

Evolution and Genomics

immersive training opportunities

[Home](#)

[Workshops](#)

[Learning](#)

[Resources](#)

[People](#)

[Apply](#)

[Information](#)

[Blog](#)

Blog

You are here: [Home](#) > [Blog](#)



Our sequencing Guru takes the podium...

Posted on January 15, 2014 by Dr. Mel in [Genomics](#), [Workshops](#)

Dr. Mel Zody Broad Institute Topic: Genomics Study Design Dr. Zody is a jack of all sequencing trades—he's been at it for several years and his career spans the Human Genome Project, vertebrate evolution and positive selection, to genetic links to viral disease. He's probably seen and heard it all, his slides are excellent so [...]

[Leave a comment](#) • [Continue Reading](#) →

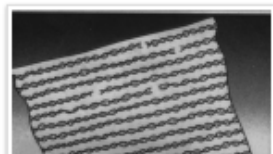


Unix 101+1: Tinkerbell has issues

Posted on January 14, 2014 by Dr. Mel in [Programming](#), [Workshops](#)

Julian Catchen University of Oregon Unix Ninja Topic: Unix Part 2: More Advance Ninja-ry So unfortunately as with this blog, I am or will be unable to give you files that we practiced on but Julian's slides are quite good. Remember, we learned about pipes and added on to our current knowledge of command line. [...]

[Leave a comment](#) • [Continue Reading](#) →



Sequencing Technology: Where's my Minlon!?

Posted on January 14, 2014 by Dr. Mel in [Genomics](#), [Workshops](#)

So this morning started off with a lecture from Dr. Konrad Paszkiewicz on the 'state of the union' with respect to Sequencing

New Illumina sequencer launched



Mi-seq



Next-seq 500



HiSeq 2500

<http://biomickwatson.wordpress.com/>

HiSeq X10

Population power. Extreme throughput. \$1,000 human genome.

The HiSeq X Ten is a set of ten ultra-high-throughput sequencers, purpose-built for large-scale human whole-genome sequencing.



General queries

- Technical replicates
 - <http://genomebiology.com/2011/12/3/R22#B18>
 - <http://genomebiology.com/2011/12/3/R22#B19>
- Allele drop out for double-digest RAD
 - <http://www.ncbi.nlm.nih.gov/pubmed/23110526>



Workshop on Genomics

Short read alignment:

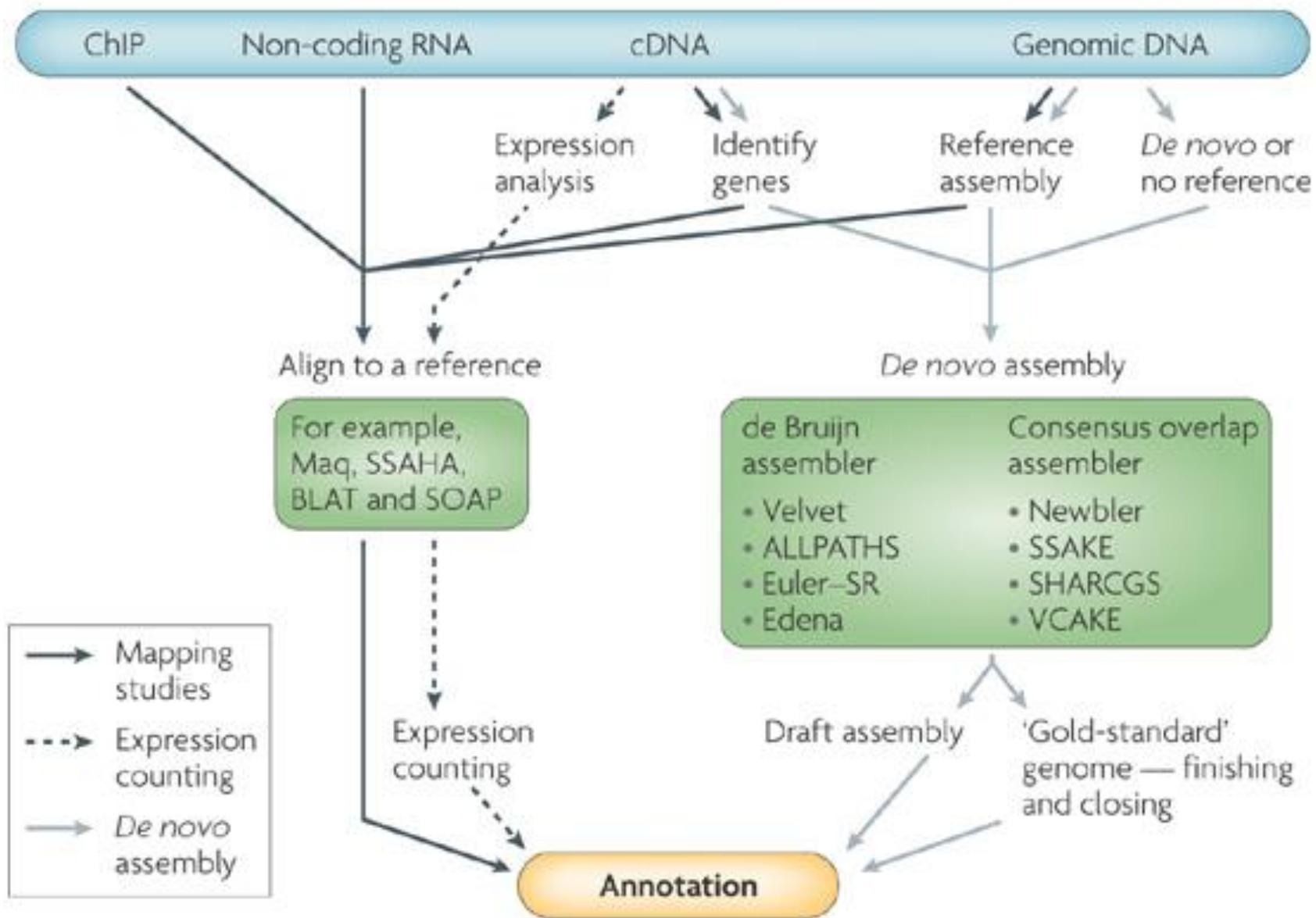
An introduction

Dr Konrad Paszkiewicz
University of Exeter, UK
k.h.paszkiewicz@exeter.ac.uk

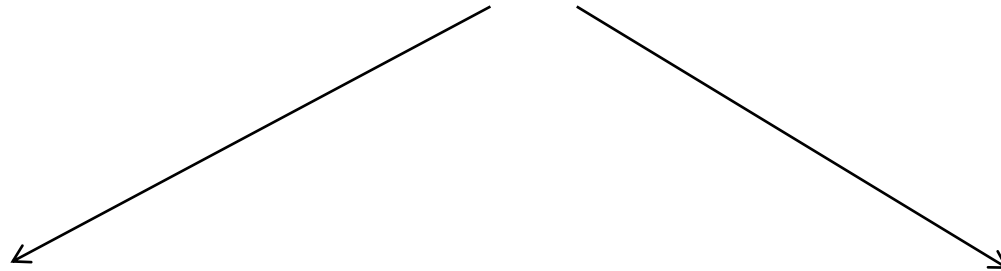
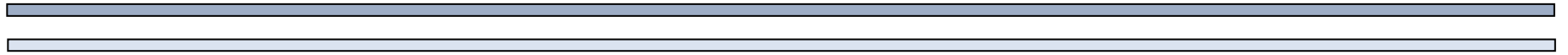
Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Suffix/Prefix Tries
 - Indels
 - Other alignment considerations
 - Typical alignment pipeline
 - SNP calling

Raw sequence source

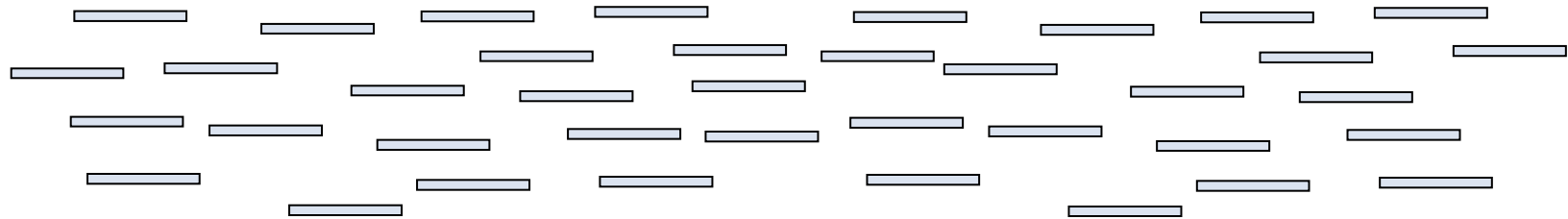


Alignment of reads to a reference



..ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA.. Reference

..ACTGGGTCATCGTACGATCGATAGATCGATCGATCGCTAGCTAGCTA.. Sample



Why is short read alignment hard?

The shorter a read, the less likely it is to have a unique match to a reference sequence

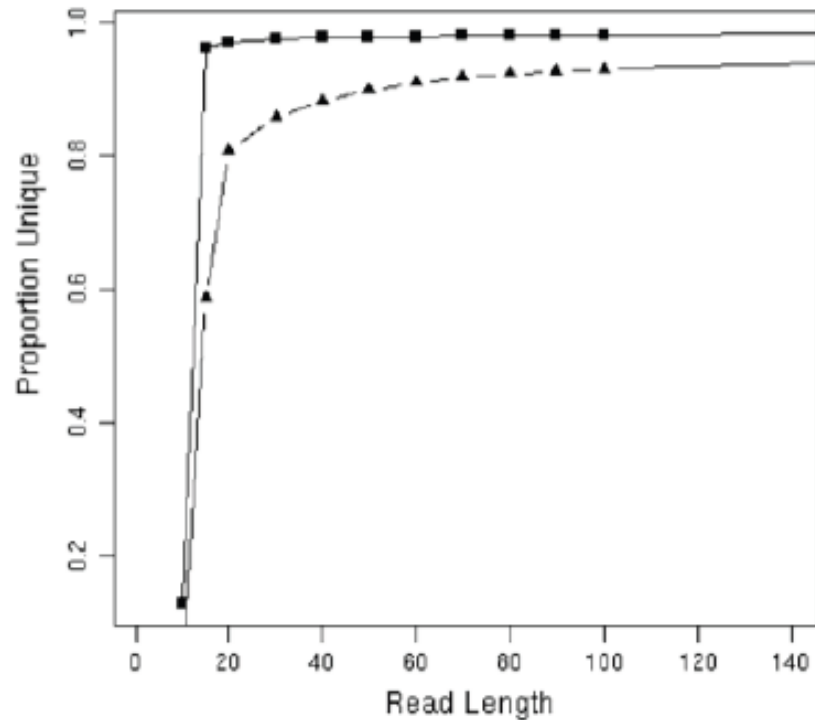


Fig. 1 The proportion of unique sequence in the *Streptococcus suis* (squares) and *Mus musculus* (triangles) genomes for varying read lengths. This graph indicates that read length has a critical affect on the ability to place reads uniquely to the genome

Why do we generate short reads?

- Sanger reads lengths ~ 800-2000bp
- Generally we define short reads as anything below 200bp
 - Illumina (50bp – 300bp)
 - SoLID (80bp max)
 - Ion Torrent (200-400bp max...)
 - Roche 454 – 400-800bp
- Even with these platforms it is cheaper to produce short reads (e.g. 50bp) rather than 100 or 200bp reads
- Diminishing returns:
 - For some applications 50bp is more than sufficient
 - Resequencing of smaller organisms
 - ChIP-Seq
 - Digital Gene Expression profiling
 - Bacterial RNA-seq

Short read alignment applications

Genotyping:

Methylation

SNPs

Indels

```
...CCATAG      TATGCGCCC   CGGA AATT  CGGTATAC...
...CCAT      CTATATGCG   TCGGA AATT  CGGTATAC
...CCAT GGCTATATG   CTATCGG AAA   GCGGTATA
...CCA AGGCTATAT   CCTATCGG A   TTGCGGT  C...
...CCA AGGCTATAT   GCCCTATCG   TTTGCGGT  C...
...CC  AGGCTATAT   GCCCTATCG   A AATTTGC  ATAC...
...CC  TAGGCTATA   GCGCCCTA   A AATTTGC  GTATAC...
...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...
```

Classify and measure peaks:

ChIP-Seq

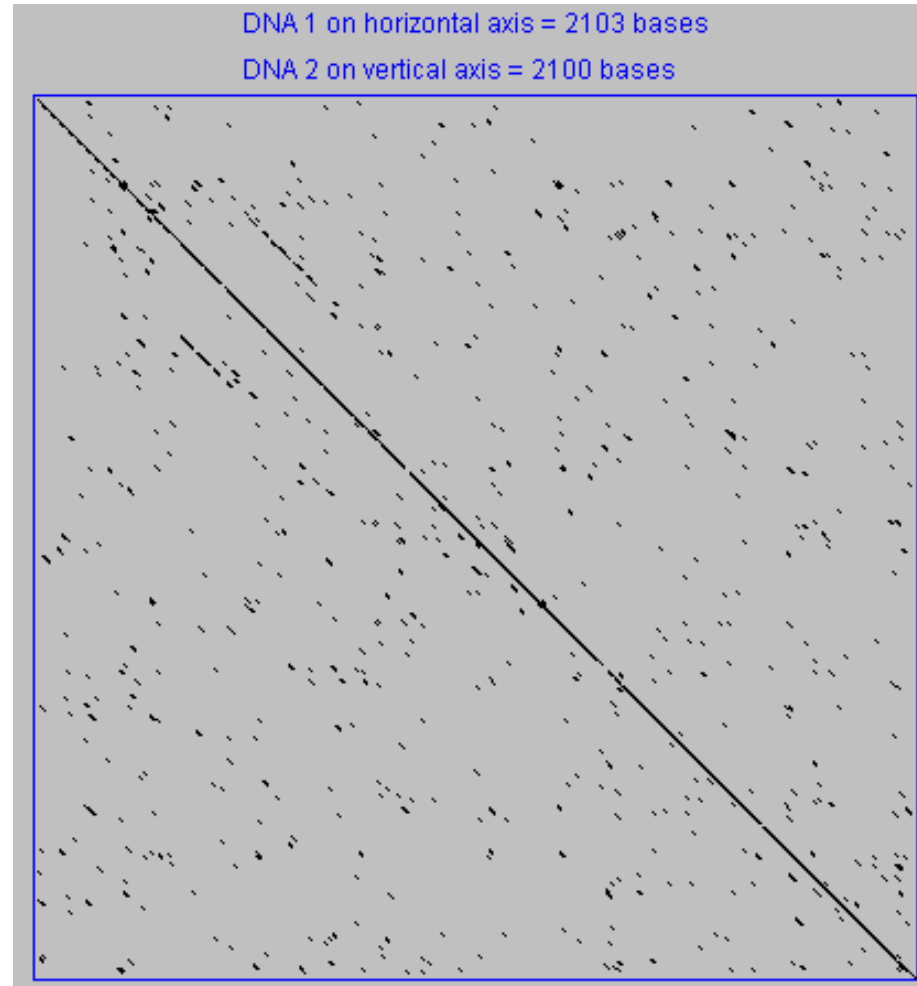
RNA-Seq

```
...CC
...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...
GAAATTTGC
GGAAATTTG
CGGAAATTT
CGGAAATTT
TCGGAAATT
CTATCGGAAA
CCTATCGGA TTTGCGGT
GCCCTATCG AAATTTGC
GCCCTATCG AAATTTGC ATAC...
```

Contents

- **Alignment algorithms for short-reads**
 - **Background – Blast (why can't we use it?)**
 - Global alignment
 - Local alignment
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Indel detection
 - Suffix/Prefix Tries
 - Other alignment considerations
 - Typical alignment pipeline
 - New methods of SNP calling

Dot Matrix Method - Aligning by eye



Sequence Alignment

~~AATCCGATA~~ CCG
AATCGGATTTACCG

3 possibilities

Match

...A...
|
...A...

Mismatch

...C...
...G...

Indel

...-...
...T...

A very simple alignment scoring system

Points for a matching letter: 1

Points for a non-matching letter: 0

Points for inserting a gap: 0

Global Pair-wise Alignment

ATCGATACG, ATGGATTACG

A T C G A T - A C G
| | | | | | | | |
A T G G A T T A C G

Matches:	+1	+1		+1	+1	+1		+1	+1	+1	= +8
Mismatches:			0								= 0
Gaps:							0				= 0
											<hr/>
											Total score = +8

But, what does this score mean??
Could we get a better alignment?

How to choose the best alignment?

- Sequence 1: ACTGAGC
- Sequence 2: ATGATGC
- Some possible alignments:

ACTGAGC-- ACTGA-GC A-----CTGAGC
A-TGA-TGC A-TGATGC ATGAT-----GC

Global alignment – Needleman-Wunsch

A **global** alignment covers the **entire lengths of the sequences** involved

The Needleman-Wunsch algorithm finds the best global alignment between 2 sequences across their whole length

Step 1: Initialise

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T							0
G							0
C	0	1	0	0	0	0	1

Fill in far-right column and bottom row with:

0 for a mis-match

1 for a match

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T							0
G							0
C	0	1	0	0	0	0	1

For each box, find the highest number out of the blue boxes

Step 3:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T							0
G						1+1=2	0
C	0	1	0	0	0	0	1

If there is a match in the yellow box as, take the highest value from the blue boxes and add 1 to it

G matches G in the yellow box, so add 1 to the 1 in the blue box

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T							0
G					0+0=0	2	0
C	0	1	0	0	0	0	1

A does not match G. So add zero to the zero in the blue box.

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T							0
G				0+1=1	1	2	0
C	0	1	0	0	0	0	1

If there is a match as here, take the highest value and add 1 to it

G matches G so add 1 to zero in the blue box

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T							0
G			0+0=0	1	0	2	0
C	0	1	0	0	0	0	1

If there is a match as here, take the highest value and add 1 to it

T does not match G. So add zero.

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T						0+1=1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Highest out of the blue boxes is zero

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T					2+0=2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Highest out of the blue boxes is 2

A does not match T

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T				2+0=2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Highest out of the blue boxes is 2

G does not match T

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T			3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Highest out of the blue boxes is 2

T does match **T**

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T		2+0=2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Highest out of the blue boxes is 2

C does not match T

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A							0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Do the same for all remaining rows

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A						1+0=0	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Do the same for all remaining rows

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A					2+1=3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Do the same for all remaining rows

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A				2+0=2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Do the same for all remaining rows

Step 2:

	A	C	T	G	A	G	C
A							0
T							0
G							0
A			2+0=2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Do the same for all remaining rows

Step 2:

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Do the same for all remaining rows

Step 3: Backtracking

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Follow largest numbers starting from top-left going down and to the right

Step 3: Backtracking

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Follow largest numbers starting from top-left going down and to the right

Step 3: Backtracking

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Follow largest numbers starting from top-left going down and to the right

Step 3: Backtracking

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Follow largest numbers starting from top-left going down and to the right

Step 4: Generate alignment

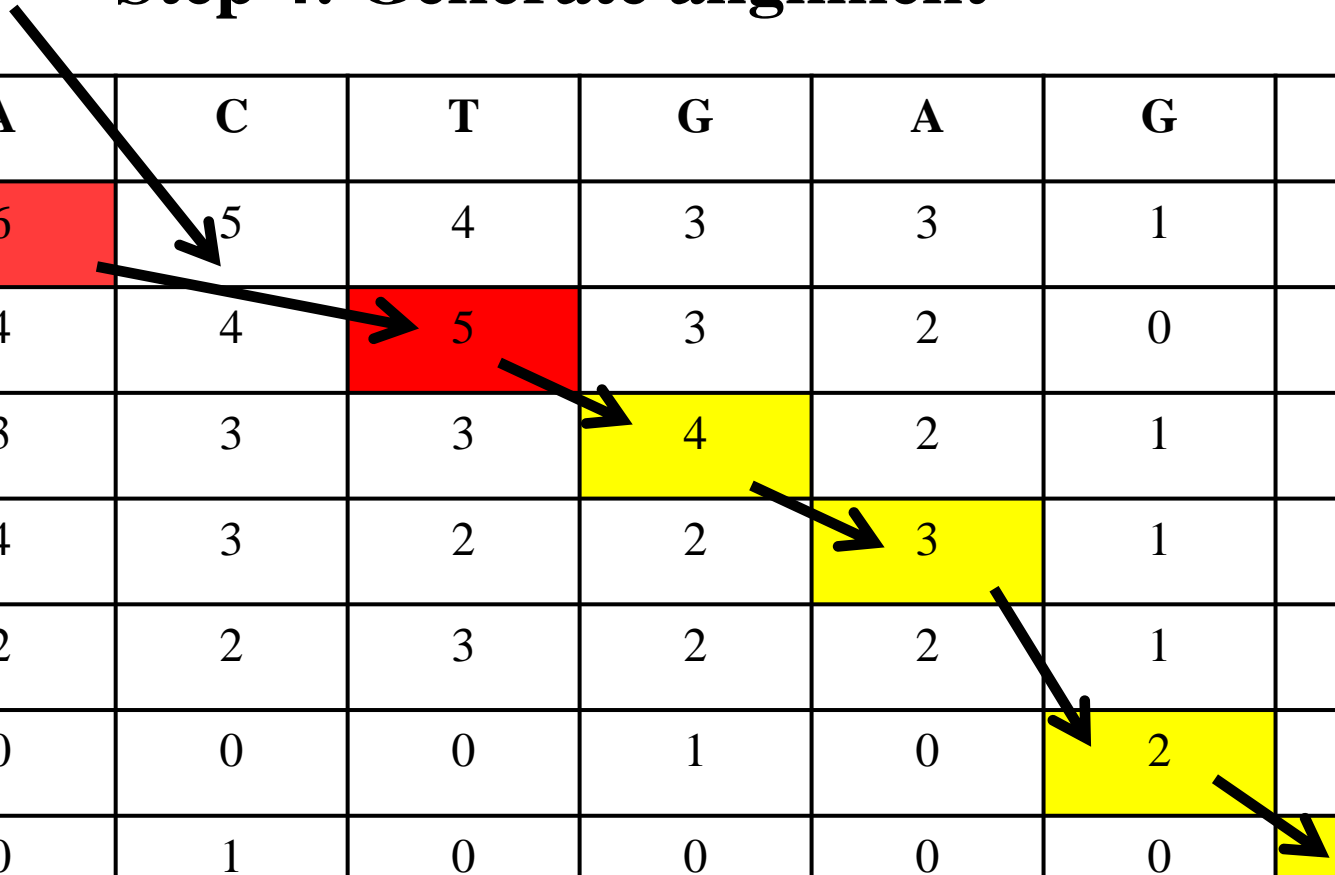
	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Horizontal seq A
Vertical seq A

Step 4: Generate alignment

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Gap



Horizontal seq

Vertical seq

ACT

A-T

Step 4: Generate alignment

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Horizontal seq
Vertical seq

ACTG
A-TG

Step 4: Generate alignment

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Horizontal seq
Vertical seq

ACTGA
A-TGA

Step 4: Generate alignment

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Horizontal seq
Vertical seq

ACTGA-
A-TGAG

Step 4: Generate alignment

	A	C	T	G	A	G	C
A	6	5	4	3	3	1	0
T	4	4	5	3	2	0	0
G	3	3	3	4	2	1	0
A	4	3	2	2	3	1	0
T	2	2	3	2	2	1	0
G	0	0	0	1	0	2	0
C	0	1	0	0	0	0	1

Horizontal seq
Vertical seq

ACTGA-C
A-TGAGC

Optimal global alignment

ACTGA-C

| | | | |

A-TGAGC

Local alignment

A **global** alignment is often not appropriate as only parts of sequences may be conserved

A **local** alignment only covers **parts of the sequences**

The **Smith-Waterman** algorithm finds the **best local alignment** between 2 sequences

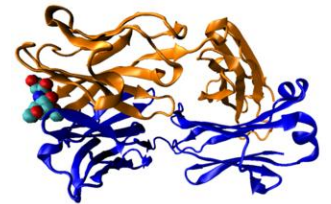
		Q	K	E	S	G	P	S	S	S	Y	C
Global alignment	V	Q	Q	E	S	G	L	V	R	T	T	C

				E	S	G						
Local alignment				E	S	G						

Local alignment

A **local alignment** of 2 sequences is an alignment between **parts** of the 2 sequences

E.g. Two proteins may be very similar in a functional site, but be very dissimilar outside that region



A global alignment of such sequences would have:

- (i) lots of matches in the region of high sequence similarity
- (ii) lots of mismatches & gaps (insertions/deletions) outside the region of similarity

It makes sense to find the **best local alignment** instead


```

human/1-422  .....
fly/1-898  1 MFTLQPTPTAIGTVVPPWSAGTLIERLPSLEDMAHKDNIAMRNLPCLGT 50

human/1-422  1 .....MQNSHSG 7
fly/1-898  51 AGGSGLGGIAGKPSPTEAVEASTASHPHSTSSYFATTYYHLTDDEC HSG 100

human/1-422  8 VNQLGGVFVNGRPLPDSTRQKIVELAHSGARPCDISRILQVSNGCVSKIL 57
fly/1-898  101 VNQLGGVFVNGRPLPDSTRQKIVELAHSGARPCDISRILQVSNGCVSKIL 150

human/1-422  58 GRYYETGSIRPRAIGGSKPRVATPEVVSKIADYKRECPSIFAWEIRDRL 107
fly/1-898  151 GRYYETGSIRPRAIGGSKPRVATPEVVSKISDYKRECPSIFAWEIRDRL 200

human/1-422  108 SEG VCTNDNIPSVSSINRVLRLNLA SEKQDM ..... 137
fly/1-898  201 QEN VCTNDNIPSVSSINRVLRLNLA AQKEQDSTGSGSSSTSAAGNSISAKVS 250

human/1-422  138 .....GA ..... DG 141
fly/1-898  251 VSIGGNVSNVASGSRGTLSSSTDLMQTATPLNSSSESGGASNSGEGSEGEA 300

human/1-422  142 MYDKLRMLNGQTG ..... 154
fly/1-898  301 IYEKLRLLNTQHAAGPGPLEPARAAPLVGQSPNHLGTRSSHPQLVHGNHQ 350

human/1-422  155 .....SWGTR...PGWYPTSVPGQPTQ ..... 174
fly/1-898  351 ALQHQHQQSWPPRHYSGSWYP-TSLSEIPISSAPNIASVTAYASGPSLAH 399

human/1-422  175 .....DGCQQQE...GGGENTN 188
fly/1-898  400 SLSPPNDIESLASIGHQRNCPVATEDIHLKKELDGHQSDETGSGEGENS 449

human/1-422  189 SISISNGEDSDEADMRLQLKRRKLRNRTSFTQEQIEALEKEFERTHYPDVF 238
fly/1-898  450 GGA SNIGNTEDDQARLILKRRKLRNRTSFTNDQIDSLEKEFERTHYPDVF 499

human/1-422  239 ARERLAAKIDLPEARIQVWFSNRRAKWRREEKLRNQRRTQASNTPSHIPIS 288
fly/1-898  500 ARERLAGKIGLPEARIQVWFSNRRAKWRREEKLRNQRRTPNSTGASATSS 549

human/1-422  289 SFSTSVYQPI PQPTTPV SSFTSGSMLGRD TALNTYSALPPMPSTMA 338
fly/1-898  550 STSATASLTDS PNL SACSLLSGSAGGPSVSTINGLSS ..... PSTLST 594

human/1-422  339 N-NLP .....MQPVPVPSQTSSYSQMLPTSPSVNGRSYD ..... TYT 373
fly/1-898  595 NVNAPT LGAGIDSSSESTPIPHIRPSC...TSDNDNGRQSEDCRRVCSPC 641

human/1-422  374 PPHMQTHMNSQPMGTS GTTSTGLISFGVSVVQVPGSEPDMSQYW PRLQ 422
fly/1-898  642 PLGVGGHQNTHHIQSNQHAQGHALVPAIS ..... PRLNF 675

human/1-422  .....
fly/1-898  676 NSGSGFAMYSNMHHTALSMSDSYGAVTPIPSFNHSAVGPLAPPSPIPQGG 725

human/1-422  .....
fly/1-898  726 DLTPSSLYPCHMTLRPPMAPAHHHIVPGDGGRPAGVGLGSGQSANLGAS 775

human/1-422  .....
fly/1-898  776 CSGSGYEVLSAYALPPPMASSAADSSSFAASSASANVTPHHTIAQESC 825

human/1-422  .....
fly/1-898  826 PSPCSSASHFGVAHSSGFSDDPISPAVSSYAHMSYNYASSANTMT PSSAS 875

human/1-422  .....
fly/1-898  876 G TSAHVAPGKQQFFASC FYSPWV 898

```

Alignment of an orthologous protein in *D.melanogaster* vs *H.sapiens*

Not suitable for global alignment

2 main regions of similarity

Better to use local alignment

Local alignment – Smith-Waterman algorithm

Example – align GATC to GAC

0	-	G	A	T	C
-	0				
G					
A					
C					

Local alignment – Smith-Waterman algorithm

Local alignment algorithm

GATC

||||

—

	-	G	A	T	C
-	0	-2	-4	-6	-8
G					
A					
C					

Points for match = +1

Points for mismatch = -1

Points for a gap insertion = -2

Local alignment – Smith-Waterman algorithm

Local alignment algorithm

GAC

|||

—



	-	G	A	T	C
-	0	-2	-4	-6	-8
G	-2				
A	-4				
C	-6				

Points for match = +1

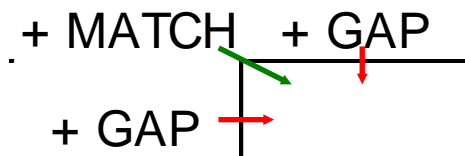
Points for mismatch = -1

Points for a gap insertion = -2

Local alignment – Smith-Waterman algorithm

GATC
|

	-	G	A	T	C
-	0 ⁺¹	-2	-4	-6	-8
G	-2	Max=1			
A	-4				
C	-6				



Points for match = +1
 Points for mismatch = -1
 Points for a gap insertion = -2

Local alignment – Smith-Waterman algorithm

	-	G	A	T	C
-	0	-2	-4	-6	-8
G	-2	1			
A	-4	-1			
C	-6				

Points for match = +1
 Points for mismatch = -1
 Points for a gap insertion = -2

Dynamic Programming

Local alignment algorithm

	-	G	A	T	C
-	0	-2	-4	-6	-8
G	-2	1	-1	-3	-5
A	-4	-1	2	0	-2
C	-6	-3	0	1	1

Points for match = +1

Points for mismatch = -1

Points for a gap insertion = -2

Backtracking and final alignment

	-	G	A	T	C
-	0	-2	-4	-6	-8
G	-2	1	-1	-3	-5
A	-4	-1	2	0	-2
C	-6	-3	0	1	1

GATC
| | |
GA-C

Smith-Waterman – more details

<http://www.youtube.com/watch?v=IVRSFaGCGeE>

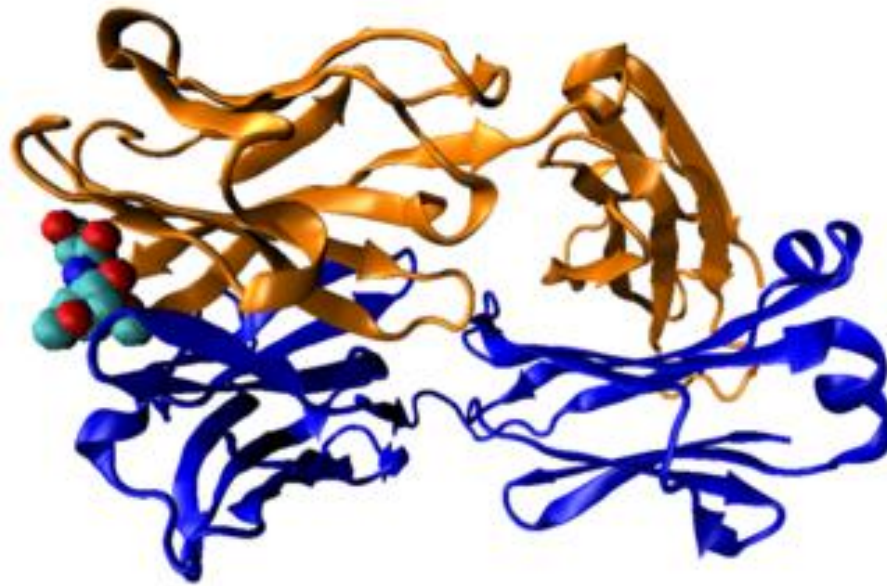
Dynamic programming

- Needleman-Wunsch and Smith-Waterman are a class of methods known as ‘Dynamic Programming’
- Guaranteed to give you the best possible alignment
- In biology, this algorithm is very inefficient because most sequences are not similar to each other
- Therefore it takes a long time to run

**BLAST –
Basic Local Alignment Search Tool**

Background – BLAST

- Primarily designed to identify homologous sequences
 - Blast is a hashed seed-extend algorithm
 - Negative selection
 - Only some parts of a sequence are usually constrained



BLAST - Original version

Example:

Seed size = 4,

No mismatches in seed

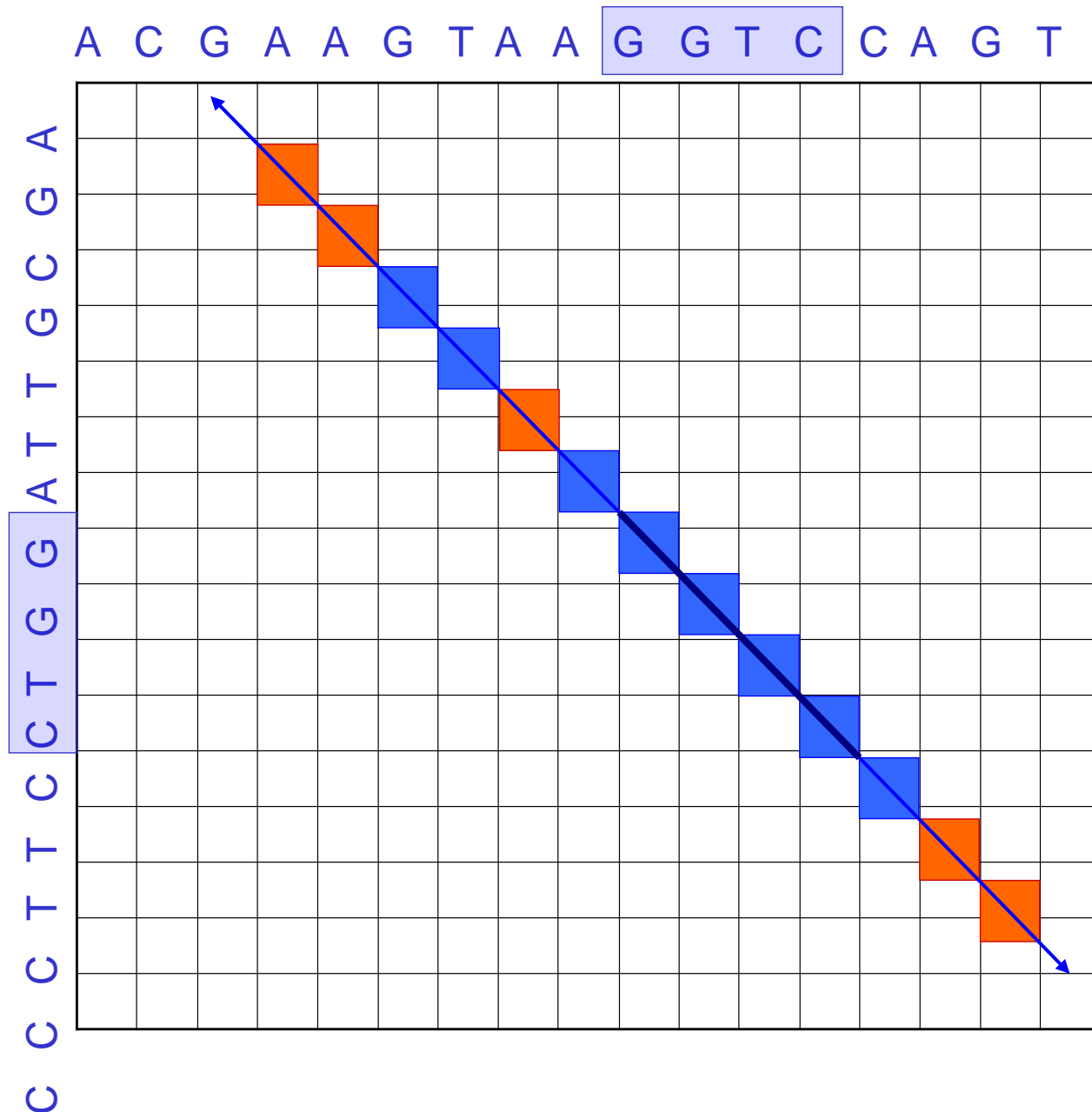
The matching word GGTC
initiates an alignment

Extension to the left and right
with no gaps until alignment
score falls below 50%

Output:

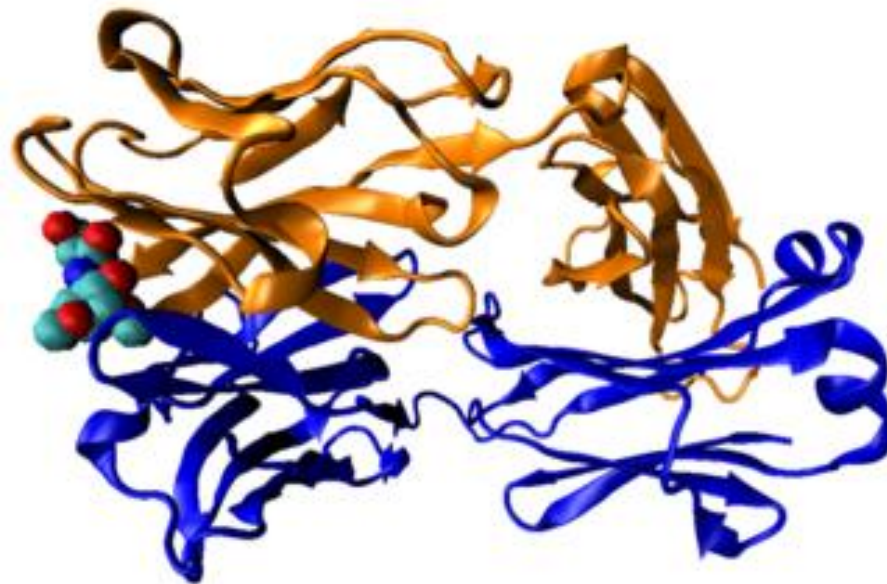
GTAAGGTCC

GTTAGGTCC



BLAST - Original algorithm

- Finding seeds significantly increases the speed of BLAST compared to doing a full local alignment over a whole sequence
- Will not guarantee the best solution
- BLAST first finds highly conserved or identical sequences which are then extended with a local alignment.



BLAST – Speed (or lack thereof)

- Typically BLAST will take approximately 0.1 – 1 second to search 1 sequence against a database
- Depends on size of database, e-value cutoff and number of hits to report selected
- 60 million reads equates to 70 CPU days!
- Even on multi-core systems this is too long!
- Especially if you have multiple samples!
- This is still true of FPGA and SIMD (vectorised) implementations of BLAST

When **NOT** to use *BLAST*

- A typical situation: you have lots DNA sequences and want to extend it or find where on a genome it maps.
- In other words, you want an **exact** or **near-exact** match to a sequence that is part of an **assembled genome**.
- Short reads require very fast algorithms for finding near-exact matches in genomic sequences:
 - BLAT
 - Highly recommended: the BLAT paper (Kent WJ (2003) *Genome Res* 12:656-64) – very well written
 - SOAP
 - Bowtie/Bowtie 2
 - MAQ
 - BWA
 - Shrimp2

Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - **Adapting hashed seed-extend algorithms to work with shorter reads**
 - Indel detection
 - Suffix/Prefix Tries
 - Other alignment considerations
 - Typical alignment pipeline
 - New methods of SNP calling

Adapting hashed seed-extend algorithms to work with shorter reads

- Improve seed matching sensitivity
 - Allow mismatches within seed
 - BLAST
 - Allow mismatches + Adopt spaced-seed approach
 - ELAND, SOAP, MAQ, RMAP, ZOOM
 - Allow mismatches + Spaced-seeds + Multi-seeds
 - SSAHA2, BLAT, ELAND2
- Above and/or Improve speed of local alignment for seed extension
 - Single Instruction Multiple Data
 - Shrimp2, CLCBio
 - Reduce search space to region around seed

Hashed seed-extend algorithms

- These are most similar to BLAST
- Are not designed to work with large databases
- **2 step process**
 - Identify a match to the seed sequence in the reference
 - Extend match using sensitive (but slow) Smith-Waterman algorithm (dynamic programming)

Seed-extend algorithm

Reference sequence:

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

Short read:

GTCATCGTACGATCGATAGATCGATCGATCGGCTA

Note that the short read has 1 difference wrt to reference

Seed-extend algorithm

Reference sequence:

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

Short read:

GTCATCGTACG ATCGATAGATCG ATCGATCGGCTA

11bp word

11bp word

11bp word

The algorithm will try to match each word to the reference. If there is a match at with any single word it will perform a local alignment to extend the match

Seed-extend algorithm

Reference sequence:

Seed Extend with Smith Waterman

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

GTCATCGTACGATCGAACGATCGATCGATCGGCTA

Short read:

GTCATCGTACG ATCGATAGATCG ATCGATCGGCTA

Here the algorithm is able to match the short read with a word length of 11bp

Seed-extend algorithm

Reference sequence:

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

Short read:

GTCATCGTACGATCGATCGATCGATCGGCA

Note that the short read has 3 differences
Possibly sequencing errors, possibly SNPs

Seed-extend algorithm

Reference sequence:

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

Short read:

GTCATCGTACG

11bp word

ATCGATCGATCG

11bp word

ATCGATCGGCAA

11bp word

Note that the short read has 3 differences

Seed-extend algorithm

Reference sequence:

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

Short read:

GTCATCGTACG ATCGATCGATCG ATCGATCGGCAA

No seeds match

Therefore the algorithm would find no hits at all!

Adapting hashed seed-extend algorithms to work with shorter reads

- Improve seed matching sensitivity
 - **Allow mismatches within seed**
 - **BLAST**
 - Allow mismatches + Adopt spaced-seed approach
 - ELAND, SOAP, MAQ, RMAP, ZOOM
 - Allow mismatches + Spaced-seeds + Multi-seeds
 - SSAHA2, BLAT, ELAND2
- Above and/or Improve speed of local alignment for seed extension
 - Single Instruction Multiple Data
 - Shrimp2, CLCBio
 - Reduce search space to region around seed

Adapting hashed seed-extend algorithms to work with shorter reads

- Improve seed matching sensitivity
 - **Allow mismatches within seed**
 - **BLAST**
 - **Allow mismatches + Adopt spaced-seed approach**
 - **ELAND, MAQ, RMAP, ZOOM**
 - Allow mismatches + Spaced-seeds + Multi-seeds
 - SSAHA2, BLAT, ELAND2
- Above and/or Improve speed of local alignment for seed extension
 - Single Instruction Multiple Data
 - Shrimp2, CLCBio
 - Reduce search space to region around seed

Consecutive seed

Consecutive seed 9bp with no mismatches:

ACTCCCATCGTCATCGTACTAGGGATCGTAACA

CCACTGTCCTCTACATAGGAACGA

TCATCGTAC
TCCTCTAC

Cannot find seed match due to A->C SNP
and G->C SNP

Reference sequence
SNP 'heavy' read

Even allowing for 2 mismatches in
the seed - no seeds match.
No hits!

Spaced seeds

To increase sensitivity we can use spaced-seeds:

1111111111

Consecutive seed template with *length* 9bp

ACTATCATCGTACACAT

Reference

TCATCGTAC

Query

11001100110011001

Spaced-seed template with *weight* 9bp

ACTATCATCGTACACAT

Reference

ACTCTCACC GTACACAT

Query

Spaced seeds

Spaced seed with weight 9bp and no mismatches:

ACTCCCATTTGTCATCGTACTTGGGATCGTAACA

Reference sequence

CCACTGTAAATCGTACATGGGAACGA

SNP 'heavy' read



CCATTGTCATCGTACAT

CCXXTGXXATXXTAXXT

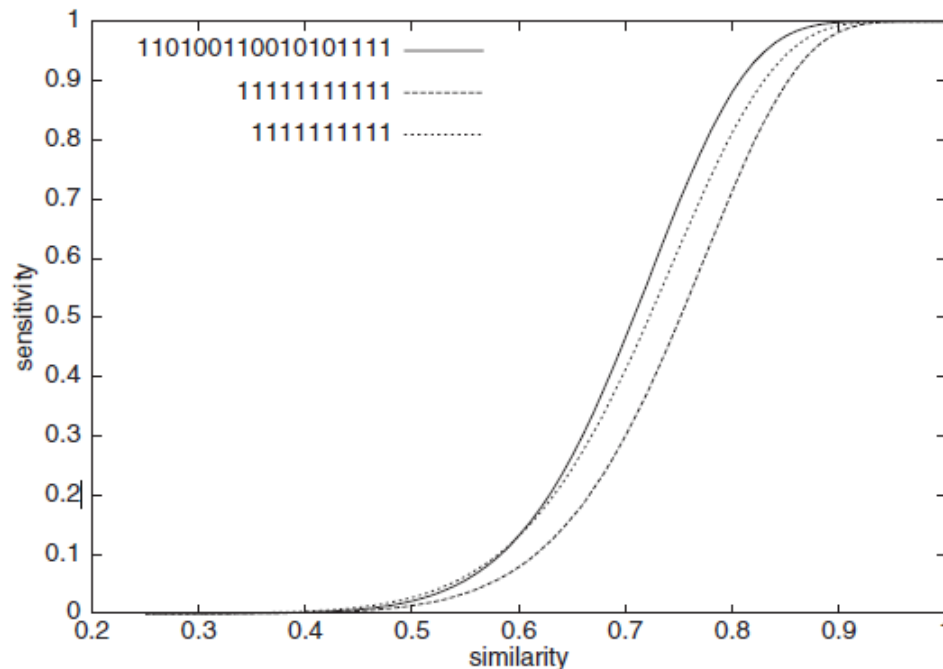
Despite SNPs – seed matched with 0 mismatches

Can now extend with Smith-Waterman or other local alignment

Spaced seeds

Spaced seeds:

- A seed template '111010010100110111' is 55% more sensitive than BLAST's default template '1111111111' for two sequences of 70% similarity
- Typically seeds of length ~30bp and allow up to 2 mismatches in short read datasets



Contents

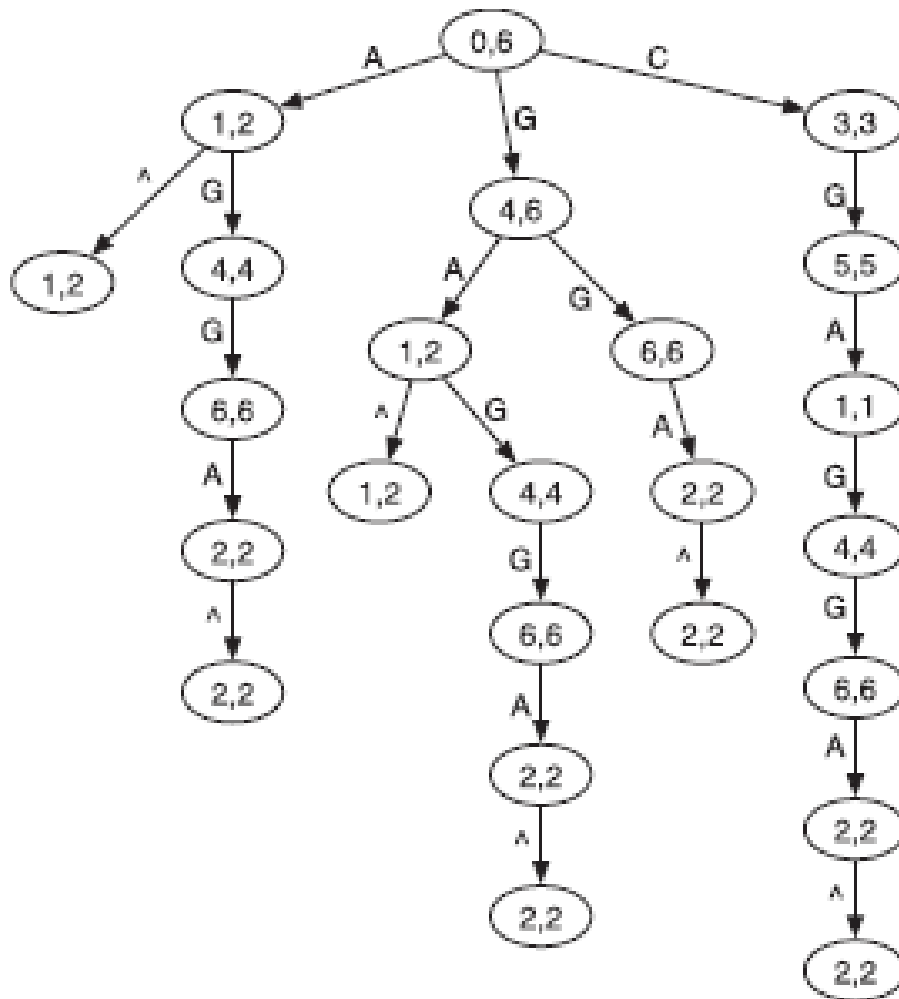
- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - **Suffix/Prefix Tries**
 - Indel detection
 - Other alignment considerations
 - Typical alignment pipeline
 - New methods of SNP calling

Suffix-Prefix Trie

- Trie – data structure which stores the suffixes (i.e. ends of a sequence)
- A family of methods which uses a Trie structure to search a reference sequence
 - Bowtie
 - BWA aln (<70bp reads) and MEM algorithm (>70bp reads)
 - SOAP version 2
- Key advantages:
 - Alignment of multiple copies of an identical sequence in the reference only needs to be done once
 - Use of an FM-Index to store Trie can drastically reduce memory requirements (e.g. Human genome can be stored in 2Gb of RAM)
 - Burrows Wheeler Transform to perform fast lookups

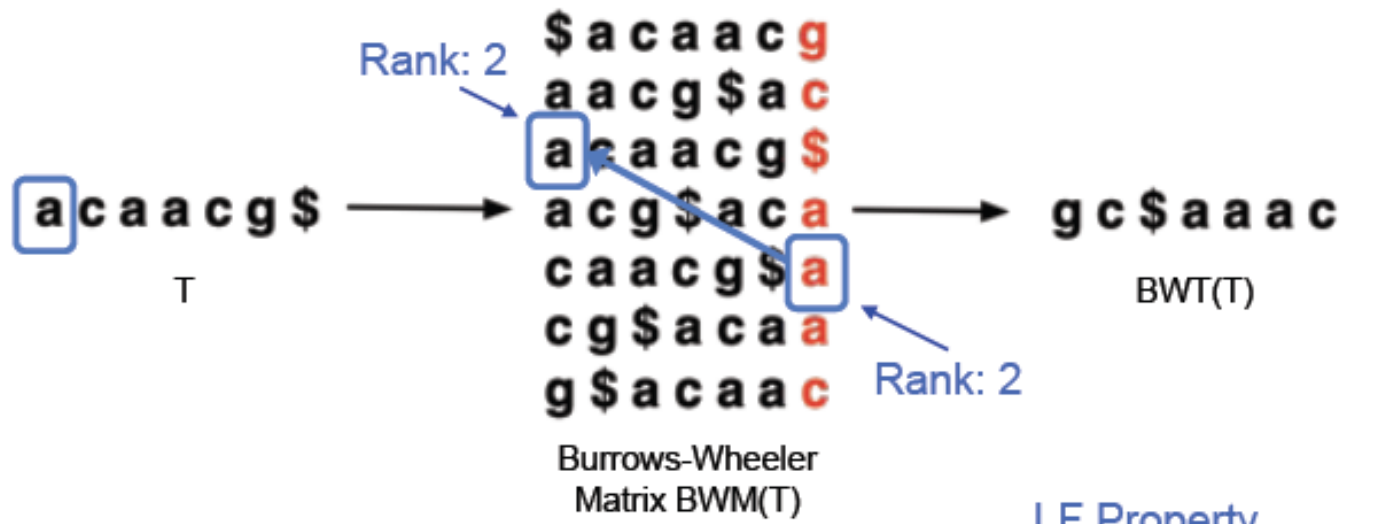
Suffix Trie

Read
AGGAGC



Heng Li & Nils Homer.
Sequence alignment
algorithms for next-
generation sequencing.
Briefings in
Bioinformatics. Vol 11.
No 5. 473 483, 2010

Suffix Trie



- $BWT(T)$ is the index for T

A block sorting lossless data compression algorithm.

Burrows M, Wheeler DJ (1994) *Digital Equipment Corporation*. Technical Report 124

Burrows-Wheeler Algorithm

- Encodes data so that it is easier to compress
- Burrows-Wheeler transform of the word BANANA
- Can later be reversed to recover the original word

Transformation				
Input	All Rotations	Sorting All Rows in Alphabetical Order by their first letters	Taking Last Column	Output Last Column
<div style="border: 1px dashed blue; padding: 5px; width: fit-content;"> <code>^BANANA </code> </div>	<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>	<div style="border: 1px dashed blue; padding: 5px; width: fit-content;"> <code>BNN^AA A</code> </div>

More Burrows-Wheeler

Input	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES
Burrows-Wheeler Output	TEXYDST.E.IXIXIXSSMPPS.B..E.S.EUSFXDIIIOIIT

Repeated characters mean that it is easier to compress

Suffix Trie for a bacterial genome would be $> 1\text{Tb}$

We have to compress it

Use FM-Index/BW transform to do this compression

Bowtie/BWA example

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATCTA

Bowtie/BWA example

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATCTA

Bowtie/BWA example

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATC **TA**

Bowtie/BWA example

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATCTA

Bowtie/BWA example

Reference



BWT(Reference)

Query:

AATGATACGGCGAC **CACCGAGATCTA**

Bowtie/BWA example

Reference



BWT(Reference)



Query:

AATGATACGGCGACCAACGAGATCTA

Bowtie/BWA example

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATCTA

Bowtie/BWA example

Reference



BWT(Reference)



Query:

AATG**T**TACGGCGACCACCGAGATCTA

Bowtie/BWA example

Reference



BWT(Reference)



Query:

AATGTTACGGCGACCACCGAGATCTA


Bowtie/Soap2 vs. BWA

- Bowtie 1 and Soap2 cannot handle gapped alignments
 - No indel detection => Many false SNP calls

Bowtie/Soap2:

```
ACTCCATTGTCATCGTACTTGGGATCGTAACA   Reference
      CCATTGTCATCGTACTTGGGATCTA
          TCATCGTACTTGGGATCTA
              TTGGGATCTA
```

False SNPs



N.B. Bowtie2 can handle gapped alignments

Bowtie/Soap2 vs. BWA

- Bowtie 1 and Soap2 cannot handle gapped alignments
 - No indel detection => Many false SNP calls

BWA:

```
ACTCCCATTGTCATCGTACTTGGGATCGTAACA   Reference
      CCATTGTCATCGTACTTGGGATC-TA
          TCATCGTACTTGGGATC-TA
              TTGGGATC-TA
```

N.B. Bowtie2 can handle gapped alignments

Comparison

Hash referenced spaced seeds

- Requires ~50Gb of memory
- Runs 30-fold slower
- Is much simpler to program
- Most sensitive

Indexed Suffix/Prefix Trie

- Requires <2Gb of memory
- Runs 30-fold faster
- Is much more complicated to program
- Least sensitive

There are limits however

With longer 100-300 bp reads, multiple indels or variable regions longer than a few bp are likely to be missed

```
ACTCCCATTGTCATCGTACTTGGGATCGTAACA   Reference
      CCATTGTCAACCATCTAGTAGCT-TA
          TCAACCATCTAGTAGCT-TA
              ACCATCTA-TA
```

You only find what you are looking for

- What happens if there are SNPs and Indels in the same region?

Let's assume that the SNP caller made this call of a single SNP:

ATGT**A**TGTA
ATGT**G**TGTA

and the indel caller produced this call of a 3 base deletion:

ATGTATGTA
ATGT- - - TA

Should we assume this is a heterozygous SNP opposite a heterozygous Indel or a more complex locus?

Comparison

- Bowtie's reported 30-fold speed increase over hash-based methods with small loss in sensitivity
- Limitations to Trie-based approaches:
 - Only able to find alignments within a certain 'edit distance'
 - Important to quality clip reads (-q in BWA)
 - Non-A/C/G/T bases on reads are often treated as mismatches
 - Make sure Ns are removed!

Hash based approaches are more suitable for divergent alignments

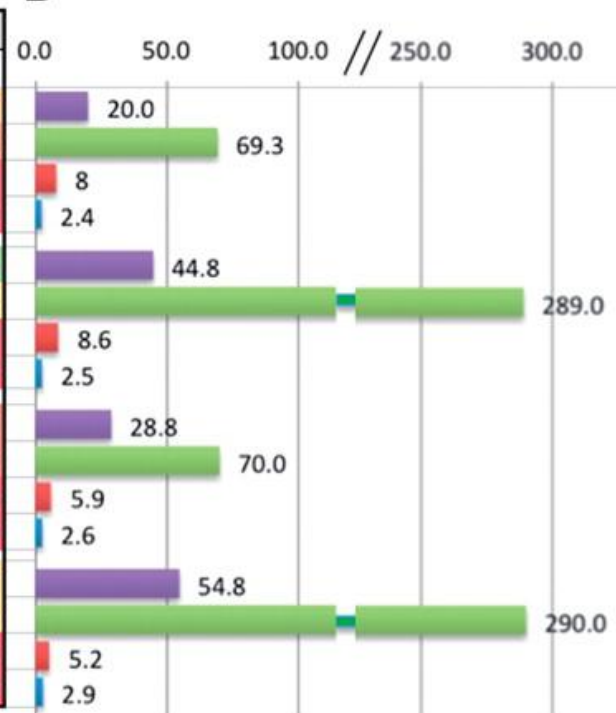
- Rule of thumb:
 - <2% divergence -> Trie-based
 - E.g. human alignments
 - >2% divergence -> Seed-extend based approach
 - E.g. wild mouse strain alignments

Precision and recall by amount of variation for 4 datasets, by polymorphism: (number of SNPs, Indel size)

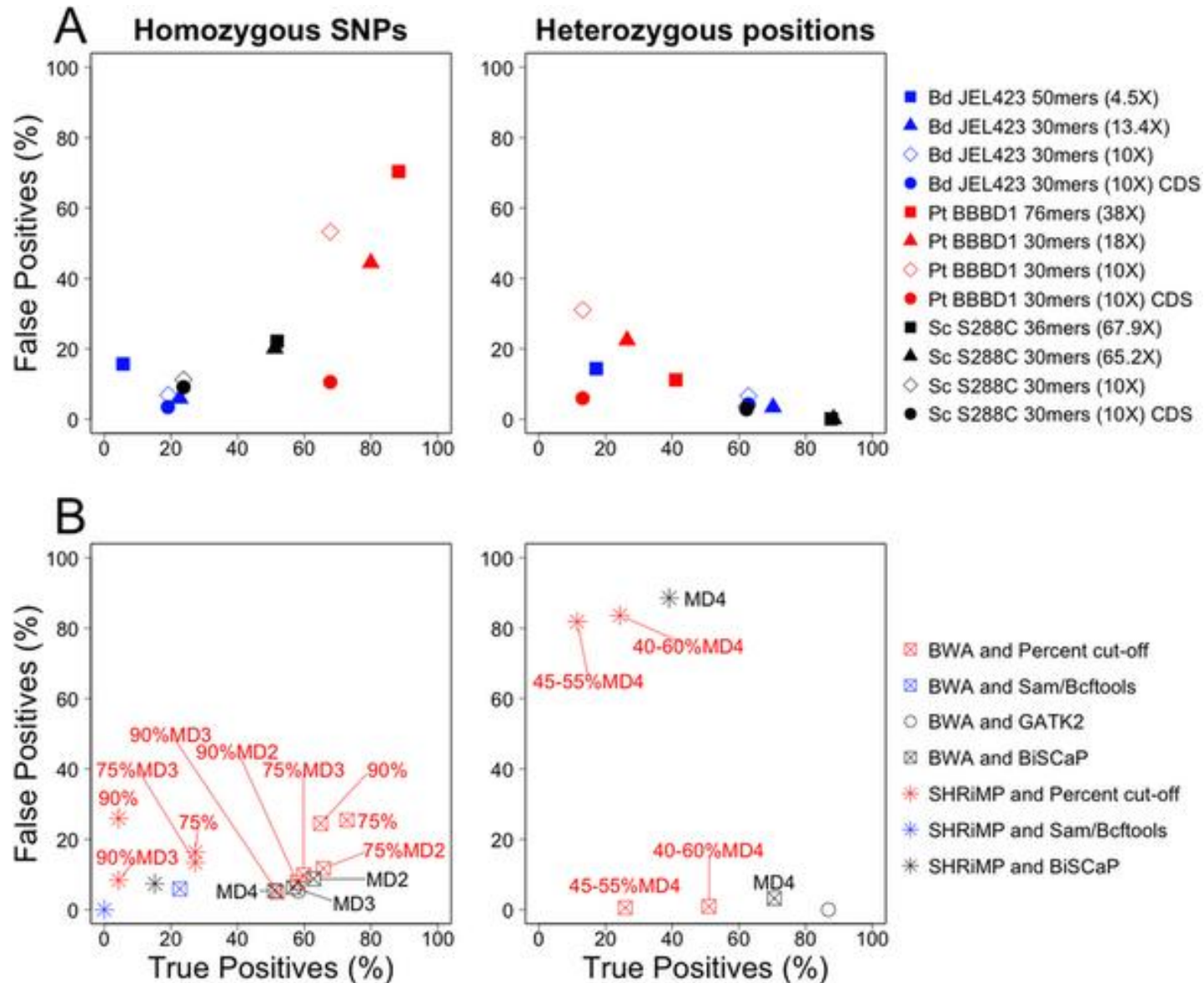
A

	Program	(0,0)		(1,0)		(2,0)		(4,0)		(0,3)		(1,3)		(2,3)		(4,3)	
		Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.
50 paired	SHRiMP	99.7	96.6	99.6	96.4	99.6	95.7	99.3	89.3	99.3	93.5	99.3	90.6	98.6	85.7	97.6	69.7
	BFAST	95.4	93.8	94.3	91.6	92.6	86.2	87.0	63.5	91.6	78.8	89.3	71.8	86.8	61.9	80.7	38.8
	BWA	91.1	65.2	85.4	27.7	64.7	5.4	17.7	0.3	62.0	4.4	49.2	1.5	29.6	0.4	11.9	0.1
	Bowtie	97.5	46.6	97.5	11.1	96.9	1.0	0.0	0.0	97.1	1.3	100	0.2	100	0.0	0.0	0.0
75 paired	SHRiMP	99.6	97.5	99.6	97.2	99.6	97.3	99.6	96.9	99.3	96.6	99.5	96.9	99.4	96.5	99.2	94.5
	BFAST	97.4	97.1	97.1	96.8	96.8	96.5	95.9	94.5	96.4	96.0	96.0	95.5	95.9	94.8	94.1	89.5
	BWA	93.2	62.3	86.5	30.2	68.2	8.8	14.7	0.4	65.0	7.5	41.5	2.2	22.4	0.6	11.7	0.1
	Bowtie	98.1	18.1	98.4	2.6	96.2	0.1	100	0.0	97.1	0.5	100	0.0	0.0	0.0	0.0	0.0
50 single	SHRiMP	99.7	93.3	98.9	92.6	98.0	91.1	94.8	72.5	97.0	89.5	95.3	83.5	93.0	69.6	83.4	25.6
	BFAST	98.9	93.0	97.9	90.5	96.2	83.7	87.7	50.7	95.2	80.4	92.8	68.7	89.0	53.5	78.0	24.6
	BWA	95.3	79.7	93.0	33.7	71.8	2.1	15.2	0.0	89.5	5.6	83.7	1.1	61.9	0.1	0.0	0.0
	Bowtie	95.2	65.5	92.1	15.7	49.1	0.3	2.5	0.0	92.1	2.2	85.4	0.4	36.8	0.0	0.0	0.0
75 single	SHRiMP	99.7	96.0	99.6	95.8	99.4	95.6	98.9	94.4	99.2	95.5	98.8	94.9	98.5	93.7	97.2	79.7
	BFAST	99.3	96.0	99.1	95.6	98.8	95.1	97.4	91.6	98.5	95.1	98.0	94.1	97.4	92.1	94.3	81.6
	BWA	97.5	78.2	97.0	38.0	95.1	6.5	56.4	0.0	96.7	9.4	94.6	1.2	90.4	0.2	100	0.0
	Bowtie	97.4	42.0	96.2	6.0	75.7	0.1	0.0	0.0	95.8	0.8	96.3	0.1	100	0.0	0.0	0.0

B



False discovery rates for variants were ascertained using cFDR for three fungal NGS datasets



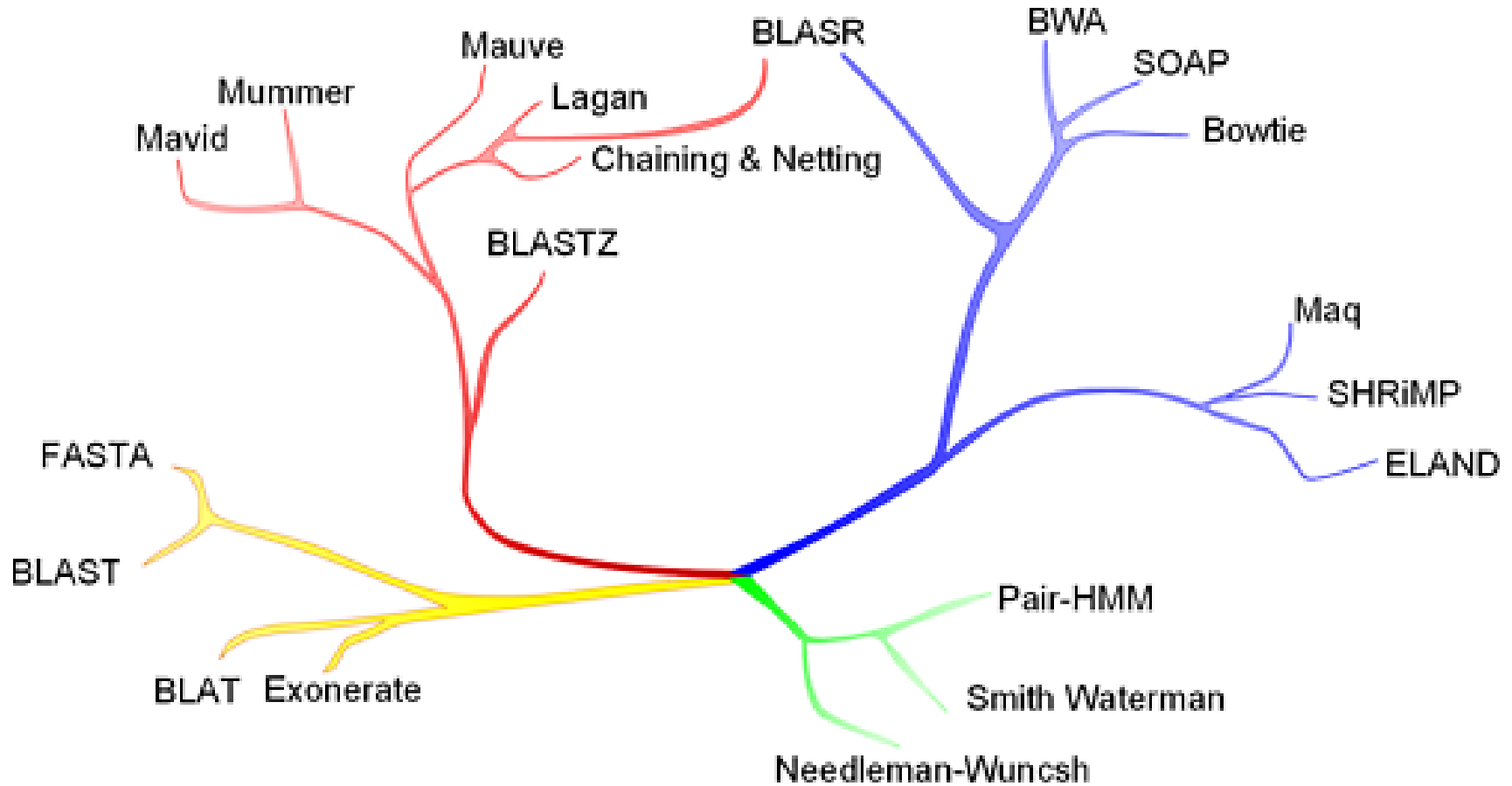
Summary of open-source short read alignment programs

Program	Algorithm	SoLID	Long reads	Gapped alignment	Paired-end	Quality scores used?
Bfast	Hashing ref	Yes	No	Yes	Yes	No
Bowtie2*	FM-Index	Yes	Yes	Yes	Yes	Yes
Blat	Hashing ref	No	Yes	Yes	No	No
BWA	FM-Index	Yes	Yes	Yes	Yes	No
MAQ	Hashing reads	Yes	No	Yes	Yes	Yes
Mosaik	Hashing ref	Yes	Yes	Yes	Yes	No
Novoalign	Hashing ref	No	No	Yes	Yes	Yes
Shrimp2	Hashing ref	Yes	Yes	Yes	Yes	Yes
SOAP2	FM-Index	No	No	No	Yes	Yes
SSAHA2	Hashing ref.	No	No	No	Yes	Yes

Heng Li & Nils Homer. Sequence alignment algorithms for next-generation sequencing. Briefings in Bioinformatics. Vol 11. No 5. 473-483, 2010

* Bowtie1 does not support gapped alignments

Aligner phylogeny



Whole genome

Short read

Pairwise heuristic

Sensitive global aligners

Alignment format for short reads – Sequence AlignMent (SAM format)

- Plain text format – human readable (sort-of)
- Eleven mandatory fields and a variable amount of optional fields.
- The optional fields are a key-value pair of TAG:TYPE:VALUE. These store extra information
- Can be converted to Binary AlignMent format (BAM) to save space and speed up look-up operations using SAMTools

Alignment format for short reads – Sequence Alignment (SAM format)

Table 1. Mandatory fields in the SAM format

No.	Name	Description
1	QNAME	Query NAME of the read or the read pair
2	FLAG	Bitwise FLAG (pairing, strand, mate strand, etc.)
3	RNAME	Reference sequence NAME
4	POS	1-Based leftmost POSition of clipped alignment
5	MAPQ	MAPping Quality (Phred-scaled)
6	CIGAR	Extended CIGAR string (operations: MIDNSHP)
7	MRNM	Mate Reference NaMe ('=' if same as RNAME)
8	MPOS	1-Based leftmost Mate POSition
9	ISIZE	Inferred Insert SIZE
10	SEQ	Query SEQUENCE on the same strand as the reference
11	QUAL	Query QUALity (ASCII-33=Phred base quality)

SAM format – Optional fields

Tag	Type	Description
X?	?	Reserved fields for end users (together with Y? and Z?)
AM	i	The smallest template-independent mapping quality of fragments in the rest
AS	i	Alignment score generated by aligner
BQ	Z	Offset to base alignment quality (BAQ), of the same length as the read sequence. At the i -th read base, $BAQ_i = Q_i - (BQ_i - 64)$ where Q_i is the i -th base quality.
CM	i	Edit distance between the color sequence and the color reference (see also NM)
CQ	Z	Color read quality on the original strand of the read. Same encoding as QUAL; same length as CS.
CS	Z	Color read sequence on the original strand of the read. The primer base must be included.
E2	Z	The 2nd most likely base calls. Same encoding and same length as QUAL.
FI	i	The index of fragment in the template.
FS	Z	Fragment suffix.
LB	Z	Library. Value to be consistent with the header RG-LB tag if @RG is present.
H0	i	Number of perfect hits
H1	i	Number of 1-difference hits (see also NM)
H2	i	Number of 2-difference hits
HI	i	Query hit index, indicating the alignment record is the i -th one stored in SAM
IH	i	Number of stored alignments in SAM that contains the query in the current record
MD	Z	String for mismatching positions. <i>Regex</i> : $[0-9]+((([ACGTN] \^[ACGTN]+)[0-9]+)^*)^1$
MQ	i	Mapping quality of the mate/next fragment
NH	i	Number of reported alignments that contains the query in the current record
NM	i	Edit distance to the reference, including ambiguous bases but excluding clipping
OQ	Z	Original base quality (usually before recalibration). Same encoding as QUAL.
OP	i	Original mapping position (usually before realignment)
OC	Z	Original CIGAR (usually before realignment)
PG	Z	Program. Value matches the header PG-ID tag if @PG is present.
PQ	i	Phred likelihood of the template, conditional on both the mapping being correct
PU	Z	Platform unit. Value to be consistent with the header RG-PU tag if @RG is present.
Q2	Z	Phred quality of the mate/next fragment. Same encoding as QUAL.
R2	Z	Sequence of the mate/next fragment in the template.
RG	Z	Read group. Value matches the header RG-ID tag if @RG is present in the header.
SM	i	Template-independent mapping quality
TC	i	The number of fragments in the template.
U2	Z	Phred probability of the 2nd call being wrong conditional on the best being wrong. The same encoding as QUAL.
UQ	i	Phred likelihood of the fragment, conditional on the mapping being correct

SAM output

```
amaxwell@pinfish:~/ecoli/illum/bwamem/i4m1_1/seqAssist1/result
[amaxwell@pinfish result]$ head -10 aln.sam
@SQ      SN:gi|49175990|ref|NC_000913.2| LN:4639675
HWI-EAS397:8:1:1067:18713#CTTGTA      16      gi|49175990|ref|NC_000913.2|      2909788 35      13S36M *      0      0      TACAATAACCCCCGCCCCTGGGTAATAAAGCCGACAATCTCATCTCCA      BBBBBBBBBBBBBBBB^aa\^^\V_dkwdfad
affccffcYadcL^Y      NM:i:1 AS:i:31 XS:i:0
HWI-EAS397:8:1:1070:11813#CTTGTA      16      gi|49175990|ref|NC_000913.2|      4516890 24      49M      *      0      0      CGTCGGTGCTAAAGCAGGTGACGCTGGCTGTTTACGCGTATGACAGG      BaYa[\Kaa0GVHLZLIXH`]VaJR]V]V_J]
Qca\aaS[J_aacccc      NM:i:2 AS:i:39 XS:i:30
HWI-EAS397:8:1:1070:11813#CTTGTA      256      gi|49175990|ref|NC_000913.2|      278711 0      30M19H *      0      0      *      *      NM:i:0 AS:i:30
HWI-EAS397:8:1:1070:11813#CTTGTA      256      gi|49175990|ref|NC_000913.2|      290182 0      30M19H *      0      0      *      *      NM:i:0 AS:i:30
HWI-EAS397:8:1:1070:11813#CTTGTA      272      gi|49175990|ref|NC_000913.2|      1049415 0      19H30M *      0      0      *      *      NM:i:0 AS:i:30
HWI-EAS397:8:1:1121:19907#CTTGTA      0      gi|49175990|ref|NC_000913.2|      953460 60      49M      *      0      0      AGAGAGAAGCAAATGCCGCCAACCAAGTTTTGCCATGCCGAAGGGCATTG      SIaJaL^ZT`[da^WcacfK^acSacfffff[
cY^dcdd^f]\cffff      NM:i:0 AS:i:49 XS:i:0
HWI-EAS397:8:1:1138:20219#CTTGTA      4      *      0      0      *      *      0      0      TAAAAAAGGCCATAACCTGCCGCTTTGTATAATAAAAAAGGGCCCGG      Yd[cK^__wdc\ddada]df[fccff_f_ffffcfdcBBBBBBBBBB A
S:i:0 XS:i:0
HWI-EAS397:8:1:1149:21173#GTTGTA      4      *      0      0      *      *      0      0      GTTCGTTTTTTTTGTACAGTGCCAACATCTGCTCCCGCCAATGGTGGC      Scc_ZccccaVZR^[W_K_M_\^VSQ]\aS^aVaXL]wY^^IaaaB A
S:i:0 XS:i:0
HWI-EAS397:8:1:1152:20684#CTTGTA      4      *      0      0      *      *      0      0      GCTGATGCGTACGTGAAAGGGCGGGACACCTTGGATGTATGGTTGAGT      aQ^aZ^M0aXOX_V\_cc_BBBBBBBBBBBBBBBBBBBBBBBBBBB A
S:i:0 XS:i:0
[amaxwell@pinfish result]$
```

Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Indel detection
 - Suffix/Prefix Tries
 - **Other alignment considerations**
 - Typical alignment pipeline
 - New methods of SNP calling

Other alignment considerations

- Indel detection
- Effect of paired-end alignments
- Using base quality to inform alignments
- PCR duplicates
- Methylation experiments – bisulfite treated reads
- Multi-mapping reads
- Aligning spliced-reads from RNA-seq experiments
- Local realignment to improve SNP/Indel detection
- Platform specific errors
- Unmapped reads

Indel detection

Spaced seed with weight 9bp and no mismatches:

ACTCCCATTTGTCATCGTACTTGGGATCGTAACA

Reference sequence

CCATTGTCATGTACTTGGGATCGT

Read containing a
deletion



CCATTGTCATCGTACAT

CCXXTGXXATXXACXXG

Seed not matched due to frame shift caused
by gap

No seed match. No alignment!

Indel detection

Reference sequence:

Seed Extend with Smith Waterman

...ACTGGGTCATCGTACGATCGATCGATCGATCGATCGGCTAGCTAGCTA...

GTCATCGTACGATCGA-CGATCGATCGATCGGCTA

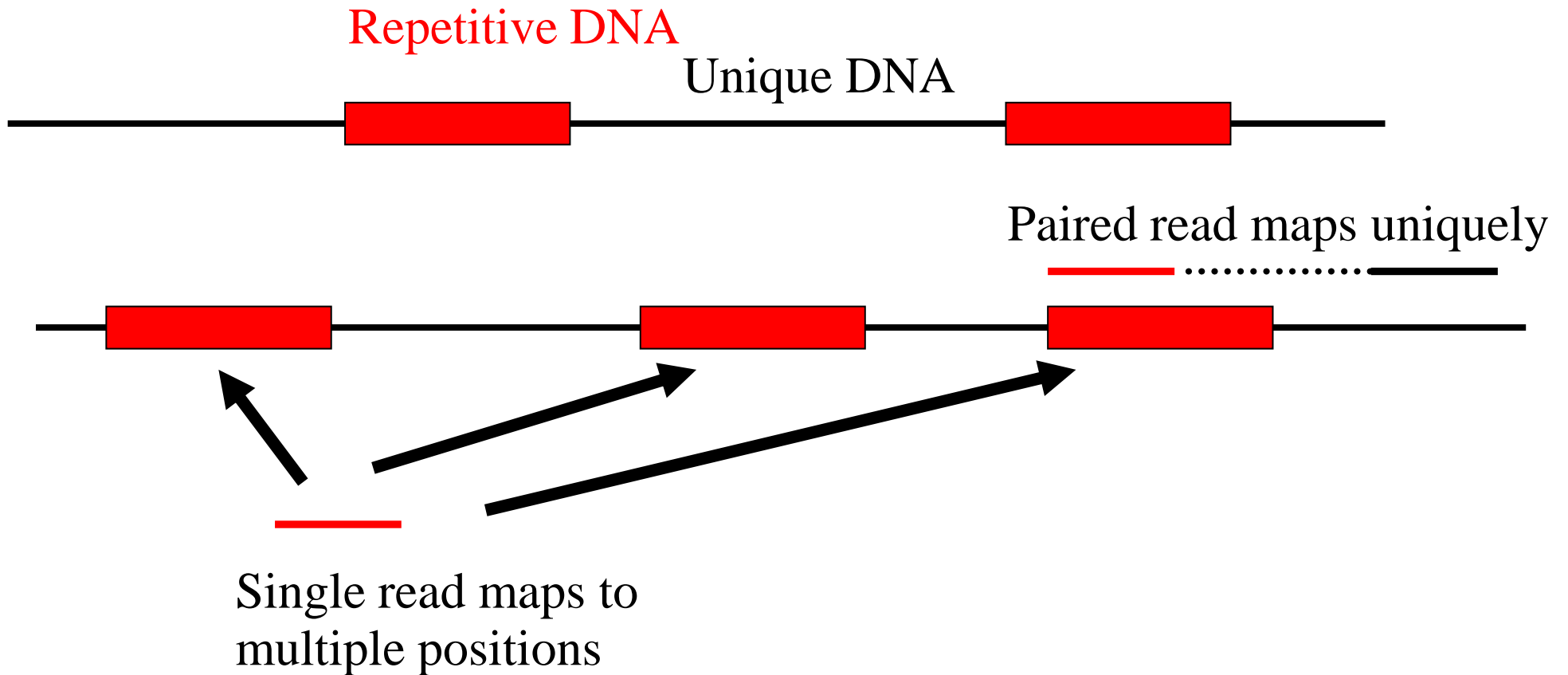
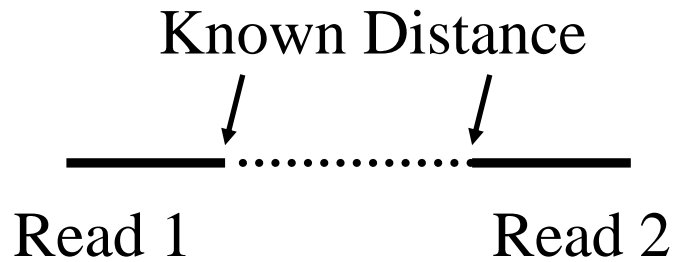
Most alignment programs can only detect gaps in
Smith-Waterman phase
once a seed has been identified. Some algorithms (e.g.
Bowtie) do not allow gaps at this stage to improve
speed

**This reduces sensitivity especially with multiple
insertions in a small region**

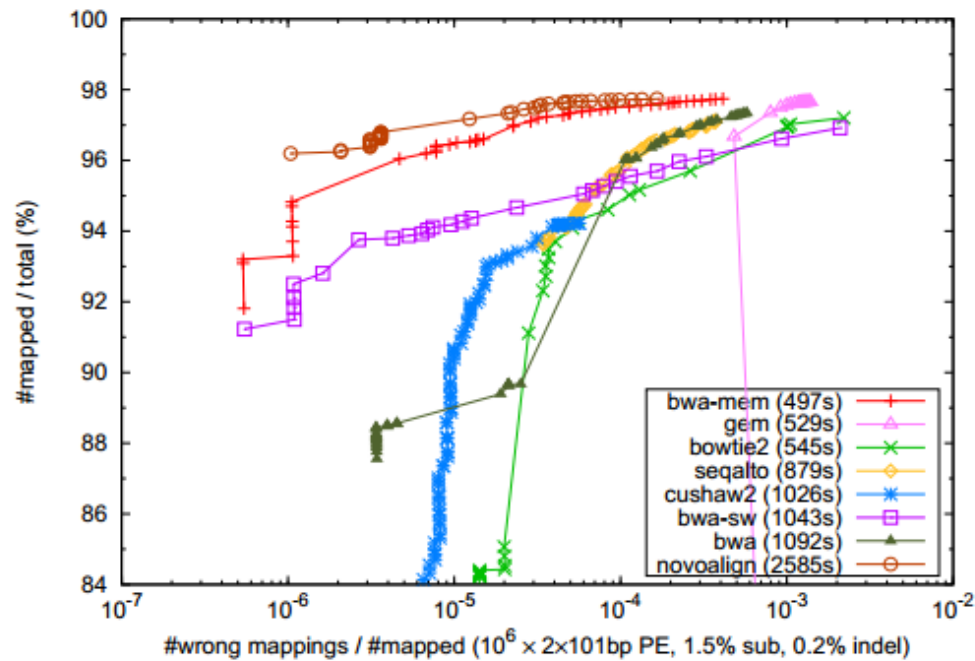
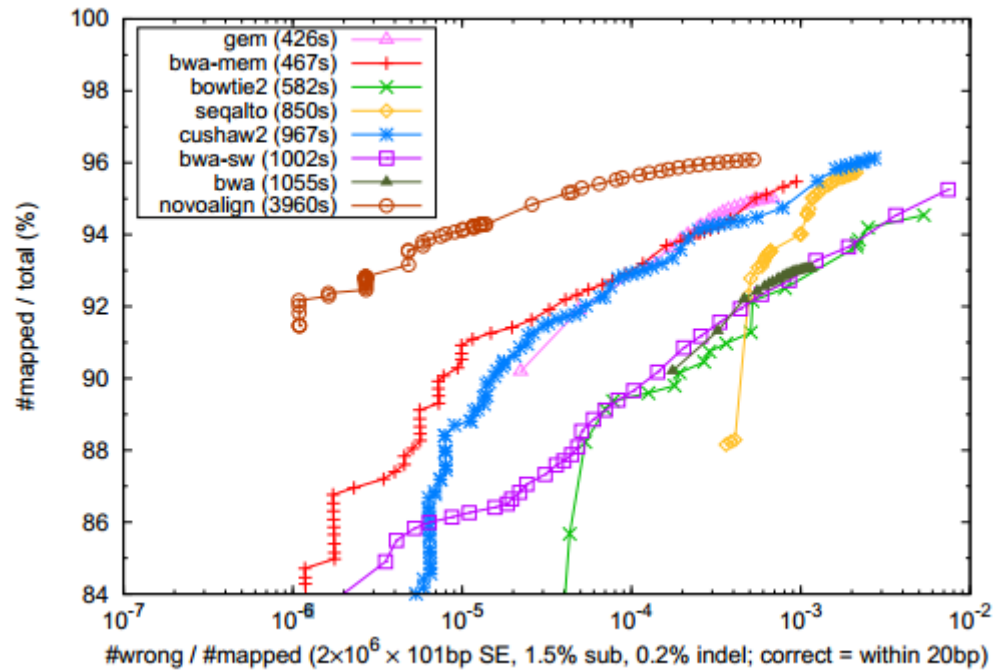
Indel detection

- Some algorithms do allow gaps within seed
 - Indel seeds for homology search *Bioinformatics* (2006) 22(14): e341-e349
doi:10.1093/bioinformatics/btl263
 - Weese D, Emde AK, Rausch T, et al. RazerS—fast read mapping with sensitivity control. *Genome Res* 2009;19:1646–54
 - Rumble SM, Lacroute P, Dalca AV, et al. SHRiMP: accurate mapping of short color-space reads. *PLoS Comput Biol* 2009;5:e1000386
- Use of multiple seeds
 - Especially useful for longer reads (>50bp)
 - Li R, Li Y, Kristiansen K, et al. SOAP: short oligonucleotide alignment program. *Bioinformatics* 2008;24:713–4
 - Jiang H, Wong WH. SeqMap: mapping massive amount of oligonucleotides to the genome. *Bioinformatics* 2008;24: 2395–6

Paired-end reads are important



Effect of paired-end alignments

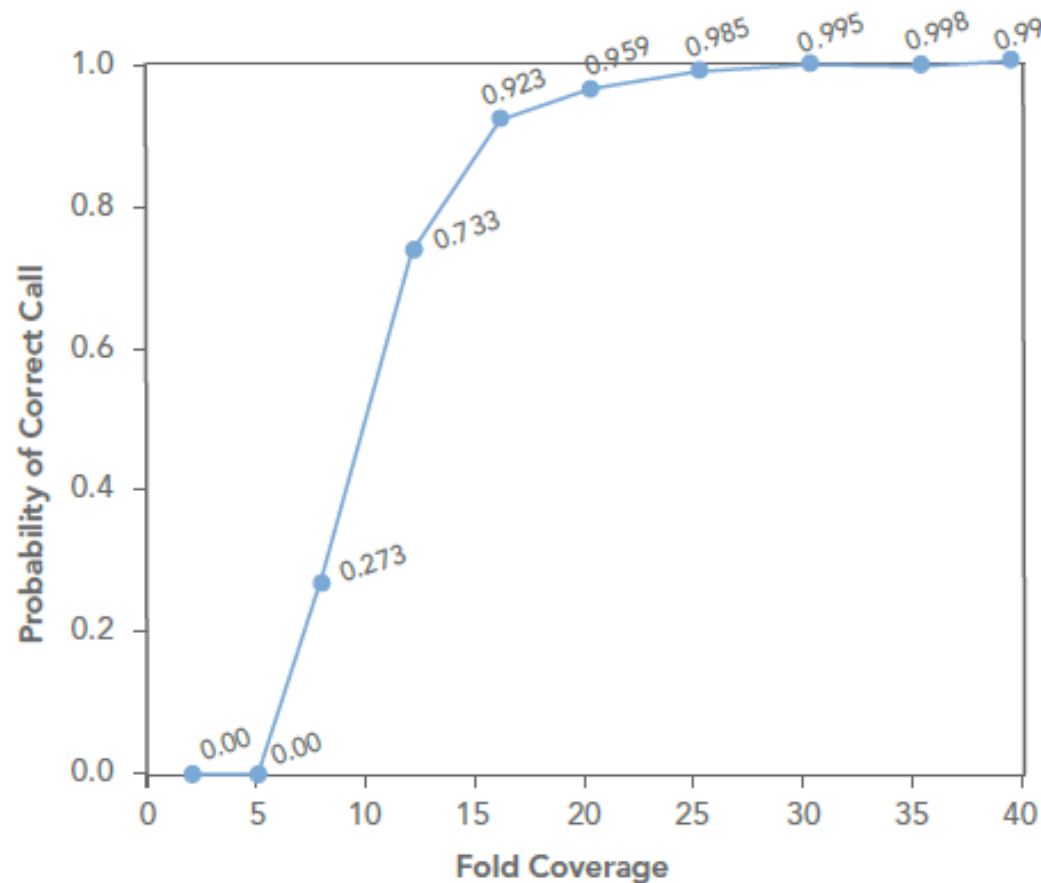


BWA-MEM

<http://arxiv.org/pdf/1303.3997v2.pdf>

Effect of coverage on SNP call accuracy

- Depends crucially on ploidy
- Bacterial genomes can get away with 10-20x
- For human genomes and other diploids 20-30x is regarded as standard
- Poly-ploids (e.g wheat) may need much higher coverage



Source – Illumina Tech Note
Human diploid sample

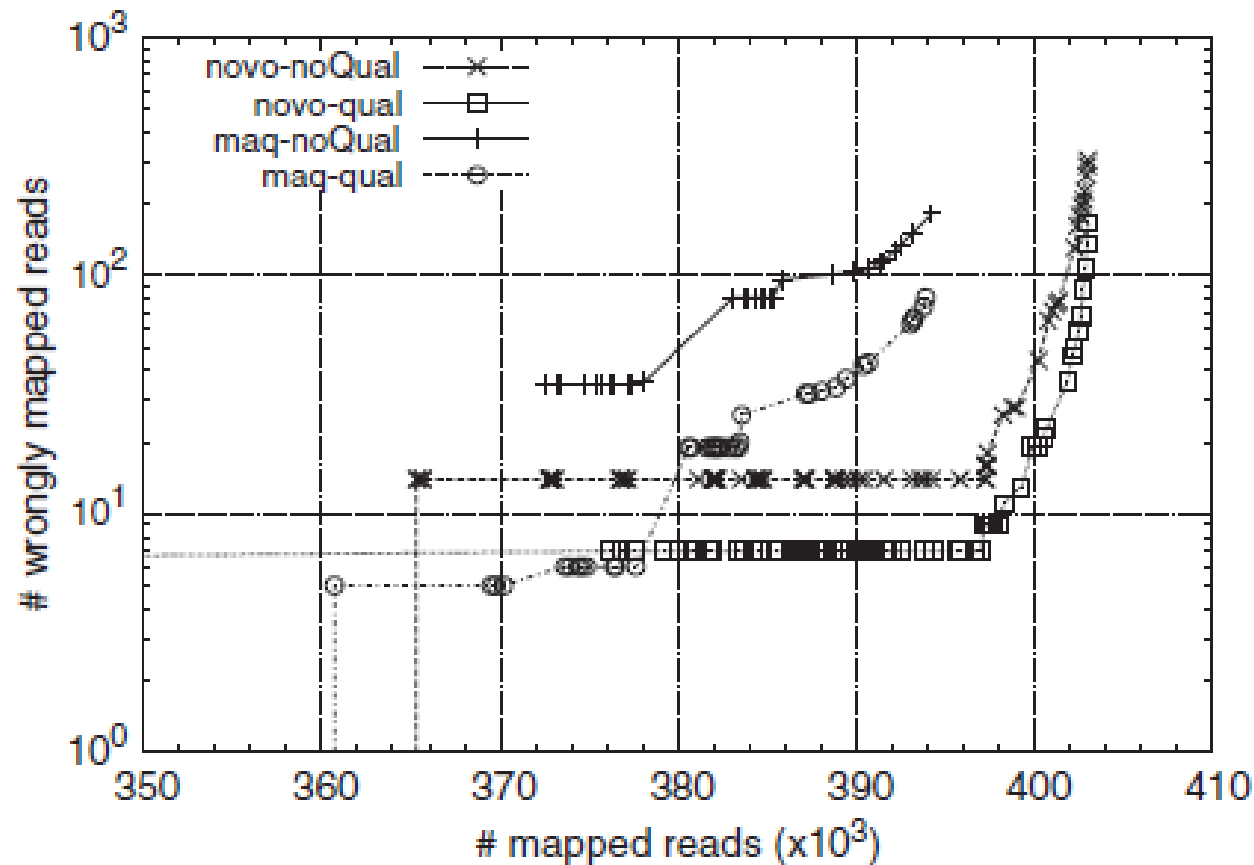
PCR duplicates

- **2nd generation sequencers are not single-molecule sequencers**
 - All have at least one PCR amplification step
 - Can result in duplicate DNA fragments
 - This can bias SNP calls or introduce false SNPs
- **Generally duplicates only make up a small fraction of the results**
 - Good libraries have < 2-3% of duplicates
 - SAMtools and Picard can identify and remove these when aligned against a reference genome
 - Do NOT do this for RNA and ChIP-seq data!

PCR duplicates



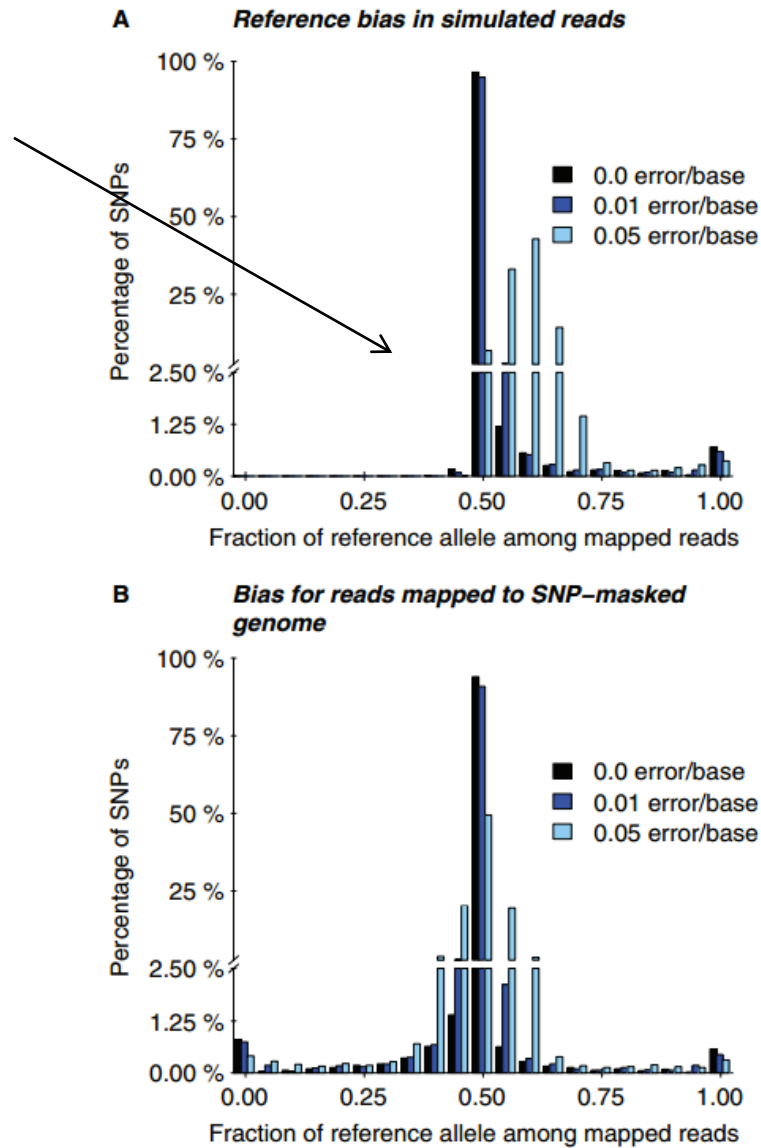
Base quality impacts on read mapping



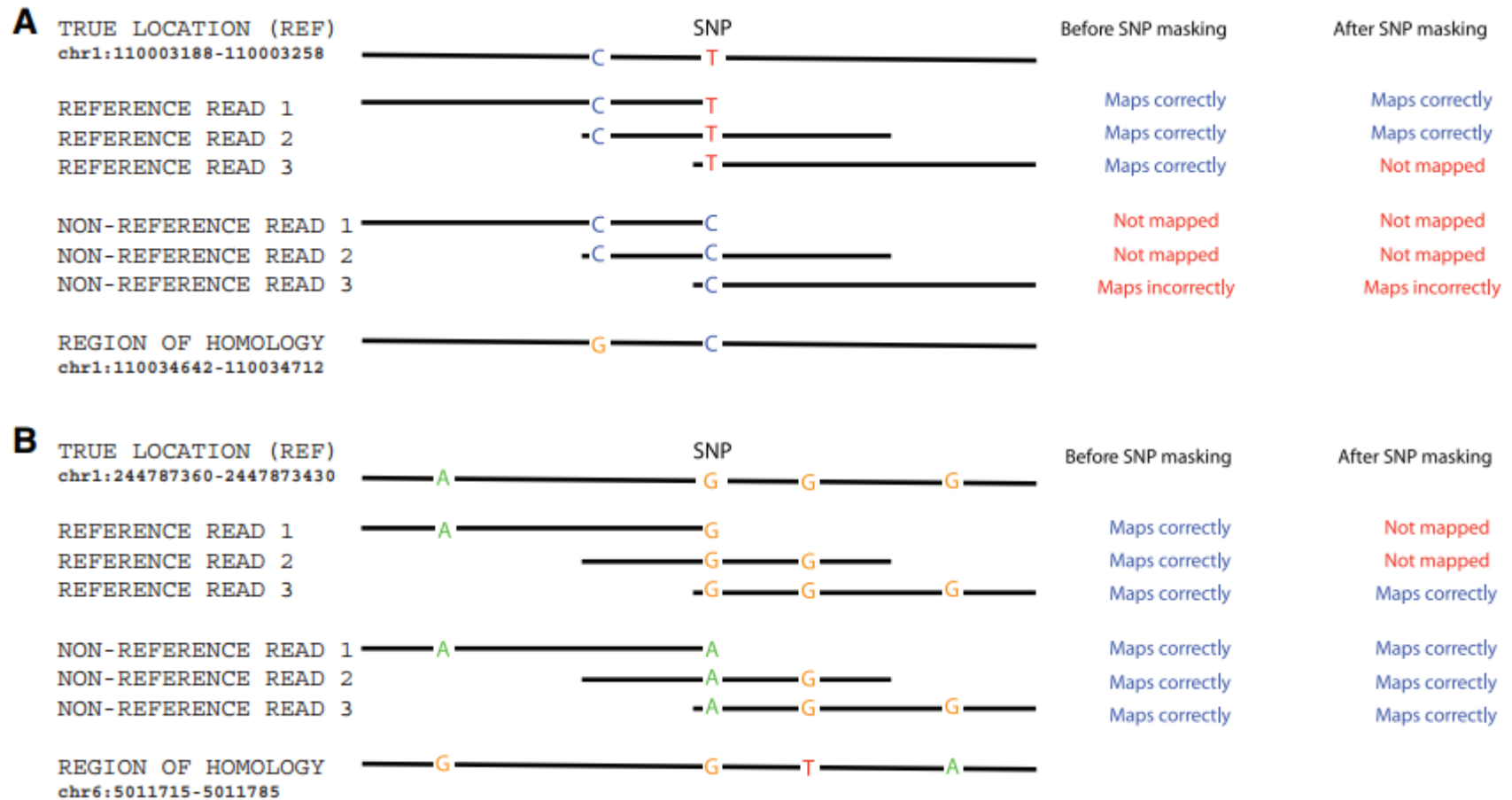
Heng Li & Nils Homer.
Sequence alignment
algorithms for next-
generation sequencing.
Briefings in
Bioinformatics. Vol 11.
No 5. 473-483, 2010

Allele-specific sequencing

Missing alternate allele



Biasing towards and against the reference allele



Methylation experiments

Unmethylated cytosine

5' - atcgCCcgataCga - 3'
3' - tagcgggCtatgct - 5'

Bisulfite
treatment

5' - atcgUUcgataUga - 3'

3' - tagcgggUtatgct - 5'

5' - atcgTTcgataIga - 3' (1)

3' - tagcAAgCtataAct - 5' (2)

Amplification

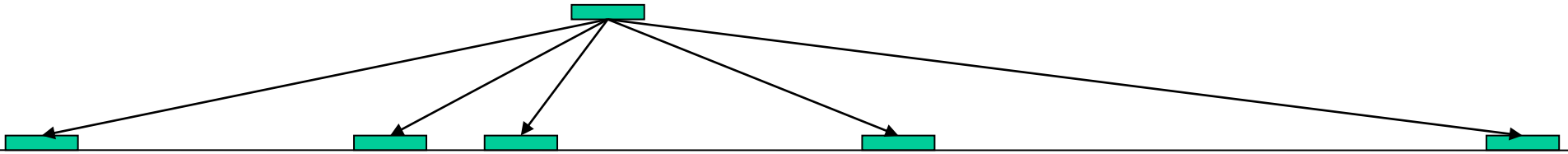
5' - atcgcccAatacga - 3' (3)

3' - tagcgggItatgct - 5' (4)

Methylation experiments

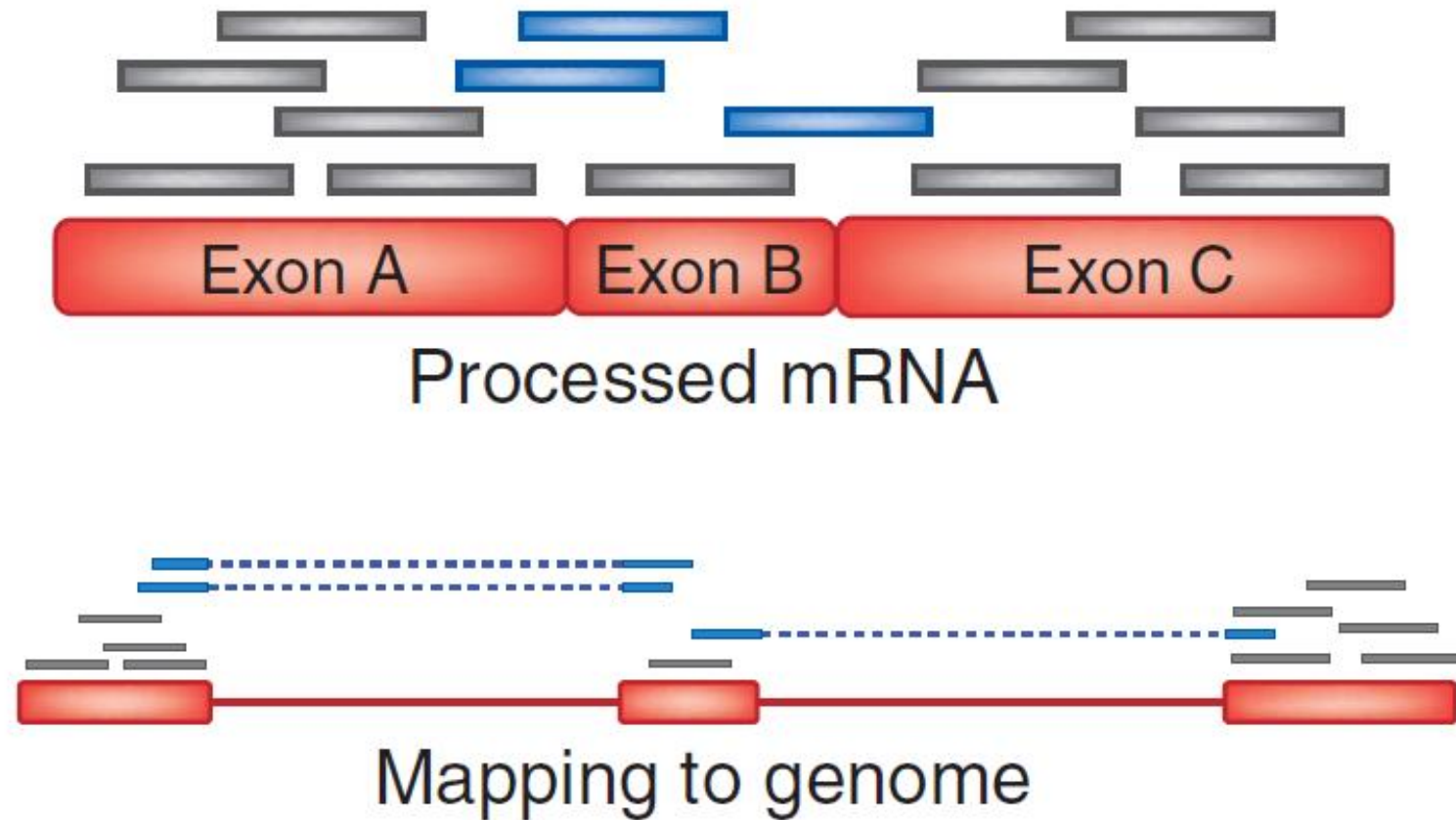
- Directly aligning reads against a reference will fail due to excessive mismatches in non-methylated regions
- Most aligners deal with this by creating 2 reference sequences
 - One has all Cs converted to Ts
 - One has all Gs converted to As
- Convert Cs to Ts in all reads aligned against C-T reference
- Convert Gs to As in all reads aligned against G-A reference
- If there are no mutations or sequencing errors the reads will always map to one of the two references

Multiple mapping reads



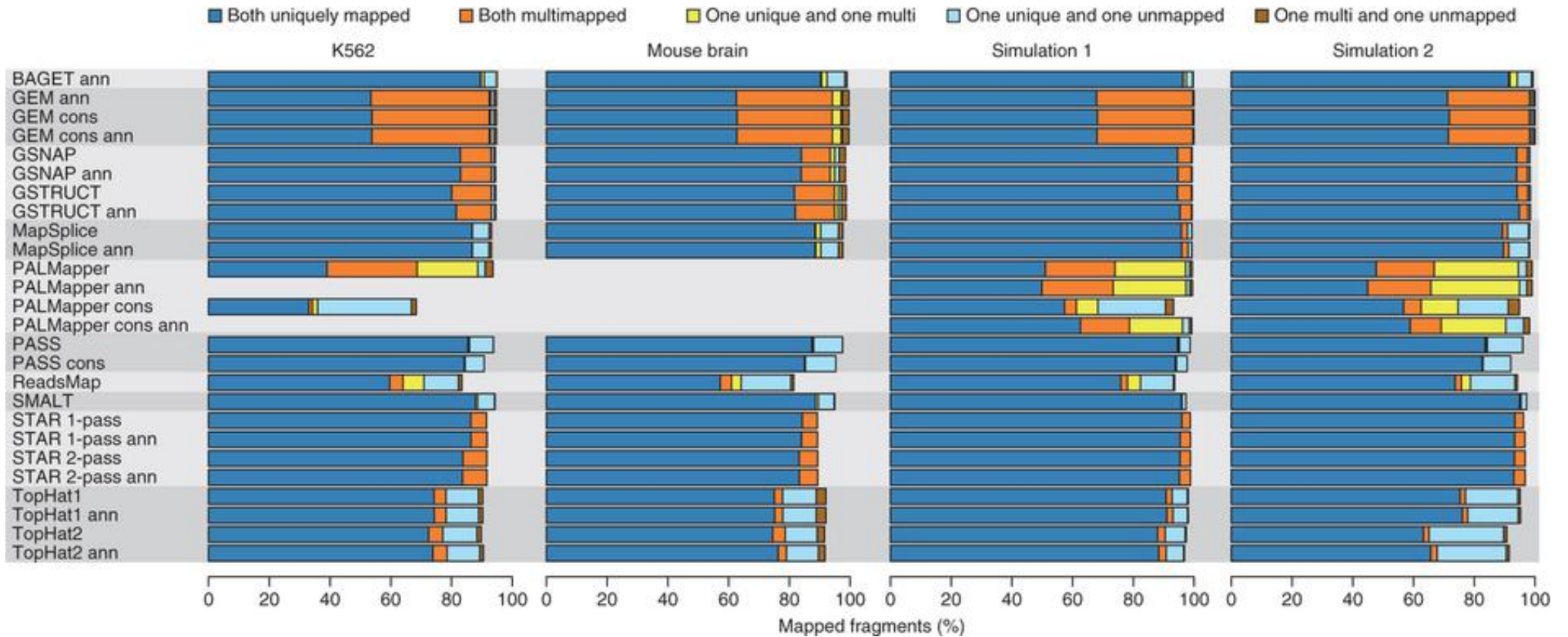
- A single read may occur more than once in the reference genome.
- Could be due to:
 - Paralogs (duplicated genes).
 - Transcripts which share exons.
 - Mutations in genotype relative to the reference.
 - Transposons and other common repetitive sequences
- Some aligners automatically assign a multi-mapping read to one of the locations at random (e.g. Tophat)
- Aligners may allow you to choose how these are dealt with – others may not

Spliced-read mapping



- Need packages which can account for splice variants
- Examples: TopHat, STAR, GSNAP, MapSplice

Spliced-read mapper evaluation



Local realignment to improve SNP/Indel detection

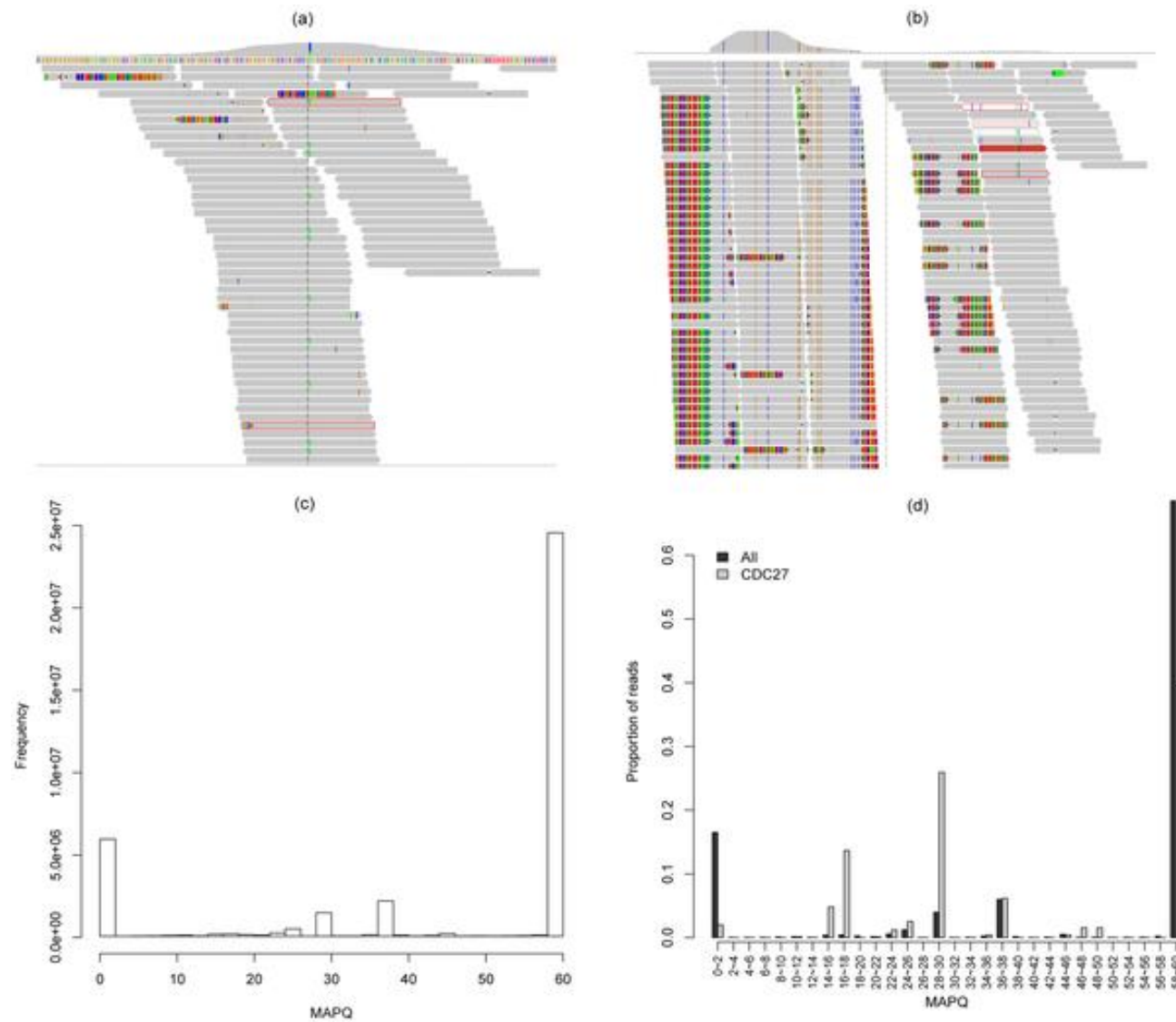
- Read aligners map each read (or read pair) independently of all other reads
- Around indels and other variants it can be helpful to make use of other metrics
 - e.g. Global median coverage for multi-mapping reads
- Tools such as GATK, SAMtools, Pindel and Breakdancer realign reads in the vicinity of variants to improve calls

http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit

Chen, K. BreakDancer: an algorithm for high-resolution mapping of genomic structural variation *Nature Methods* 6, 677 - 681 (2009)

Li H.*, Handsaker B.*, Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R. and 1000 Genome Project Data Processing Subgroup (2009) The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics*, 25, 2078-9

Figure 6. A visual examination of a spurious gene (CDC27).



<http://www.plosone.org/article/info:doi/10.1371/journal.pone.0038470>

All platforms have errors and artefacts



Illumina



PacBio



Roche 454



Ion Torrent

1. Removal of low quality bases
2. Removal of adaptor sequences
3. Platform specific artefacts (e.g homopolymers)

Table 2. Spurious genes having mutations detected in 30 samples.

CCDS ID	Gene symbol	Exon	# samples
CCDS11509.1	<i>CDC27</i>	13 th	36
CCDS12749.1	<i>CGB</i>	3 rd	36
CCDS12752.1	<i>CGB5</i>	1 st	36
CCDS41378.1	<i>NBPF11</i>	19 th	36
CCDS43407.1	<i>FAM153C</i>	4 th	36
CCDS5931.1	<i>MLL3</i>	42 nd	36
CCDS34703.1	<i>STAG3</i>	33 rd	34
CCDS5590.1	<i>POMZP3</i>	1 st	34
CCDS10638.1	<i>EIF3C</i>	8 th	32
CCDS30836.1	<i>NBPF14</i>	22 nd	31

CCDS: Consensus coding sequence. Exon: the specific exon in which the variants are detected.

doi:10.1371/journal.pone.0038470.t002

Illumina artefacts

Sequence-specific error profile of Illumina sequencers

Kensuke Nakamura^{1,*}, Taku Oshima², Takuya Morimoto^{2,3}, Shun Ikeda¹, Hirofumi Yoshikawa^{4,5}, Yuh Shiwa⁵, Shu Ishikawa², Margaret C. Linak⁶, Aki Hirai¹, Hiroki Takahashi¹, Md. Altaf-Ul-Amin¹, Naotake Ogasawara² and Shigehiko Kanaya¹

¹Graduate School of Information Science, ²Graduate School of Biological Sciences, Nara Institute of Science and Technology, 8916-5 Takayama-cho, Ikoma, Nara 630-0192, Japan, ³Biological Science Laboratories, Kao Corporation, 2606 Akabane, Ichikai, Haga, Tochigi 321-3497, ⁴Department of Bioscience, Tokyo University of Agriculture, ⁵Genome Research Center, NODAI Research Institute, Tokyo University of Agriculture, 1-1-1 Sakuragaoka Setagaya-ku, Tokyo, 156-8502, Japan and ⁶Department of Chemical Engineering and Material Science, University of Minnesota, 223 Amundson Hall, 421 Washington Avenue S.E., Minneapolis, MN 55455, USA

Received February 3, 2011; Revised April 25, 2011; Accepted April 26, 2011

ABSTRACT

We identified the sequence-specific starting positions of consecutive miscalls in the mapping of reads obtained from the Illumina Genome Analyser (GA). Detailed analysis of the miscall pattern indicated that the underlying mechanism involves sequence-specific interference of the base elongation process during sequencing. The two major sequence patterns that trigger this sequence-specific error (SSE) are: (i) inverted repeats and (ii) GGC sequences. We speculate that these sequences favor dephasing by inhibiting single-base

platforms [Illumina/Solexa Genome Analyser (4), Life Technologies/ABI SOLiD System (5) and Roche/454 Genome Sequencer FLX (6)], the Illumina Genome Analyser (GA) is, at the moment, the most popular choice for the analysis of genomic information (7). The Illumina/Solexa sequencers are characterized by: (i) solid-phase amplification and (ii) a cyclic reversible termination (CRT) process, also termed sequencing-by-synthesis (SBS) technology (8). The sequencer can generate hundreds of millions of relatively short (30–100 bp) read sequences per run.

The application of data obtained from this NGS technology can be roughly categorized into the following three

Nakamura, K. et al. Sequence-specific error profile of Illumina sequencers
Nucl. Acids Res. (2011) May 16, 2011

Illumina artefacts

1. GC rich regions are under represented
 - a. PCR
 - b. Sequencing
2. Substitutions more common than insertions
3. GGC/GCC motif is associated with low quality and mismatches
4. Filtering low quality reads exacerbates low coverage of GC regions

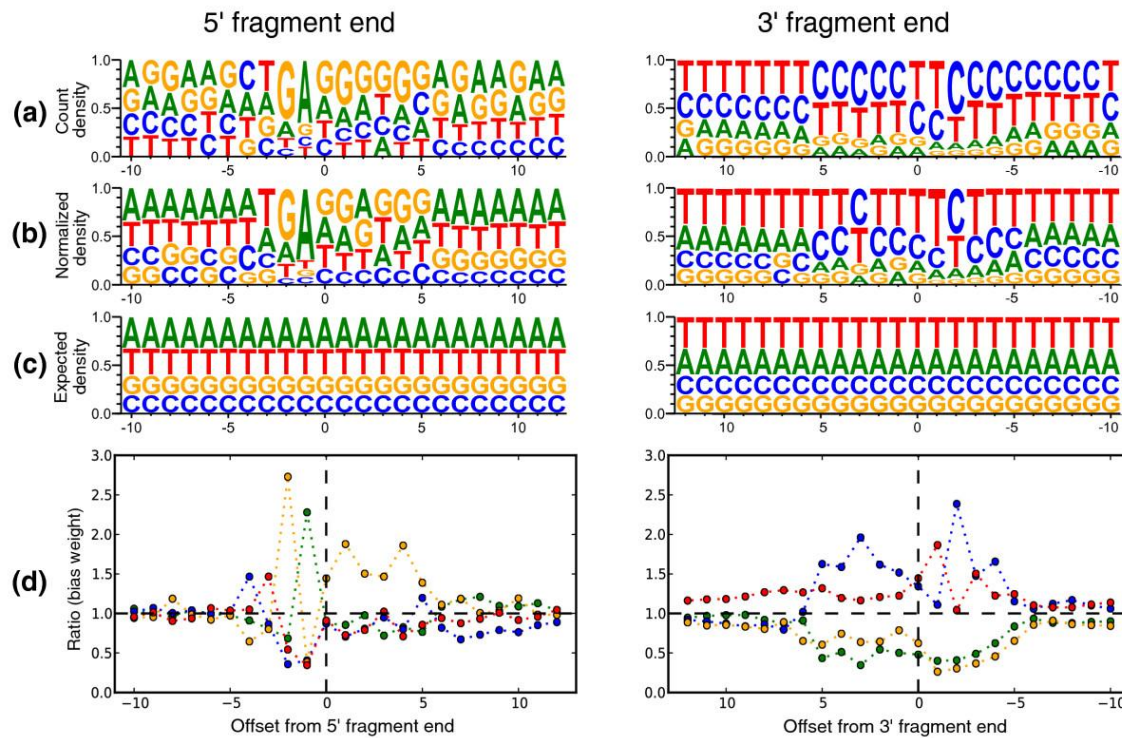
Alignment software should ideally account for technology specific bias but generally does not

Its up to you to filter before alignment

Your alignments are only as good as your library prep

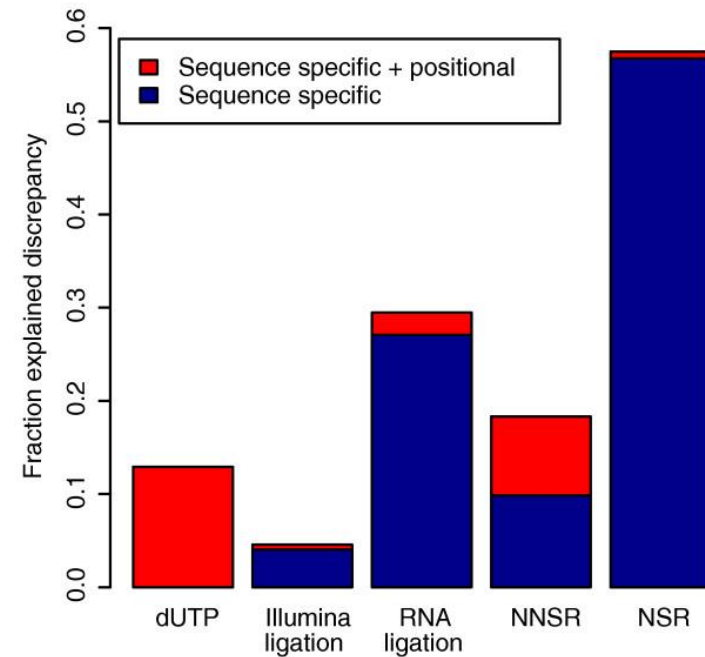
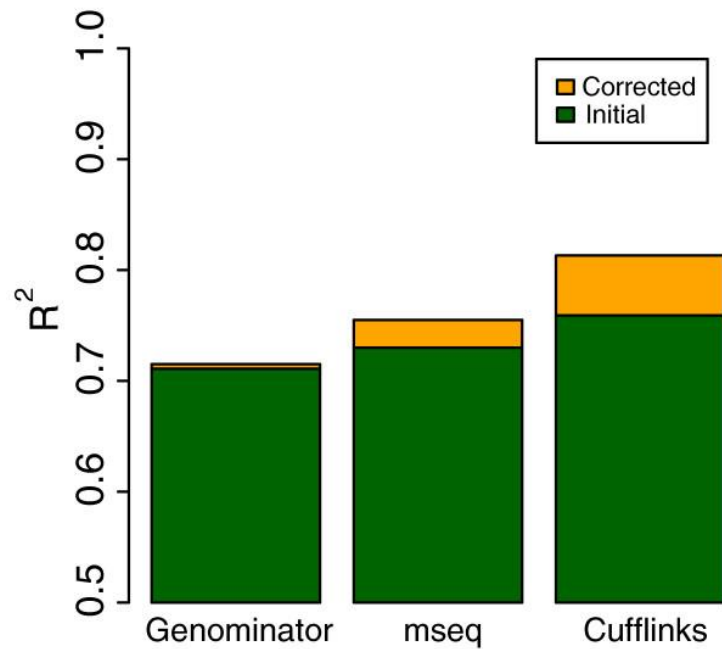
- Even if all other artefacts are removed:
- If your library prep is biased, your alignments will also reflect this bias

Tophat/Cufflinks aside



- Applies to random primed RNA-seq libraries
- Main potential biases:
 - Random hexamer priming biases
 - Fragments near 5' or 3' are likely to

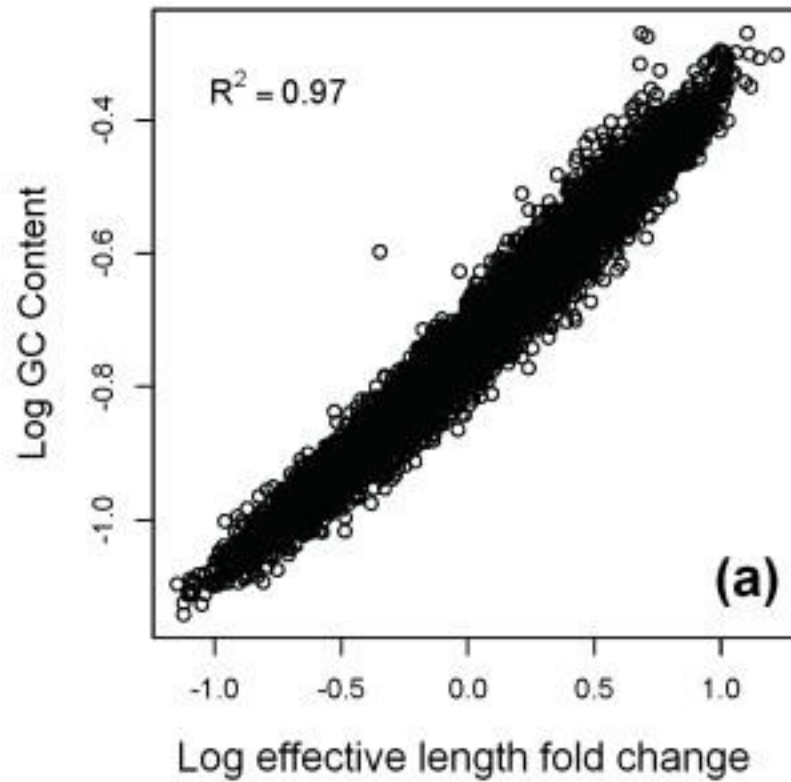
Effect of bias correction



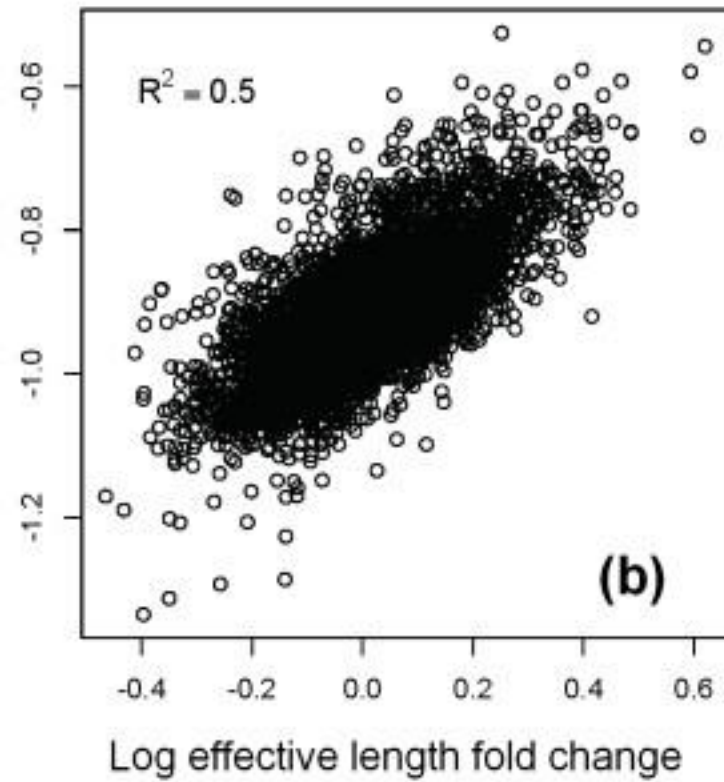
<http://genomebiology.com/2011/12/3/R22>

N.B. Out-dated version of Cufflinks used here

Correcting for GC-bias in RNA-seq

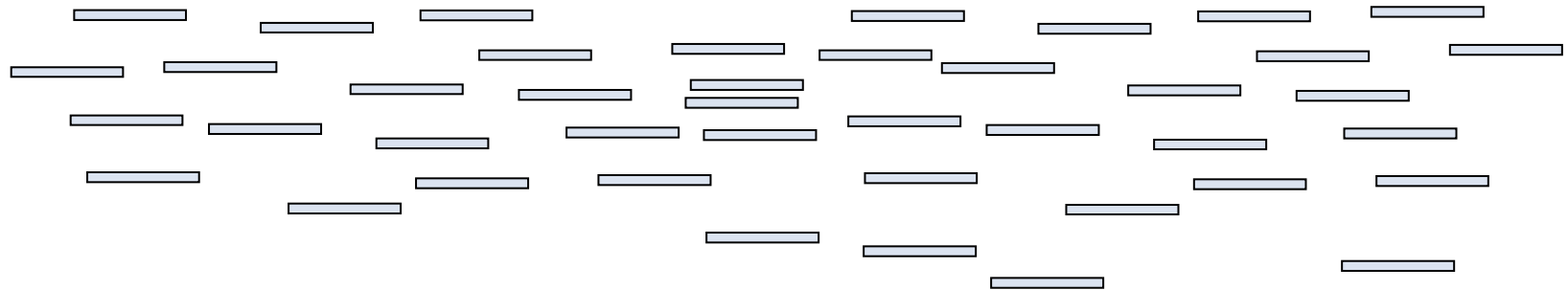
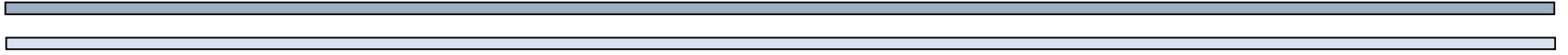


Human



Yeast

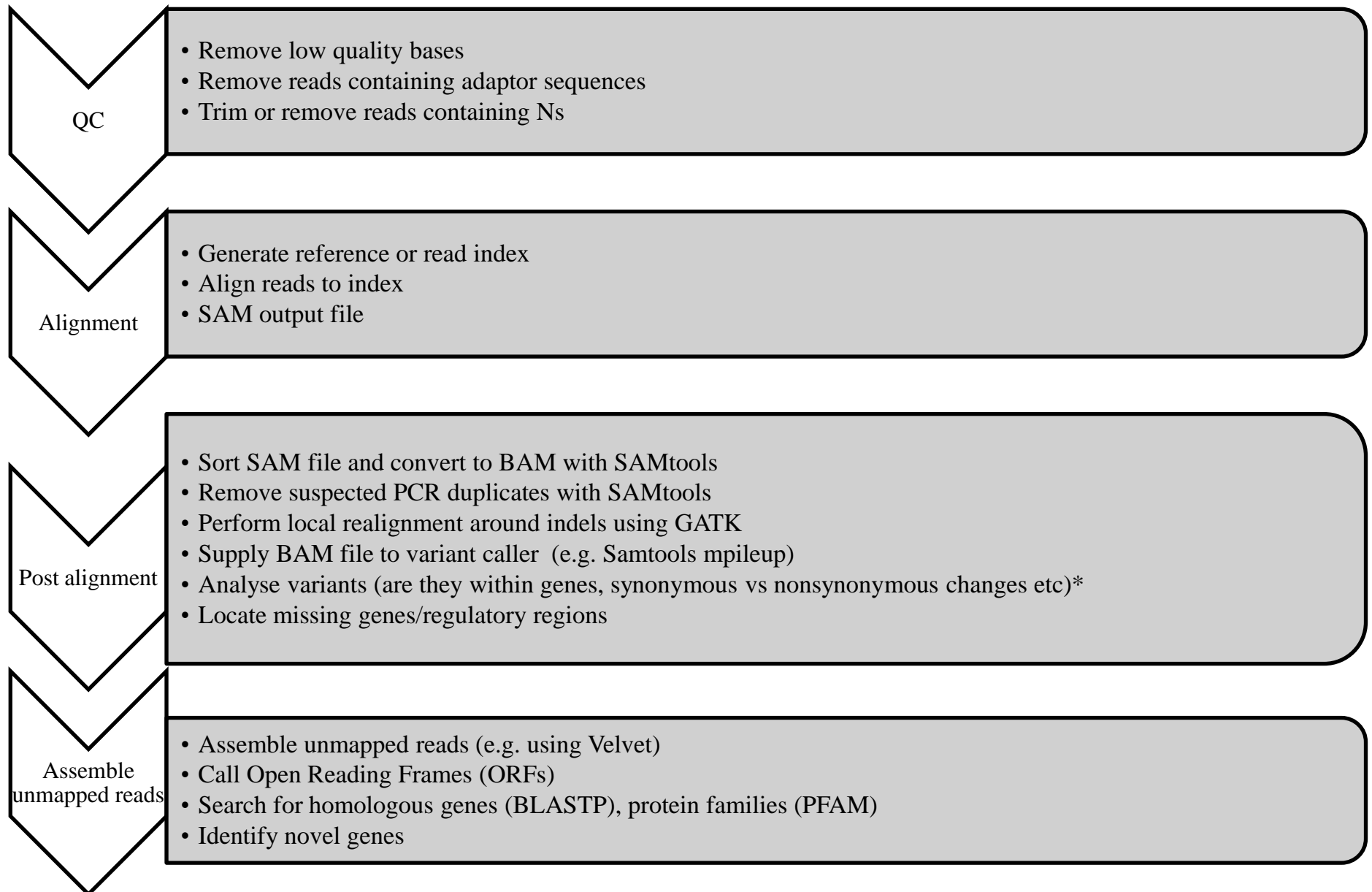
Unmapped reads



Unmapped reads

- Can be the result of:
 - Sequencing errors (should be small fraction if quality filtering applied before mapping)
 - Contamination
 - Excessive matches to repeats
 - Highly divergent regions between samples
 - Novel genetic material not present in reference
 - Plasmids
- Should be assembled de-novo with paired-end information if possible
- Resulting contigs run through MegaBlast against NCBI NT to check species
- Check against RepBase to remove repetitive contigs
- Call ORFs
- Blast ORFs using BlastP against NCBI NR or Swissprot and Blast2GO
- Run through PFAM

Typical alignment pipeline



* <http://bioinformatics.net.au/software.nesoni.shtml>

Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Indel detection
 - Suffix/Prefix Tries
 - Other alignment considerations
 - Typical alignment pipeline
 - **New methods of SNP calling**

New methods of SNP calling

- FreeBayes (<http://arxiv.org/pdf/1207.3907v2.pdf>)
- Warning - unpublished
 - Haplotype calling in polyploids

ACA Reference Genome

Assume a SNP at both 5' A->T and 3' A->G

Do we have a heterozygous?

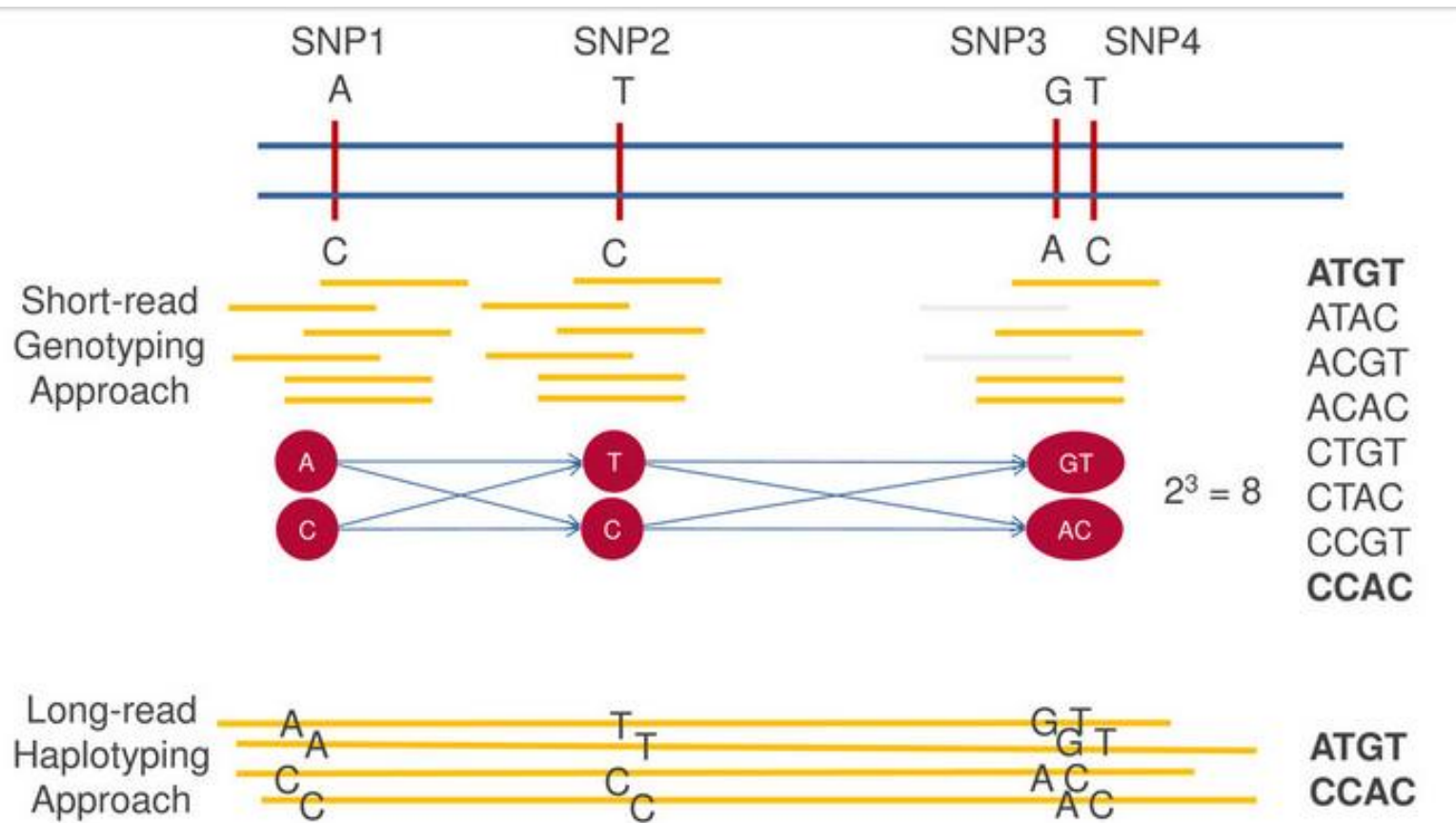
ACG

TCA

Or do we have a homozygous?

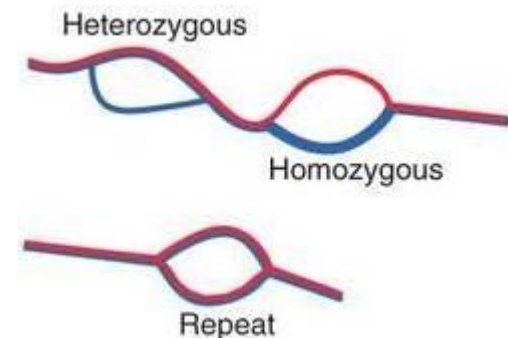
TCG

Haplotype issue calling – Long reads to the rescue



New methods of SNP calling

- Why align at all?
 - We only do this because of computational constraints
 - Ideally we want to assemble denovo and then align to reference genome
- Cortex is a step in this direction:
 - Denovo genome assembler, but keeps track of differences which could be due to SNPs/Indels



Variant calling with de-novo assembly

Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly

Heng Li^{1,*}

¹Broad Institute, 7 Cambridge Center, Cambridge, MA 02142, USA

Associate Editor: Dr. Michael Brudno

nature
genetics

ABSTRACT

Motivation: Eugene Myers in his stri suggested that in a string graph or e path spells a valid assembly. As a st every valid assembly of reads, such be constructed correctly, is in fact reads. In principle, every analysis bas sequencing (WGS) data, such as SNP calling, can also be achieved with uniti

De novo assembly and genotyping of variants using colored de Bruijn graphs

Zamin Iqbal^{1,2,5}, Mario Caccamo^{3,5}, Isaac Turner¹, Paul Flicek² & Gil McVean^{1,4}

Detecting genetic variants that are highly divergent from a reference sequence remains a major challenge in genome sequencing. We introduce *de novo* assembly algorithms using colored de Bruijn graphs for detecting and genotyping simple and complex genetic variants in an individual or population. We provide an efficient software implementation, Cortex, the first *de novo* assembler capable of assembling multiple eukaryotic genomes simultaneously. Four applications of Cortex are presented. First, we detect and validate both simple

a single suitable reference, as in ecological sequencing²¹. Fourth, methods for variant calling from mapped reads typically focus on a single variant type. However, in cases in which variants of different types cluster, focus on a single type can lead to errors, for example, through incorrect alignment around indel polymorphisms^{6,7}. Fifth, although there are methods for detecting large structural variants, such as using array comparative genomic hybridization (aCGH)²²⁻²⁵ and mapped reads^{11,12,14,26}, these cannot determine the exact location, size or allelic sequence of variants. Finally, mapping

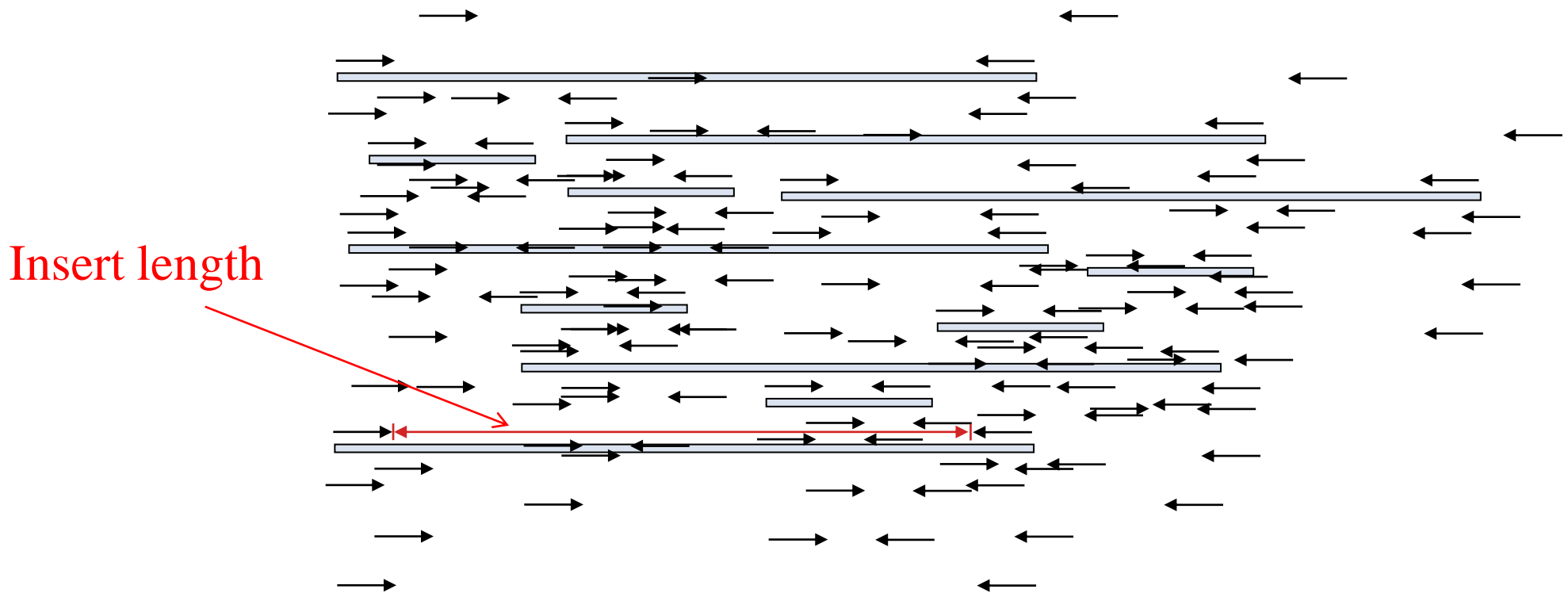
Questions!

biosciences.exeter.ac.uk/facilities/sequencing/usefulresources/

Assembly algorithms for short reads

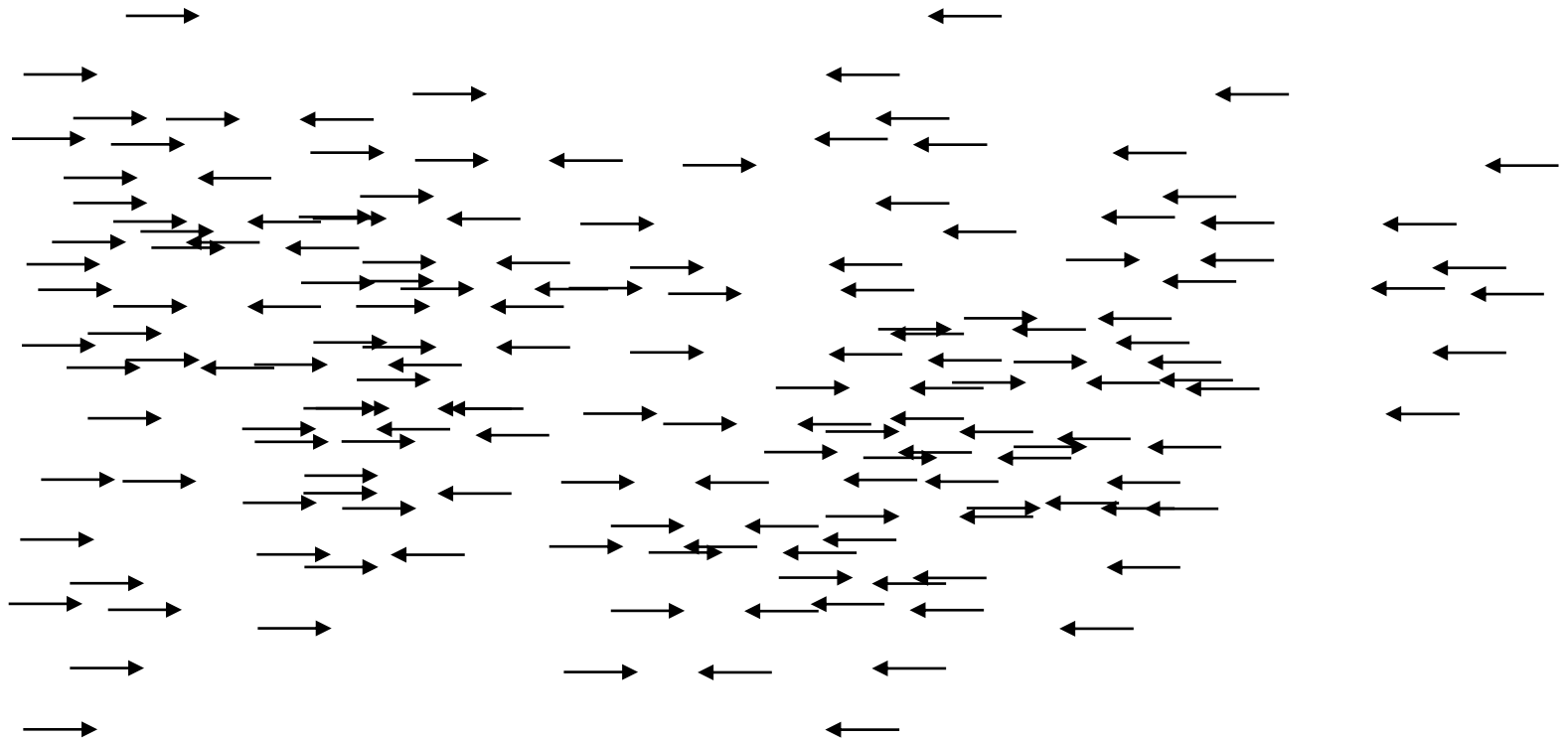
De-novo sequence assembly

1. Sequence DNA fragments from each end



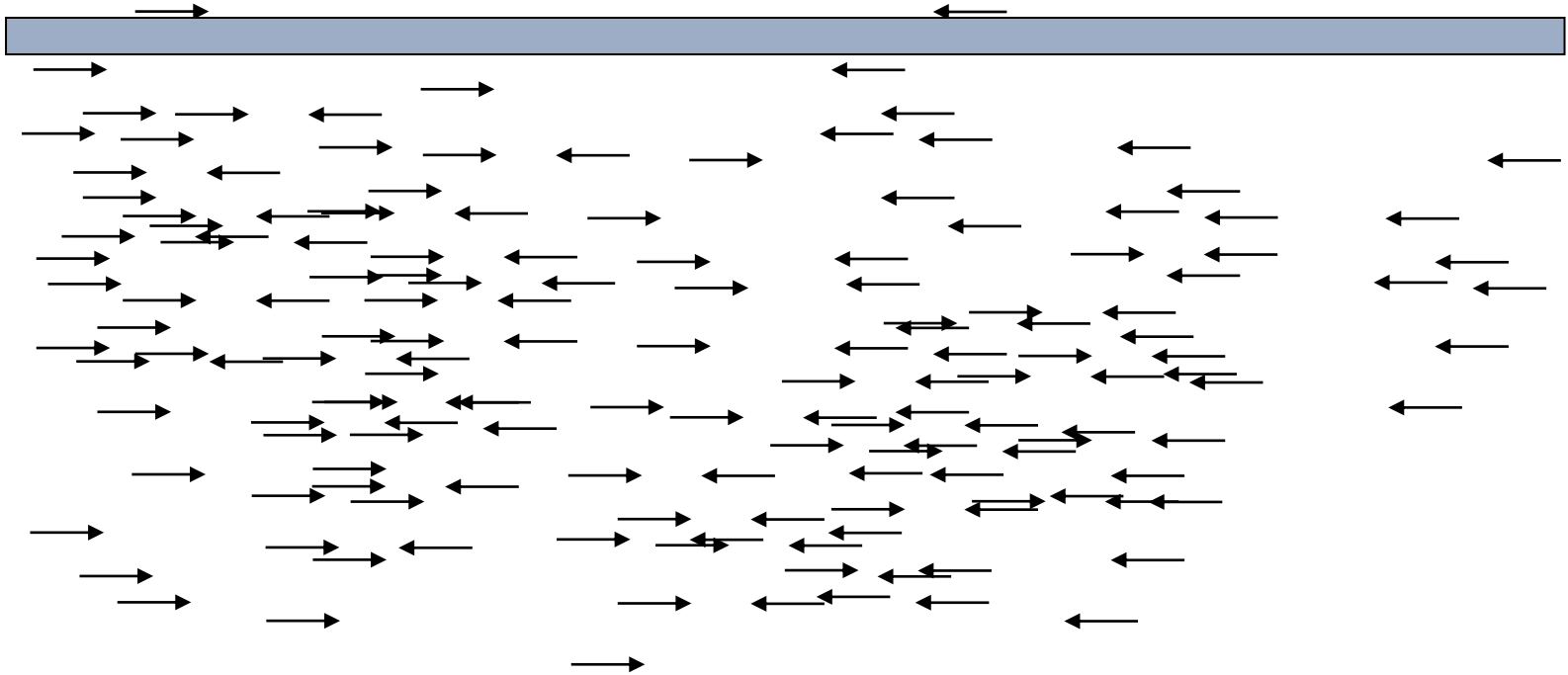
De-novo sequence assembly

1. Sequence DNA fragments from each end
2. Reads aligned to generate contigs



De-novo sequence assembly

1. Sequence DNA fragments from each end
2. Reads aligned to generate contigs



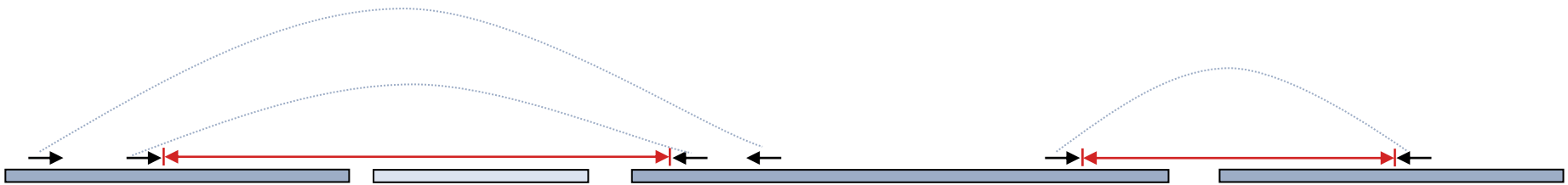
De-novo sequence assembly

1. Sequence DNA fragments from each end
2. Reads aligned to generate contigs
3. Supercontigs derived from paired reads on different contigs



De-novo sequence assembly

1. Sequence DNA fragments from each end
2. Reads aligned to generate contigs
3. Supercontigs derived from paired reads on different contigs



4. Ordering of contigs is determined
5. Different insert lengths and read lengths can resolve ambiguities
6. Insert size can be increased to 2-20kb by using mate-pair libraries (helps to span repetitive regions)

Mate-pair vs paired-end

- Often causes confusion
- Paired-end usually refers to libraries prepared for the Illumina platform with insert sizes 50-500bp.
- Mate-pair is a different library preparation protocol and usually produces insert sizes 2kb-20kb.

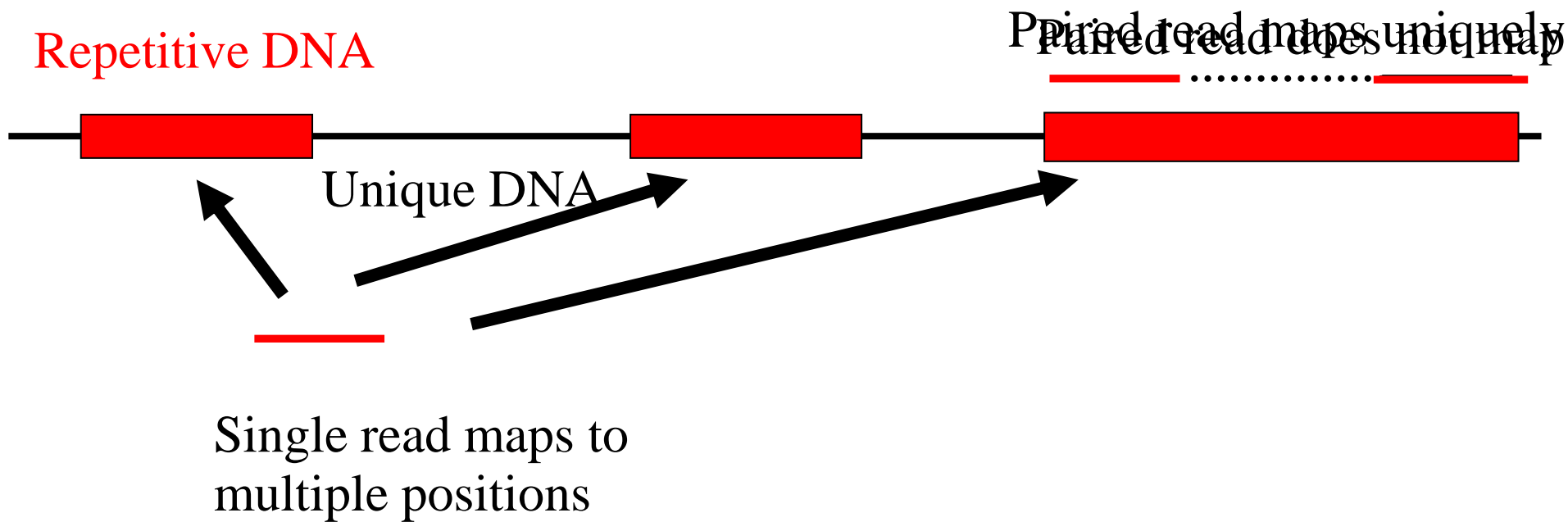
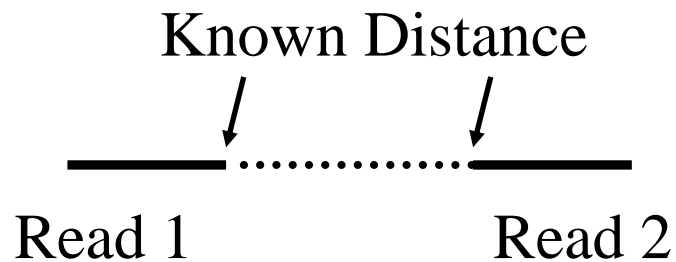
Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Suffix/Prefix Tries
 - Other alignment considerations
 - Typical alignment pipeline
- **Assembly algorithms for short reads**
 - **Effect of repeats**
 - Overlap-Consensus
 - de Bruijn graphs
 - Assembly evaluation metrics
 - Typical assembly pipeline

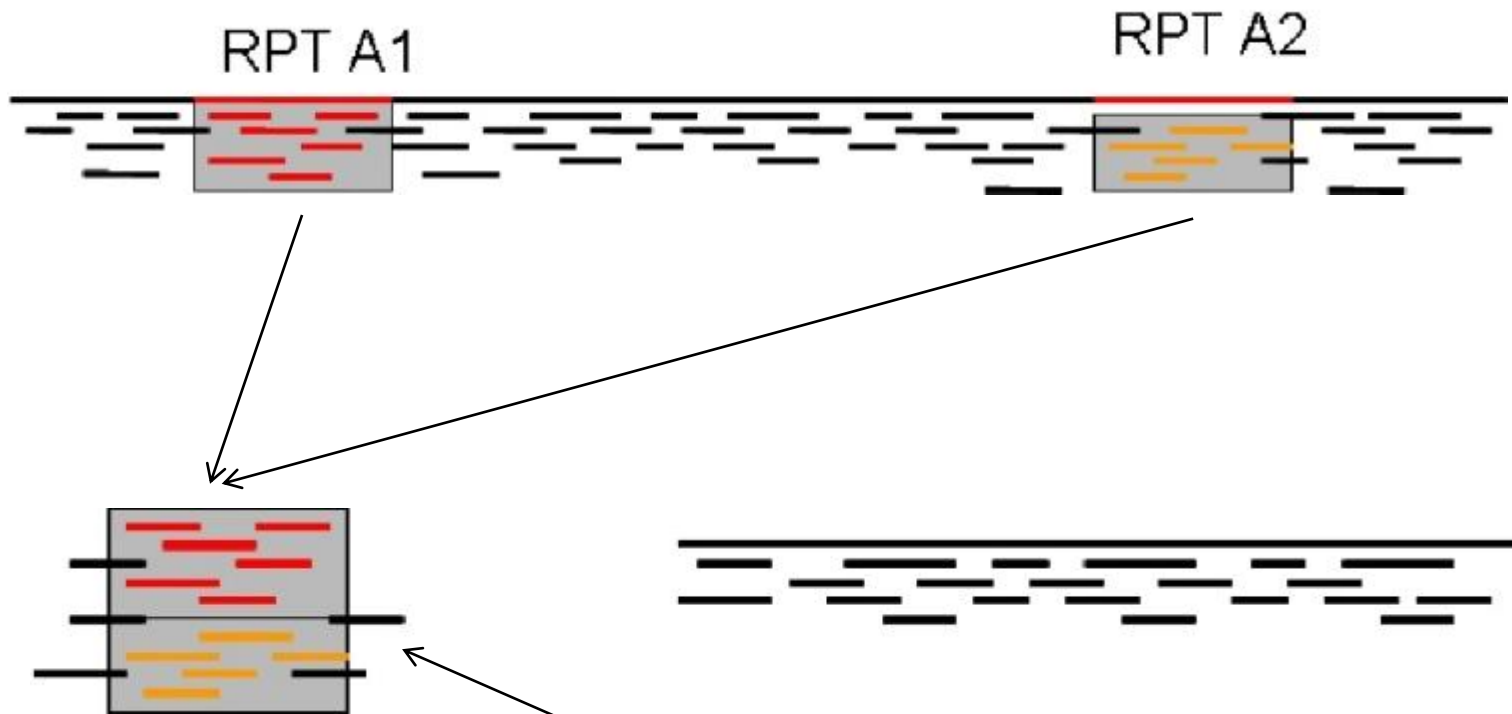
Repetitive sequence

- Main reason for fragmented genome assemblies
- Additional sequencing depth will not help overcome repeat limited assemblies

Repetitive sequence



Repetitive sequence



Can try to identify collapsed repeats by increased relative coverage

Repetitive sequence

- Main reason for fragmented genome assemblies
- Additional sequencing depth will not help overcome repeat limited assemblies
- Can estimate the number of repetitive regions, based on relative coverage
- Only longer reads or paired-end/mate-pair reads can overcome this
- PacBio reads can extend up to 10-20kb but expensive and impractical for most labs
- Large mate pair insert sizes ~20kb are possible, but library preparation is inefficient (2-3 days of trial and error). Also a significant fraction will be error-prone and/or chimeric

Assumptions made by de-novo assemblers

Based on Lander-Waterman model

Number of times a base is sequenced follows a Poisson distribution

Reads are randomly distributed throughout a genome

The ability to detect an overlap between two reads is not dependent on the base-composition of the read

$$P = 1 - \left[1 - \frac{L}{G} \right]^N .$$

L = Read length

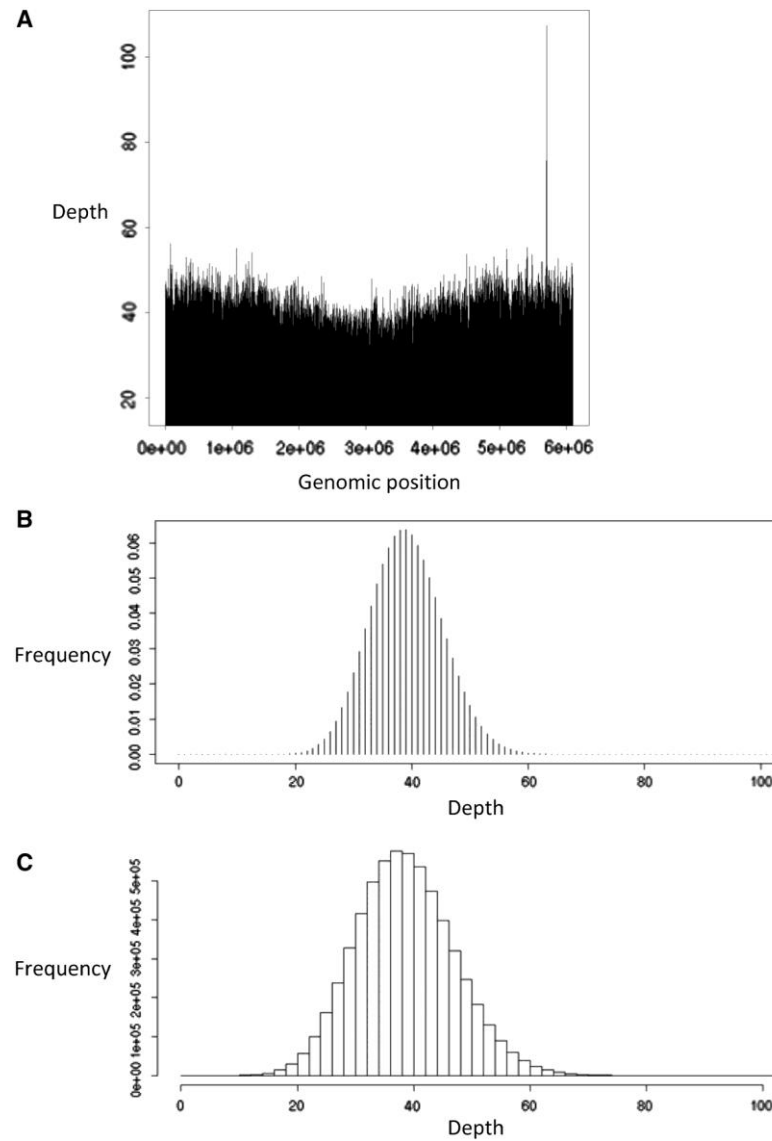
N = Number of reads

G = Genome size

P = Probability base is sequenced

Lander, E.S. and Waterman, M.S. (1988). "Genomic Mapping by Fingerprinting Random Clones: A Mathematical Analysis". *Genomics* **2** (3): 231–239

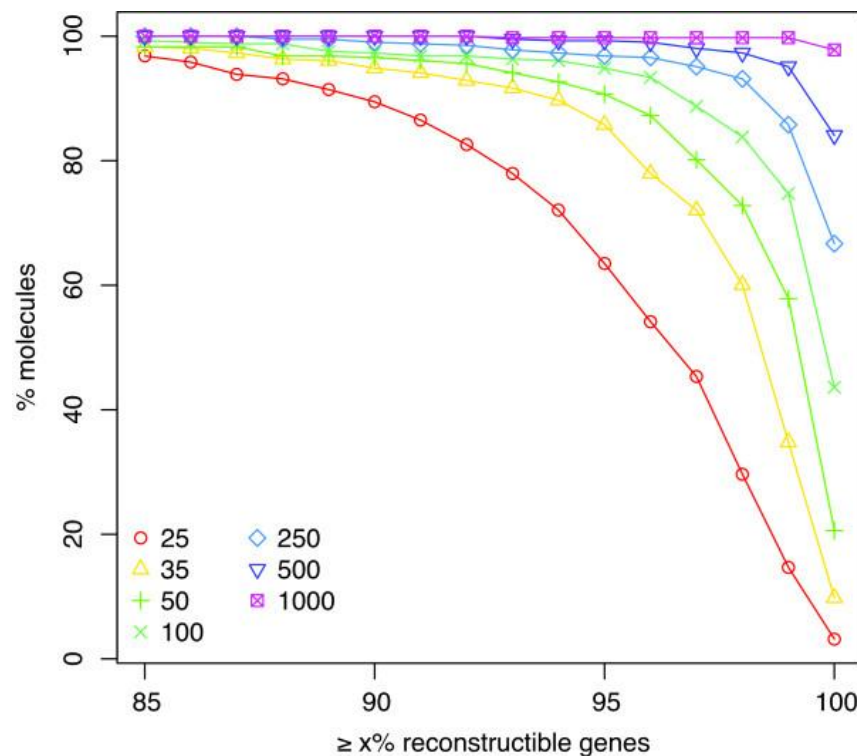
Assumptions are not true



Paszkiwicz K , Studholme D J Brief Bioinform
2010;11:457-472

NGS de-novo assemblies are draft quality at best

- 500 contigs covering most of a bacterial genome can be obtained in 1 week from genomic DNA to Genbank submission
- To get 1 contigs covering **all** genomic sequence could take many months
- Is the extra effort worth it?
- Short answer: Usually not.



Assembly complexity of
prokaryotic genomes using
short reads
Carl Kingsford , Michael C
Schatz and Mihai Pop
BMC Bioinformatics 2010, **11**:21

Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Suffix/Prefix Tries
 - Other alignment considerations
 - Typical alignment pipeline
- **Assembly algorithms for short reads**
 - Effect of repeats
 - **Overlap-Consensus**
 - **de Bruijn graphs**
 - Assembly evaluation metrics
 - Typical assembly pipeline

Overlap consensus vs. de Bruijn

- **2 main categories of assembly algorithms**
 - Overlap Consensus (OLC) and de Bruijn graph assemblers
- OLC
 - Primarily used for Sanger and hybrid assemblies
 - Memory constraints prevent its use beyond 1 million reads or so
- de Bruijn
 - Primarily used for NGS assemblies
 - Still memory hungry but possible

Original sequence
GTAGTATAGTCAGTATCA

Sequence reads

GTAGTA TAGTAT AGTATA
GTATAG TATAGT
ATAGTC TAGTCA AGTCAG
GTCAGT TCAGTA
CAGTAT AGTATC GTATCA

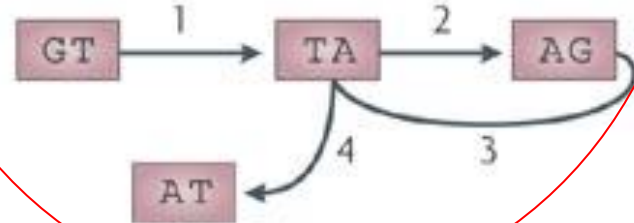
Consensus overlap assembly

GTAGTA
TAGTAT
AGTATA
GTATAG
TATAGT
ATAGTC
TAGTCA
AGTCAG
GTCAGT
TCAGTA
CAGTAT
AGTATC
GTATCA
GTAGTATAGTCAGTATCA

k-mers (2-mers)

GT TA AG AT TC CA

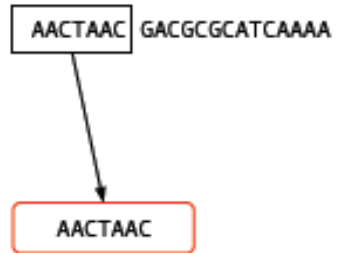
de Bruijn graph



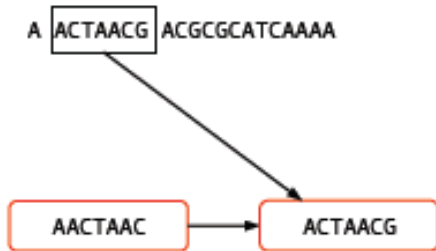
de Bruijn graph assembly

AACTAACGACGCGCATCAAAA

de Bruijn graph assembly

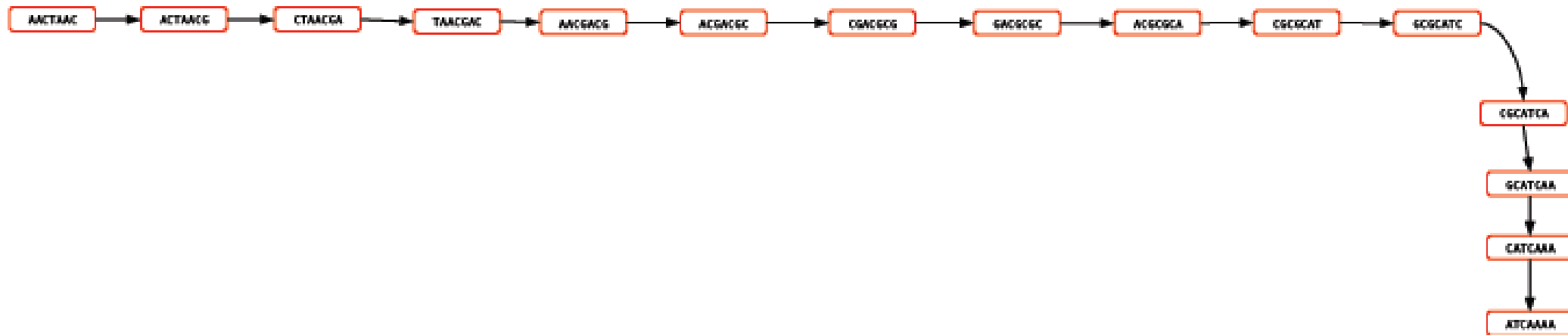


de Bruijn graph assembly



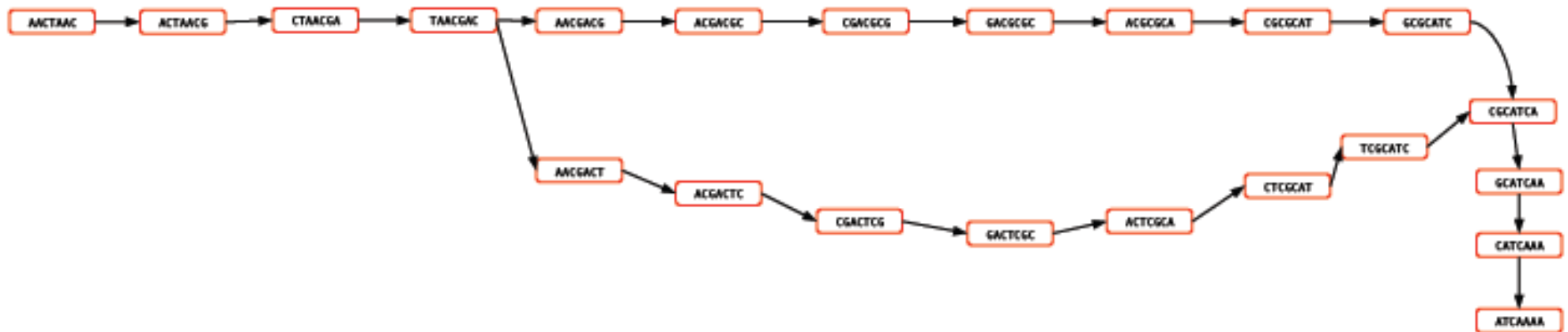
de Bruijn graph assembly

AACTAACGACGCGCATCAAAA



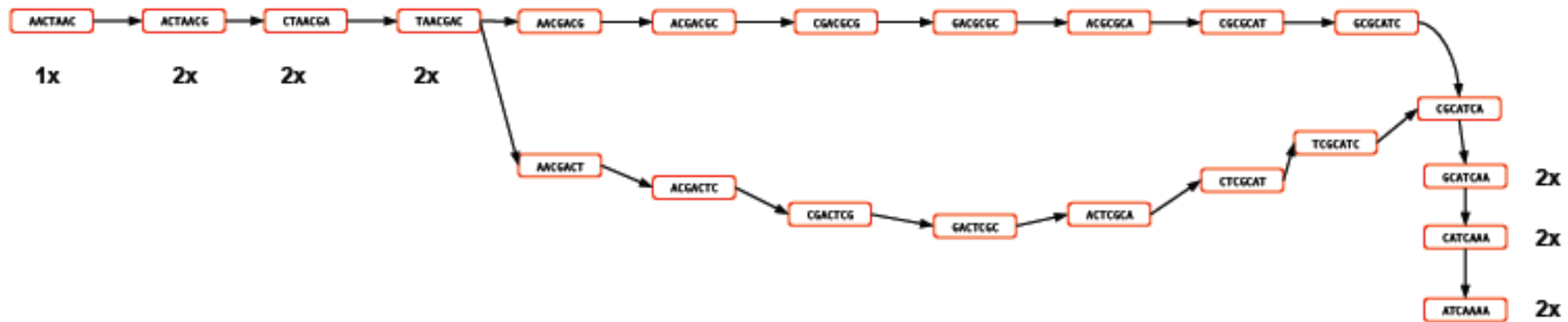
de Bruijn graph assembly

AACTAACGAC G CGCATCAAAA
ACTAACGAC T CGCATCAAAA



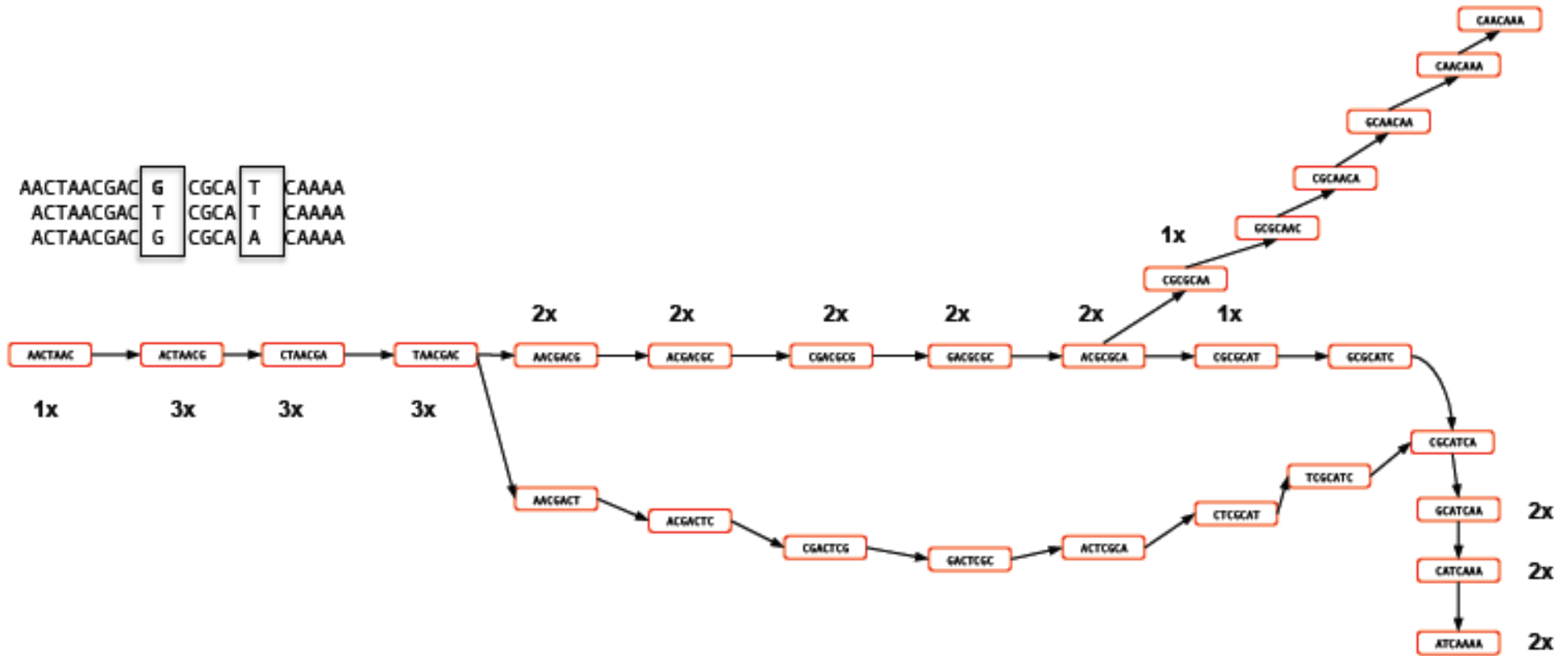
de Bruijn graph assembly

AACTAACGAC G CGCATCAAAA
ACTAACGAC T CGCATCAAAA

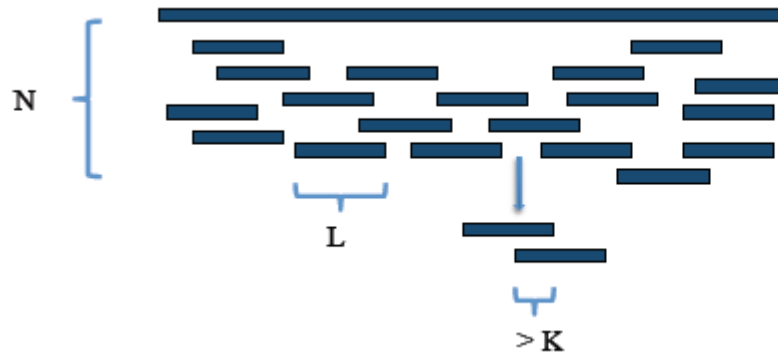
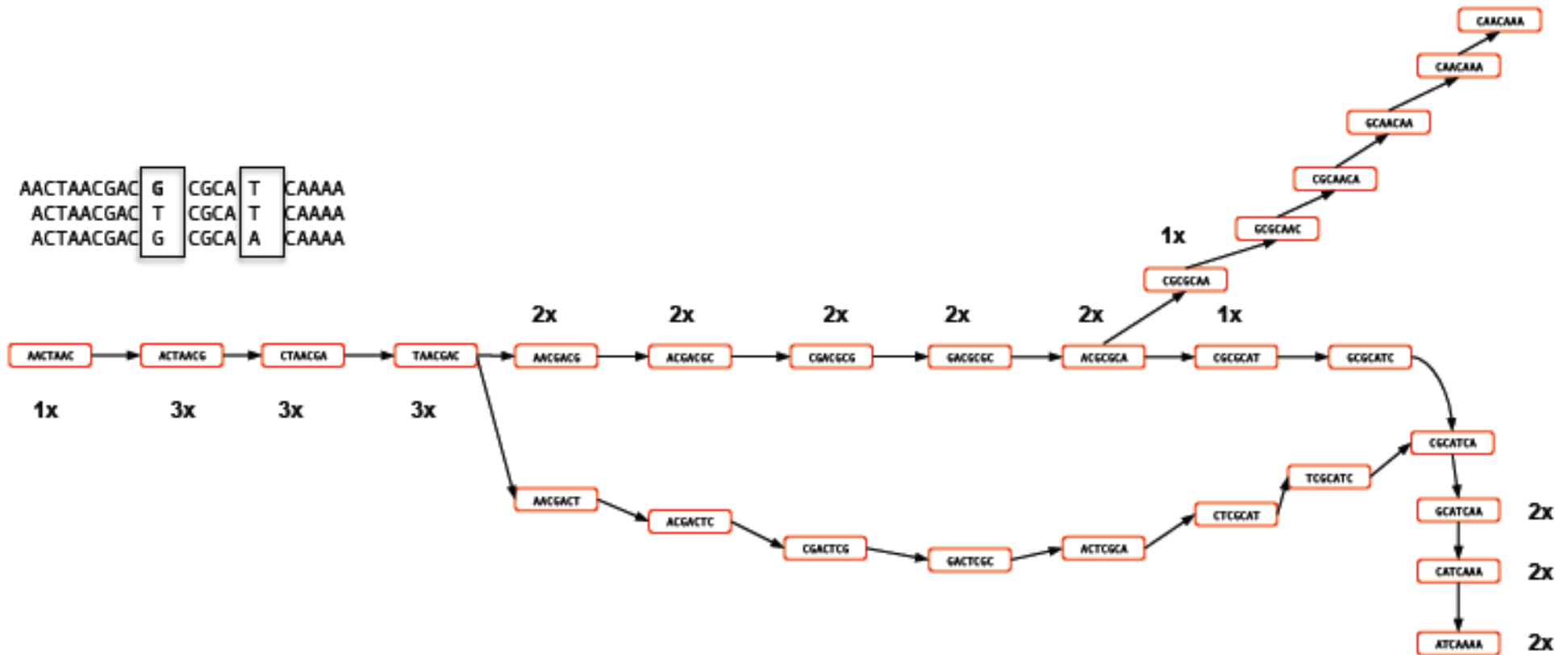


de Bruijn graph assembly

AACTAACGAC	G	CGCA	T	CAAAA
ACTAACGAC	T	CGCA	T	CAAAA
ACTAACGAC	G	CGCA	A	CAAAA



de Bruijn graph assembly



$$P(d > 0) = 1 - e^{-N(L-K)/G}$$

Diagrams courtesy M. Caccamo, TGAC

Dealing with errors

Illumina sequencing error rate 1-2% depending on read length

many of the 25-mers will contain errors

Error correction before assembly for small data sets is less important

- ▶ Can be removed during the graph assembly

Large datasets

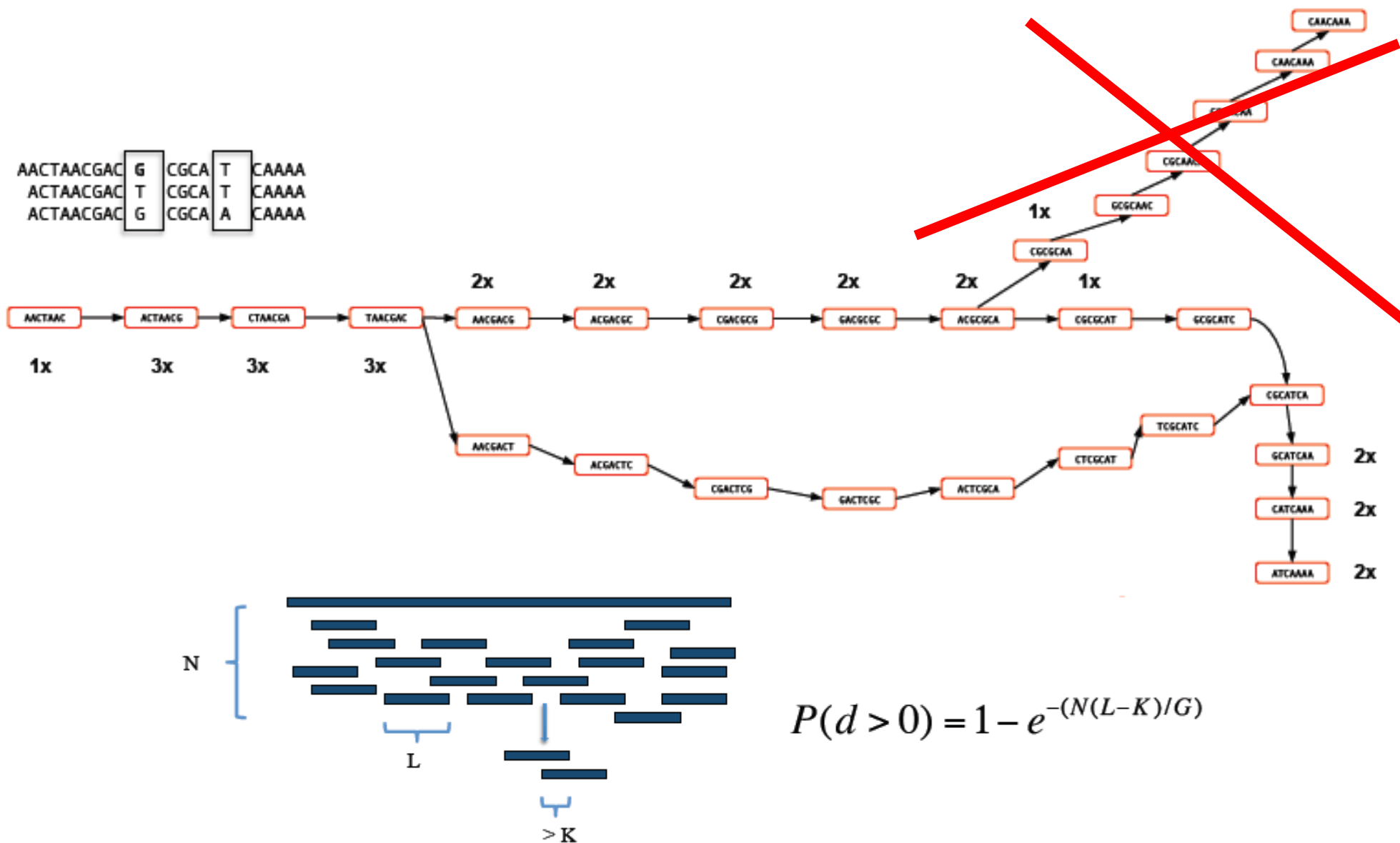
- ▶ Removal of singleton kmers is essential as will drastically reduce the memory footprint of the graph
- ▶ e.g. Asian human genome data, the total number of distinct 25-mers was reduced from 14.6 billion to 5.0 billion

Table 1. Summary of preassembly error correction in the Asian genome sequencing

	Total reads	Error-free reads (%)	25-mer no.
Original reads	4,083,271,441	60.1	14,551,534,812
After correction	3,312,495,883	74.0	4,966,416,149

Li et al (2009) Gen Res, 20

de Bruijn graph assembly error correction



Diagrams courtesy M. Caccamo, TGAC

Errors or rare sequence?

- Depends on the type of data:
 - Assumptions are probably true for single haploid genome data
 - Diploid and polyploid expect any branches to have equal coverage
 - Less clear for RNA-seq due to splicing
 - Completely false assumption for metagenomic and metatranscriptomic data!

Short read assemblers

- First de Bruijn based assembler was Newbler
 - Adapted to handle main 454 error – indels in homopolymers
- Several other de Bruijn assemblers developed subsequently
 - Velvet, Euler-SR, ABySS, ALLPATHS2
 - Most can use paired-end and mate-pair information
- Most cannot deal with mammalian sized genomes
 - ABySS – distributed genome assembly via MPI
 - SOAPde-novo (BGI) Cortex (TGAC)
 - Early removal of spurious errors
- Hybrid assemblers
 - MIRA – capable of assembling 454, Sanger and short reads
 - Memory hungry
- Other approaches
 - String graph assemblers
 - Fermi, SGA
 - Correcting PacBio reads with Illumina

Contents

- **Alignment algorithms for short-reads**
 - Background – Blast (why can't we use it?)
 - Adapting hashed seed-extend algorithms to work with shorter reads
 - Suffix/Prefix Tries
 - Other alignment considerations
 - Typical alignment pipeline
- **Assembly algorithms for short reads**
 - Effect of repeats
 - Overlap-Consensus
 - de Bruijn graphs
 - **Assembly evaluation metrics**
 - Typical assembly pipeline

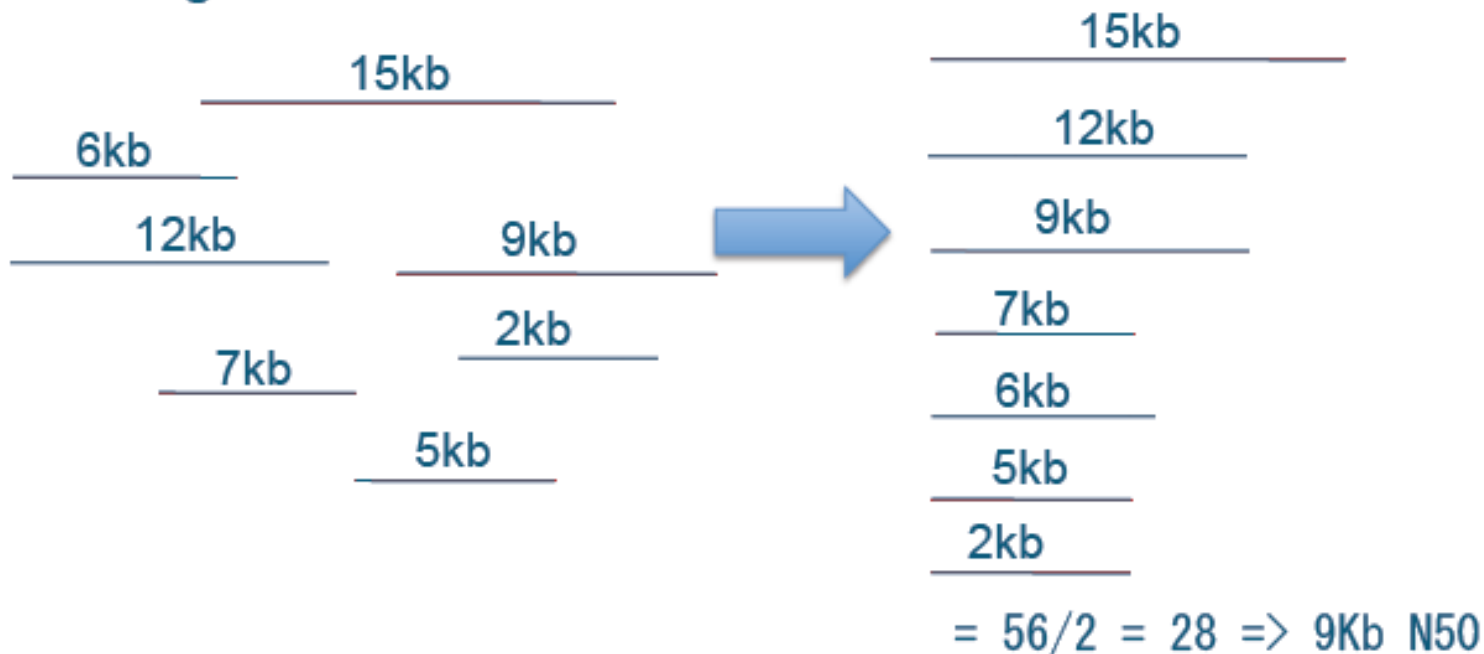
Assembly evaluation – N50

N50 has traditionally been used to compare assemblies

If you order the set of contigs produced by the assembler by size

- ▶ N50 is the size of the contig such that 50% of the total bases are in contigs of equal or greater size

E.g.



Assembly length vs. N50

Another informative measure is total length of the assembly

- ▶ Most genomes have an expected size prior to running assembly
- ▶ Assemblers assume diploid genome

Contig total length less than scaffold total length

- ▶ Scaffolds are contigs with runs of N's between the contigs

If you remove smaller contigs -> N50 increases :0)

- ▶ Total length decreases i.e. less of the genome sequence in the assembly :0(

Most assemblers will remove contigs less than 100bp or less than the read length

Assembly evaluation metrics

N50 just measures the continuity of the assembly

- ▶ Larger values are generally better

However it does not assess the quality of the assembled sequence

- ▶ E.g. if there are incorrect joins in the assembly the N50 could appear to be larger

Assembly quality measures

- ▶ Methods using contigs only:

- ▶ N50
- ▶ Total contig length
- ▶ Number of contigs

- ▶ Metrics using an alignment of reads onto the contigs

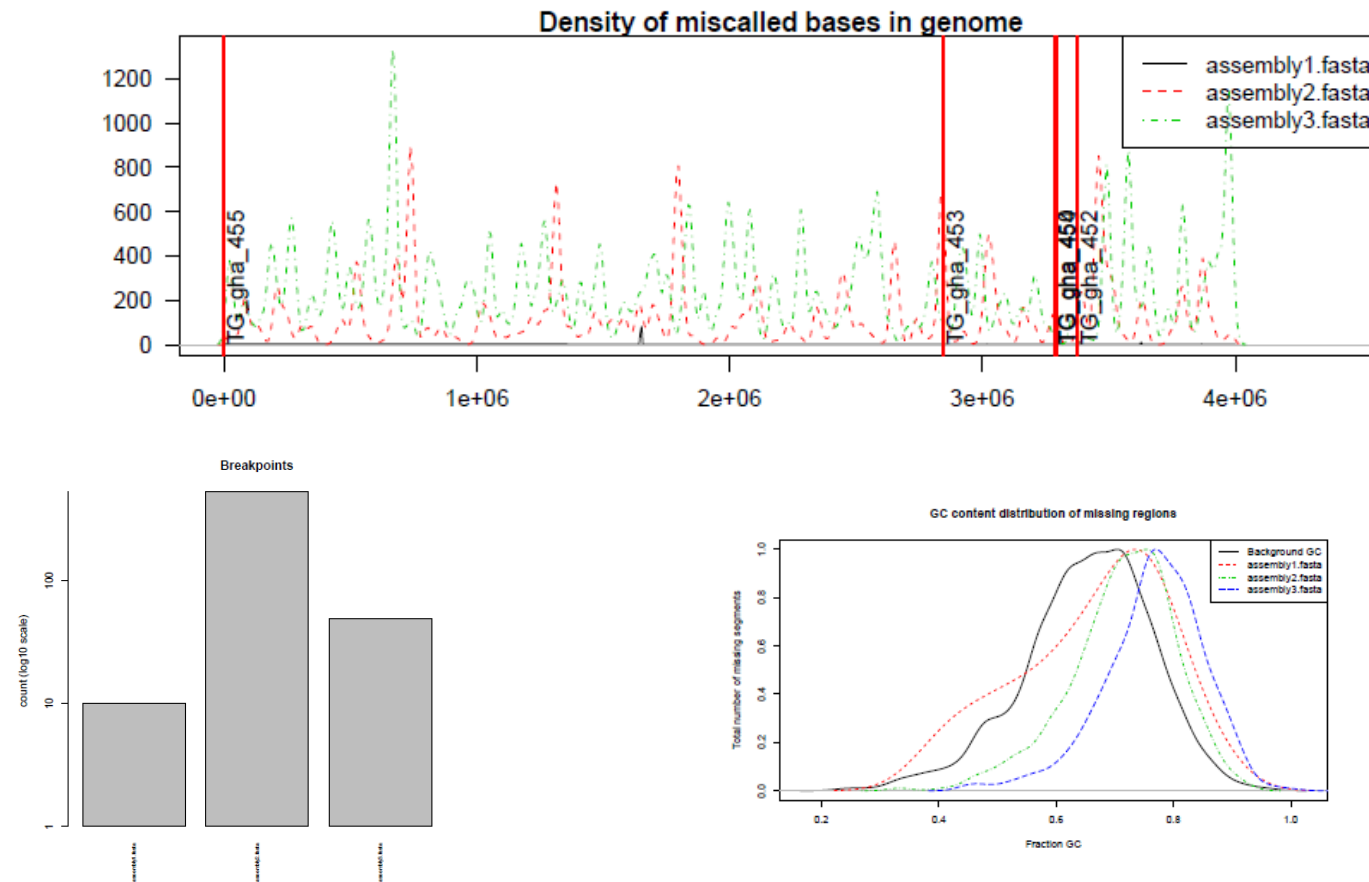
- ▶ Mapping Fraction (No. reads mapped/total reads) + pairing rate
- ▶ Count the SNPs and indels
- ▶ Misassemblies (regions not spanned by read pairs)



Which human assembly is better? Why?

	Assembly 1	Assembly 2		Assembly 1	Assembly 2
N50	51kb	42Kb		50Kb	20Kb
Total length	2.7Gb	2.69Gb		1.2Gb	2.7Gb
Avg. length	45Kb	39kb		40Kb	18Kb
Mapping rate	0.82	0.78		0.6	0.85
SNP rate	0.02	0.02		0.02	0.02
Indel rate	0.01	0.01		0.01	0.012
Pairing rate	0.8	0.9		0.9	0.88
Misassemblies	15	5		2	2

Assembly benchmarking software



Darling et al Mauve Assembly Metrics *Bioinformatics* (2011) btr451 first published online August 2, 2011
<http://t.co/BbpbTPz>

Types of assemblers

2 main categories, many variations

Each tends to have its own niche

Memory and hardware requirements can differ substantially

Typically a parameter scan is need to get the ‘best’ assembly

This means many assemblies need to be generated

Name	Read Type	Algorithm	Reference
SUTTA	long & short	B&B	(Narzisi and Mishra [25], 2010)
ARACHNE	long	OLC	(Batoglou et al. [14], 2002)
CABOG	long & short	OLC	(Miller et al. [13], 2008)
Celera	long	OLC	(Myers et al. [12], 2000)
Edena	short	OLC	(Hernandez et al. [16], 2008)
Minimus (AMOS)	long	OLC	(Sommer et al. [15], 2007)
Newbler	long	OLC	454/Roche
CAP3	long	Greedy	(Huang and Madan [7], 1999)
PCAP	long	Greedy	(Huang et al. [8], 2003)
Phrap	long	Greedy	(Green [6], 1996)
Phusion	long	Greedy	(Mullikin and Ning [9], 2003)
TIGR	long	Greedy	(Sutton et al. [5], 1995)
ABYSS	short	SBH	(Simpson et al. [19], 2009)
ALLPATHS	short	SBH	(Butler et al. [46,47], 2008/2011)
Euler	long	SBH	(Pevzner et al. [17], 2001)
Euler-SR	short	SBH	(Chalissou and Pevzner [35], 2008)
Ray	long & short	SBH	(Boisvert et al. [48], 2010)
SOAPdenovo	short	SBH	(Li et al. [20], 2010)
Velvet	long & short	SBH	(Zerbino and Birney [18,49], 2008/2009)
PE-Assembler	short	Seed-and-Extend	(Ariyaratne and Sung [50], 2011)
QSORA	short	Seed-and-Extend	(Bryant et al. [23], 2009)
SHARCGS	short	Seed-and-Extend	(Dohm et al. [21], 2007)
SHORTY	short	Seed-and-Extend	(Hossain et al. [51], 2009)
SSAKE	short	Seed-and-Extend	(Warren et al. [22], 2007)
Taipan	short	Seed-and-Extend	(Schmidt et al. [24], 2009)
VCAKE	short	Seed-and-Extend	(Jeck et al. [52], 2007)

Reads are defined as "long" if produced by Sanger technology and "short" if produced by Illumina technology . Note that Velvet was designed for micro-reads (e.g. Illumina) but long reads can be given in input as additional data to resolve repeats in a greedy fashion.
doi:10.1371/journal.pone.0019175.t001

Narzisi G, Mishra B, Comparing De Novo Genome Assembly:
The Long and Short of It. 2011 PLoS ONE 6(4):

De novo assembly of short sequence reads
Paszkievicz, K. Studholme, D.
Briefings in Bioinformatics
August 2010 11(5): 457-472

Which assembler is best?

- Depends on:
 - Type of reads (Illumina, SoLID, 454, Ion Torrent, PacBio, Sanger etc)
 - Paired/mate-pair data?
 - Genome
 - Repeat content
 - Available hardware
- Prokaryote genomes – Velvet
- Larger genomes ABySS or Soapdenovo

Merging assemblies

- Often assemblies are produced from 454 or Sanger data and need to be merged with Illumina data
- In order of preference:
 1. Attempt to assemble 454/Sanger reads with Illumina reads using MIRA
 2. Merge assemblies separately using minimus2 or SSPACE
 3. Input 454/Sanger contigs as part of a reference guided assembly (e.g. Velvet/Columbus)

Transcriptome assembly

- de-novo transcriptome assembly is also possible
- RNA-seq reads can be assembled and isoform abundance estimated
- Much harder as Lander-Waterman assumptions of randomly distributed reads are not true
- Also complicated by splice-variants and the need to statistically model isoform abundance based on read distributions

- Oases/Velvet
- Trans-ABYSS
- SOAPde-novo
- Trinity

Good experimental option for vertebrates and other non-model organisms where a reference genome is not available

Typical assembly pipeline



Optimal de-novo sequencing strategy and review papers

Assessing the benefits of using mate-pairs to resolve repeats in de novo short-read prokaryotic assemblies

Joshua Wetzel , Carl Kingsford and Mihai Pop

BMC Bioinformatics 2011, **12**:95

**Comparing De Novo Genome Assembly:
The Long and Short of It.**

Narzisi, G. Mishra B.

2011 PLoS ONE 6(4)

De novo assembly of short sequence reads

Paszkiewicz, K. Studholme, D.

Briefings in Bioinformatics

August 2010 11(5): 457-472

A new strategy for genome assembly using short sequence reads and reduced representation libraries

Young A.L., Abaan H.O., Zerbino D, et al.

Genome Research 2010;20:249–56.

Variant calling with de-novo assembly

Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly

Heng Li^{1,*}

¹Broad Institute, 7 Cambridge Center, Cambridge, MA 02142, USA

Associate Editor: Dr. Michael Brudno

nature
genetics

ABSTRACT

Motivation: Eugene Myers in his stri suggested that in a string graph or e path spells a valid assembly. As a st every valid assembly of reads, such be constructed correctly, is in fact reads. In principle, every analysis bas sequencing (WGS) data, such as SNP calling, can also be achieved with uniti

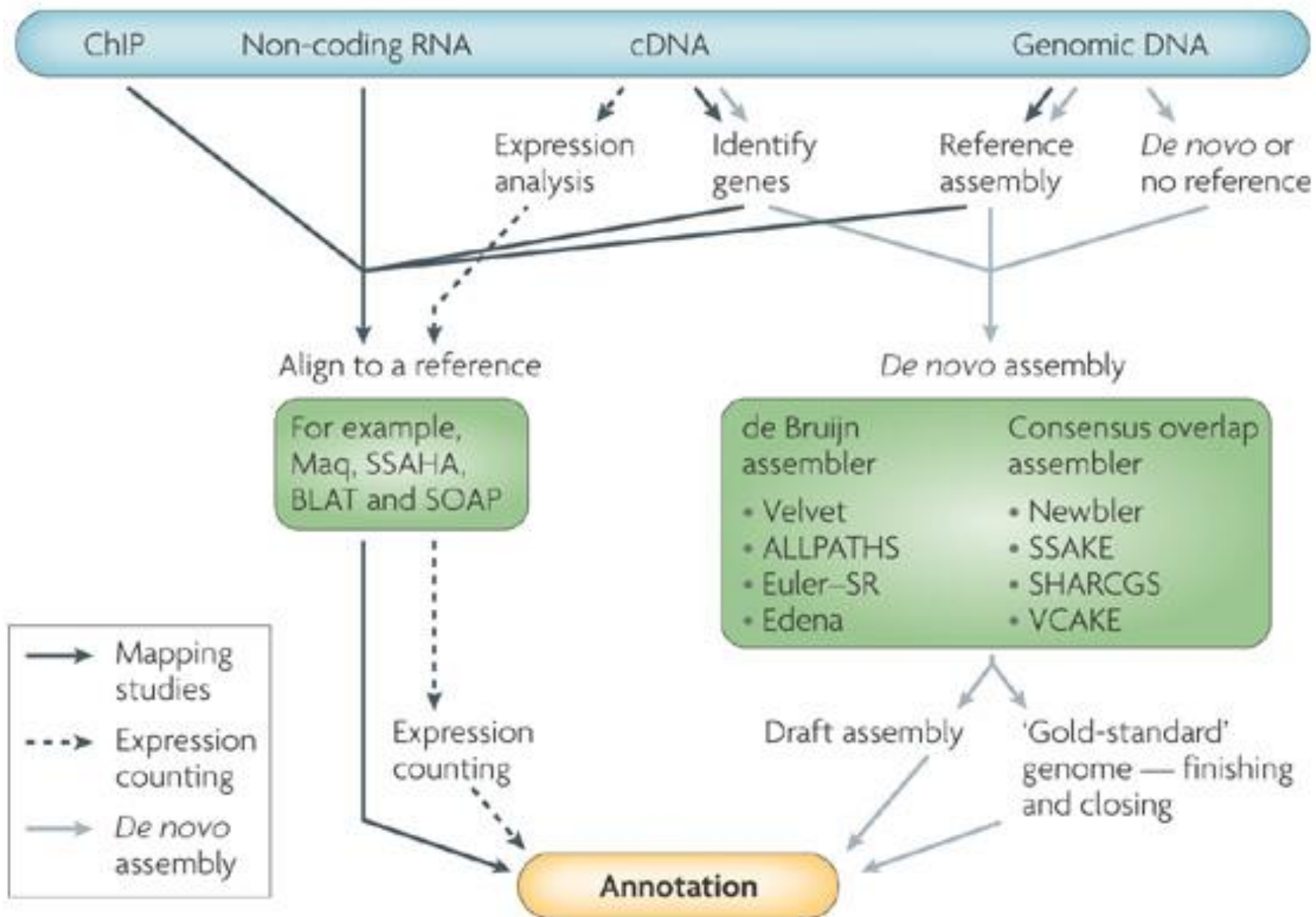
De novo assembly and genotyping of variants using colored de Bruijn graphs

Zamin Iqbal^{1,2,5}, Mario Caccamo^{3,5}, Isaac Turner¹, Paul Flicek² & Gil McVean^{1,4}

Detecting genetic variants that are highly divergent from a reference sequence remains a major challenge in genome sequencing. We introduce *de novo* assembly algorithms using colored de Bruijn graphs for detecting and genotyping simple and complex genetic variants in an individual or population. We provide an efficient software implementation, Cortex, the first *de novo* assembler capable of assembling multiple eukaryotic genomes simultaneously. Four applications of Cortex are presented. First, we detect and validate both simple

a single suitable reference, as in ecological sequencing²¹. Fourth, methods for variant calling from mapped reads typically focus on a single variant type. However, in cases in which variants of different types cluster, focus on a single type can lead to errors, for example, through incorrect alignment around indel polymorphisms^{6,7}. Fifth, although there are methods for detecting large structural variants, such as using array comparative genomic hybridization (aCGH)²²⁻²⁵ and mapped reads^{11,12,14,26}, these cannot determine the exact location, size or allelic sequence of variants. Finally, mapping

Raw sequence source



Questions!

biosciences.exeter.ac.uk/facilities/sequencing/usefulresources/

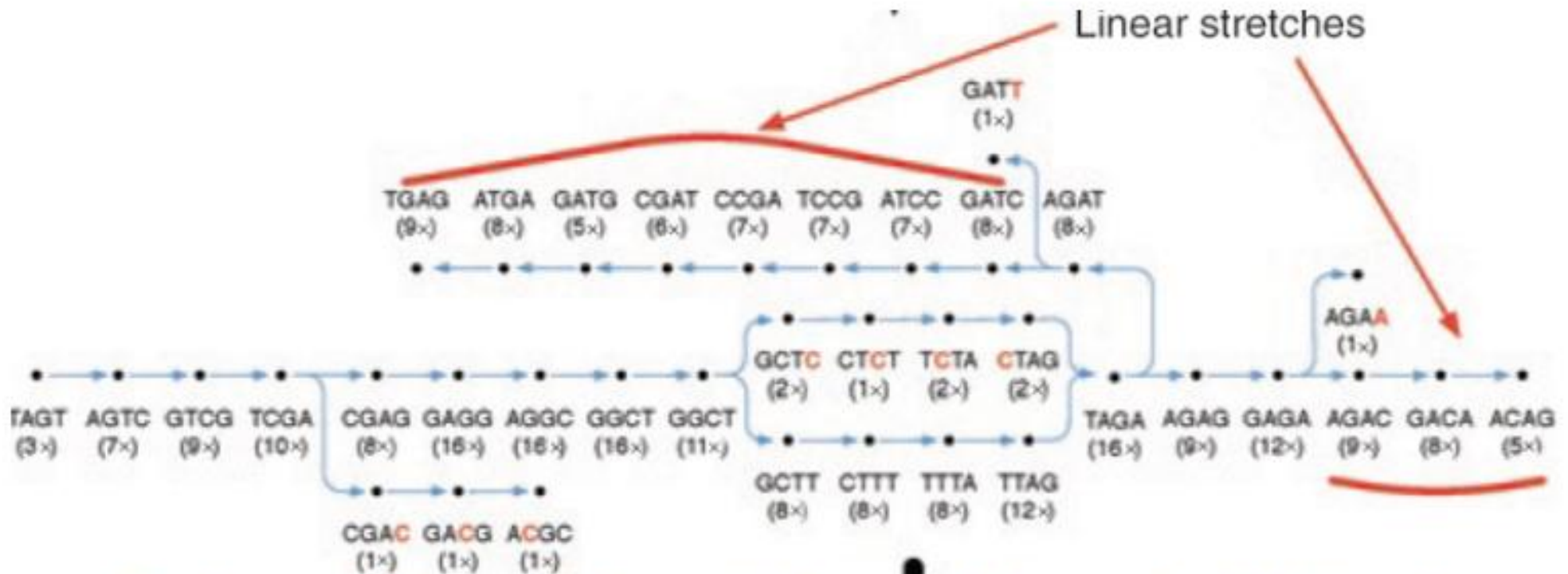
de-Bruijn graph assembly 1

TAGTCGAGGCTTTAGATCCGATGAGGCTTTAGAGACAG

```
AGTCGAG   CTTTAGA   CGATGAG   CTTTAGA
GTCGGG    TTAGATC    ATGAGGC   GAGACAG
GAGGCTC   ATCCGAT    AGGCTTT   GAGACAG
AGTCGAG   TAGATCC    ATGAGGC   TAGAGA
TAGTCGA   CTTTAGA    CCGATGA   TTAGAGA
CGAGGCT   AGATCCG    TGAGGCT   AGAGACA
TAGTCGA   GCTTTAG    TCOGATG   GCTTAG
TOGACGC   GATCCGA    GAGGCTT   AGAGACA
TAGTCGA   TTAGATC    GATGAGG   TTTAGAG
GTCGAGG   TCTAGAT    ATGAGGC   TAGAGAC
AGGCTTT   ATCCGAT    AGGCTTT   GAGACAG
AGTCGAG   TTAGATT    ATGAGGC   AGAGACA
GGCTTTA   TCOGATG    TTTAGAG
CGAGGCT   TAGATCC    TGAGGCT   GAGACAG
AGTCGAG   TTTAGATC   ATGAGGC   TTAGAGA
GAGGCTT   GATCCGA    GAGGCTT   GAGACAG
```

Genome is sampled with random sequencing 7bp reads (e.g. Illumina or 454)
Note errors in the reads are represented in red

de-Bruijn graph assembly 2

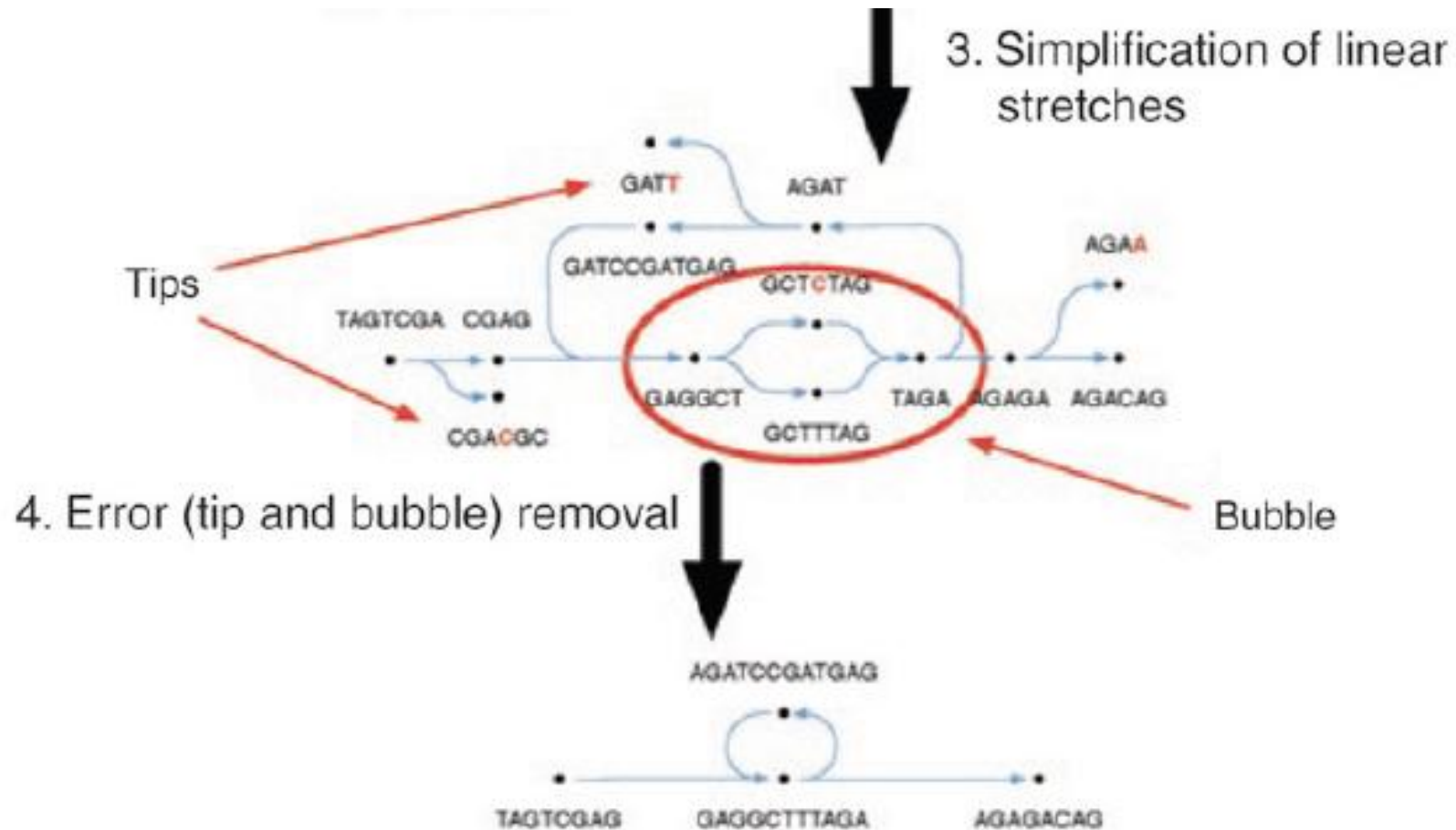


The k -mers in the reads (4-mers in this example) are collected into nodes and the coverage at each node is recorded (numbers at nodes)

Features

- ▶ continuous linear stretches within the graph
- ▶ Sequencing errors are low frequency tips in the graph

de-Bruijn graph assembly 3



Graph is simplified to combine nodes that are associated with the continuous linear stretches into single, larger nodes of various k -mer sizes

Error correction removes the tips and bubbles that result from sequencing errors

Final graph structure that accurately and completely describes in the original genome sequence