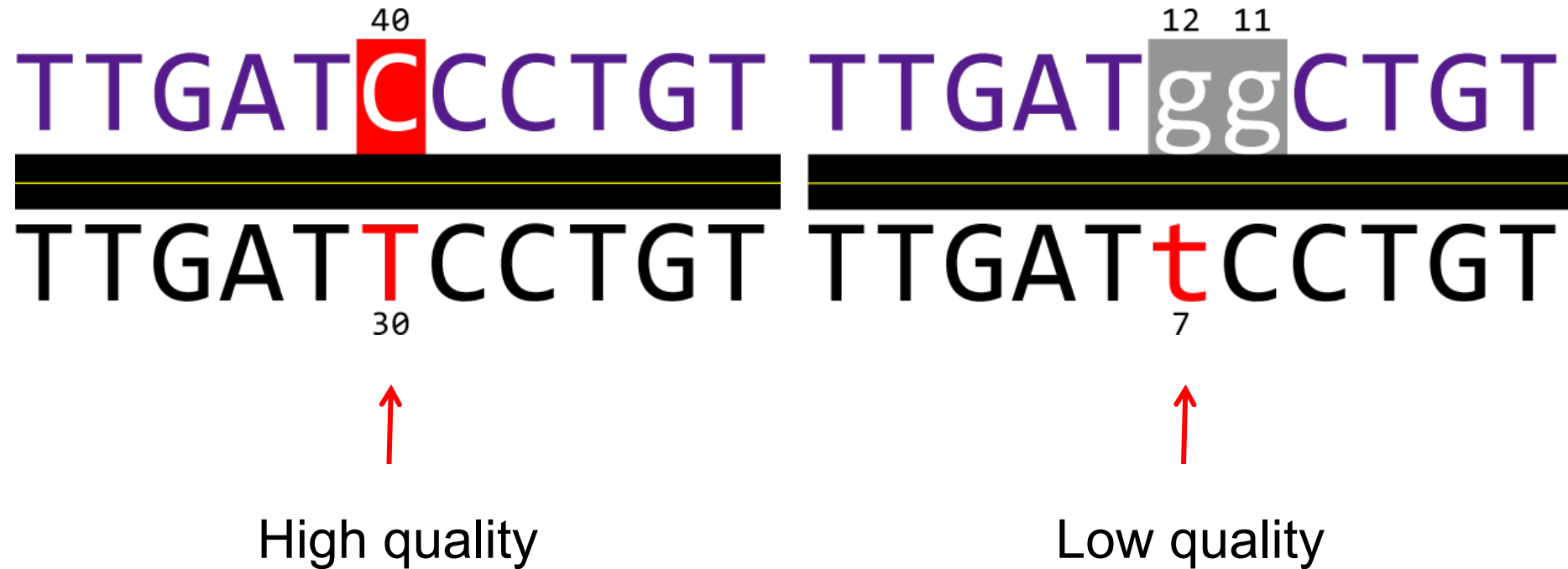
```
ATCTTCAGCCCAAAGATGAAGTT
```

4 bp insertion (orange)

```
ATCTTCAGCCCATATGTGAAAAGATGAAGTT
```



# SNP Discovery: Base Qualities



# SNPs & Bayesian Statistics

# of individuals      base quality      allele call in read

$$\Pr(G_1, G_2, \dots, G_n | B) = \frac{\prod_{i=1}^n \left[ \sum_{\forall T^k} \Pr(B_i | T_i^k) \Pr(T_i^k | G_i) \right] \Pr(G_1, G_2, \dots, G_n)}{\sum_{\forall G^l} \left\{ \prod_{i=1}^n \left[ \sum_{\forall T^k} \Pr(B_i | T_i^k) \Pr(T_i^k | G_i^l) \right] \Pr(G_1^l, G_2^l, \dots, G_n^l) \right\}}$$

# Strategies that improve variant calling

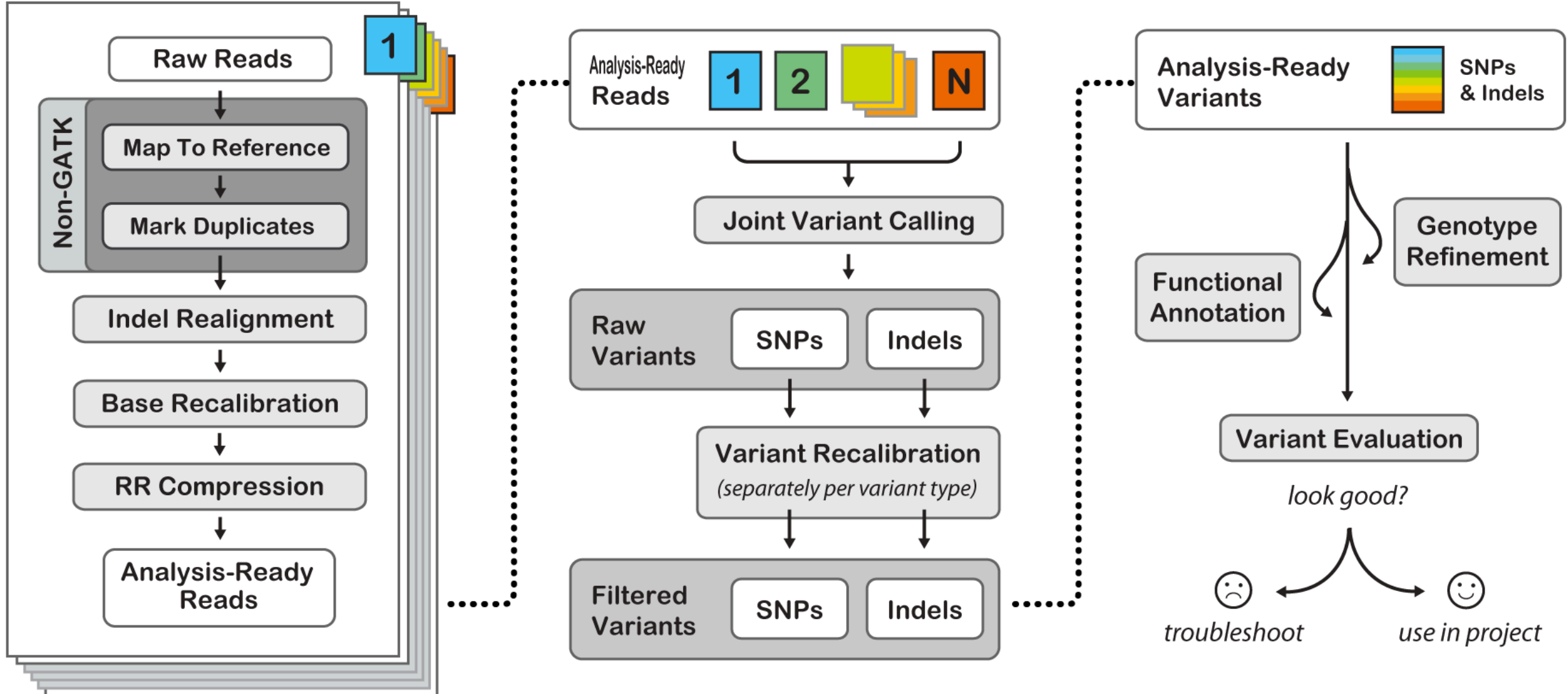
Data Pre-processing

>>

Variant Discovery

>>

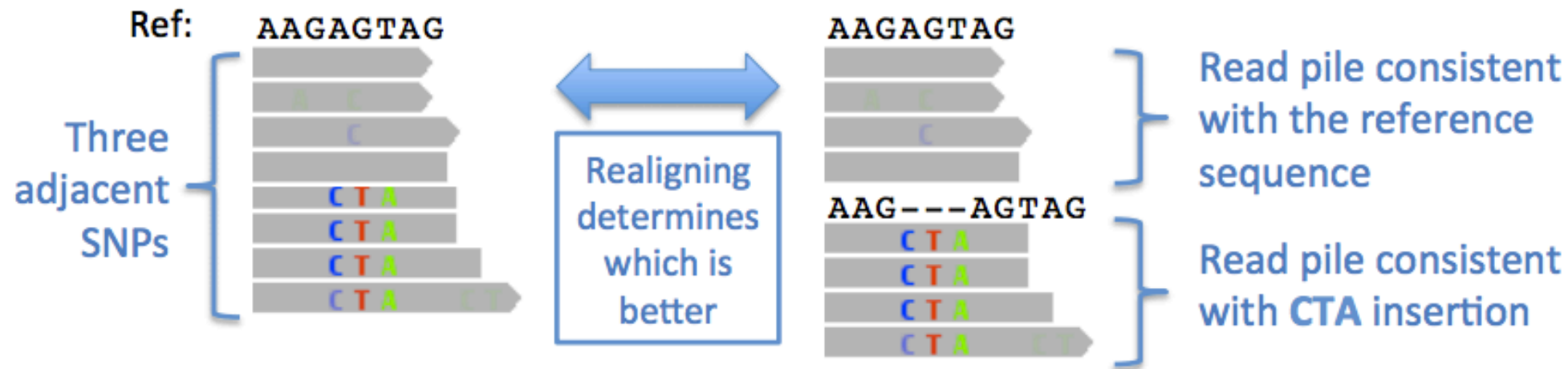
Preliminary Analyses





# Local realignment - principle

1. Find the best alternate consensus sequence that, together with the reference, best fits the reads in a pile (maximum of 1 indel)

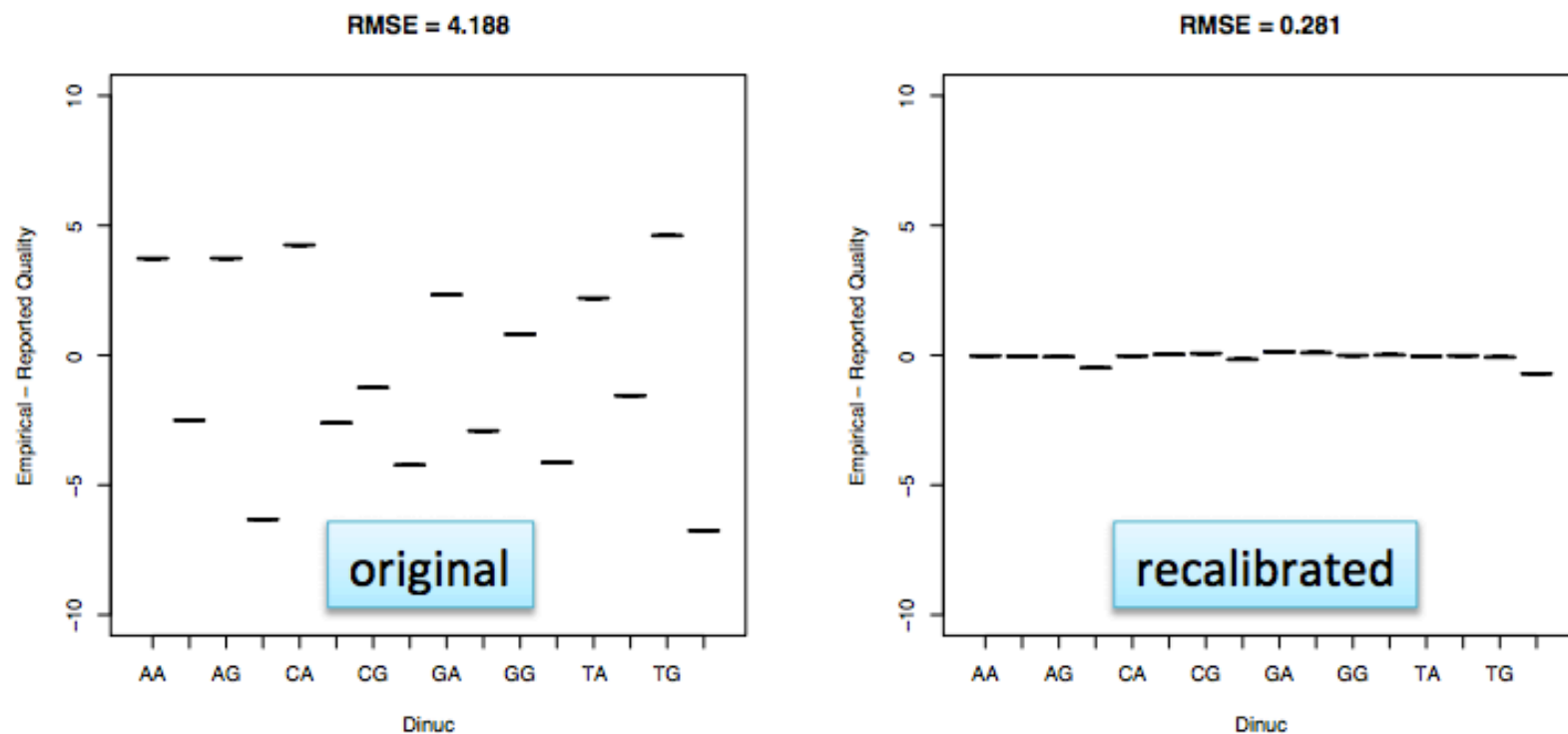


2. The score for an alternate consensus is the total sum of the quality scores of mismatching bases
3. If the score of the best alternate consensus is sufficiently better than the original alignments (using a LOD score), then we accept the proposed realignment of the reads

# Base quality recalibration

- Quality scores are critical for all downstream analysis
- Systematic biases are a major contributor to bad calls

Example of bias: qualities reported depending on nucleotide context

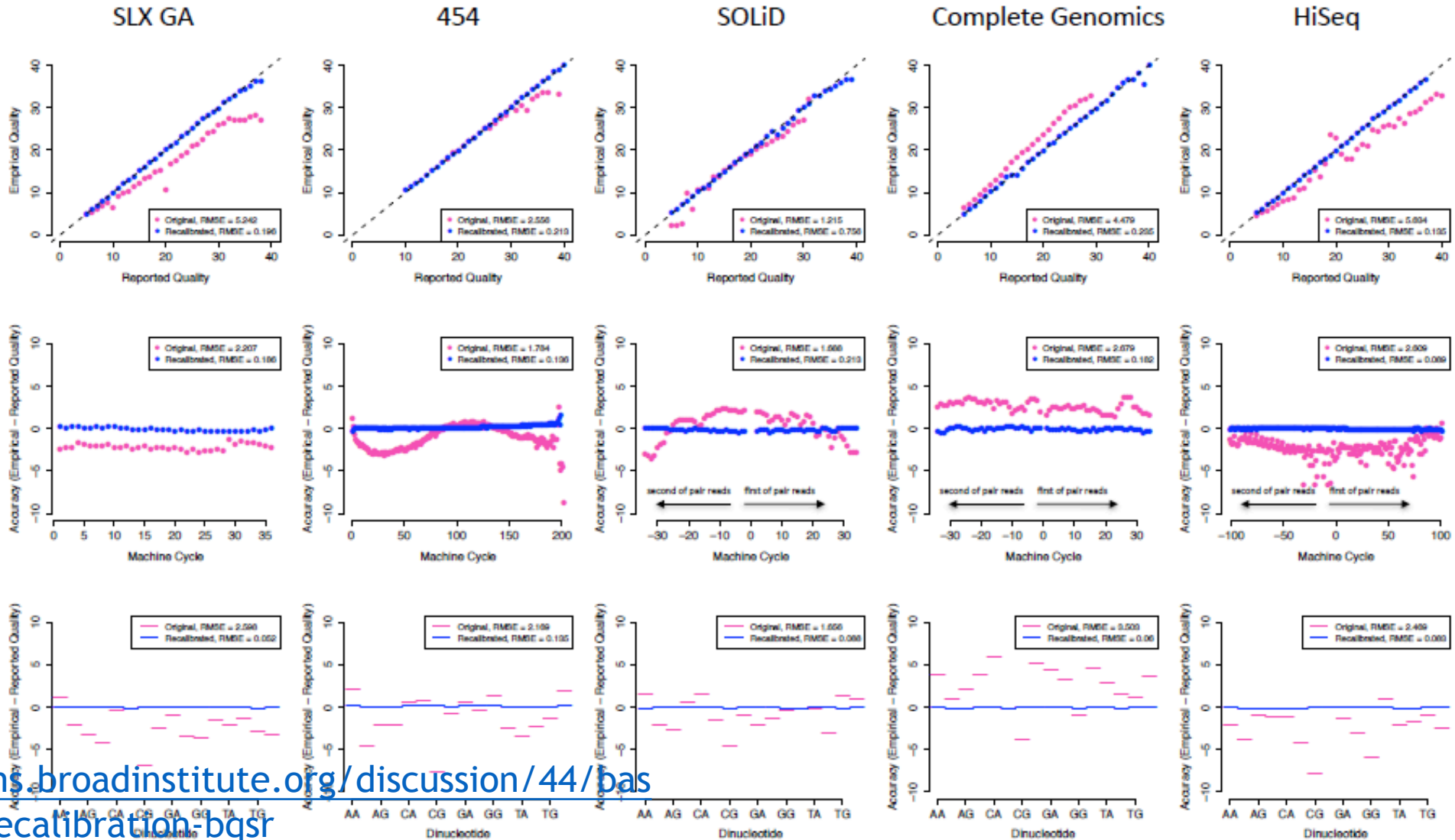


**BQSR method identifies bias and applies correction**



Highlighted as one of the major methodological advances of the 1000 Genomes Pilot Project!

# Base Quality Score Recalibration provides a calibrated error model from which to make mutation calls





# Improve beyond analysis-ready reads

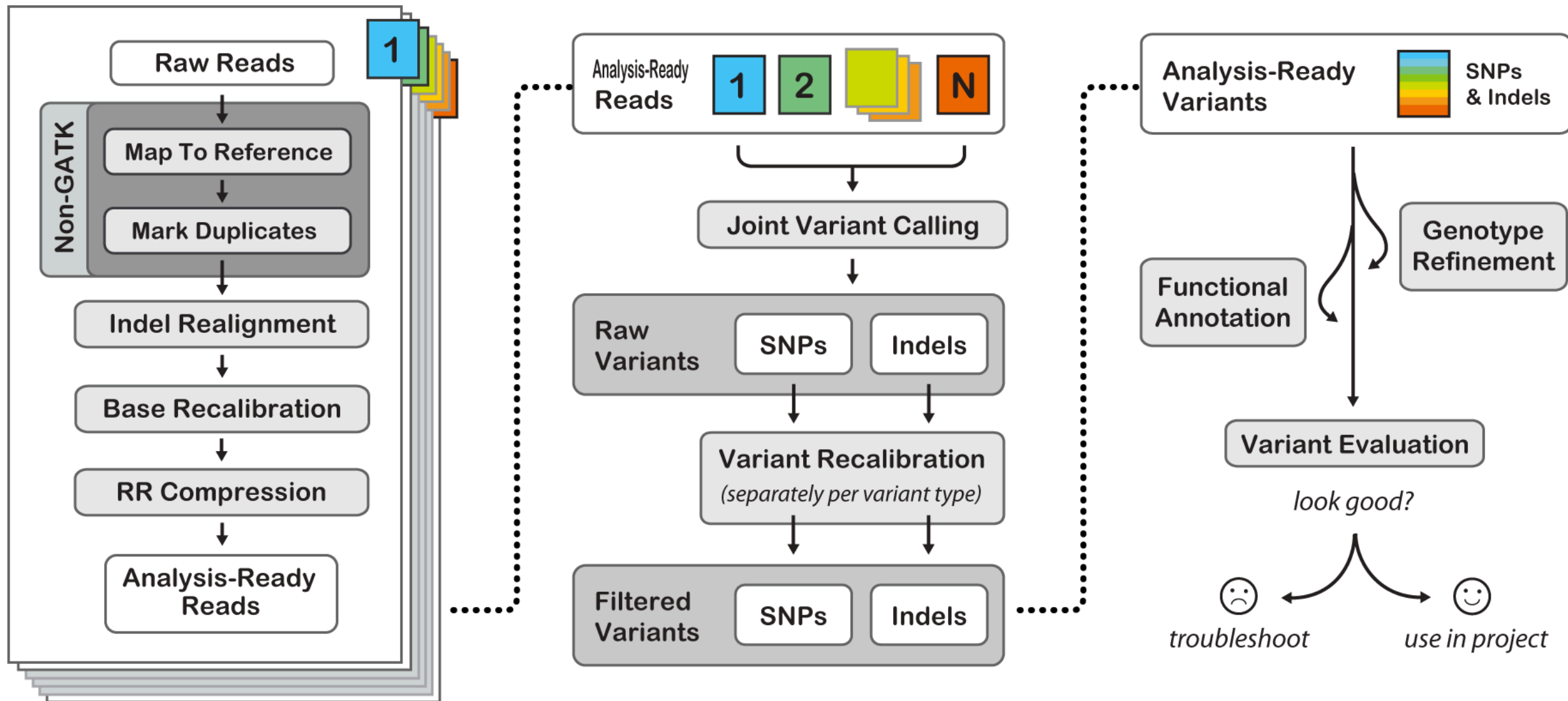
Data Pre-processing

>>

Variant Discovery

>>

Preliminary Analyses



# Using haplotypes for base calling

- Suppose that only 2 haplotypes have been observed in a population:

Chr1: .....A....T.....G.....  
Chr1: .....C....G.....A.....

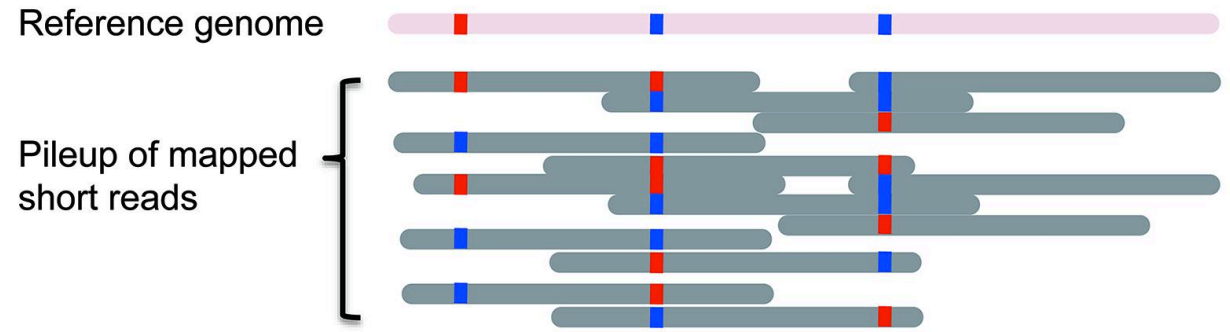
- And that you observe the following reads:

.....A....N.....G..  
..A....N.....G.....  
...A....N.....G...

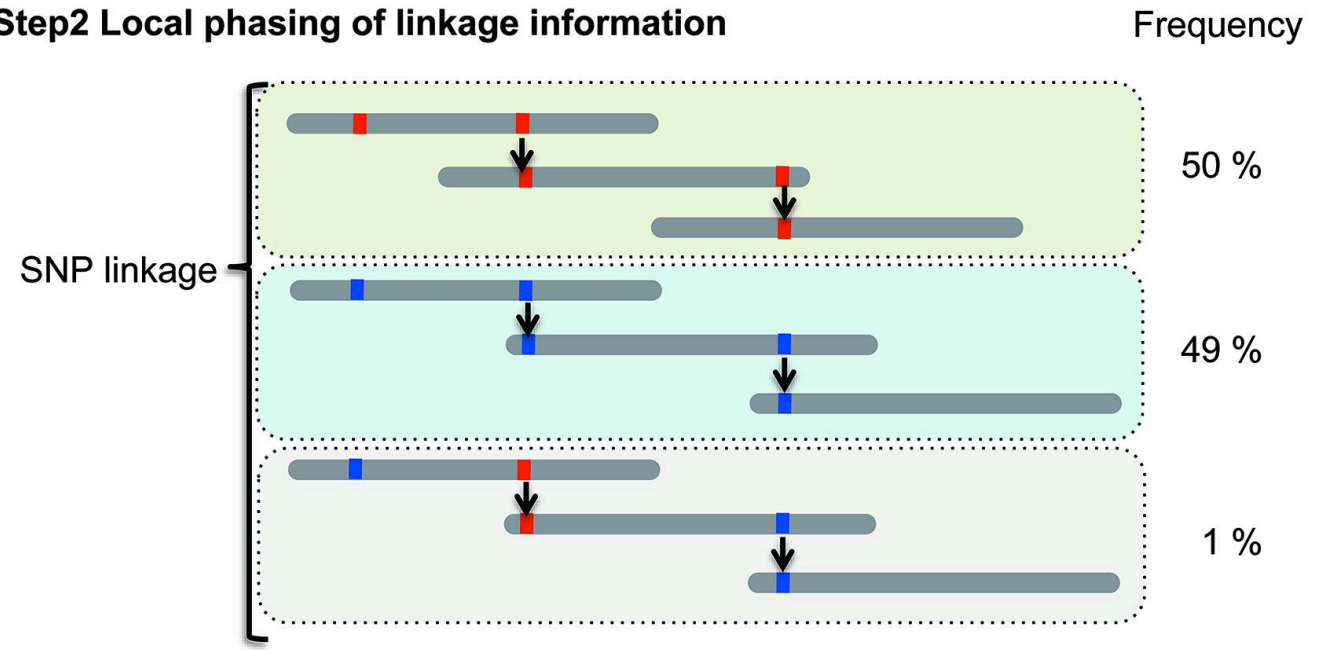
- Can you guess the value of **N** ?

# Building haplotypes

## Step1 Alignments

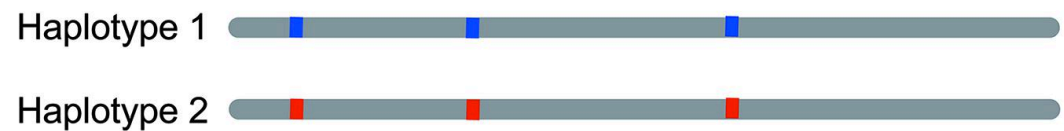


## Step2 Local phasing of linkage information

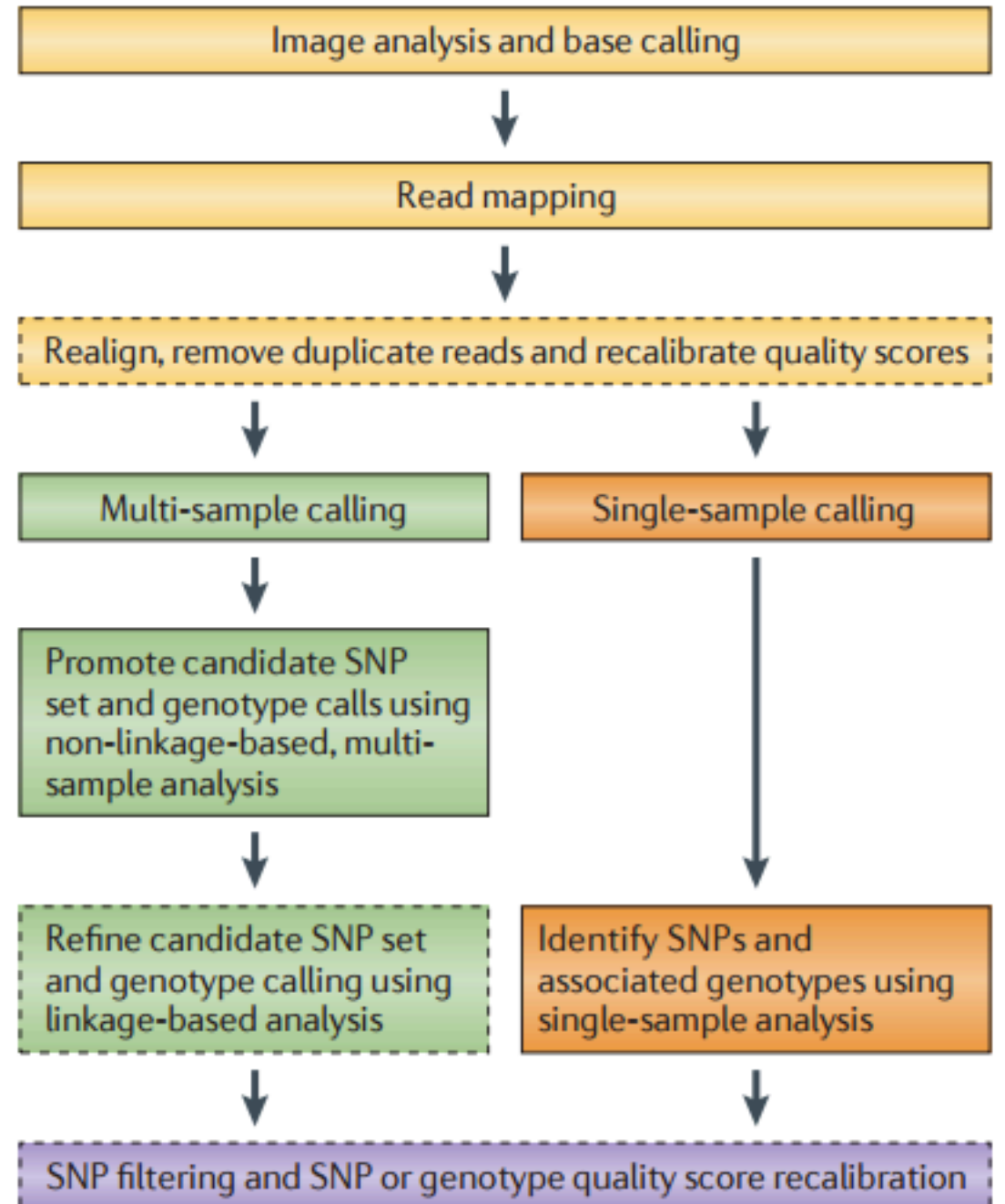


## Step3 Filtering

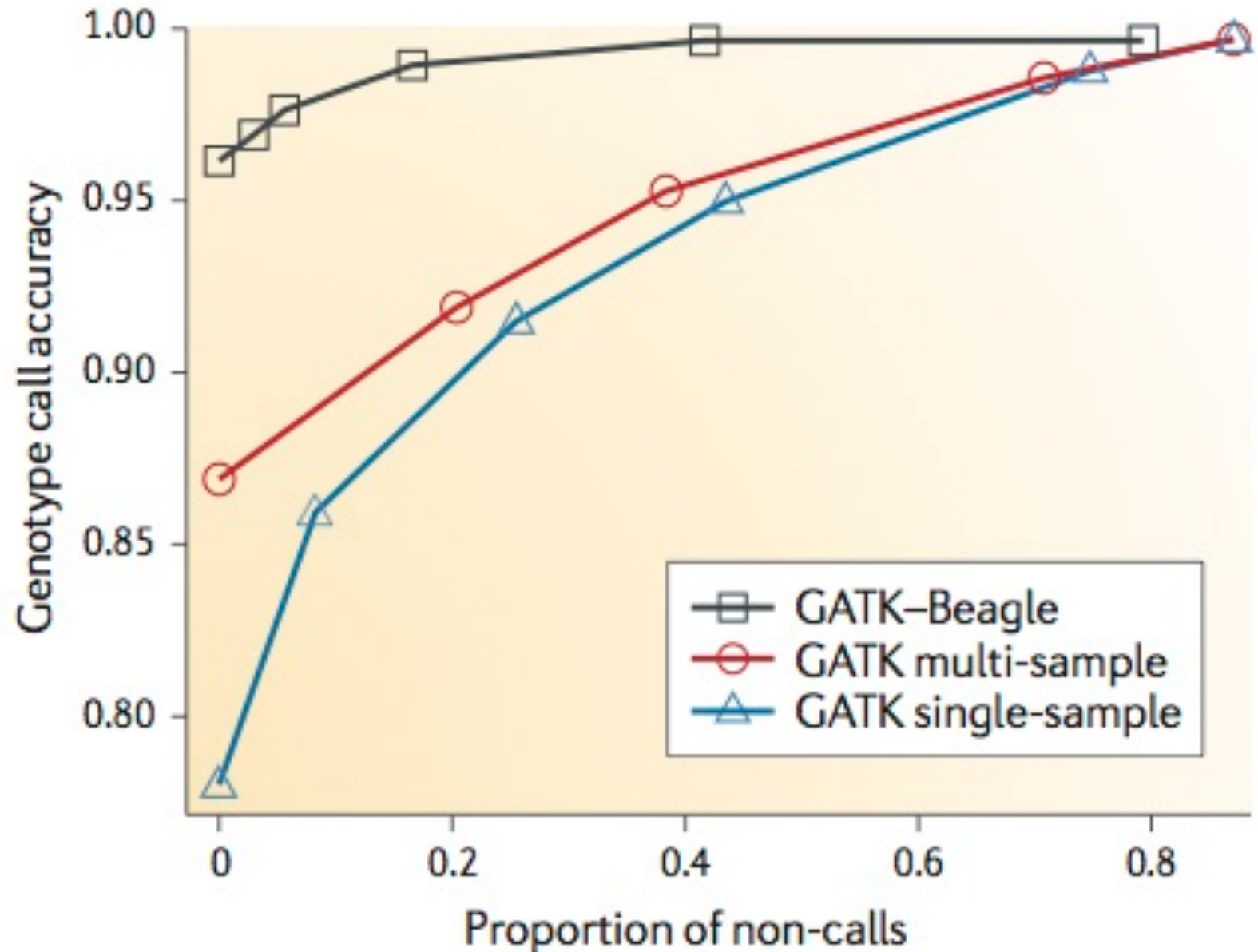
Minor haplotype is excluded.



# Use multiple samples



# Haplotype imputation increase genotype accuracy



**File size**

**File format**

**Tools**

**Time**

200 GB

**BAM files**

Recalibrated BQ, duplicates removed



**GATK**  
samtools  
freeBayes  
cortex\_var

10 hours

1 GB

**Raw variants (VCF)**

Sites with non-reference bases are  
genotyped

# VCF format

```
##fileformat=VCFv4.2
```

Mandatory header line

```
##filedate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
```

Mandatory header line

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00003
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3	0/0:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2
20	1234567	microsat1	GTC	G,GTCT	50	PASS	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4	0/2:17:2	1/1:40:3

Reference base

Alternative base

Quality score

Allele frequency, read depth, etc.

<https://samtools.github.io/hts-specs/VCFv4.2.pdf>

# Variant filtering

Raw variant calls have a lot of false positives.

How to filter?

Which one do you look at first?

Manual filtering based on different parameters

allele frequency, quality score, depth of coverage...

Location (contig ends SNPs are usually inaccurate)

Case by case

**look at the strongest effect filter**



# Annotating variants

- Annotations using reference genomes

Programs available: SNP-eff, annovar

- Calculate effects:

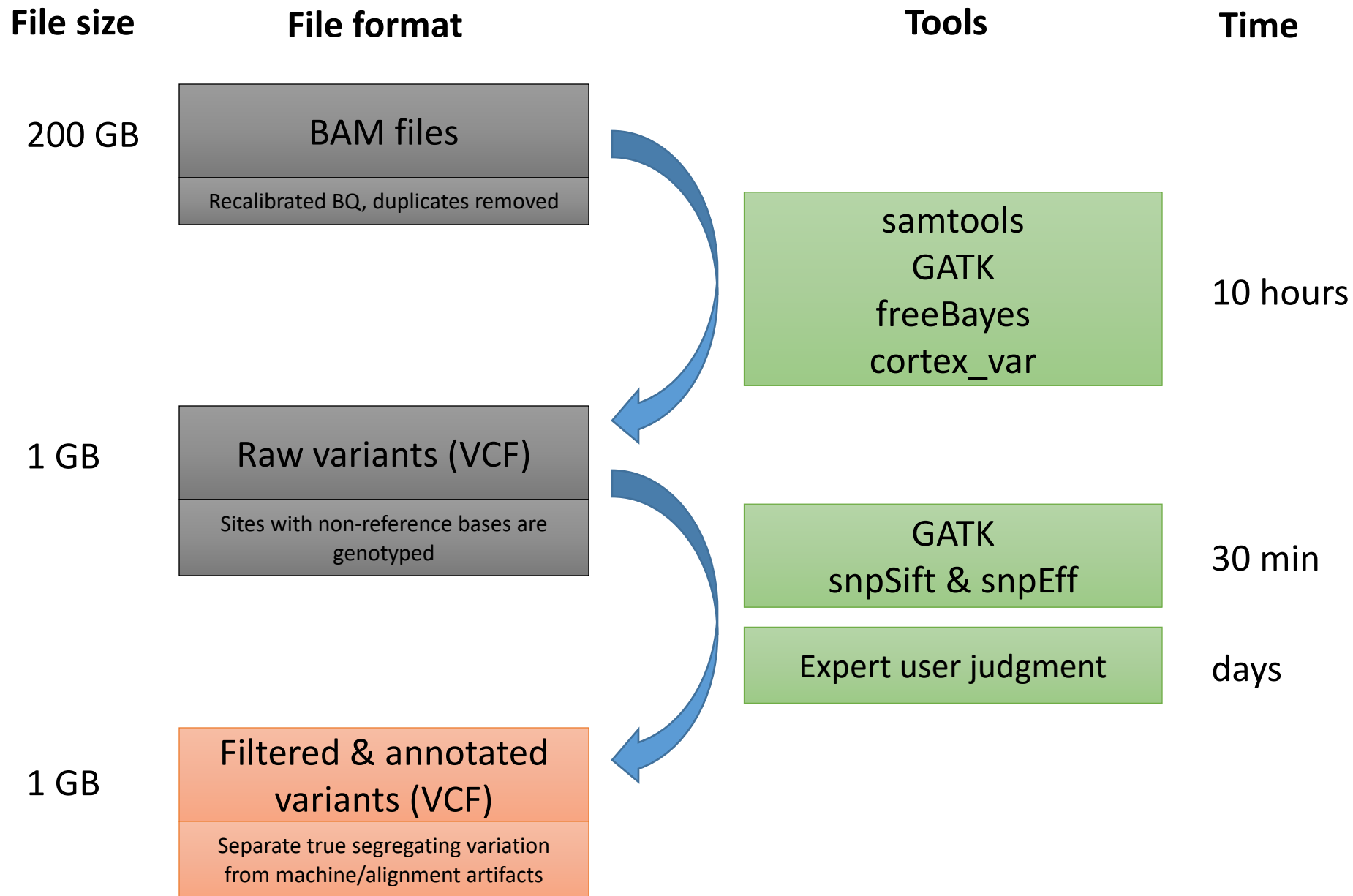
- Coding (e.g. Syn, Non-Syn, Stop gained, Splice)
- Non-coding (e.g. TFBS)

One of the mostly intensively research areas:

Linking variation to function

Unfortunately, only applicable to humans

For a new species, you have to start from scratch



Structural variation (short reads)



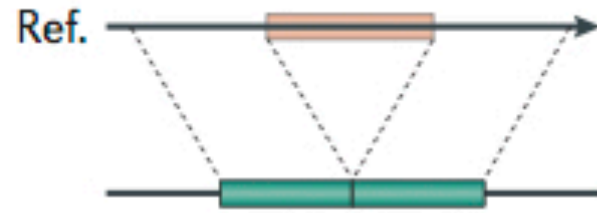
# More difficult structural variants

**Structural Variants (SVs):** Genomic rearrangements that affect **>50bp (or 100bp, or 1Kb)** of sequence, including:

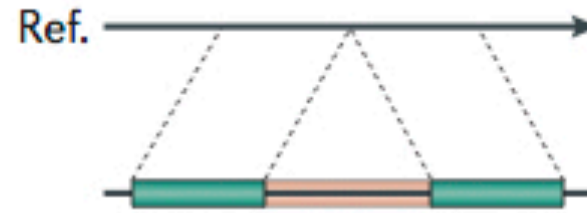
- deletions
- novel insertions
- inversions
- mobile-element transpositions
- duplications
- translocations

# SV classes

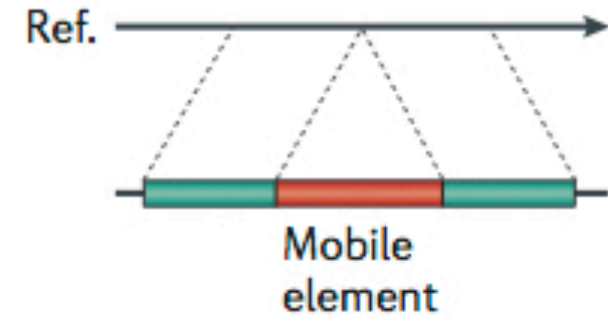
**Deletion**



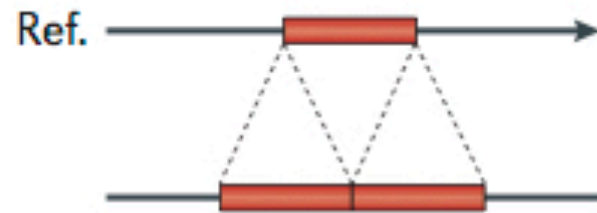
**Novel sequence insertion**



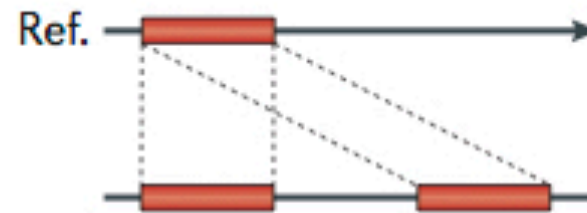
**Mobile-element insertion**



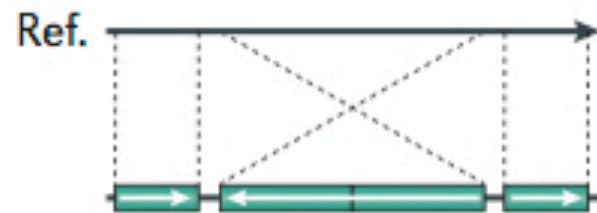
**Tandem duplication**



**Interspersed duplication**



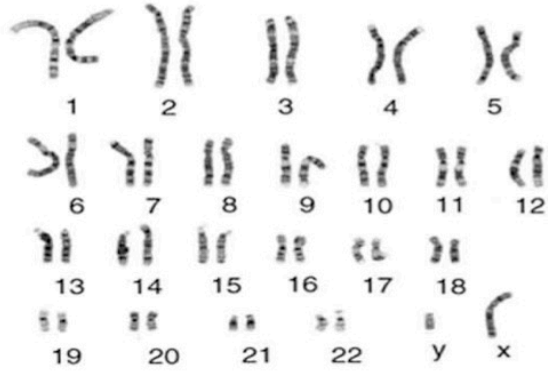
**Inversion**



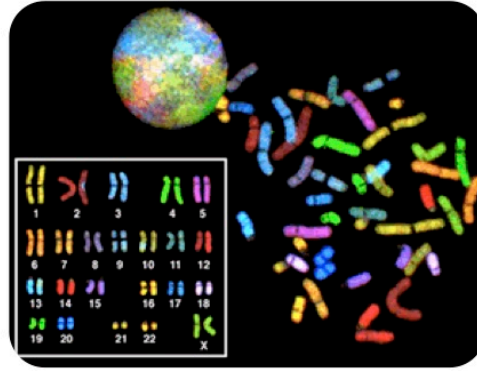
**Translocation**



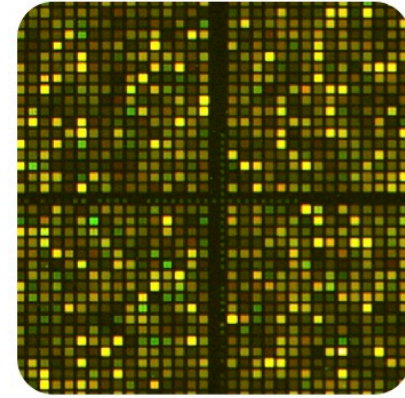
# Again, our understanding is driven by technology



1940s - 1980s  
Cytogenetics / Karyotyping



1990s  
CGH / FISH /  
SKY / COBRA



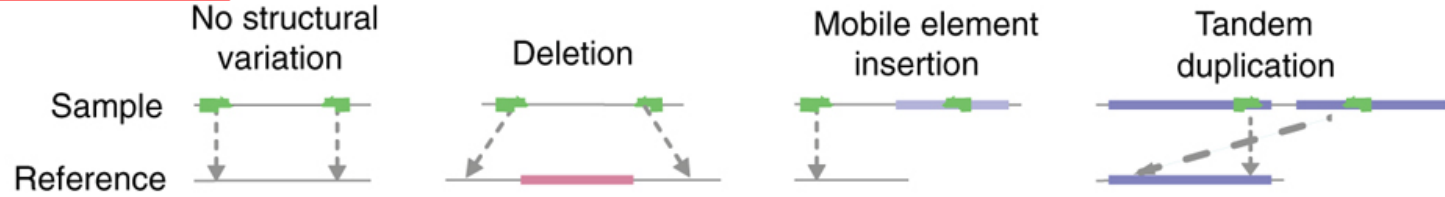
2000s  
Genomic microarrays  
BAC-aCGH / oligo-aCGH



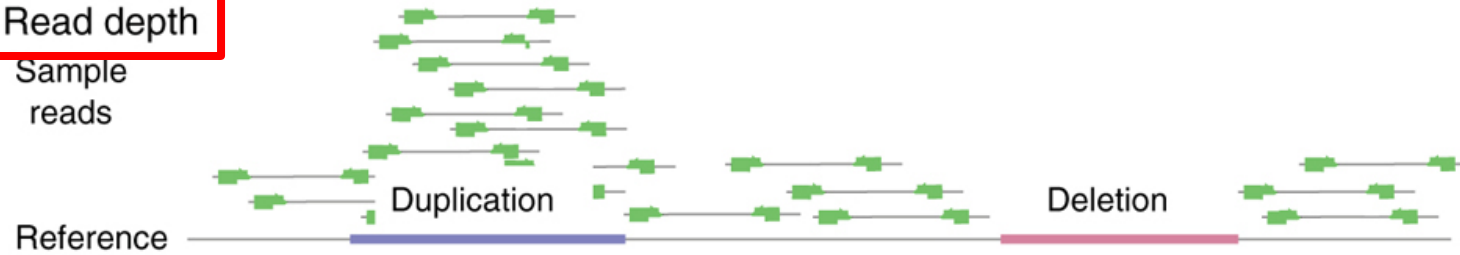
**Today**  
High throughput  
DNA sequencing

# Strategies for calling SVs from NGS data

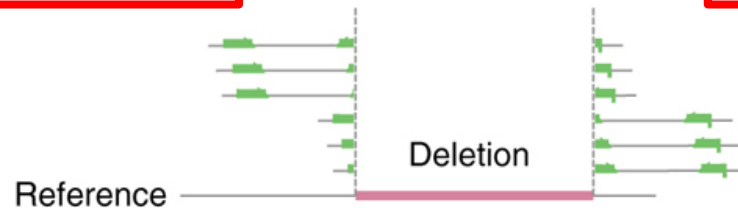
## 1. Read pairs



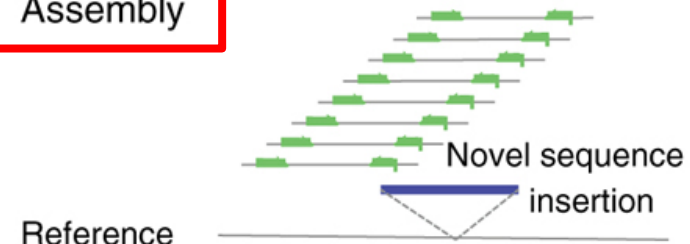
## 2. Read depth



## 3. Split reads

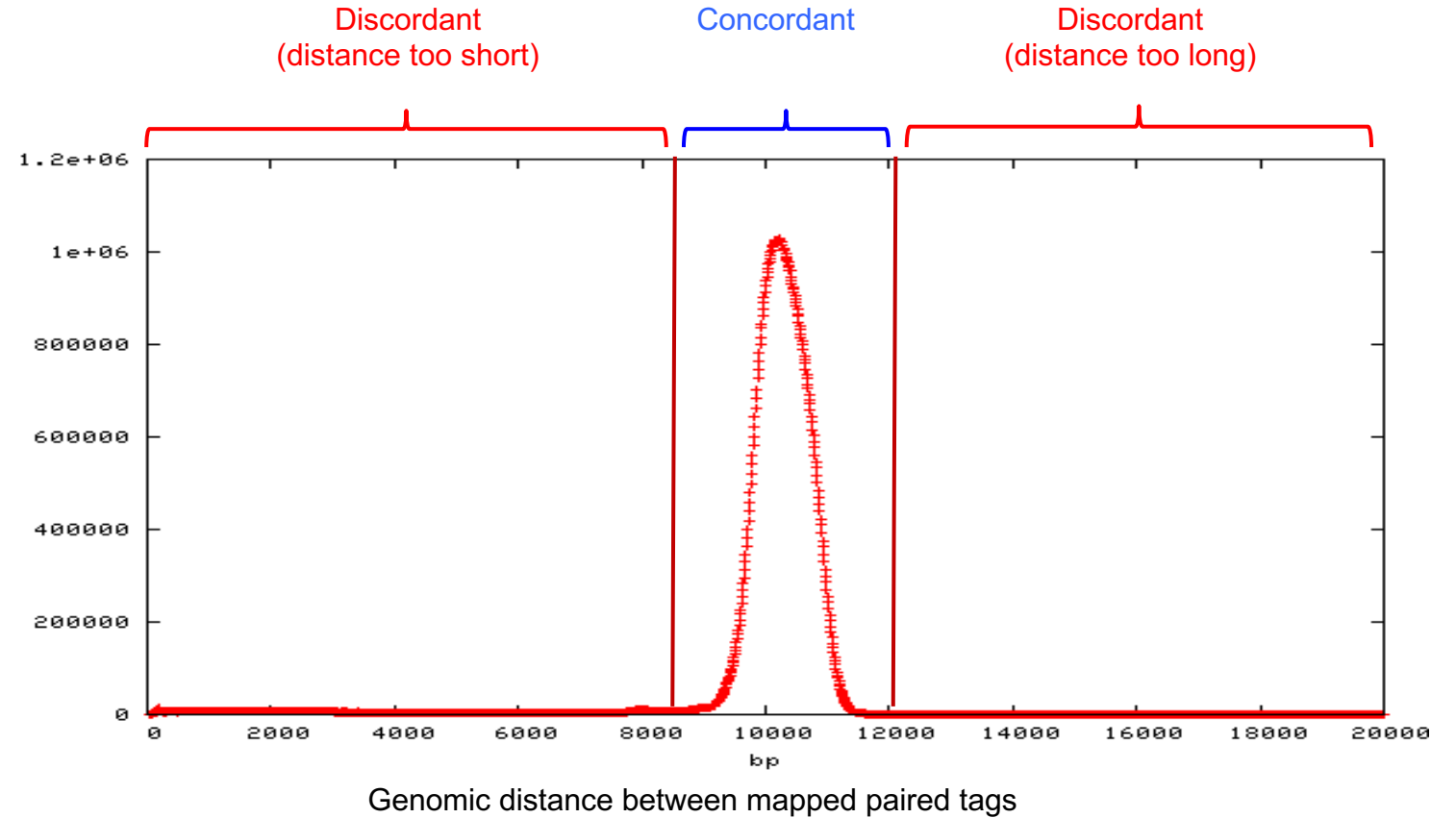
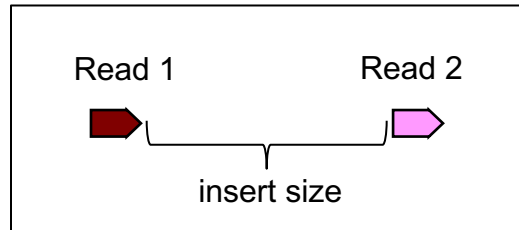


## 4. Assembly



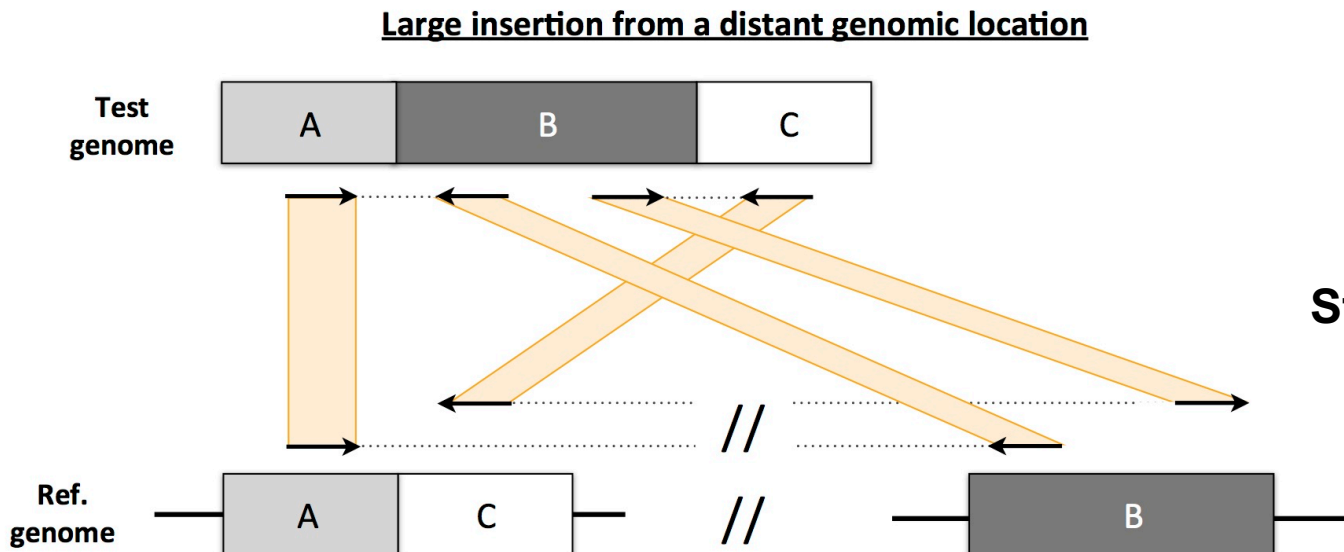
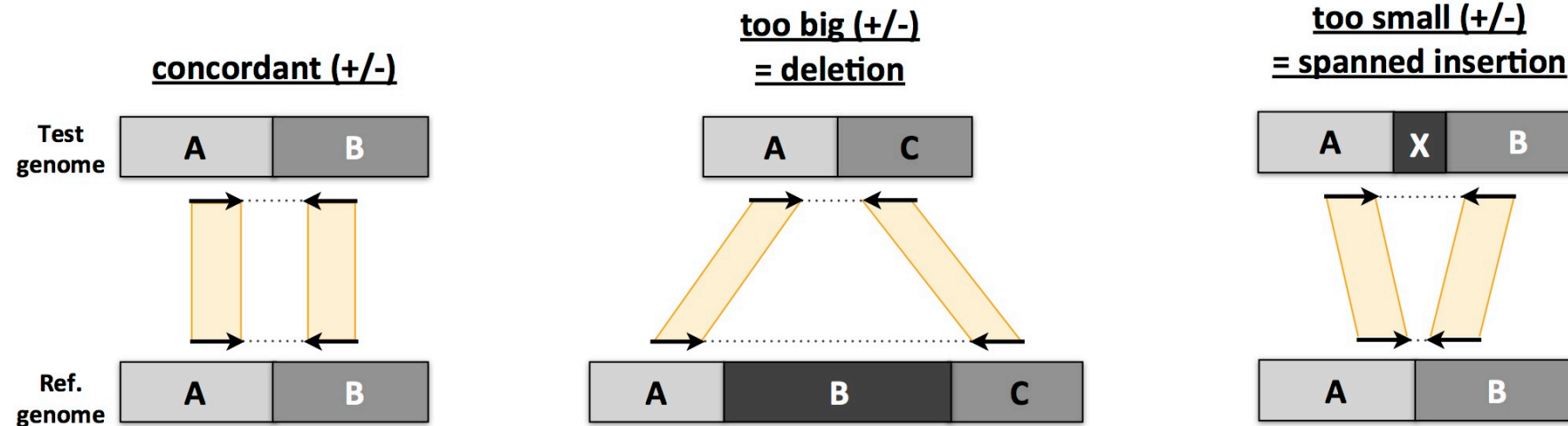


# Discordant read pairs



Reads pairs are also **Discordant** when order or orientation isn't as expected.  
Do they fall into particular region of the assembly?

# Using discordant reads to detect SVs



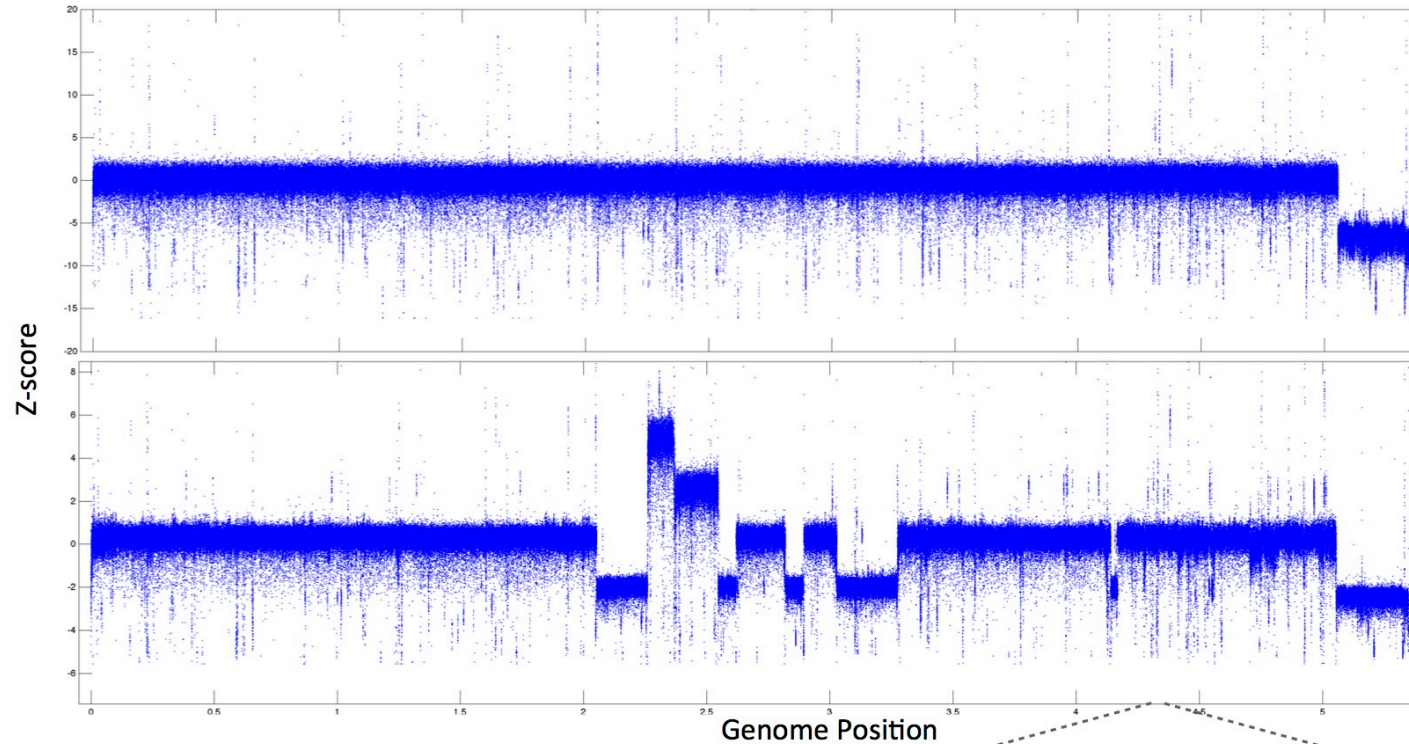
## Weaknesses

- Difficult to interpret read-pairs in repetitive regions
- Difficult to fully characterize highly rearranged regions
- High rate of false positives

## Strengths:

- Most classes of variation can, in principle, be detected

# Read-depth

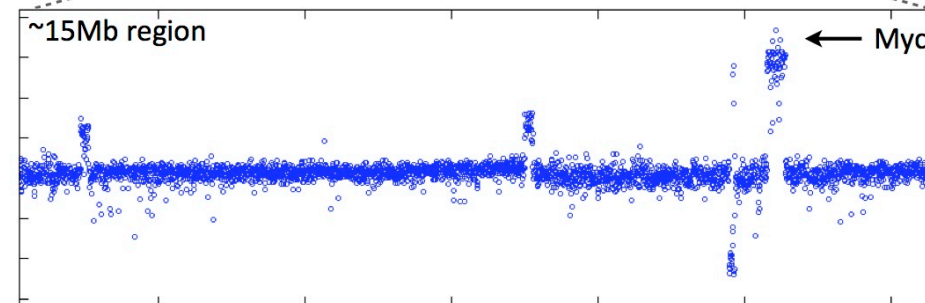


## Strengths:

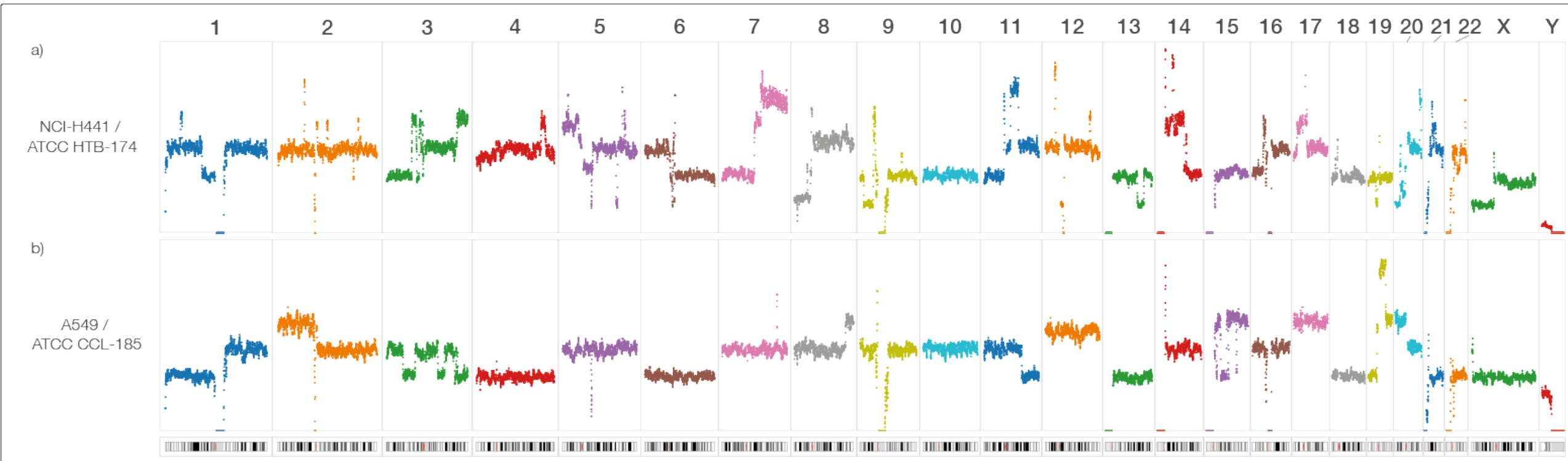
- 1) Fast and simple.
- 2) Easy to identify gene amplifications.
- 3) Relatively straightforward interpretation: is gene X amplified or deleted?

## Weaknesses:

- 1) Limited resolution (5-10kb) = imprecise boundaries
- 2) Cannot detect balanced events or reveal variant architecture.

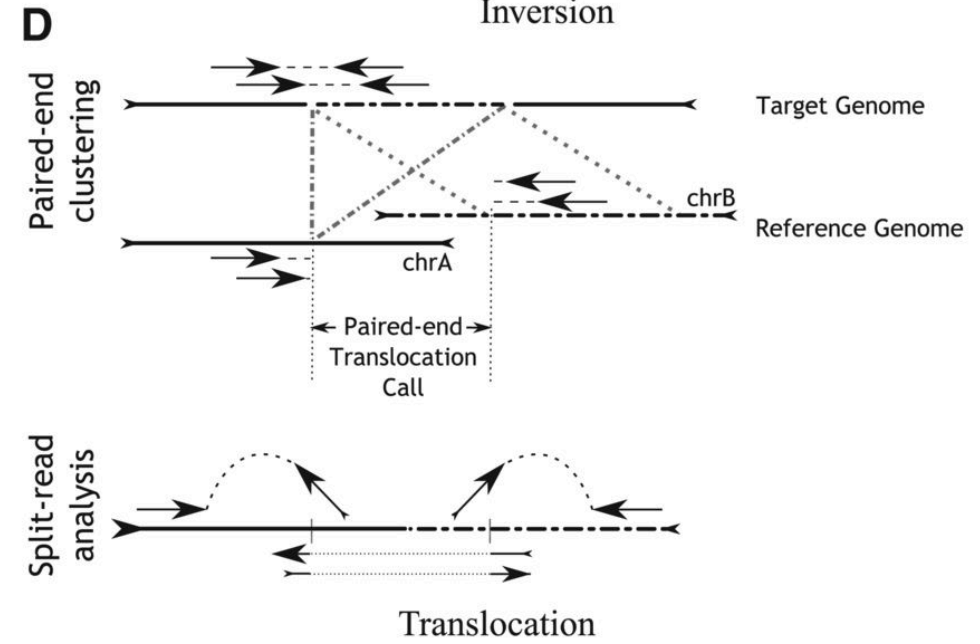
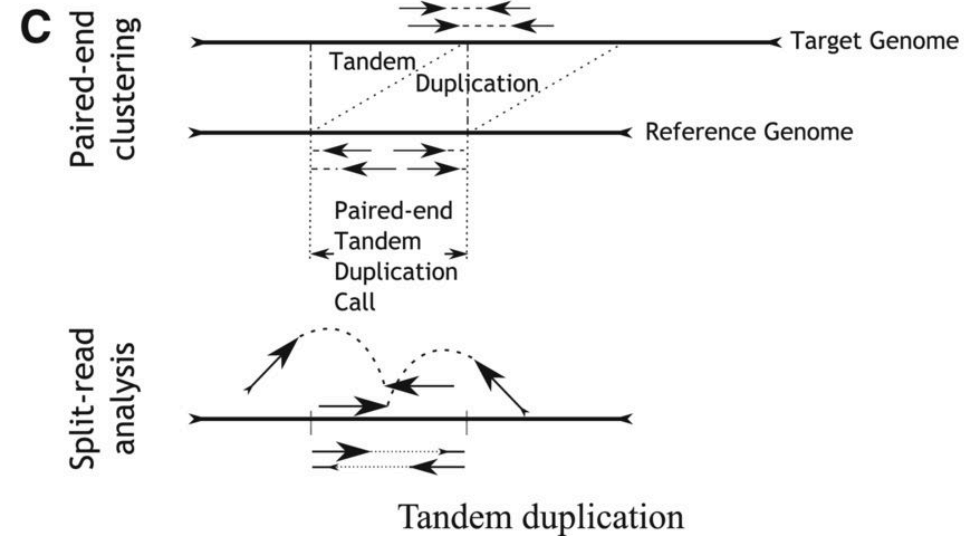
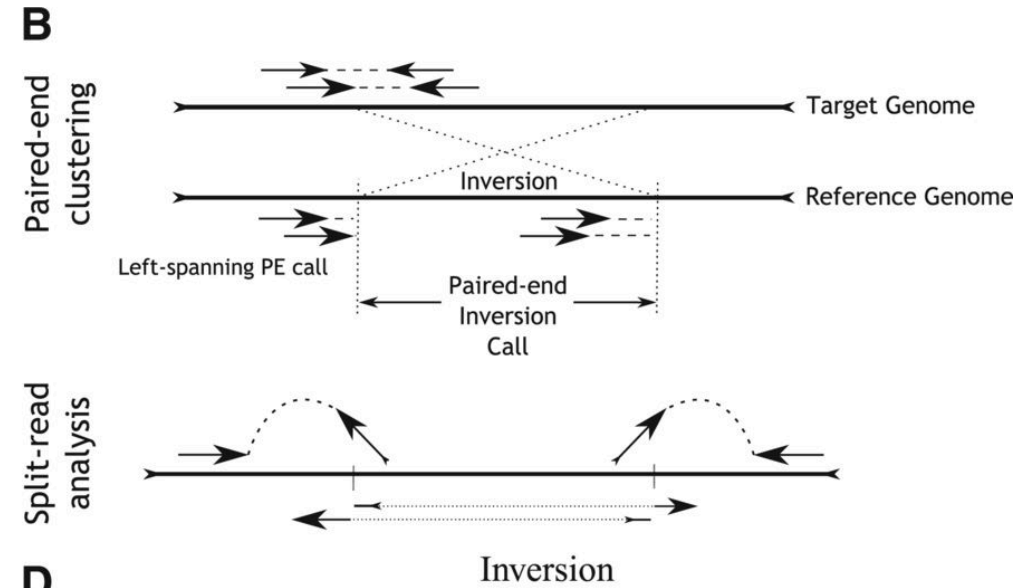
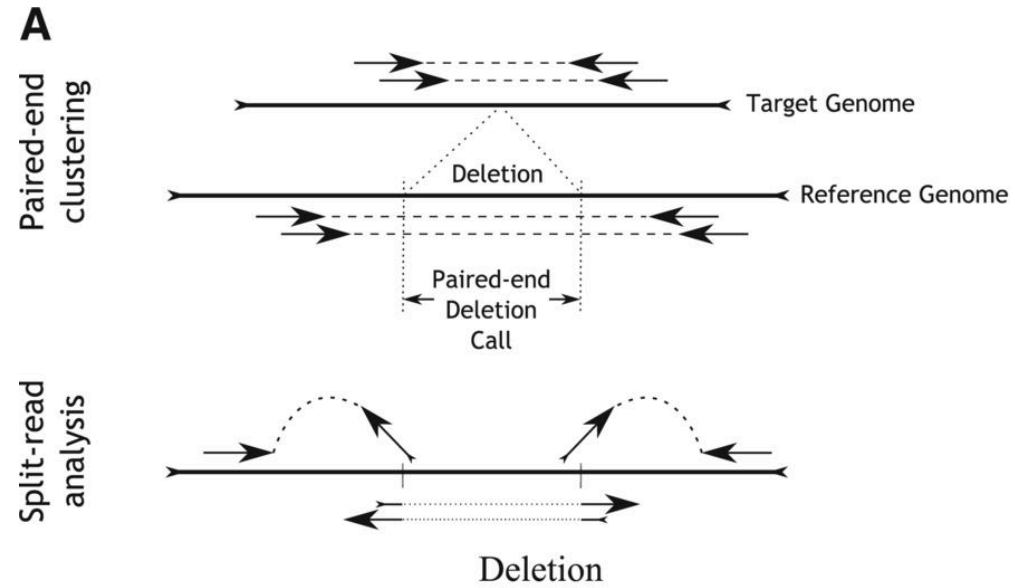


# Read-depth can be used to call aneuploidies



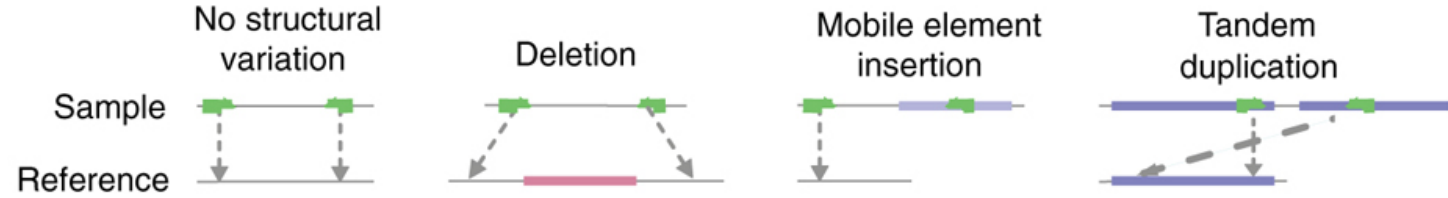
Whole-genome sequencing of two lung cancer cell lines. Each has a different pattern of duplications, deletions and translocations a) cell line H441 b) cell line A549

# Split reads

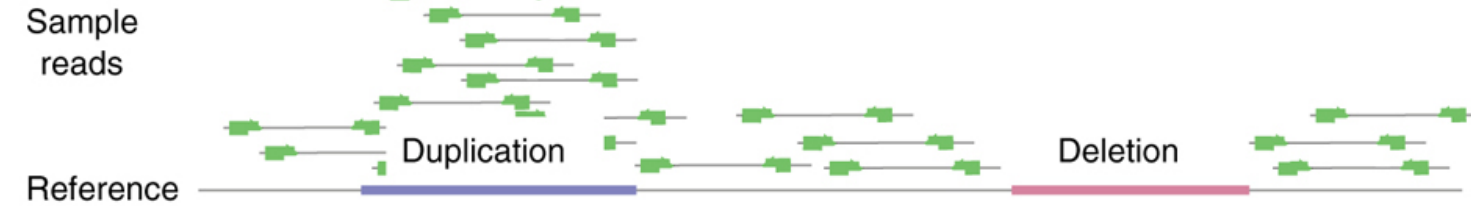


# Strategies for calling SVs from NGS data

Read pairs



Read depth

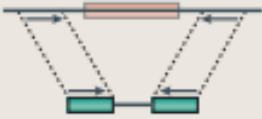

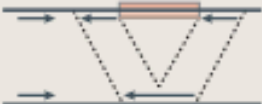
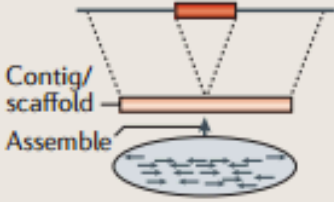
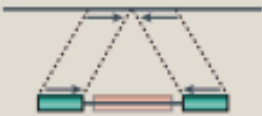
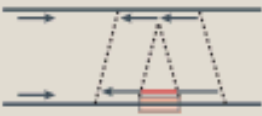
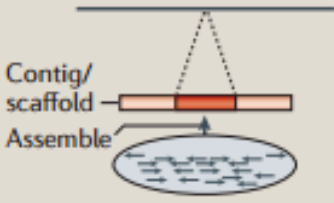
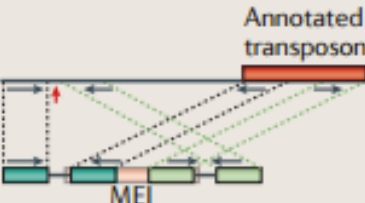
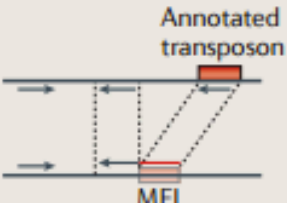
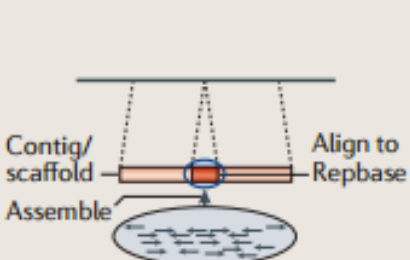
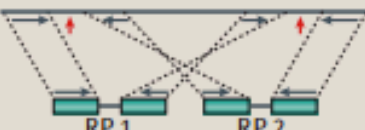
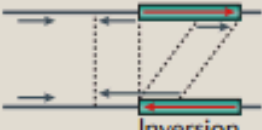
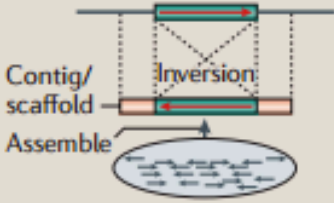


Split reads



4. Assembly

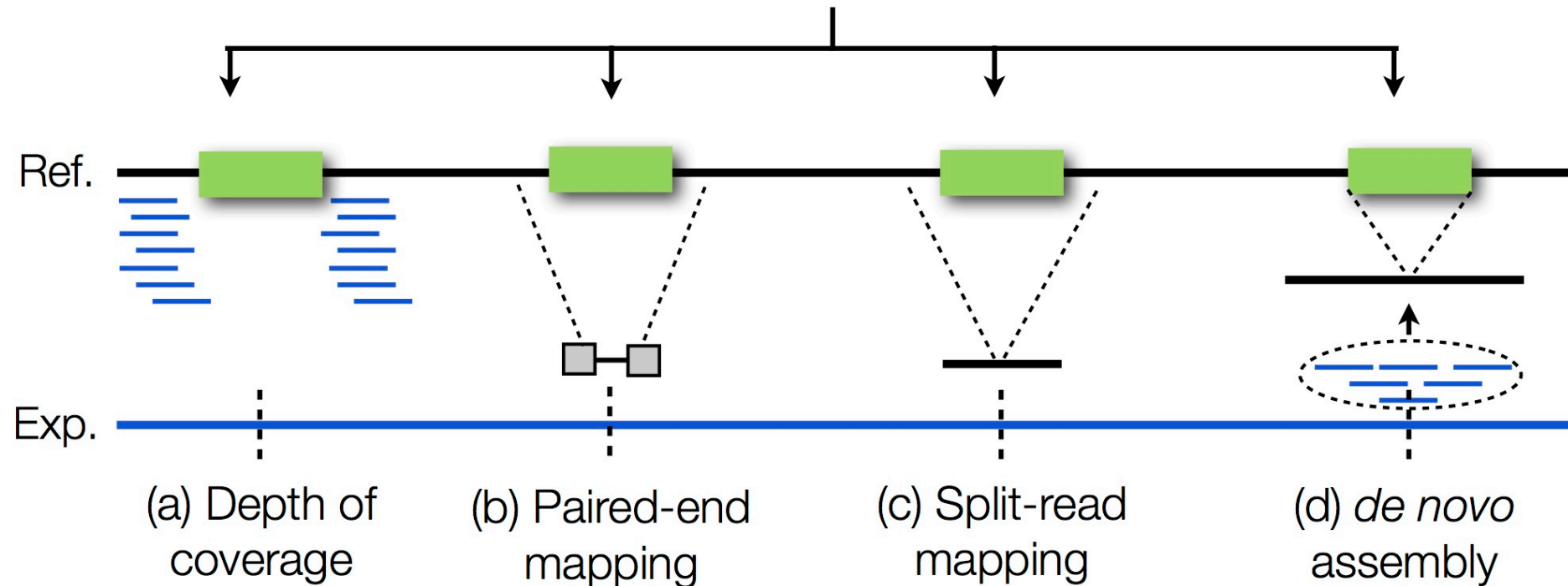
# De novo assembly for SVs

SV classes	Read pair	Read depth	Split read	Assembly
Deletion				
Novel sequence insertion		Not applicable		
Mobile-element insertion		Not applicable		
Inversion		Not applicable		

# Summary of strategies for calling SVs (short reads)

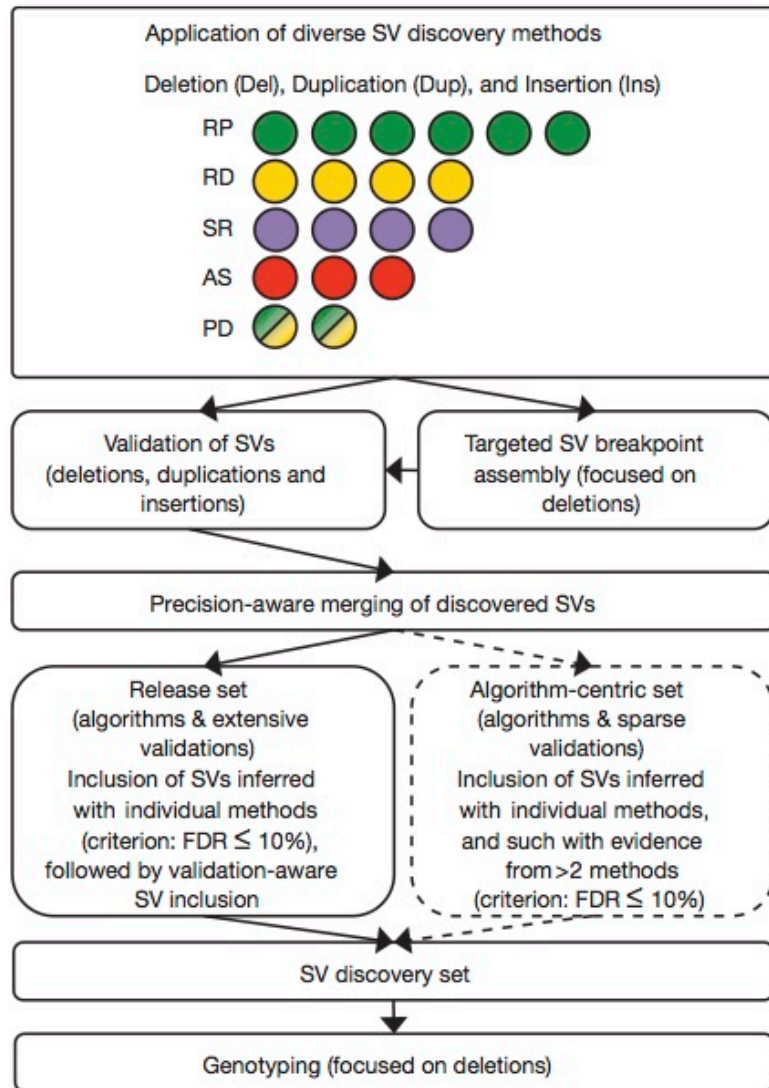
1. Align DNA sequences from sample to human reference genome

2. Look for evidence of structural differences

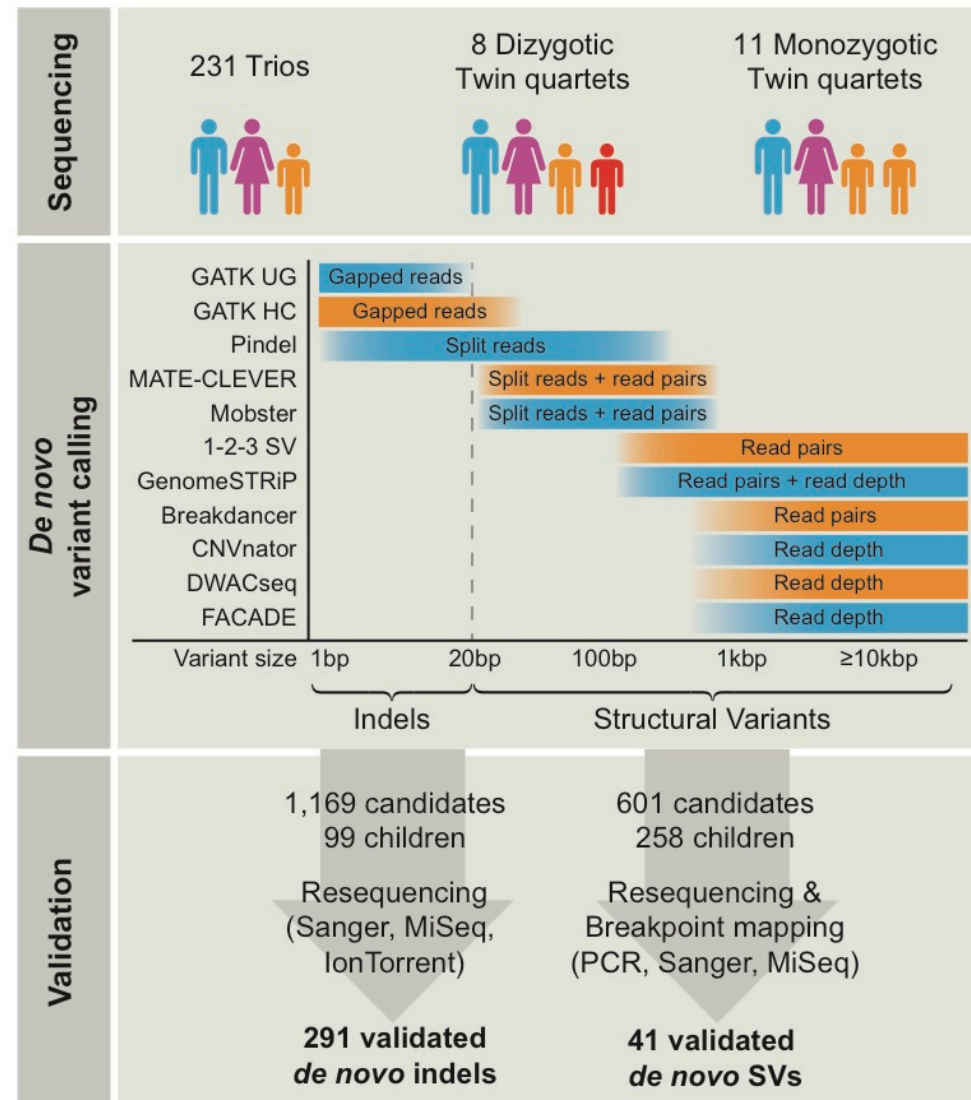




# Bottom line for short reads calling SVs : try many methods and validate

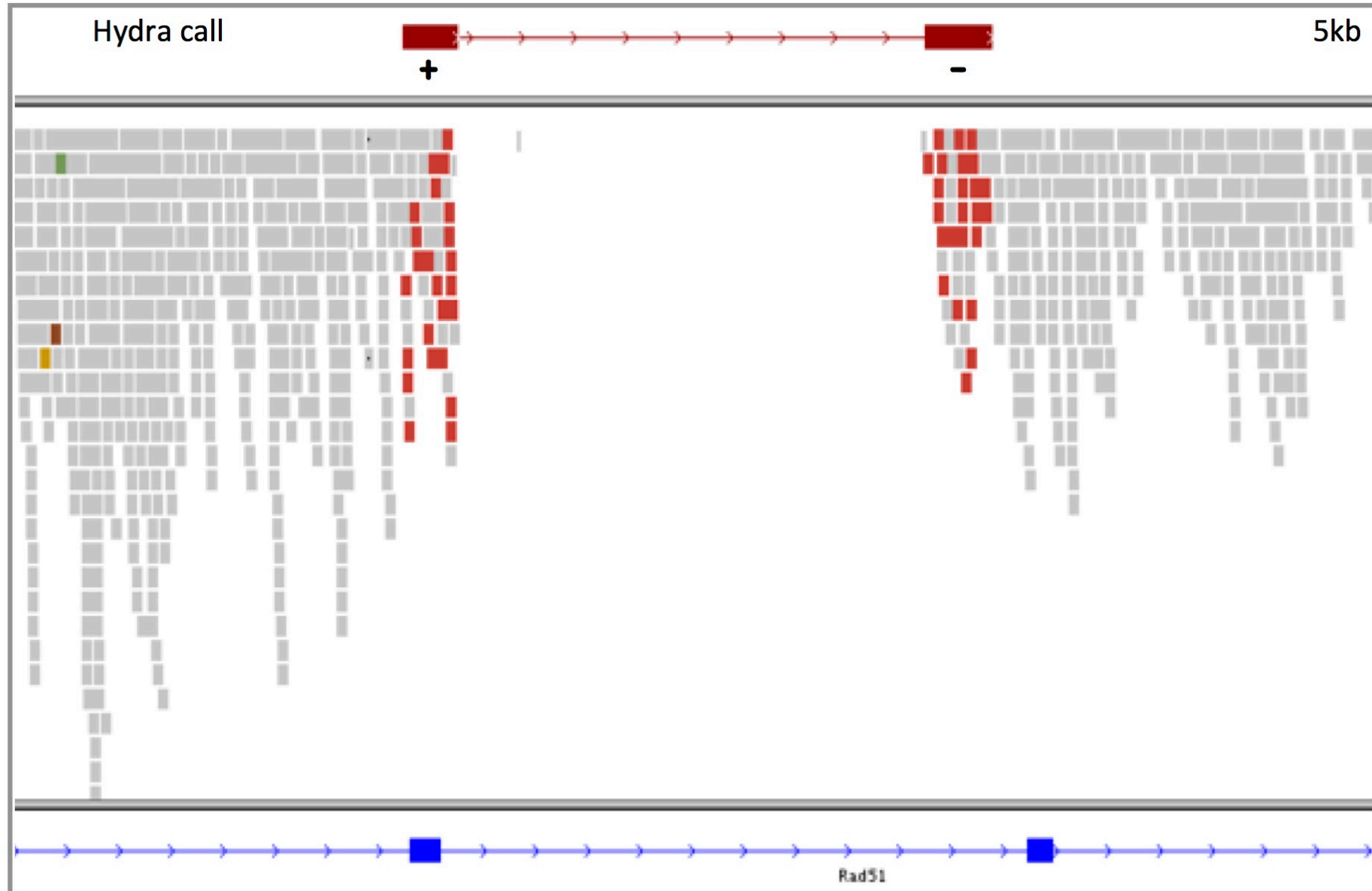


Mills et al. *Nature* 2011

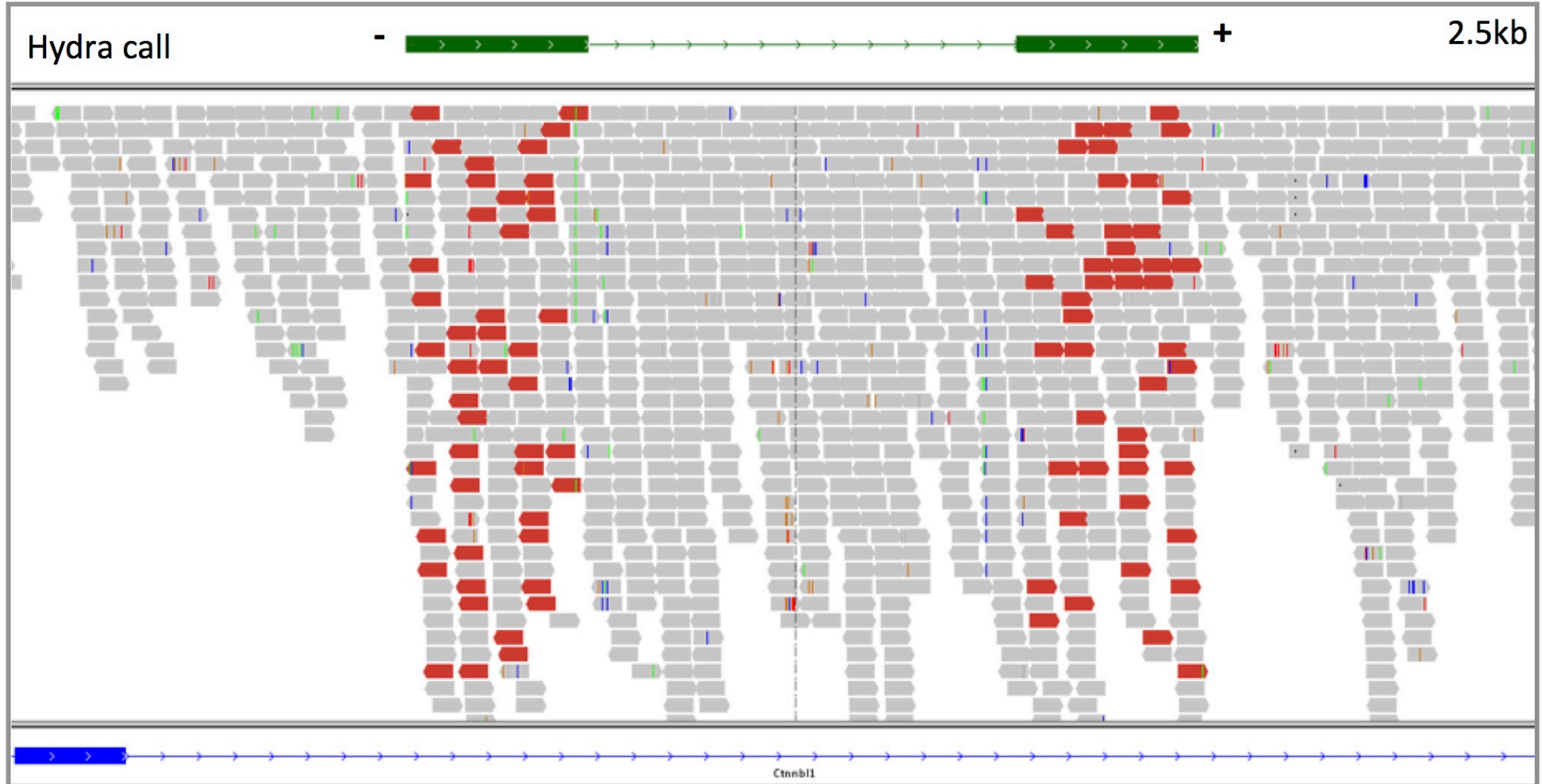


Kloosterman et al. 2015

# Visual validation: a deletion



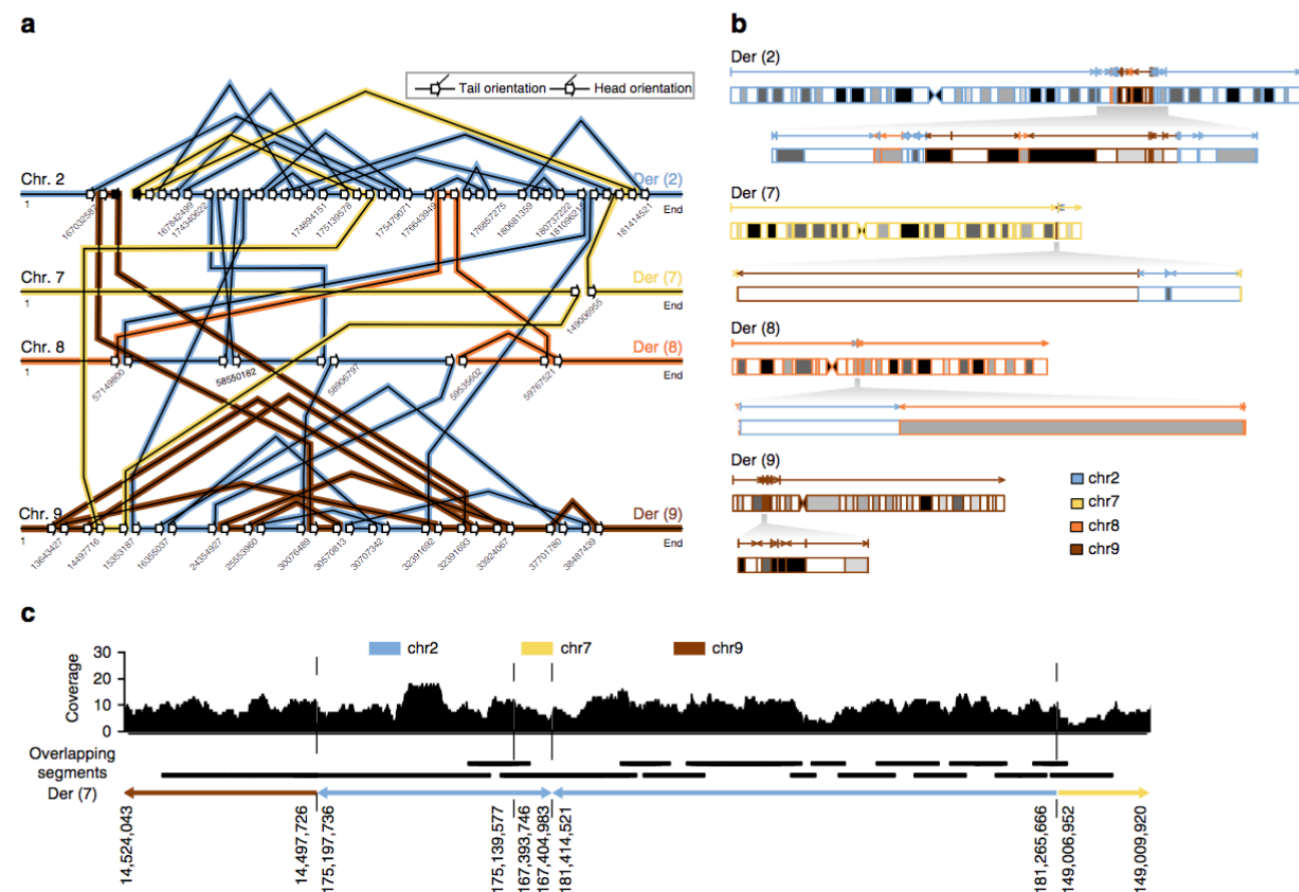
# Visual validation: a duplication



# Mapping and phasing of structural variation in patient genomes using nanopore sequencing

Mircea Cretu Stancu<sup>1</sup>, Markus J. van Roosmalen<sup>1</sup>, Ivo Renkens<sup>1</sup>, Marleen M. Nieboer<sup>1</sup>, Sjors Middelkamp<sup>1</sup>, Joep de Ligt<sup>1</sup>, Giulia Pregno<sup>2</sup>, Daniela Giachino<sup>2</sup>, Giorgia Mandrile<sup>2</sup>, Jose Espejo Valle-Inclan<sup>1</sup>, Jerome Korzelius<sup>1</sup>, Ewart de Bruijn<sup>1</sup>, Edwin Cuppen<sup>3</sup>, Michael E. Talkowski<sup>4,5,6</sup>, Tobias Marschall<sup>7,8</sup>, Jeroen de Ridder<sup>1</sup> & Wigard P. Kloosterman<sup>1</sup>

- long reads are superior to short reads with regard to detection of de novo chromothripsis rearrangements.
- long reads also enable efficient phasing of genetic variations, which we leveraged to determine the parental origin of all de novo chromothripsis breakpoints and to resolve the structure of these complex rearrangements.



# Structural variants: A summary

Actually it's all the same methods

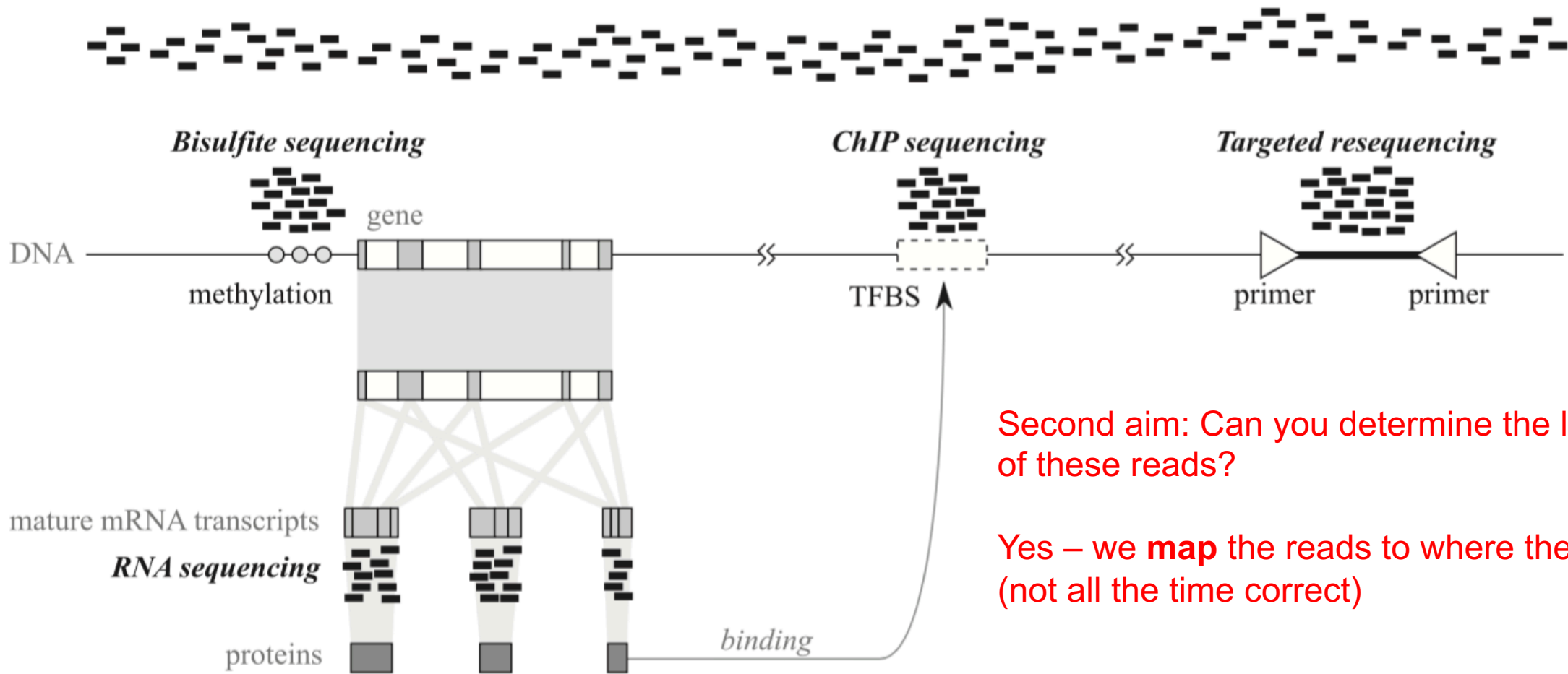
Reference assembly -> check depth -> detect duplication

New assembly -> check depth -> detect ploidy chromosomes / mis-assemblies

**Current: SV should be called using long reads**

# Other experiments requiring mapping

*De novo sequencing / Whole genome resequencing*



Second aim: Can you determine the location of these reads?

Yes – we **map** the reads to where they belong (not all the time correct)

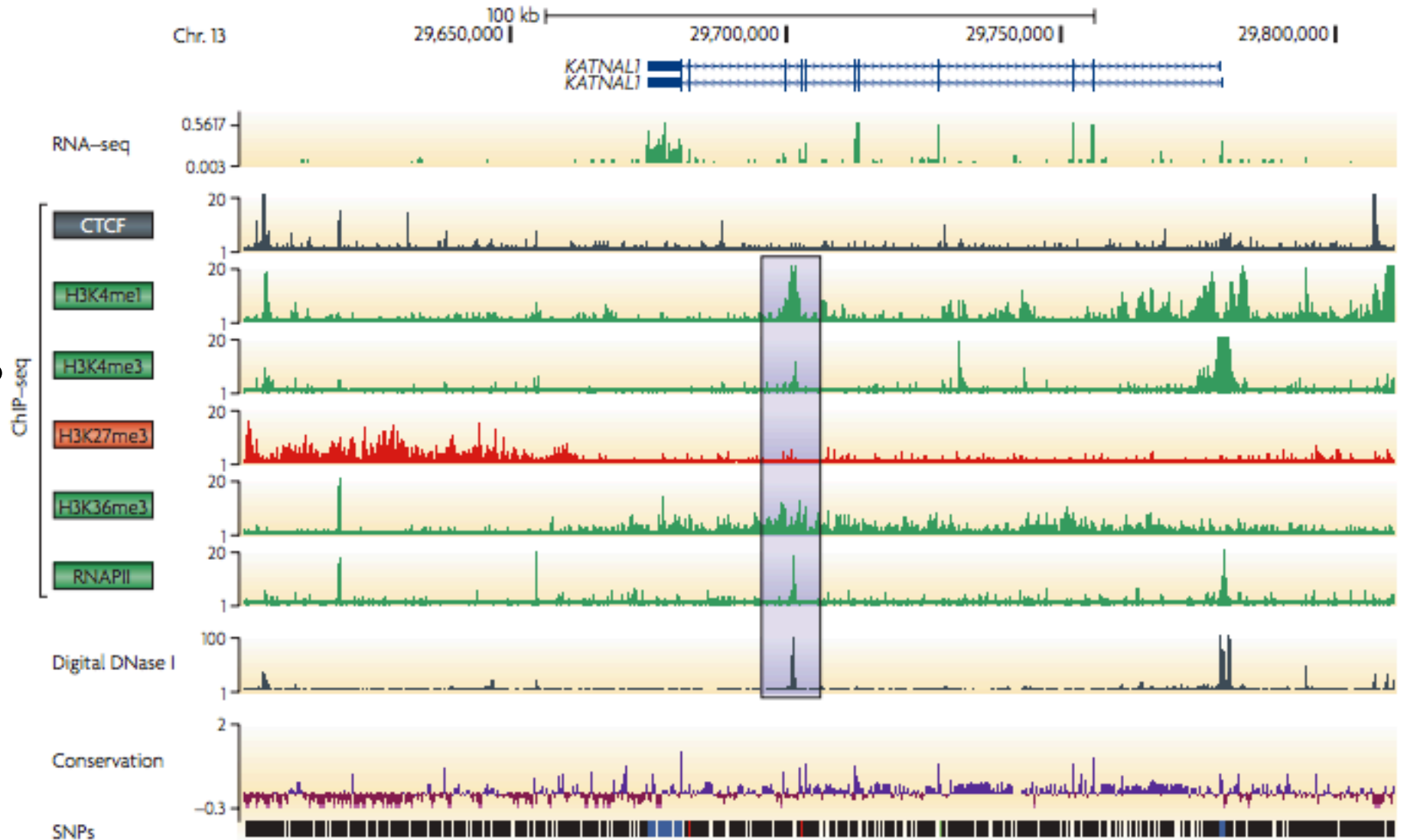
**Figure 1.2** A schematic summary of high-throughput sequencing applications. Details are described in Section 1.3.

# Other experiments requiring mapping

\*-seq methodologies

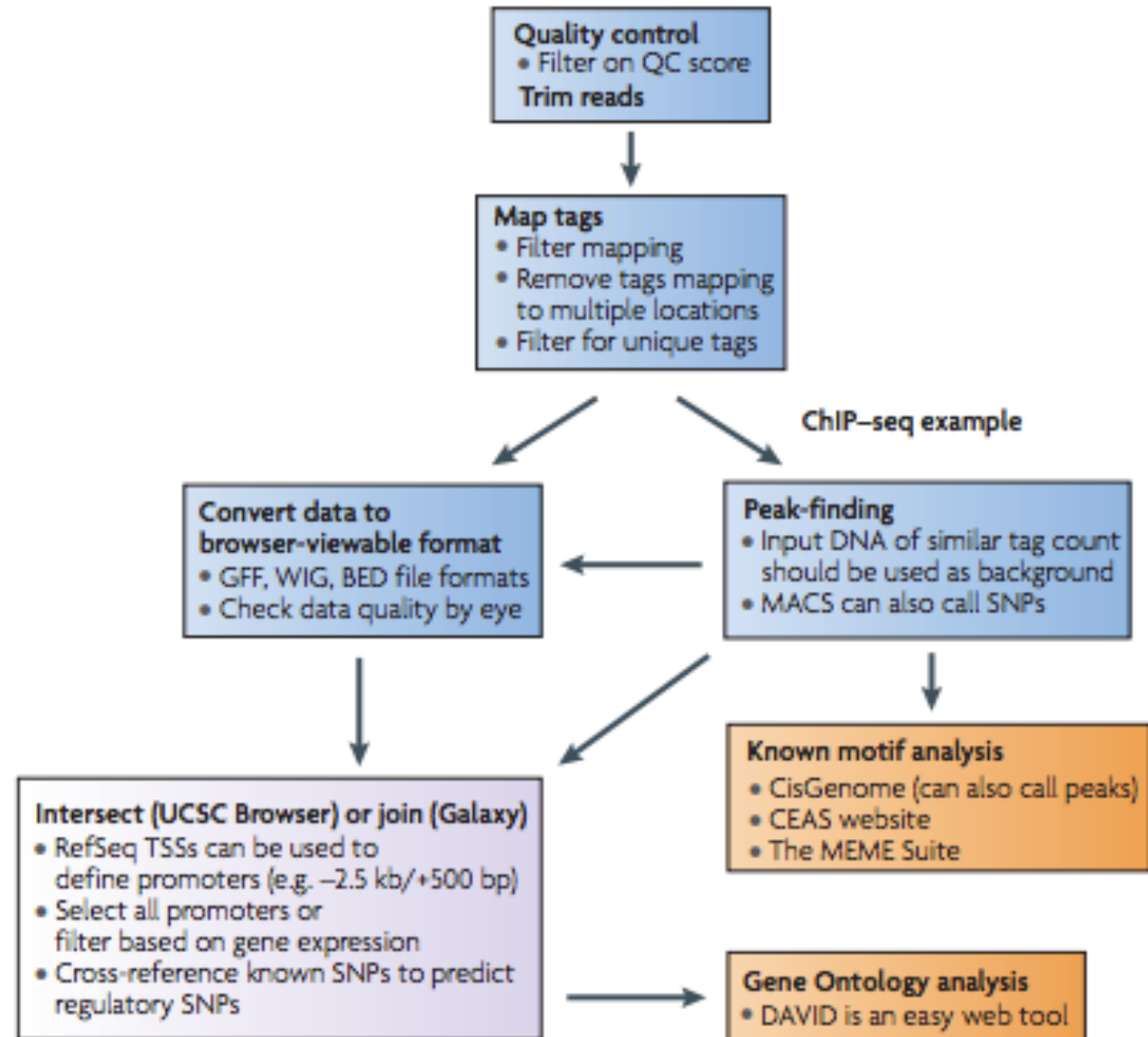
Identify peaks!

How is peak  
different to coverage?



# Other experiments requiring mapping

Similar methods  
Different analysis





# Validation and standardisation

## Genome in a Bottle Consortium

The Genome in a Bottle Consortium is a public-private-academic consortium hosted by NIST to develop the technical infrastructure (reference standards, reference methods, and reference data) to enable translation of whole human genome sequencing to clinical practice.

NA12878 cell line, sequenced many platforms, read lengths and sample preps ; A lot and lot of **Benchmarks**

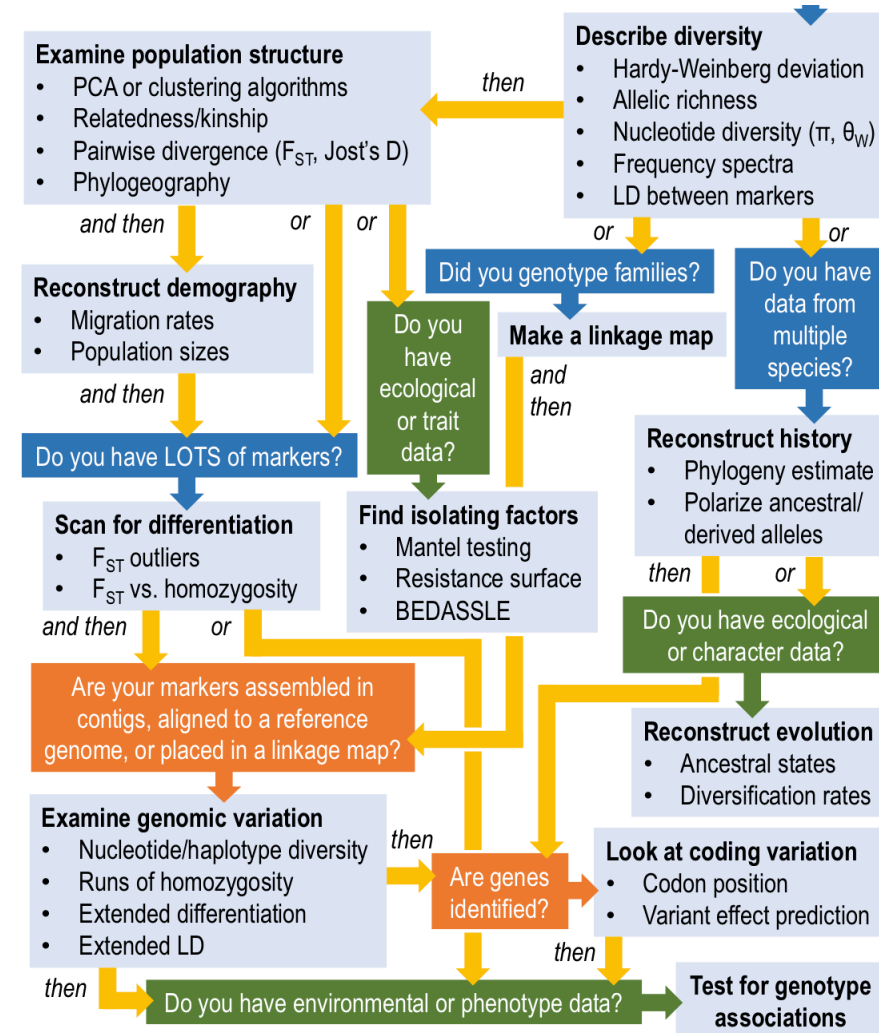
<https://sites.stanford.edu/abms/giab>

Again, only in humans...

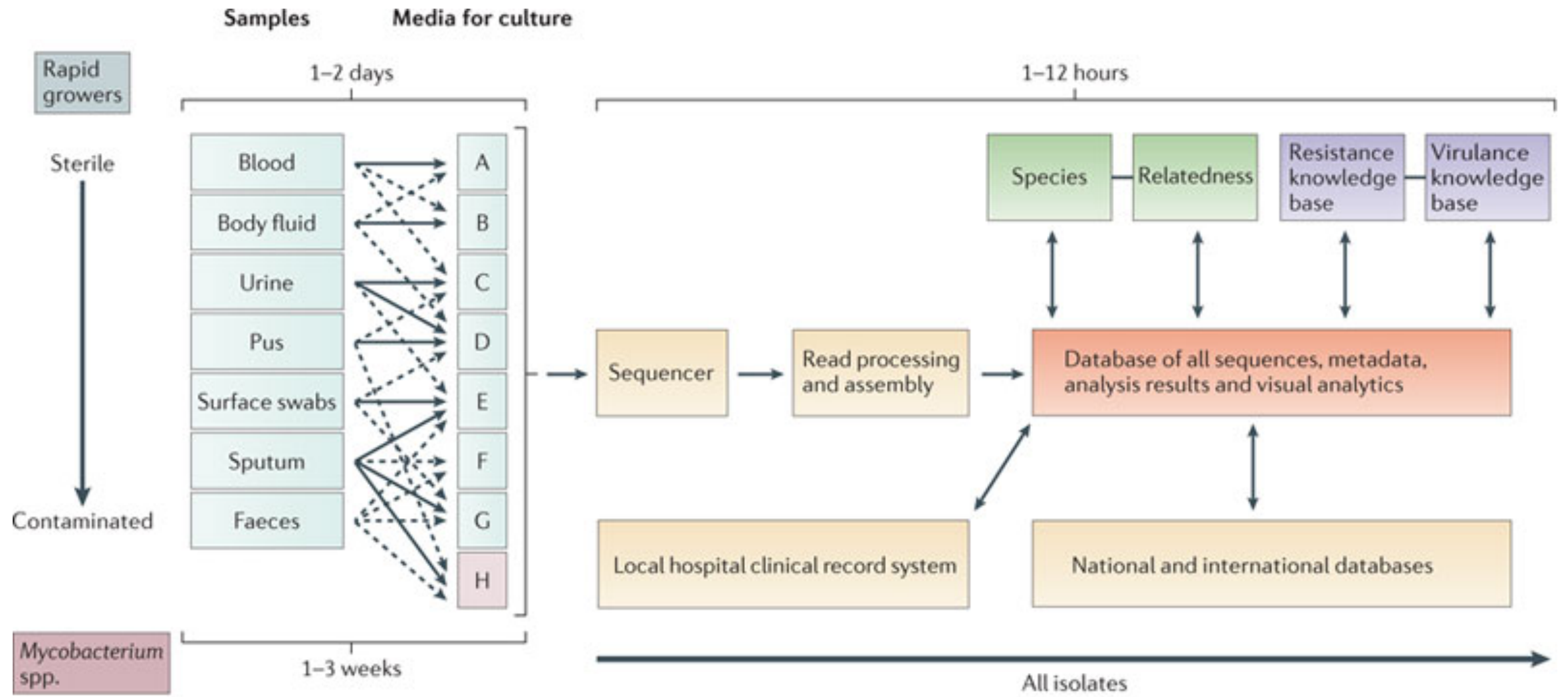


# Ultimately, mapping is to quickly identify relationship between individuals / species once reference is known

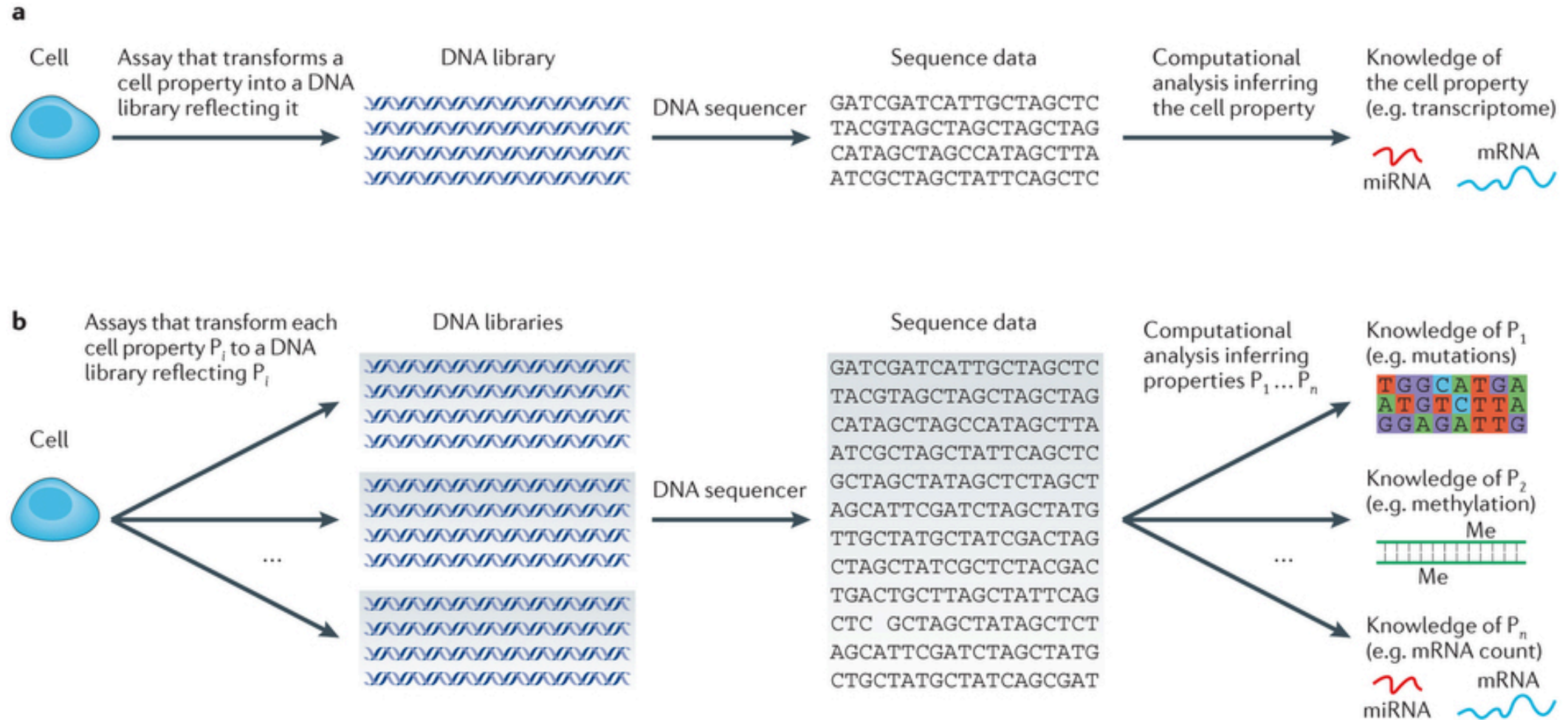
Tradeoffs between \$\$\$, sample size, sensitivity, speed



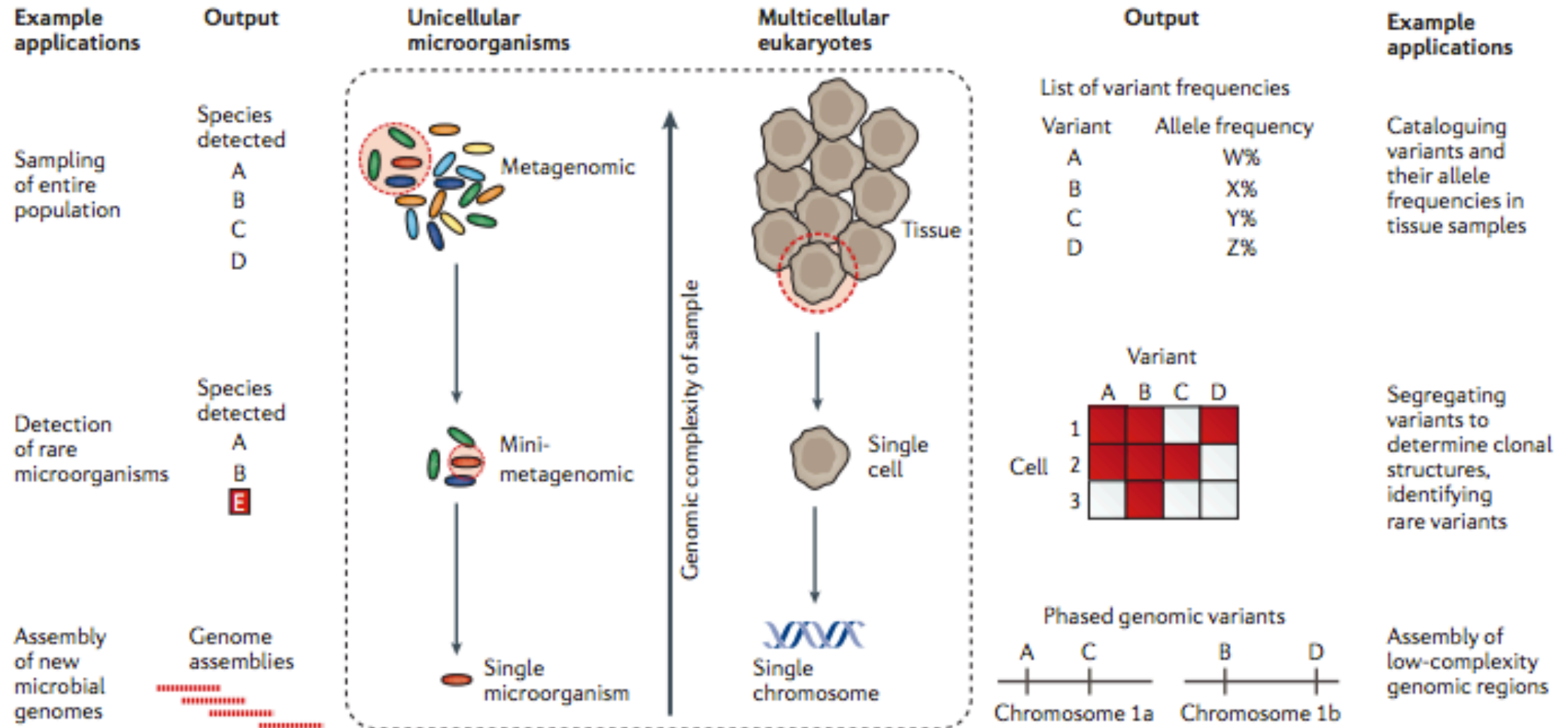
# workflow of clinical labs / ecological samples



# single cell genomics



# single cell genomics



# Example

Hoffmann *et al.* *Climate Change Responses* (2015) 2:1  
DOI 10.1186/s40665-014-0009-x



**Climate Change  
Responses**

**REVIEW**

**Open Access**

## A framework for incorporating evolutionary genomics into biodiversity conservation and management

Ary Hoffmann<sup>1\*</sup>, Philippa Griffin<sup>1</sup>, Shannon Dillon<sup>2</sup>, Renee Catullo<sup>3</sup>, Rahul Rane<sup>1</sup>, Margaret Byrne<sup>4</sup>, Rebecca Jordan<sup>1</sup>, John Oakeshott<sup>5</sup>, Andrew Weeks<sup>1</sup>, Leo Joseph<sup>6</sup>, Peter Lockhart<sup>7</sup>, Justin Borevitz<sup>3</sup> and Carla Sgrò<sup>8</sup>

# Genomics and the challenging translation into conservation practice

Aaron B.A. Shafer<sup>1</sup>, Jochen B.W. Wolf<sup>1</sup>, Paulo C. Alves<sup>2</sup>, Linnea Bergström<sup>1</sup>, Michael W. Bruford<sup>3</sup>, Ioana Brännström<sup>1</sup>, Guy Colling<sup>4</sup>, Love Dalén<sup>5</sup>, Luc De Meester<sup>6</sup>, Robert Ekblom<sup>1</sup>, Katie D. Fawcett<sup>7</sup>, Simone Fior<sup>8</sup>, Mehrdad Hajibabaei<sup>9</sup>, Jason A. Hill<sup>10</sup>, A. Rus Hoebel<sup>11</sup>, Jacob Höglund<sup>1</sup>, Evelyn L. Jensen<sup>12</sup>, Johannes Krause<sup>13</sup>, Torsten N. Kristensen<sup>14</sup>, Michael Krützen<sup>15</sup>, John K. McKay<sup>16</sup>, Anita J. Norman<sup>17</sup>, Rob Ogden<sup>18</sup>, E. Martin Österling<sup>19</sup>, N. Joop Ouborg<sup>20</sup>, John Piccolo<sup>19</sup>, Danijela Popović<sup>21</sup>, Craig R. Primmer<sup>22</sup>, Floyd A. Reed<sup>23</sup>, Marie Roumet<sup>8</sup>, Jordi Salmons<sup>24</sup>, Tamara Schenker<sup>25</sup>, Michael K. Schwartz<sup>26</sup>, Gernot Segelbacher<sup>27</sup>, Helen Senn<sup>18</sup>, Jens Thaulow<sup>28</sup>, Mia Valtonen<sup>29</sup>, Andrew Veale<sup>12</sup>, Philippine Vergeer<sup>30</sup>, Nagarjun Vijay<sup>1</sup>, Carles Vilà<sup>31</sup>, Matthias Weissensteiner<sup>1</sup>, Lovisa Wennerström<sup>10</sup>, Christopher W. Wheat<sup>10</sup>, and Piotr Zieliński<sup>32</sup>

**Table 1. Main areas traditionally addressed by conservation genetics [3], current status of genetic and genomic approaches, and the contribution that genomics can potentially make**

Category	Status of conservation genetics	Possible contribution of conservation genomics	Required for transition from basic to applied <sup>a</sup>
<i>Evolutionary genetics of natural populations</i>			
Demographic inference – population history	Regularly used Moderate resolution	Improved accuracy and precision Finer-scale population structure Less limited by sample size	Clear understanding of limitations and biases User-friendly software
Adaptive genetic variation	Minimally used Based on population correlations [77] or candidate gene approaches	Improved detection of adaptive loci Management frameworks proposed [28] Methods still emerging Interpretations unclear	In-depth validation studies Genome annotation
Quantitative genetic variation	Limited resolution Often dependent on pedigrees or targeted gene approaches	Improved detection of quantitative trait loci Active application (e.g., genome-wide association studies)	Ecological studies Genome annotation
Taxonomic identification and general diagnostics	Regularly used Moderate resolution Restricted to single individuals	Assay species simultaneously [78] Improved hybridization detection Improved detection of pathogens	Defined pipelines (Box 3) Repeatability
<i>Effects of small population size</i>			
Inbreeding detection	Regularly used Limited resolution [34]	Improved estimates of inbreeding [34,62] Novel genomic metrics [79] Assess impact on specific genomic regions or adaptive loci	User-friendly bioinformatics Genome annotation Practitioner demand
Population viability	Minimally used [80]	Improved estimates of inbreeding metrics used in viability models [80]	Practitioner demand
<i>Additional applications</i>			
Genetic monitoring	Minimally used [11]	Improved sampling regimens [63] More powerful biodiversity surveys	Practitioner demand Compliance [11]
Population census	Regularly used	Higher-throughput screening	Practitioner demand
Maternity, paternity, and kinship analysis	Regularly used	Useful when microsatellite power is limited [81]	Practitioner demand



## Genomics research and development

SNP discovery

SNP validation and selection

Genome-wide genotyping

## Marker assessment and selection

Population screening

Population genomic analysis

SNP panel selection

## Applied traceability tools

Platform selection

Method validation

Standard operating procedures

## TheFishPopTrac target species

reproduced with permission from the Scandinavian Fishing Yearbook (c)



European hake (*Merluccius merluccius* L.)



Atlantic herring (*Clupea harengus* L.)



Atlantic cod (*Gadus morhua* L.)



Common sole (*Solea solea* L.)

*TRENDS in Ecology & Evolution*

# Reading materials

Introduction to Read-Based Alignment

<http://evomicsorg.wpengine.netdna-cdn.com/wp-content/uploads/2015/06/Intro-Read-Alignment.pdf>

Ben Langmead

<http://www.langmead-lab.org/teaching-materials/>

Additional references:

<https://www.notion.so/References-papers-links-in-start-learning-genomics-b7e57b28e9194bb29a02f483e0b894ad>

# Written assignment

Construct a BWT of the following sequence:

ANNABANANA

Question:

1. What is the output of last column?
2. Write out how you searched the string ANNA  
(hint: follow wiki\*)

To be handed in **2020.04.15**

[https://en.wikipedia.org/wiki/Burrows%E2%80%93Wheeler\\_transform](https://en.wikipedia.org/wiki/Burrows%E2%80%93Wheeler_transform)