

Calling and filtering complex small variants (indels, MNPs, and haplotypes)

Erik Garrison

Wellcome Trust Sanger Institute
@ SeqShop, University of Michigan
May 20, 2015

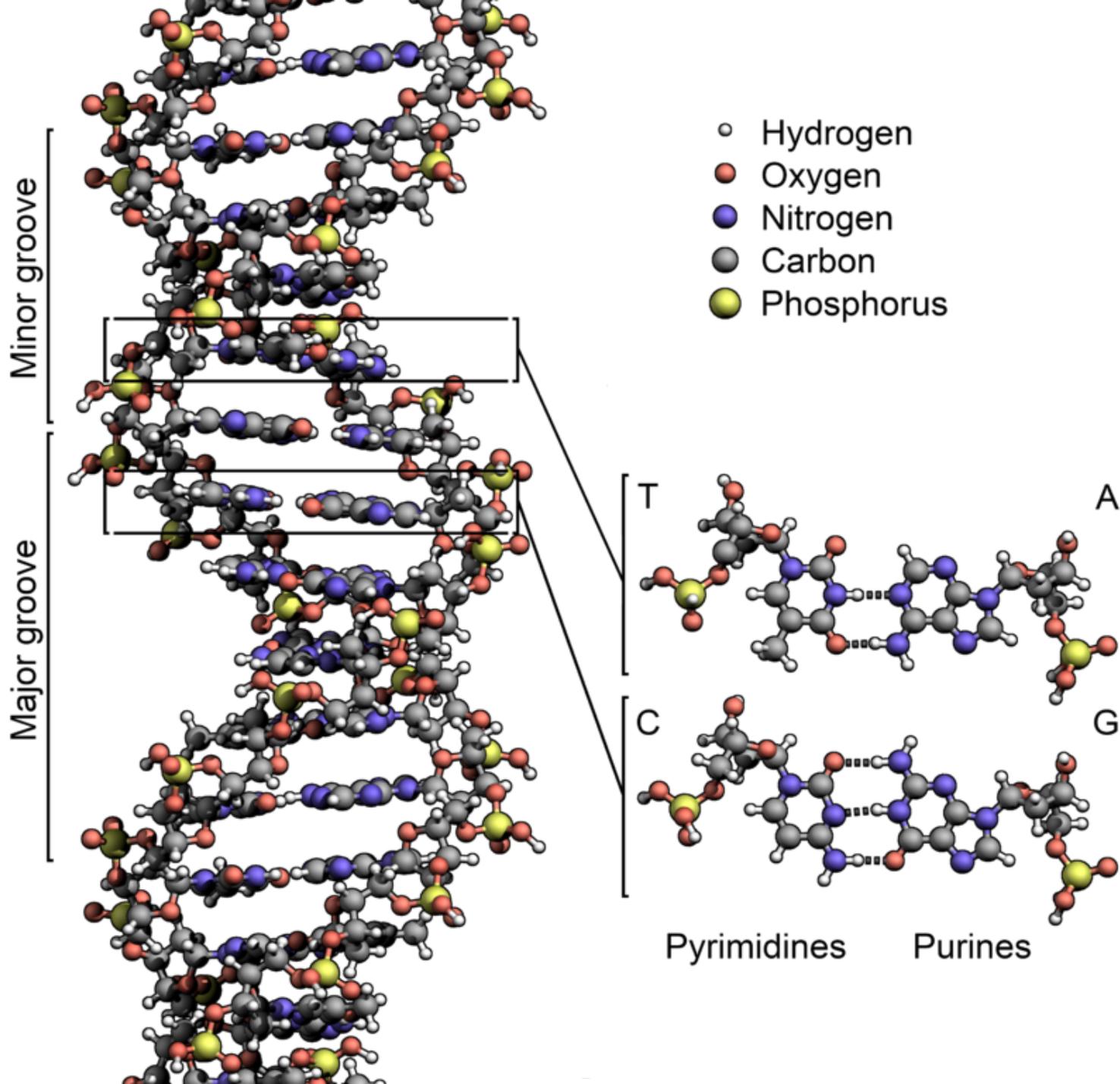


UNIVERSITY OF
CAMBRIDGE

Overview

1. **Genesis of variation (SNPs and indels)**
2. Causes of sequencing error
3. Sequence alignment
4. Alignment-based variant detection
5. Assembly-based variant detection
6. Haplotype-based variant detection
7. Graph-based methods

DNA



A SNP

A point mutation in which one base is swapped for another.

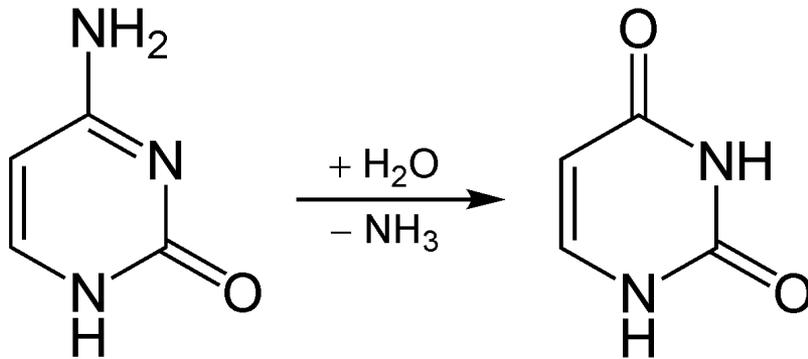
AATTAGCCATTA

AATTAGTCATTA

(some) causes of SNPs

- Deamination

- cytosine → uracil
- 5-methylcytosine → thymine
- guanine → xanthine (mispairs to A-T bp)
- adenine → hypoxanthine (mispairs to G-C bp)

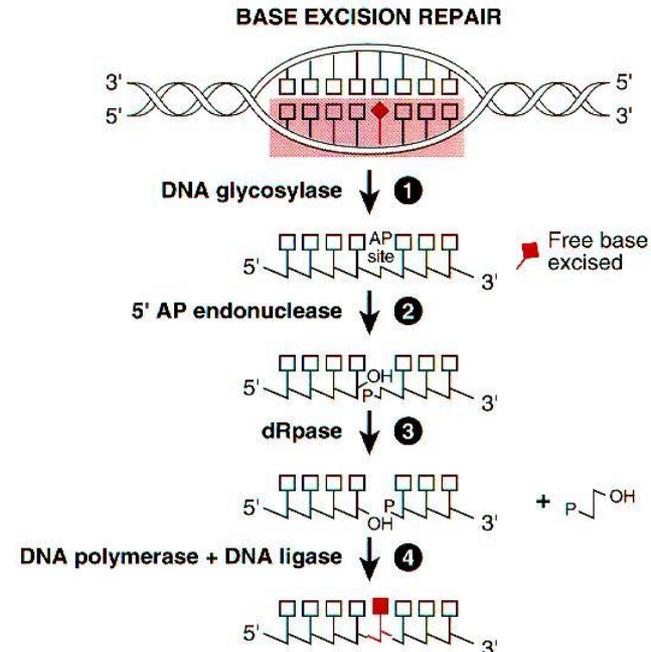
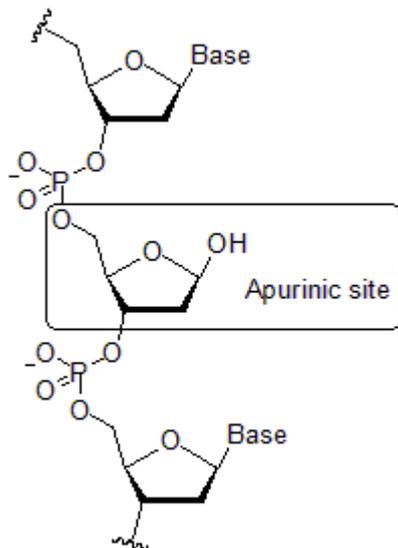


Deamination of Cytosine to Uracil
<http://en.wikipedia.org/wiki/Deamination>

(some) causes of SNPs

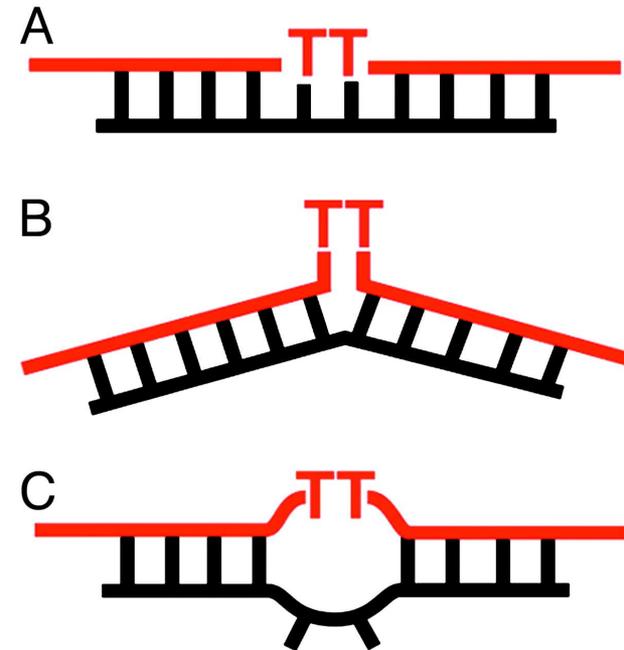
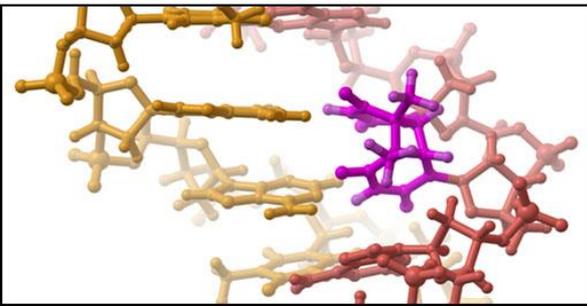
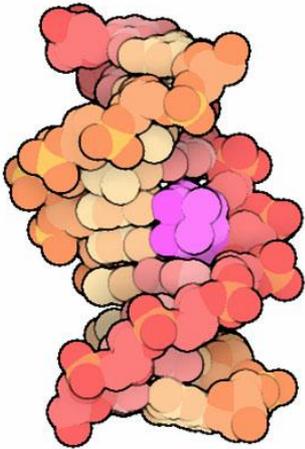
- Depurination

- purines are cleaved from DNA sugar backbone (5000/cell/day, pyrimidines at much lower rate)
- Base excision repair (BER) can fail → mutation



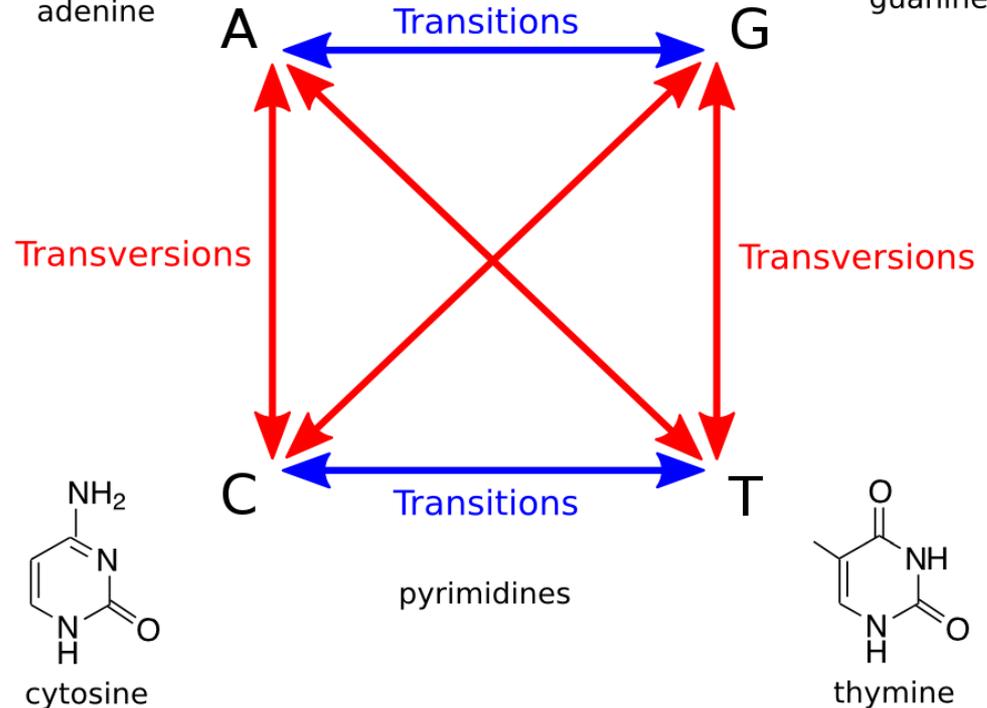
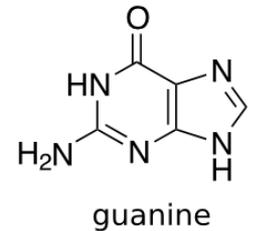
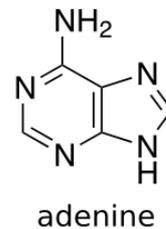
Multi-base events (MNPs)

- MNPs
 - thymine dimerization (UV induced)
 - other (e.g. oxidative stress induced)



Transitions and transversions

In general **transitions** are **2-3 times more common than transversions**. (But this depends on biological context.)



An INDEL

A mutation that results from the gain or loss of sequence.

AATTAGCCATTA

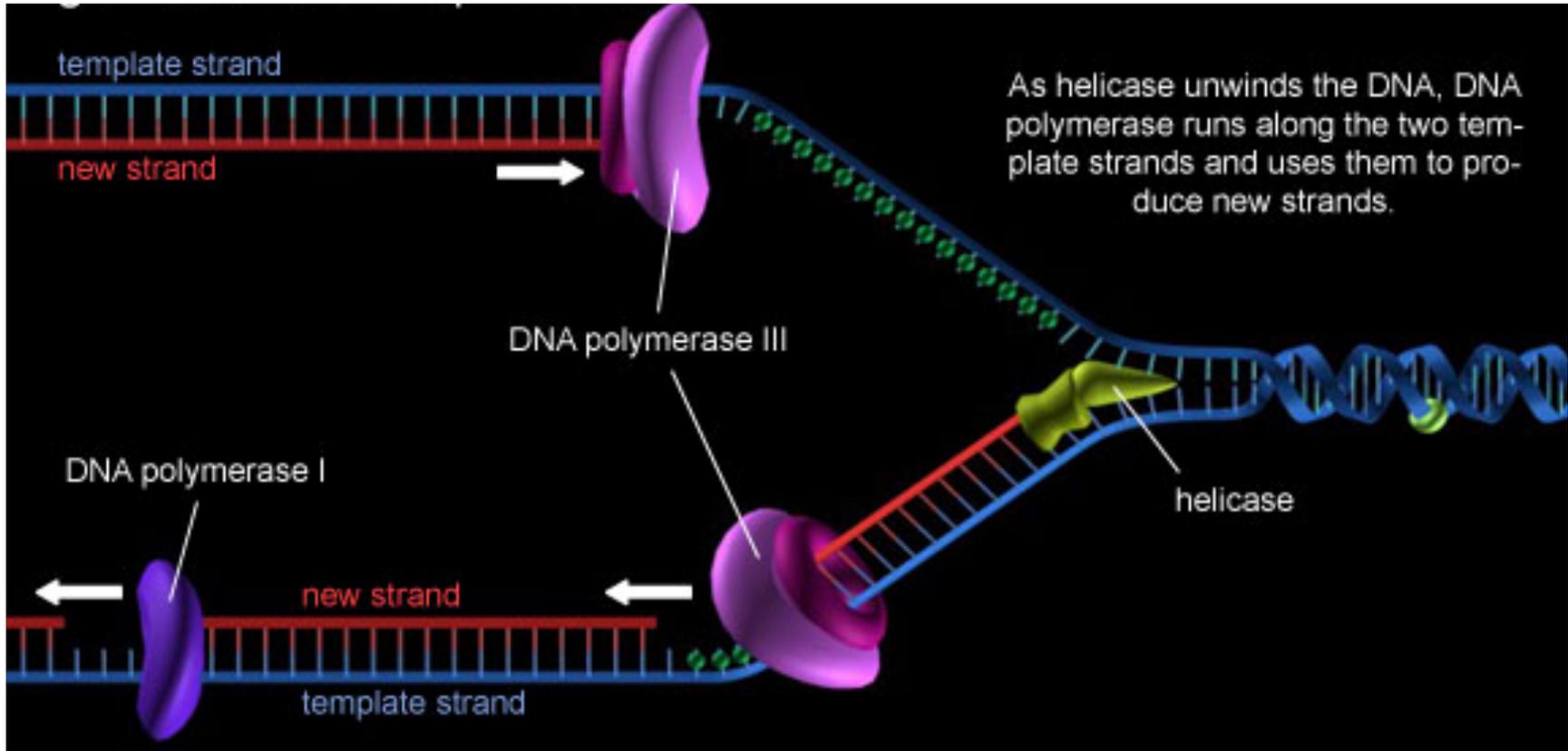
AATTA--CATTA

INDEL genesis

A number of processes are known to generate insertions and deletions in the process of DNA replication:

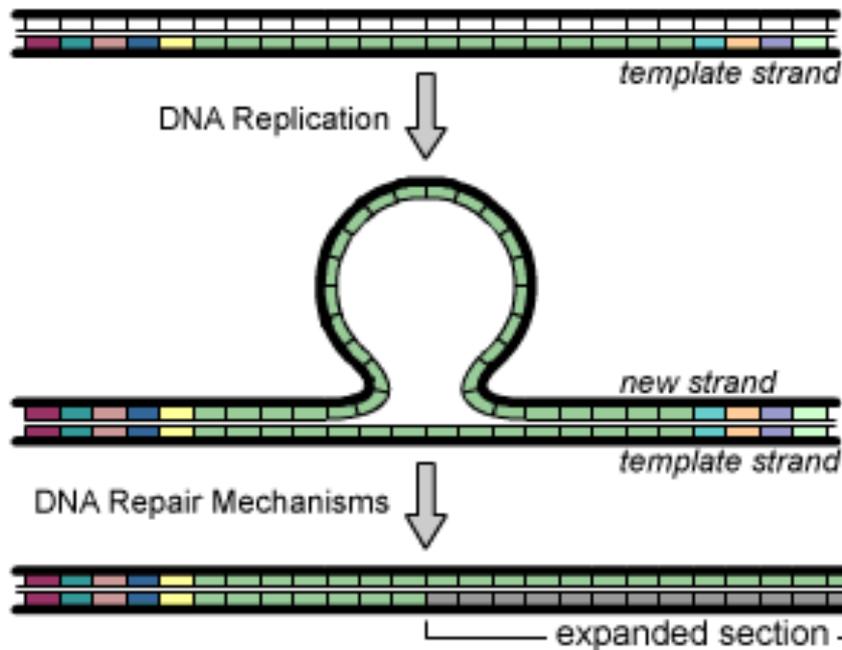
- Replication slippage
- Double-stranded break repair
- Structural variation (e.g. mobile element insertions, CNVs)

DNA replication



Polymerase *slippage*

A) Slippage Event



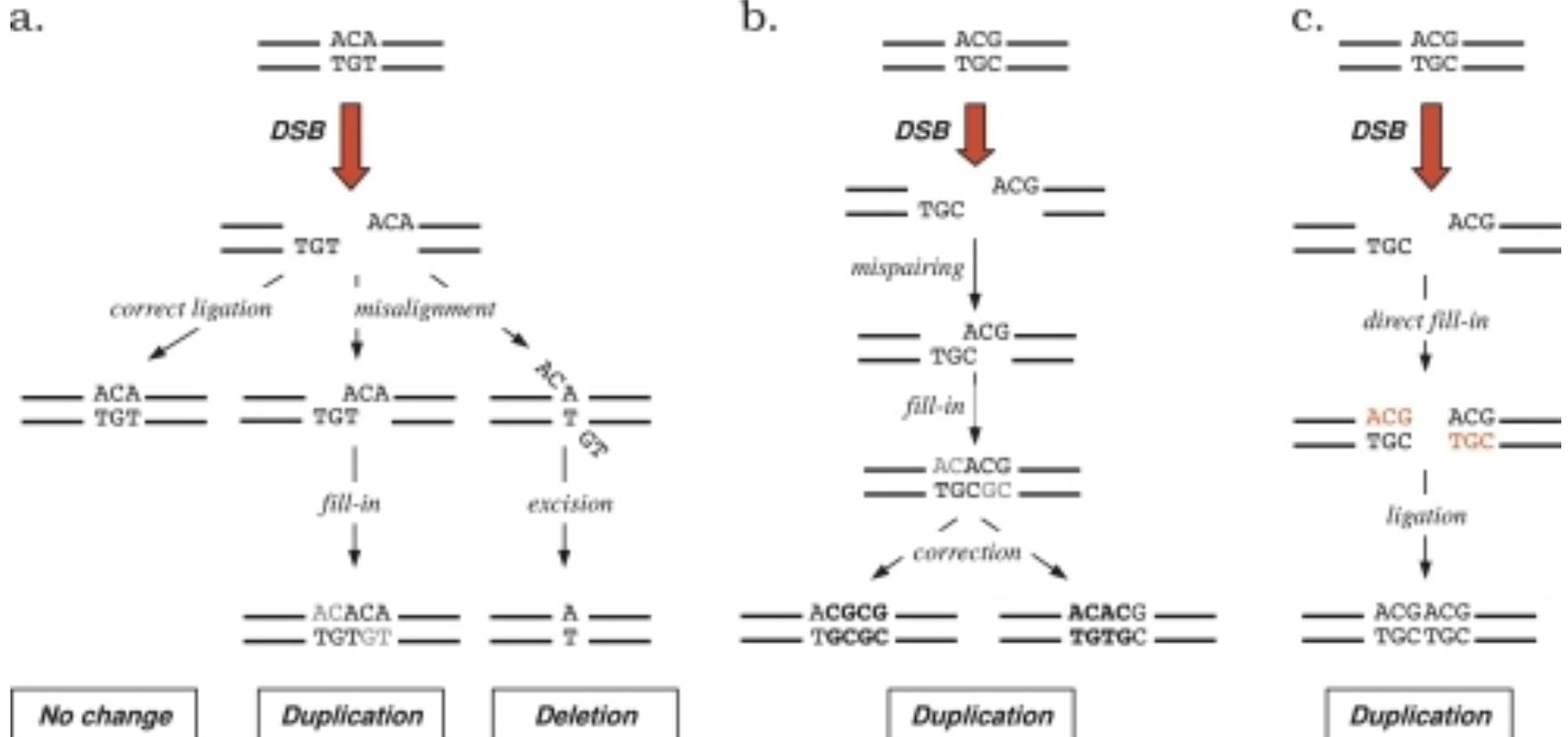
(A) During replication, polymerase slippage and subsequent reattachment may cause a bubble to form in the new strand. Slippage is thought to occur in sections of DNA with repeated patterns of bases (such as CAG), represented here by matching colors. Then, DNA repair mechanisms realign the template strand with the new strand and the bubble is straightened out. The resulting double helix is thus expanded.

B) No Slippage



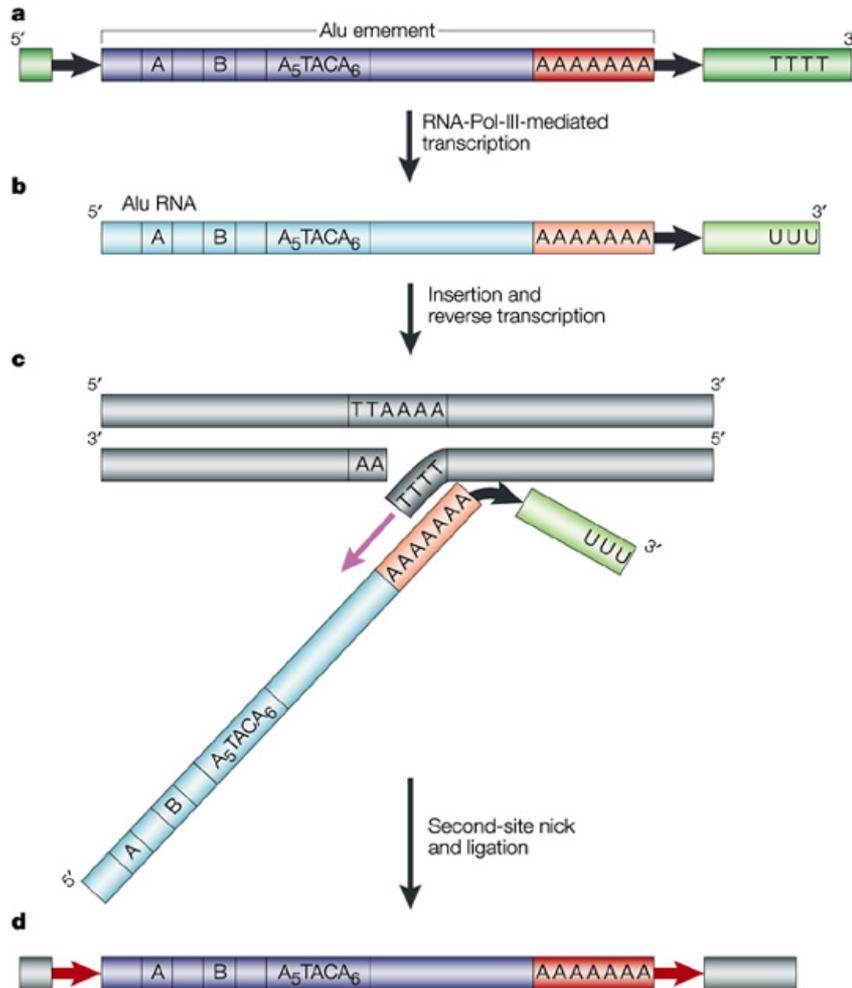
(B) Polymerase slippage, as theorized, cannot occur in DNA without repeating patterns of bases.

NHEJ-derived indels



DNA Slippage Occurs at Microsatellite Loci without Minimal Threshold Length in Humans: A Comparative Genomic Approach. Leclercq S, Rivals E, Jarne P - Genome Biol Evol (2010)

Structural variation (SV)



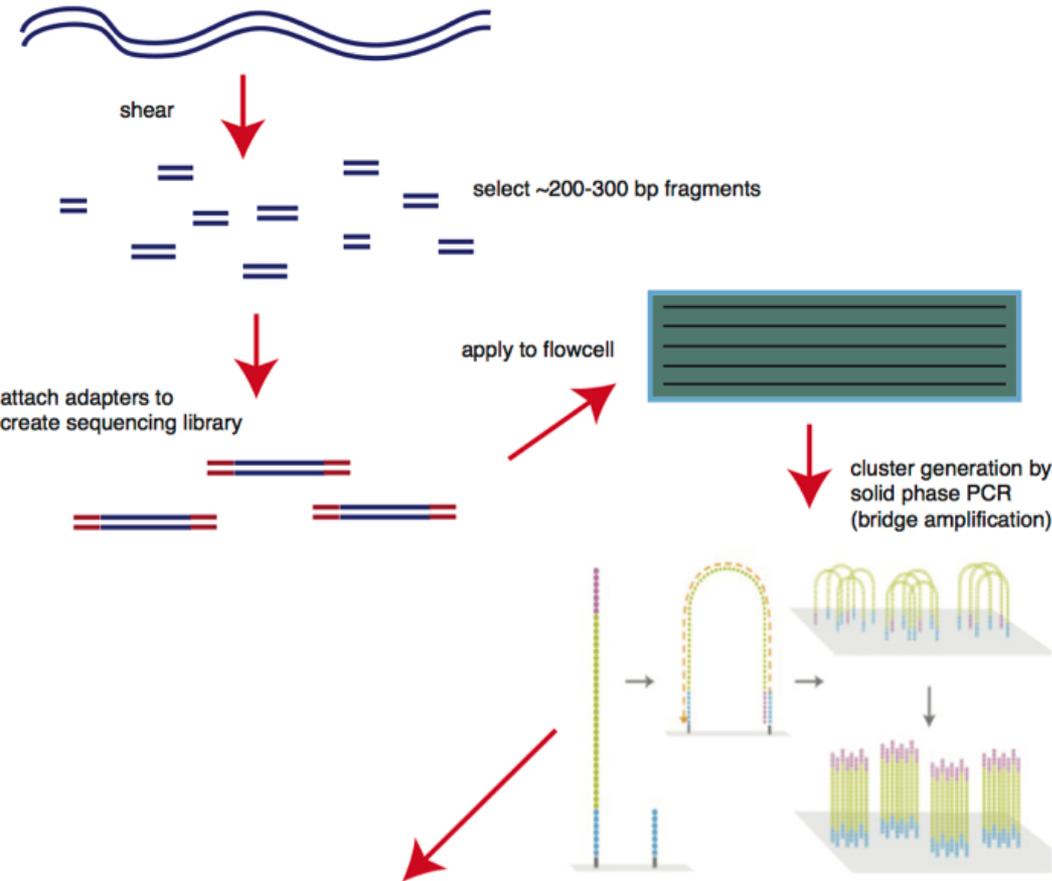
Transposable elements (in this case, an Alu) are sequences that can copy and paste themselves into genomic DNA, causing insertions.

Deletions can also be mediated by these sequences via other processes.

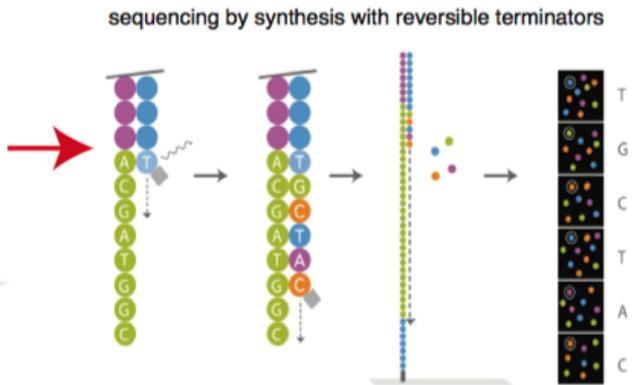
Overview

1. Genesis of variation (SNPs and indels)
2. **Causes of sequencing error**
3. Sequence alignment
4. Alignment-based variant detection
5. Assembly-based variant detection
6. Haplotype-based variant detection
7. Graph-based methods

Genomic DNA



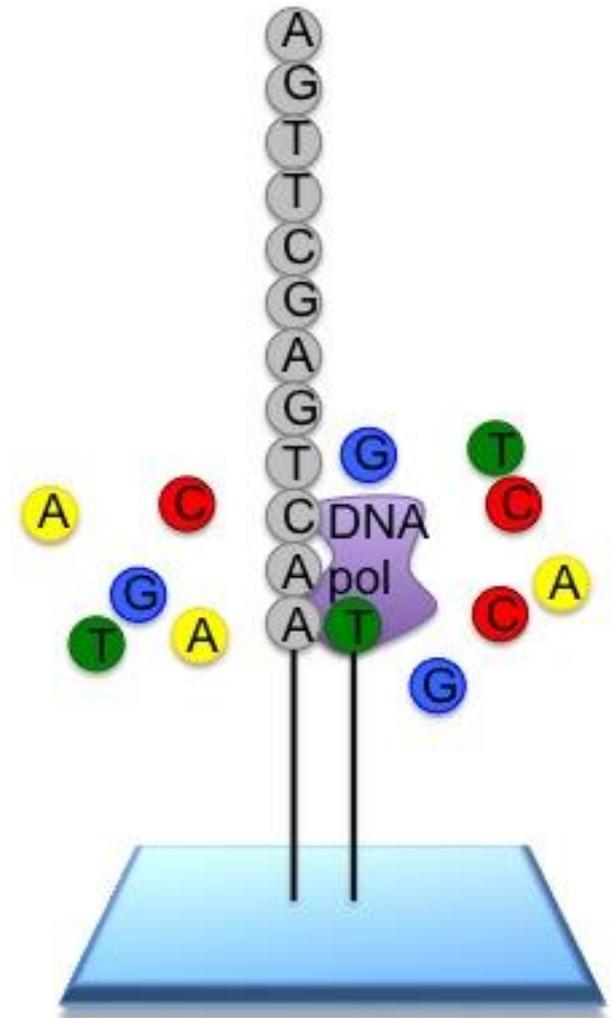
Sequencing by synthesis (Illumina)



Illumina sequencing process

(1)

Fluorescently labeled dNTPs with reversible terminators are incorporated by polymerase into growing double-stranded DNA attached to the flowcell surface.



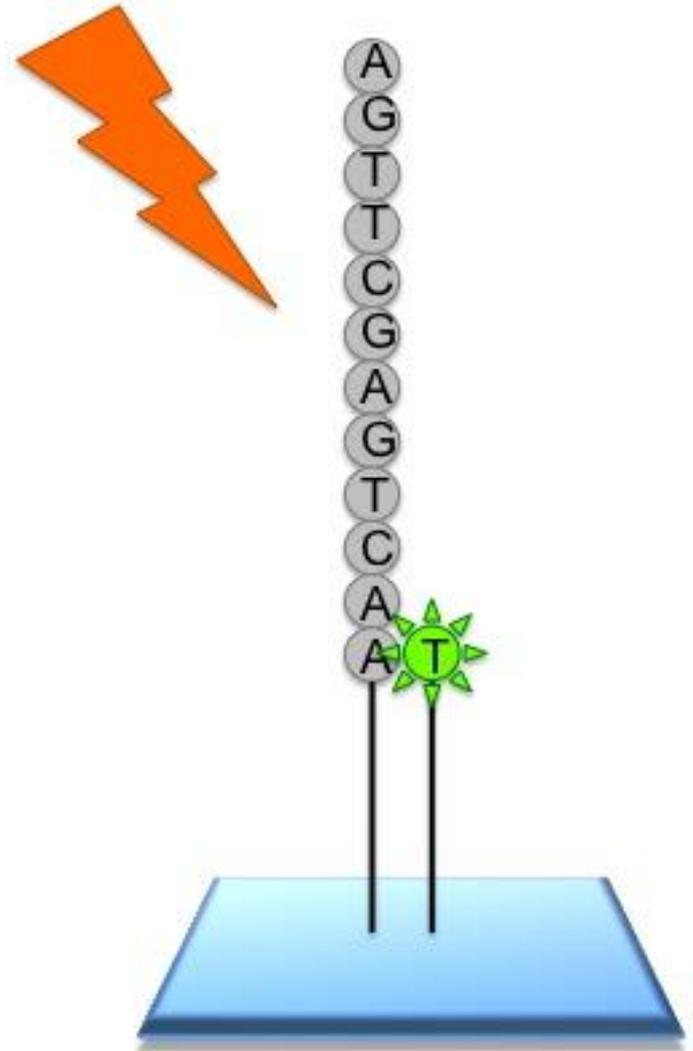
Illumina sequencing process

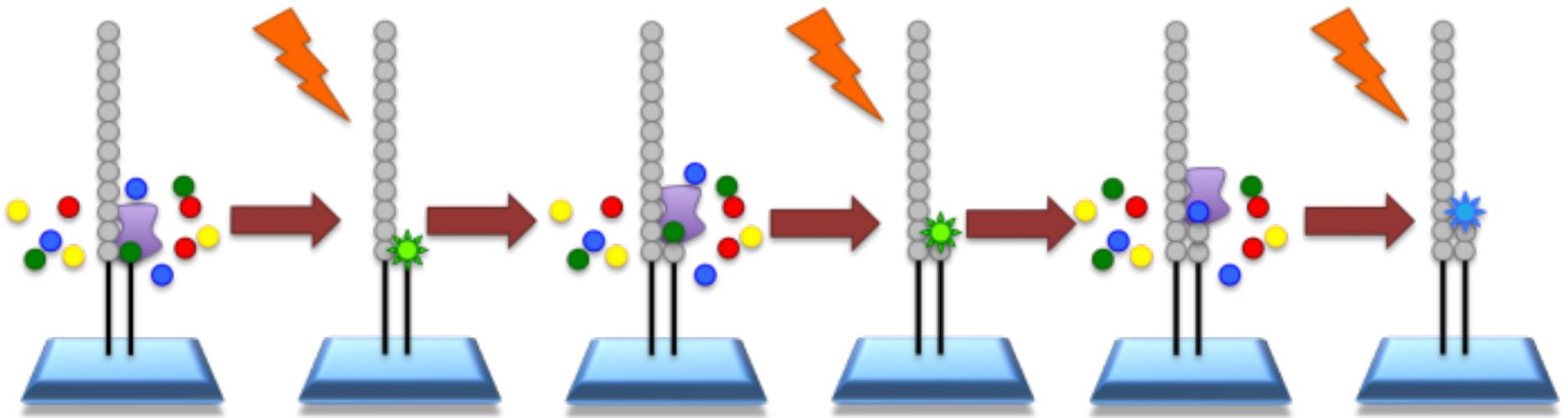
(2) A round of imaging determines the predominant base in each cluster at the given position.

(3) “reverse” terminator by washing, goto (1)

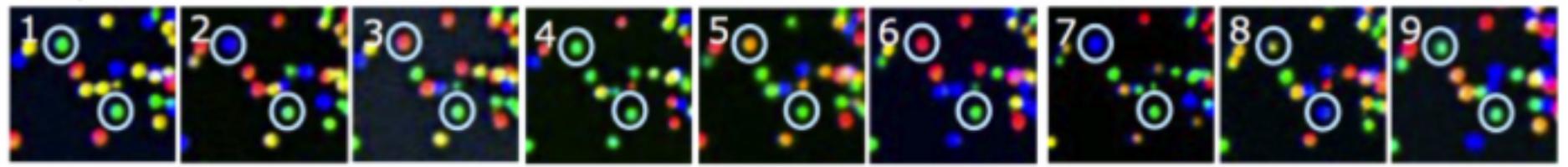
(...)

Build up images, process to determine sequence





TGCTACGAT...



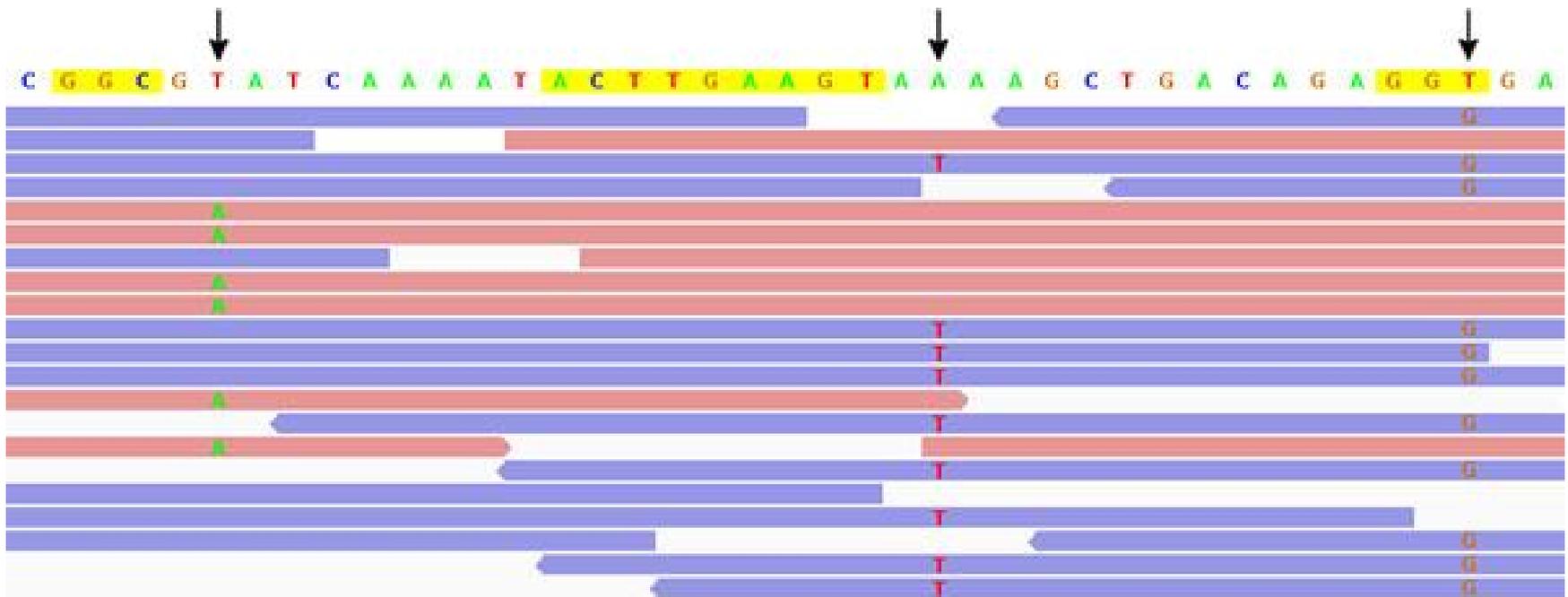
TTTTTTGT...

What can go wrong?

1. input artifacts, problems with library prep
 - a. replication in PCR has no error-correction (→ SNPs)
 - b. no quaternary structures (e.g. clamp) to prevent slippage (→ indels)
 - c. chimeras...
 - d. duplicates (worse if they are errors)
2. Sequencing-by-synthesis
 - a. phasing of step
 - i. synthesis reaction efficiency is not 100%
 - ii. particularly bad in A/T homopolymers
 - b. certain *context specific errors*
 - i. vary by sequencing protocol, device
 - ii. often strand-specific

Context specific errors

Show up as strand-specific errors:



Context specific errors (motifs)

Rank	Context	FER [%]	RER [%]	ERD [%]
1	ACGGCGGT	26.1	0.5	25.6
2	GTGGCGGT	25.1	0.7	24.4
3	GCGGCGGT	22.9	0.7	22.2
4	GTGGCTGT	22.4	0.6	21.8
5	ATGGCGGT	21.2	1.0	20.3
6	NCGGCGGT	20.0	0.7	19.3
7	GTGGCTTG	20.2	1.2	19.0
8	GNGGCGGT	19.2	0.7	18.5
9	GCGGCTGT	18.8	0.7	18.1
10	ACGGCTGT	18.6	0.8	17.7

← forward and reverse error rates for the ten most-common CSEs on a variety of illumina systems

(often GC-rich)

Overview

1. Genesis of variation (SNPs and indels)
2. Causes of sequencing error
- 3. Sequence alignment**
4. Alignment-based variant detection
5. Assembly-based variant detection
6. Haplotype-based variant detection
7. Graph-based methods

Global alignment

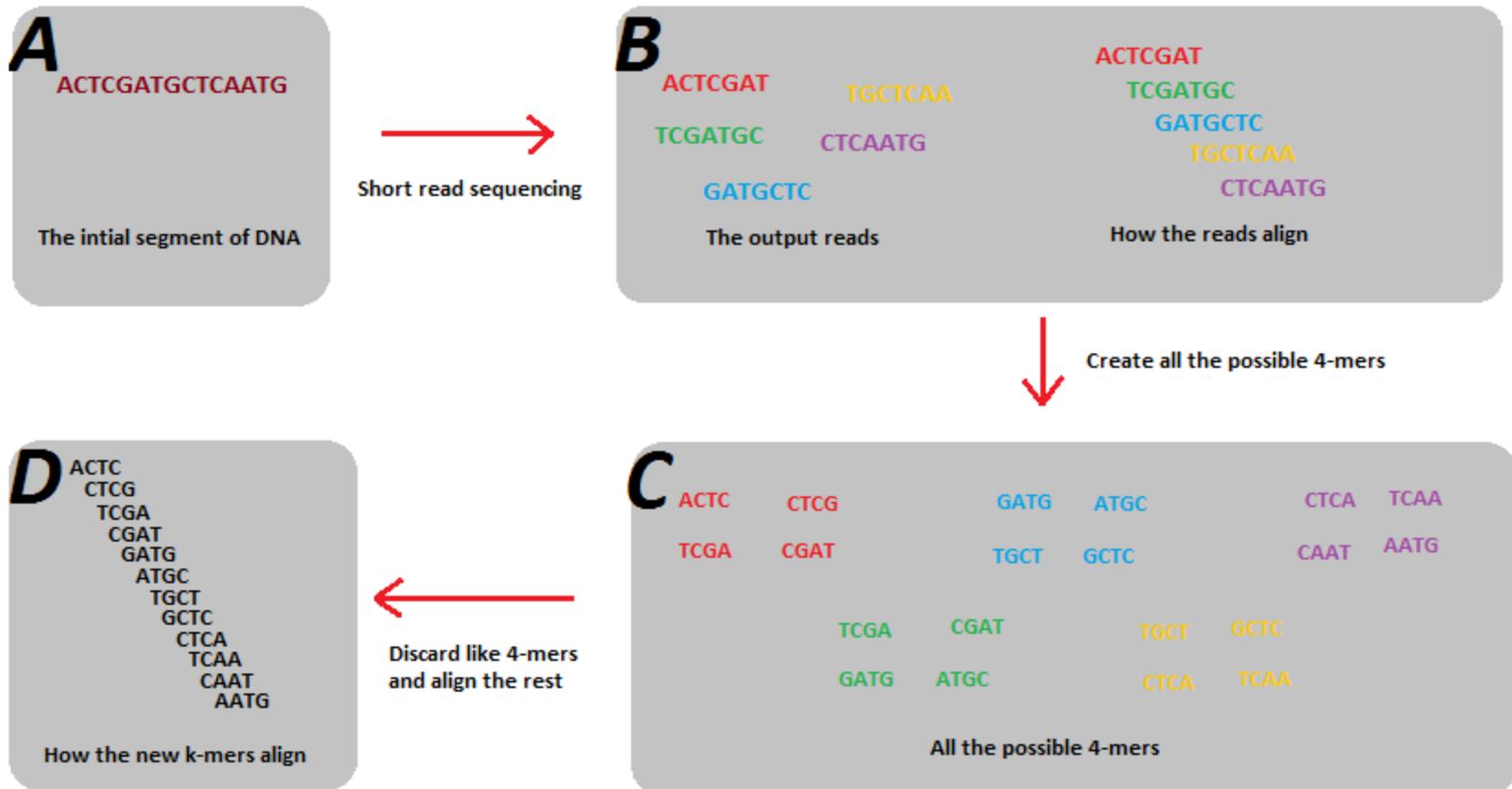
Idea: localize reads against a large reference system. Requires some cleverness.

Main approaches:

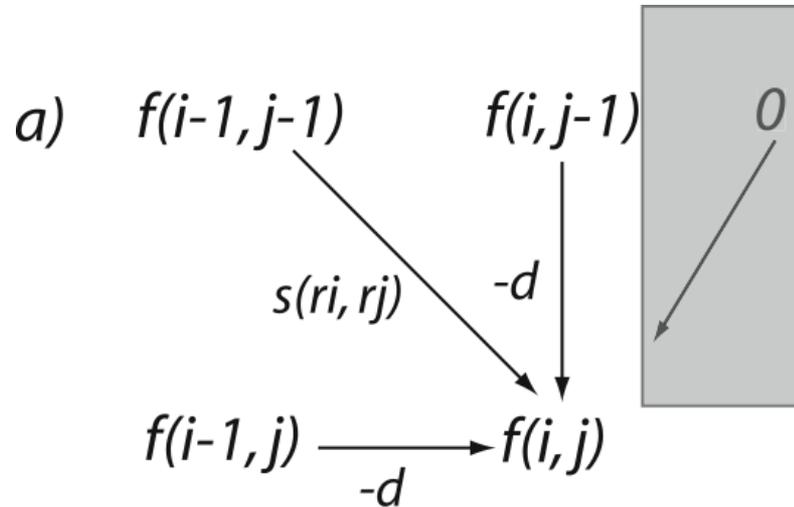
1. compressed suffix arrays
2. k-mer based seed+extend

k-mer based seed+extend

Used in novoalign, MOSAIK.



Local alignment



b)

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

Arrows in the table point from the cell (row, col) to (row-1, col-1), tracing the path of the local alignment: C(5,5) ← G(5,4) ← T(4,4) ← G(3,3) ← G(2,2) ← A(1,1).

Schematic of Needleman and Wunsch algorithm. Smith and Waterman added the 0.

[http://www.hiv.lanl.gov/content/sequence/HIV/REVIEWS/2006_7/](http://www.hiv.lanl.gov/content/sequence/HIV/REVIEWS/2006_7/ABECASIS/abecasis.html)

[ABECASIS/abecasis.html](http://www.hiv.lanl.gov/content/sequence/HIV/REVIEWS/2006_7/ABECASIS/abecasis.html)

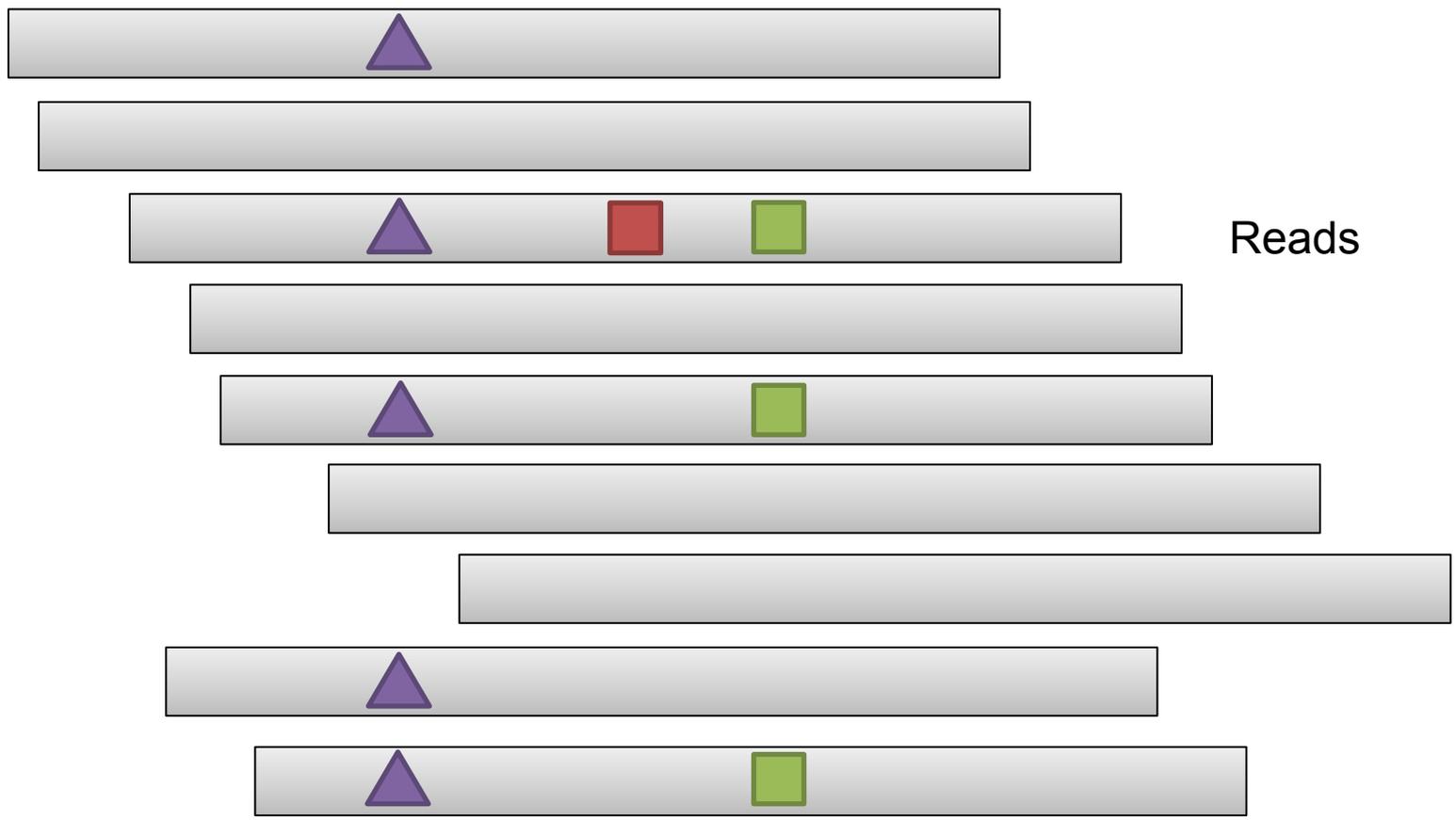
A	G	G	T	T	G	C
A	G	G	T	-	-	C

Overview

1. Genesis of variation (SNPs and indels)
2. Causes of sequencing error
3. Sequence alignment
4. **Alignment-based variant detection**
5. Assembly-based variant detection
6. Haplotype-based variant detection
7. Graph-based methods

Alignments to candidates

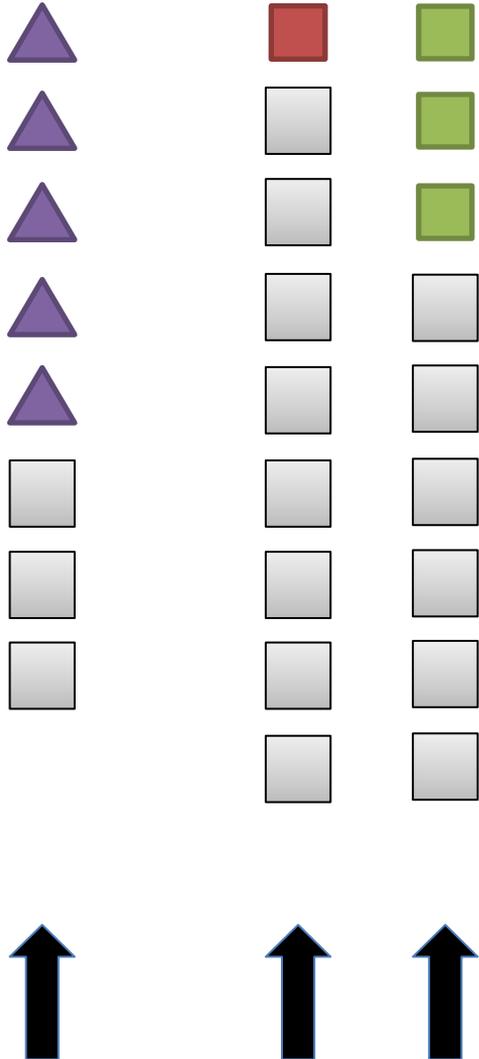
Reference



Variant observations

The data exposed to the caller

Reference



Alignment-based variant calling

For each candidate locus, collect information about support for each allele. Use heuristics (e.g. varscan) or a bayesian model (e.g. samtools, GATK, freebayes) to assign maximum-likelihood alleles and genotypes:

$$\begin{aligned} Pr(D|G = \{A_1, A_2\}) &= \prod_{i=1}^M Pr(b_i|G = \{A_1, A_2\}) \\ &= \prod_{i=1}^M \left(\frac{1}{2}p(b_i|A_1) + \frac{1}{2}p(b_i|A_2) \right), \end{aligned}$$

(more on this tomorrow...)

$$p(b|A) = \begin{cases} \frac{e}{3} & b \neq A \\ 1 - e & b = A. \end{cases}$$

But wait!

If we work directly from our alignments, we may assign too much significance to the particular way our reads were aligned.

Example: calling INDEL variation

Can we quickly design a process to detect indels from alignment data?

What are the steps you'd do to find the indel between these two sequences?

CAAATAAGGTTTGGAGTTCATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCATTATA

Indel finder

We could start by finding the long matches in both sequences at the start and end:

CAAATAAGGTTTGGAGTTCATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCATTATA

A vertical line is positioned between the 10th and 11th characters of the top sequence. The top sequence is CAAATAAGGTTTGGAGTTCATTATA and the bottom sequence is CAAATAAGGTTTGGAAATTTTCTGGAGTTCATTATA. The first 10 characters (CAAATAAGGTTT) match perfectly. At the 11th position, the top sequence has 'G' and the bottom has 'A'. From the 12th position onwards, the sequences match again: 'AGTTCATTATA'.

CAAATAAGGTTTGGAGTTCATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCATTATA

A vertical line is positioned between the 15th and 16th characters of the top sequence. The top sequence is CAAATAAGGTTTGGAGTTCATTATA and the bottom sequence is CAAATAAGGTTTGGAAATTTTCTGGAGTTCATTATA. The last 7 characters (GGAGTTCATTATA) match perfectly. At the 16th position, the top sequence has 'G' and the bottom has 'A'. From the 17th position onwards, the sequences match again: 'AGTTCATTATA'.

Indel finder

We can see this more easily like this:

CAAATAAGGTTTGGAGTTCTATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCTATTATA

CAAATAAGGTTTGGAGTTCTATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCTATTATA

CAAATAAGGTTTGGAGTTCTATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCTATTATA

Indel finder

The match structure implies that the sequence that doesn't match was inserted in one sequence, or lost from the other.

CAAATAAGGTT - - - - - TGGAGTTCTATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCTATTATA

So that's easy enough....

Something more complicated

These sequences are similar to the previous ones, but with different mutations between them.

CAAATAAGGAAATTTTCTGGAGTTCATTATA
CAAATAAGGTTTGCTATCTAGGTATTATA

They are still (kinda) homologous but it's not easy to see.

Pairwise alignment

One solution, assuming a particular set of alignment parameters, has 3 indels and a SNP:

CAAATAAGGAAATTT - - - - TCTGGAGTTCATTATA
CAAATAAGG - - - TTTGCTATCT - - AGGT - TATTATA

But if we use a higher gap-open penalty, things look different:

CAAATAAGGAAATTT - - TCTGGAGTTCATTATA
CAAATAAGG - - - TTTGCTATCTAGGT - TATTATA

Alignment as interpretation

Different parameterizations can yield different results.

Different results suggest “different” variation.

What kind of problems can this cause? (And how can we mitigate these issues?)

INDELs have multiple representations and require normalization for standard calling

Left alignment allows us to ensure that our representation is consistent across alignments and also variant calls.

CGTATGATCTAG**GCGCGC**TAGCTAGCTAGC
CGTATGATCTA - - **GCGC**TAGCTAGCTAGC

← Left aligned

CGTATGATCTAG**GCGCGC**TAGCTAGCTAGC
CGTATGATCTAG**C** - - **G**CTAGCTAGCTAGC

CGTATGATCTAG**GCGCGC**TAGCTAGCTAGC
CGTATGATCTAG**GCGC** - -TAGCTAGCTAGC

example: 1000G Phasel low coverage
chr20:708257, ref:AGC alt:CGA

ref: TATAGAGAGAGAGAGAGAGAGAGC GAGAGAGAGAGAGAGAGAGAGGGAGAGACGGAGTT
alt: TATAGAGAGAGAGAGAGAGAGAGC GAGAGAGAGAGAGAGAGAGAGAGGGAGAGACGGAGTT

ref: TATAGAGAGAGAGAGAGAGAGAGC -- GAGAGAGAGAGAGAGAGAGGGAGAGACGGAGTT
alt: TATAGAGAGAGAGAGAGAGAG -- CGAGAGAGAGAGAGAGAGAGGGAGAGACGGAGTT

Overview

1. Genesis of variation (SNPs and indels)
2. Causes of sequencing error
3. Sequence alignment
4. Alignment-based variant detection
- 5. Assembly-based variant detection**
6. Haplotype-based variant detection
7. Graph-based methods

***Problem:* inconsistent variant representation makes alignment-based variant calling difficult**

If alleles are represented in multiple ways, then to detect them correctly with a single-position based approach we need:

1. An awesome normalization method
2. Perfectly consistent filtering (so we represent our entire context correctly in the calls)
3. Highly-accurate reads

Solution: assembly and haplotype-driven detection

We can shift our focus from the specific interpretation in the alignments:

- this is a SNP
 - whereas this is a series of indels
- ... and instead focus on the underlying sequences.

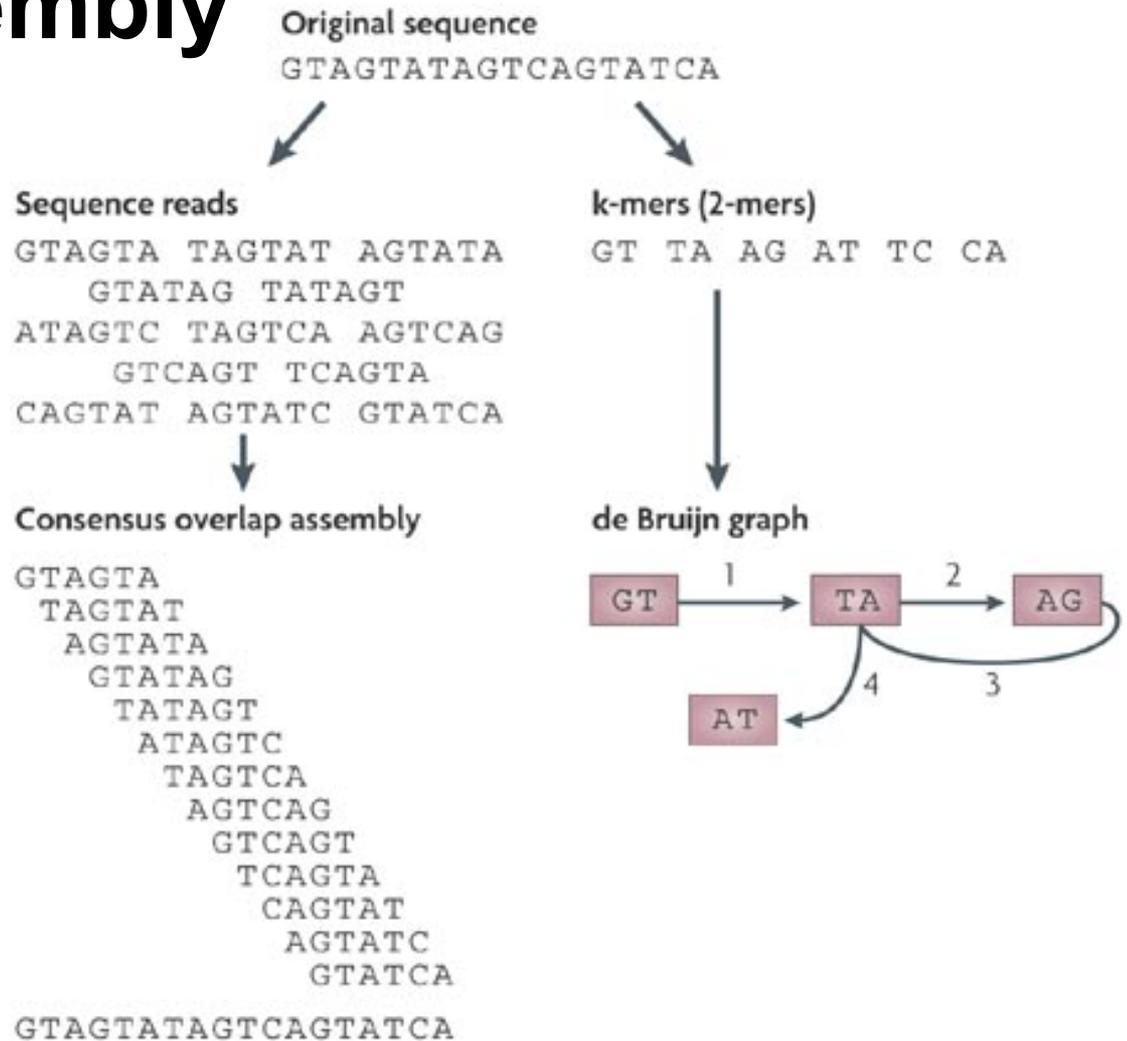
Basically, we use the alignments to localize reads, then process them again with assembly approaches to determine candidate alleles.

Variant detection by assembly

Multiple methods have been developed by members of the 1000G analysis group:

- **Global joint assembly**
 - cortex
 - SGA (localized to 5 megabase chunks)
- **Local assembly**
 - Platypus (+cortex)
 - GATK HaplotypeCaller
- **k-mer based detection**
 - FreeBayes (anchored reference-free windows)

de Bruijn Assembly

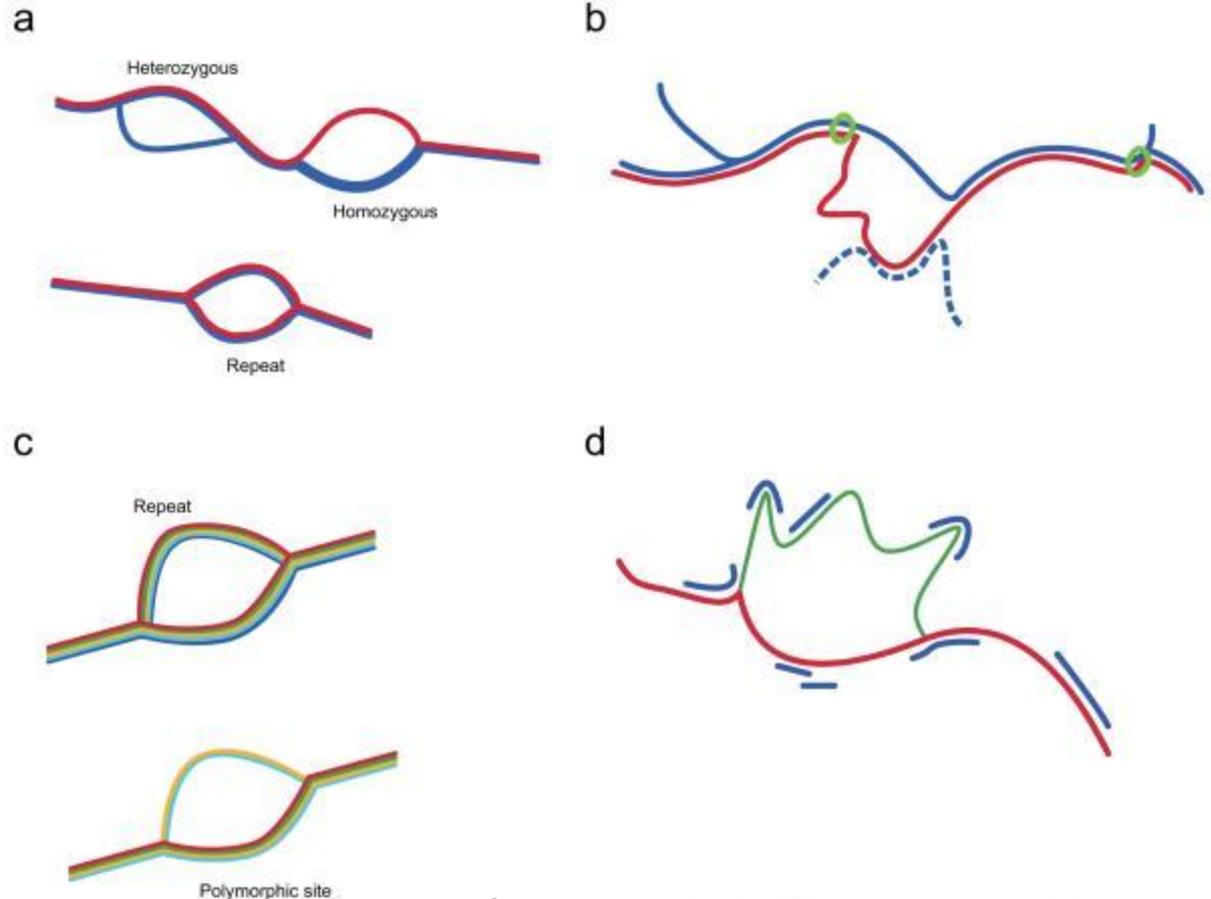


Using colored graphs (Cortex)

Variants can be called using bubbles in deBruijn graphs.

Method is completely reference-free, except for reporting of variants. The reference is threaded through the colored graph.

Many samples can be called at the same time.



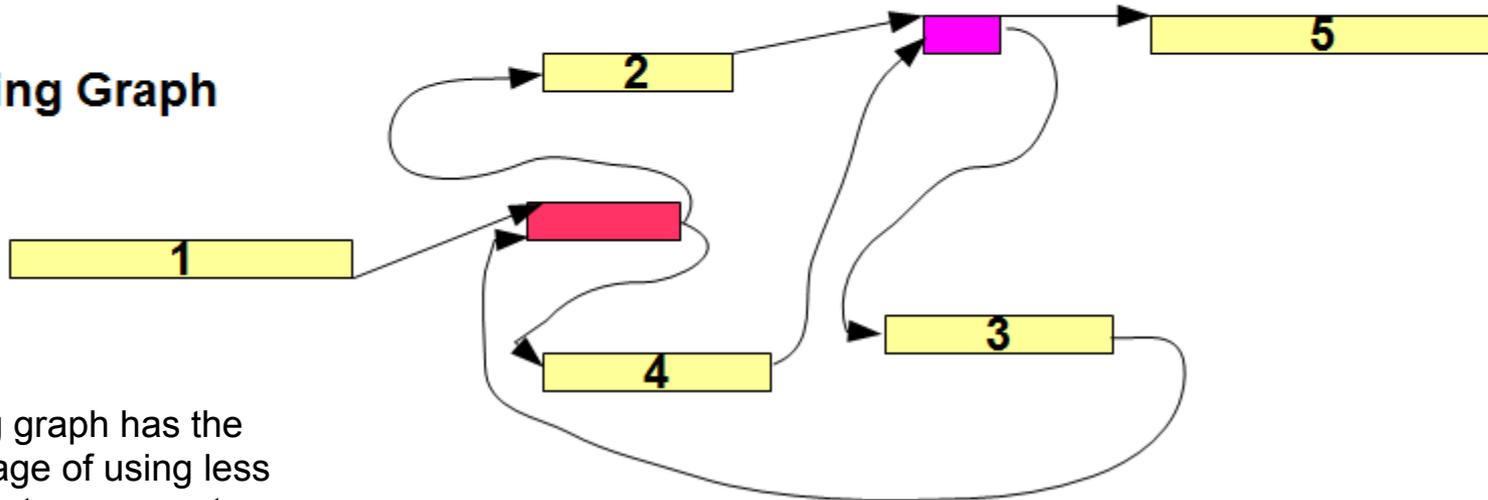
from Iqbal et. al., "De novo assembly and genotyping of variants using colored de Bruijn graphs." (2012)

String graphs (SGA)

Genome



String Graph



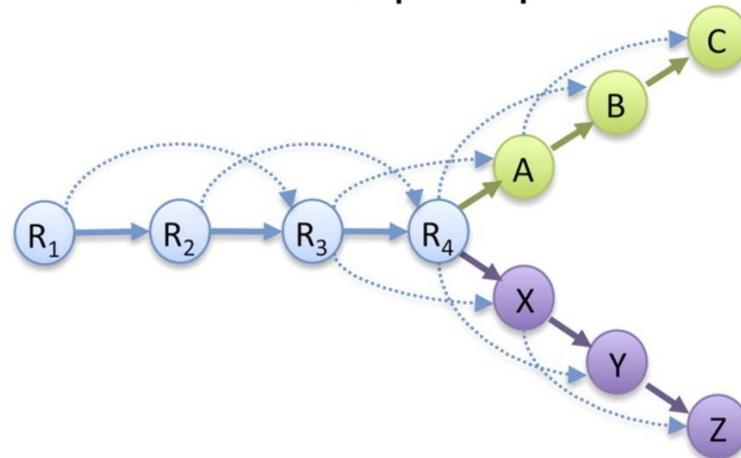
A string graph has the advantage of using less memory to represent an assembly than a de Bruijn graph. In the 1000G, SGA is run on alignments localized to ~5mb chunks.

Discovering alleles using graphs (GATK HaplotypeCaller)

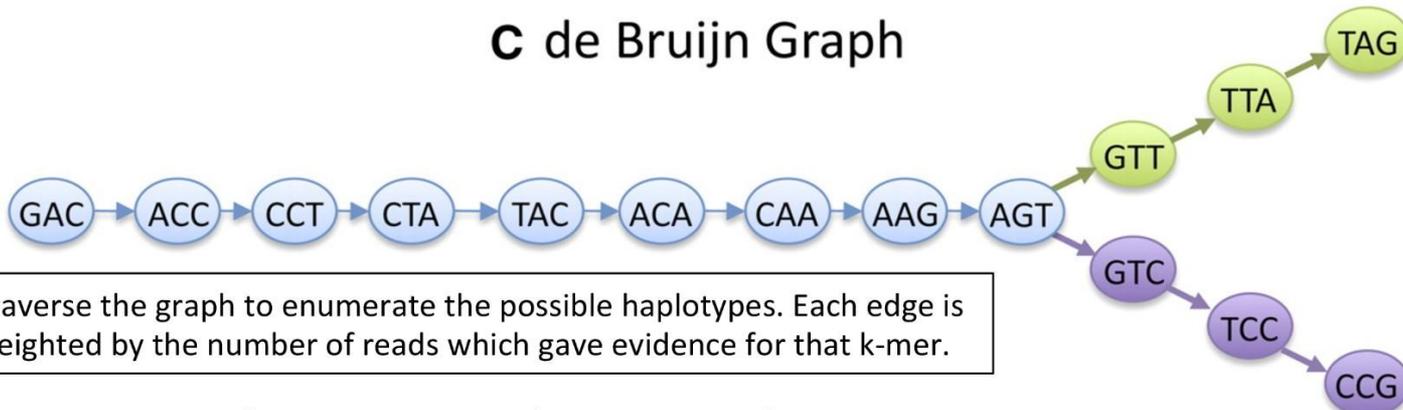
A Read Layout

R₁: GACCTACA
 R₂: ACCTACAA
 R₃: CCTACAAG
 R₄: CTACAAGT
 A: TACAAGTT
 B: ACAAGTTA
 C: CAAGTTAG
 X: TACAAGTC
 Y: ACAAGTCC
 Z: CAAGTCCG

B Overlap Graph



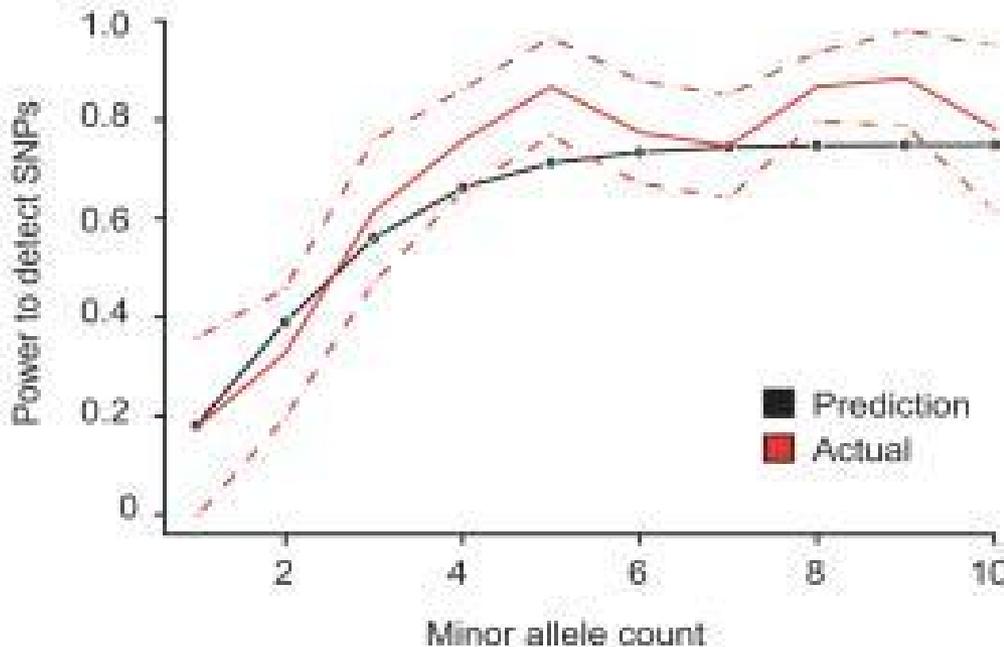
C de Bruijn Graph



Traverse the graph to enumerate the possible haplotypes. Each edge is weighted by the number of reads which gave evidence for that k-mer.

Why don't we just assemble?

Assembly-based calls tend to have high specificity, but sensitivity suffers.



The requirement of exact kmer matches means that errors disrupt coverage of alleles.

Existing assembly methods don't just detect point mutations--- they detect haplotypes.

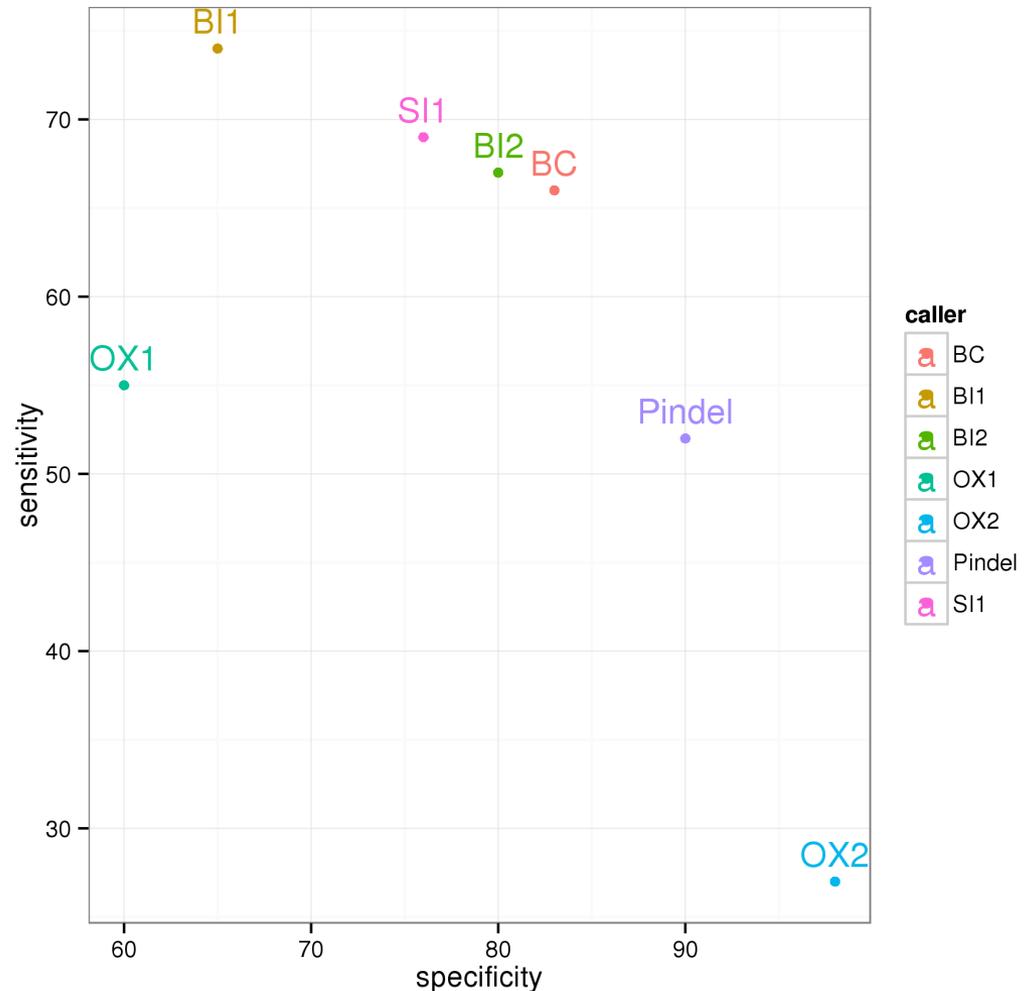
from Iqbal et. al., "De novo assembly and genotyping of variants using colored de Bruijn graphs." (2012)

Indel validation, 191 AFR samples

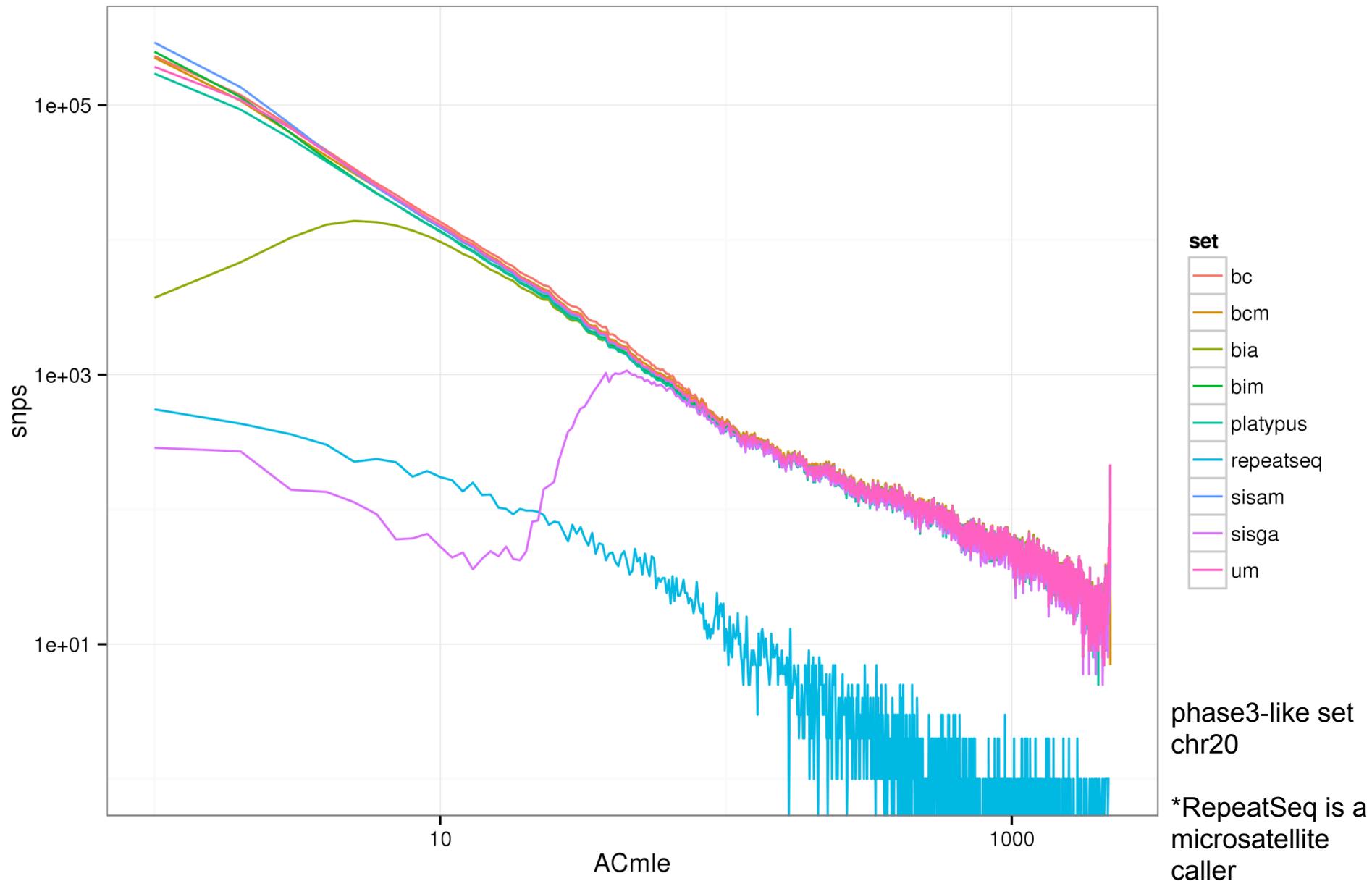
High-depth miSeq sequencing-based validation on 4 samples.

Local assembly methods (BI2, BC, SI1)* have higher specificity than baseline mapping-based calls (BI1), but lower sensitivity. Global assembly (OX2) yielded very low error, but also low sensitivity.

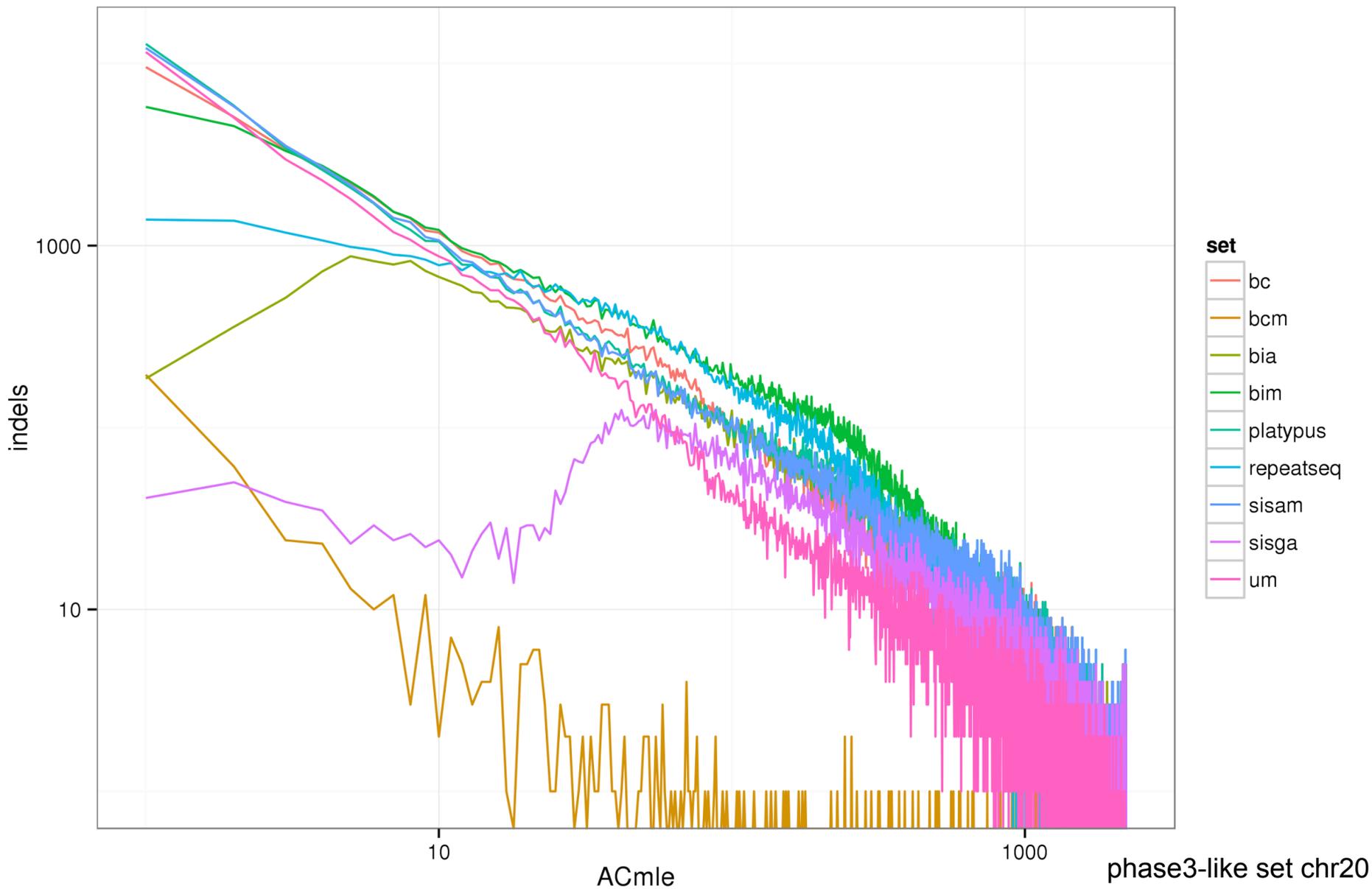
*The local assembly-based method Platypus (OX1) had a genotyping bug which caused poor performance.



Site-frequency spectrum, SNPs



Site-frequency spectrum, indels



Overview

1. Genesis of variation (SNPs and indels)
2. Causes of sequencing error
3. Sequence alignment
4. Alignment-based variant detection
5. Assembly-based variant detection
- 6. Haplotype-based variant detection**
7. Graph-based methods

Finding haplotype polymorphisms

Two
reads

AGAACCCAGTGCTCTTTCTGCT
AGAACCCAGTGGTCTTTCTGCT

a SNP

AGAACCCAGTGCCTCTTTCTGCT
AGAACCCAGTGGTCTTTCTGCT

Their
alignmen
t

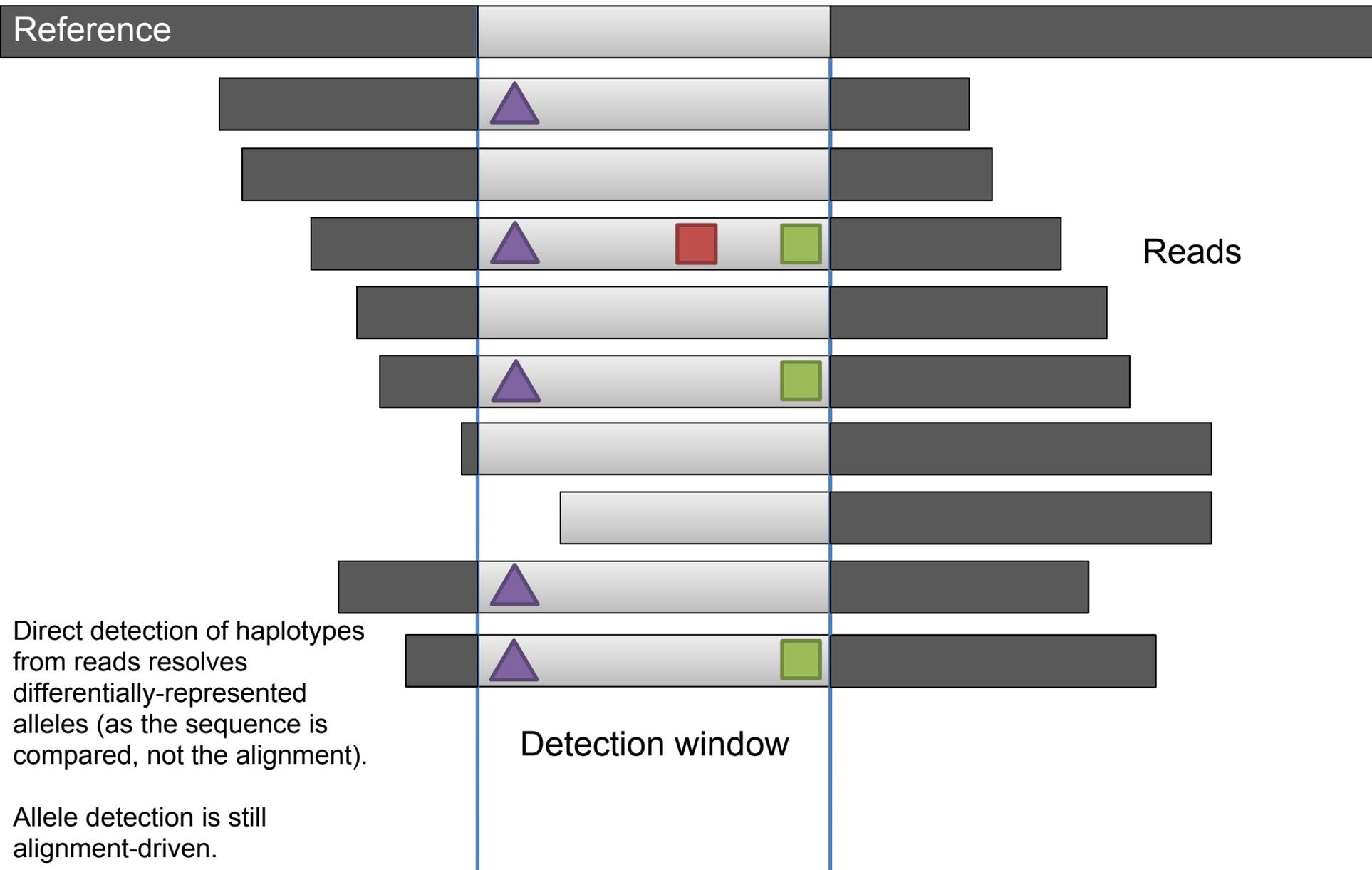
Another read
showing a SNP
on the same
haplotype as the
first

AGAACCCAGTGCTCTATCTGCT

AGAACCCAGTGCCTCTATCTGCT
AGAACCCAGTGGTCTTTCTGCT

A variant
locus implied
by
alignments

Direct detection of haplotypes



Reads

Observed Haplotypes

Ref	Variant Region	Variant Region
TACCGAT	CATTGGATCA	CGATTCC...GCATTGC
TACCGAT	CATTGGATCA	CGATTCC...GCATTGC
ACCGAT	TATTGCATCG	CGATTCC...GCATTGC
ACCGAT	CATTGGATCA	CGATTCC...GCATTGC
ACCGAT	TATTGGATCG	CGATTCC...GCATTGC
CCGAT	C-TTGGATCA	CGATTCC...GCATTGC
CCGAT	CAT G GGATCA	CGATTCC...GCATTGC

CATTGGATCA x8

TATTGGATCG x9

CTTGGATCA x1

CAT**G**GGATCA x1

...

(A)₇ x10

(A)₆ x7

(A)₅ x1

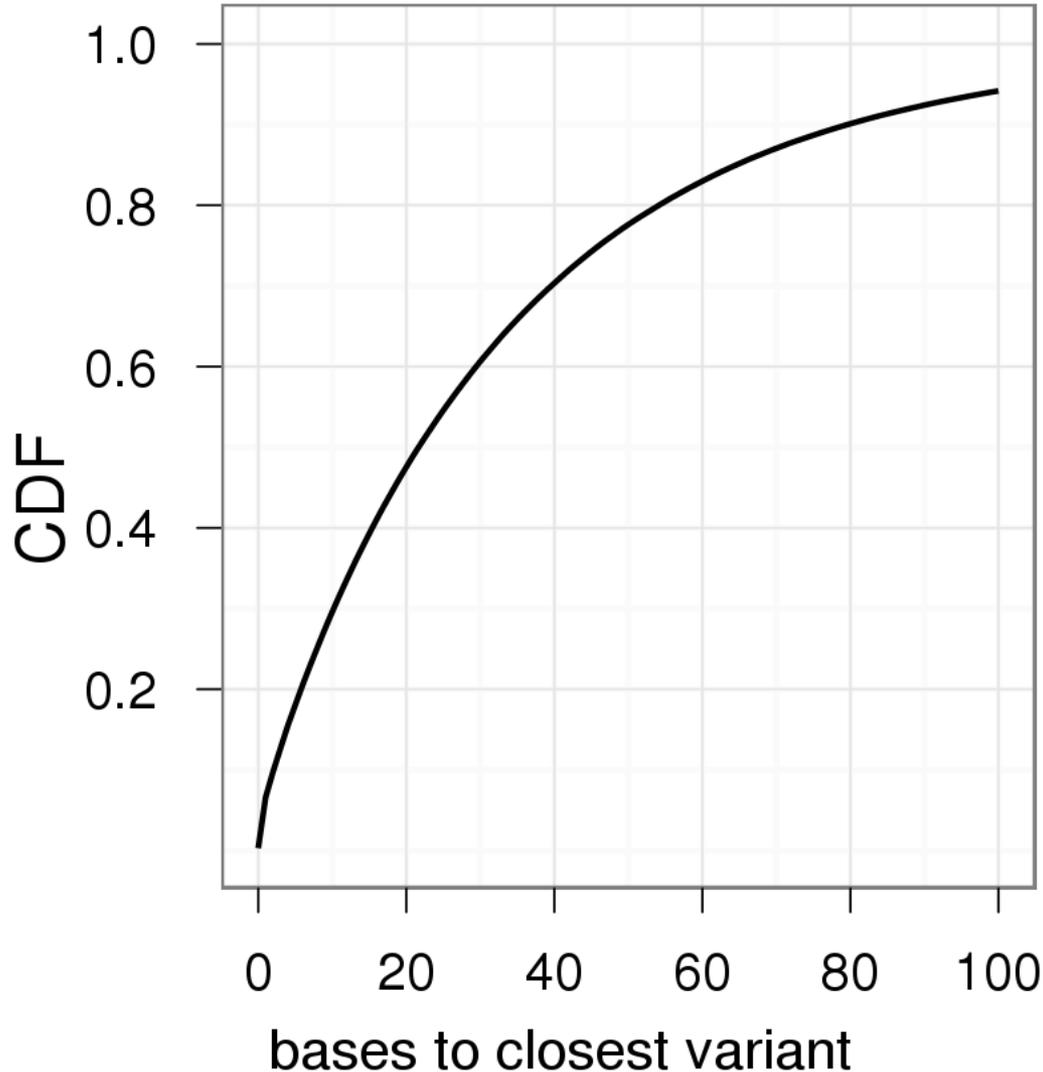
(A)₈ x1

...

Why haplotypes?

- Variants cluster.
- This has functional significance.
- Observing haplotypes lets us be more certain of the local structure of the genome.
- We can improve the detection process itself by using haplotypes rather than point mutations.
- We get the sensitivity of alignment-based approaches with the specificity of assembly-based ones.

Sequence variants cluster



In ~ 1000 individuals, $\frac{1}{2}$ of variants are within ~ 22 bp of another variant.

Variance to mean ratio (VMR) = 1.4.

The functional effect of variants depends on other nearby variants on the same haplotype

reference: AGG GAG CTG
 Arg Glu Leu

OTOF gene – mutations
cause profound recessive
deafness

apparent: AGG TAG CTG
 Arg Ter ---

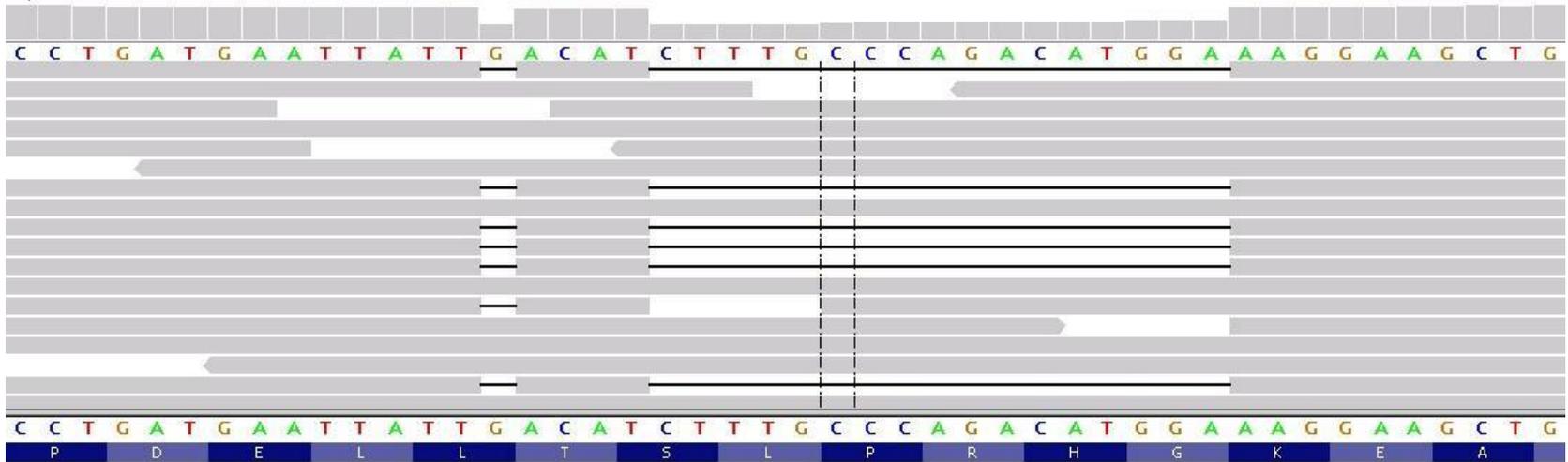
Apparent nonsense variant,
one YRI homozygote

actual: AGG TTG CTG
 Arg Leu Leu

Actually a block substitution
that results in a missense
substitution

(Daniel MacArthur)

Importance of haplotype effects: frame-restoring indels



(in NA12878)

- Two apparent frameshift deletions in the *CASP8AP2* gene (one 17 bp, one 1 bp) on the same haplotype
- Overall effect is in-frame deletion of six amino acids

(Daniel MacArthur)

Frame-restoring indels in 1000 Genomes Phase I exomes

chr6:117113761, GPRC6A (~10% AF in 1000G)

ref: A T T G T A A T T C T C A - - T A - - T T - - T G C C T T T G A A A G C

alt: A T T G T A A T T C T C A G G T A A T T T C C T G C C T T T G A A A G C

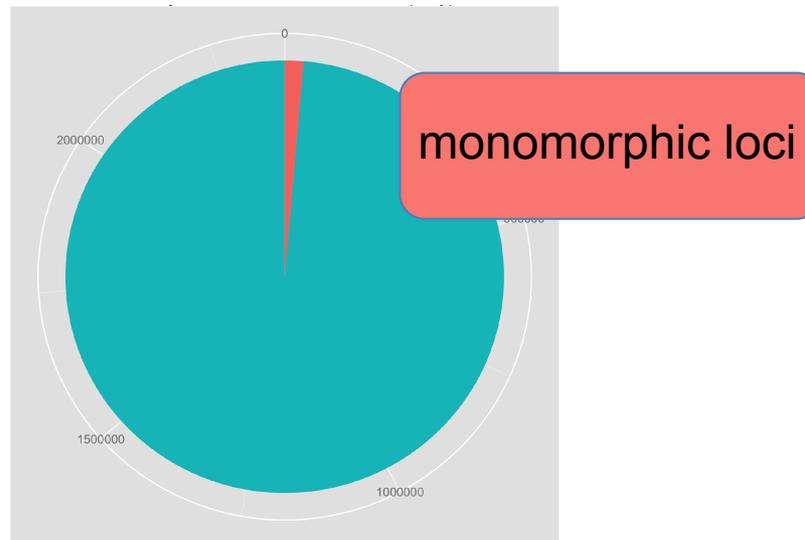
chr6:32551935, HLA-DRB1 (~11% AF in 1000G)

ref: C C A C C G C G G C C C G C G C C T G - C - T C C A G G A T G T C C

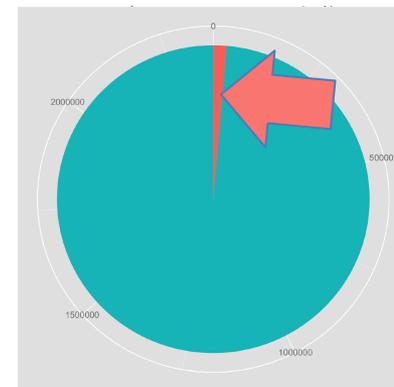
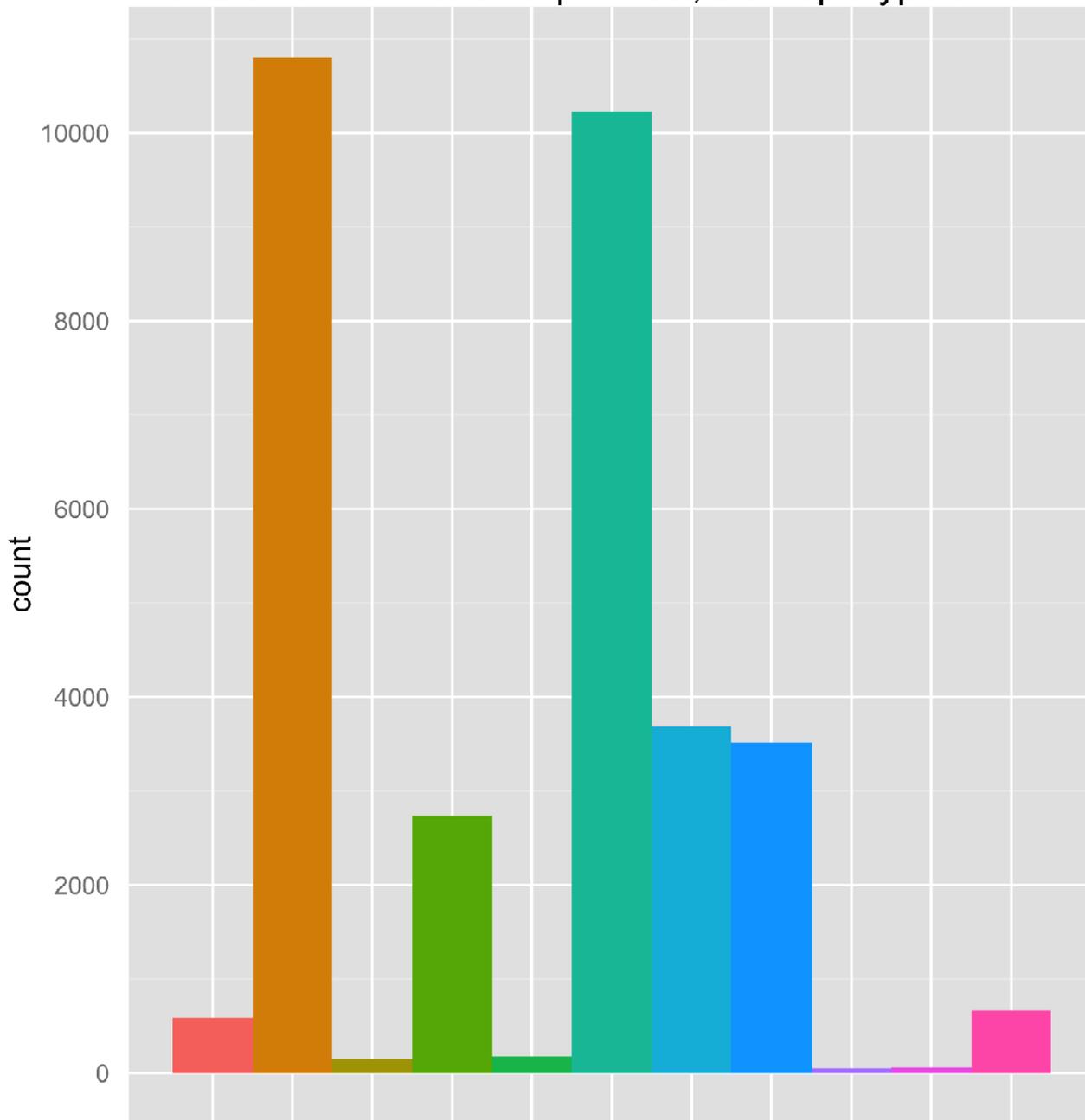
Alt: C C A C C G C G G - - C G C G C C T G T C T T C C A G G A G G T C C

Impact on genotyping chip design

- Biallelic SNPs detected during the 1000 Genomes Pilot project were used to design a genotyping microarray (Omni 2.5).
- When the 1000 Genomes samples were genotyped using the chip, 100k of the 2.5 million loci showed no polymorphism (monomorphs).



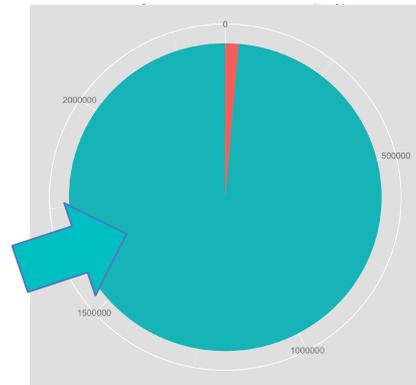
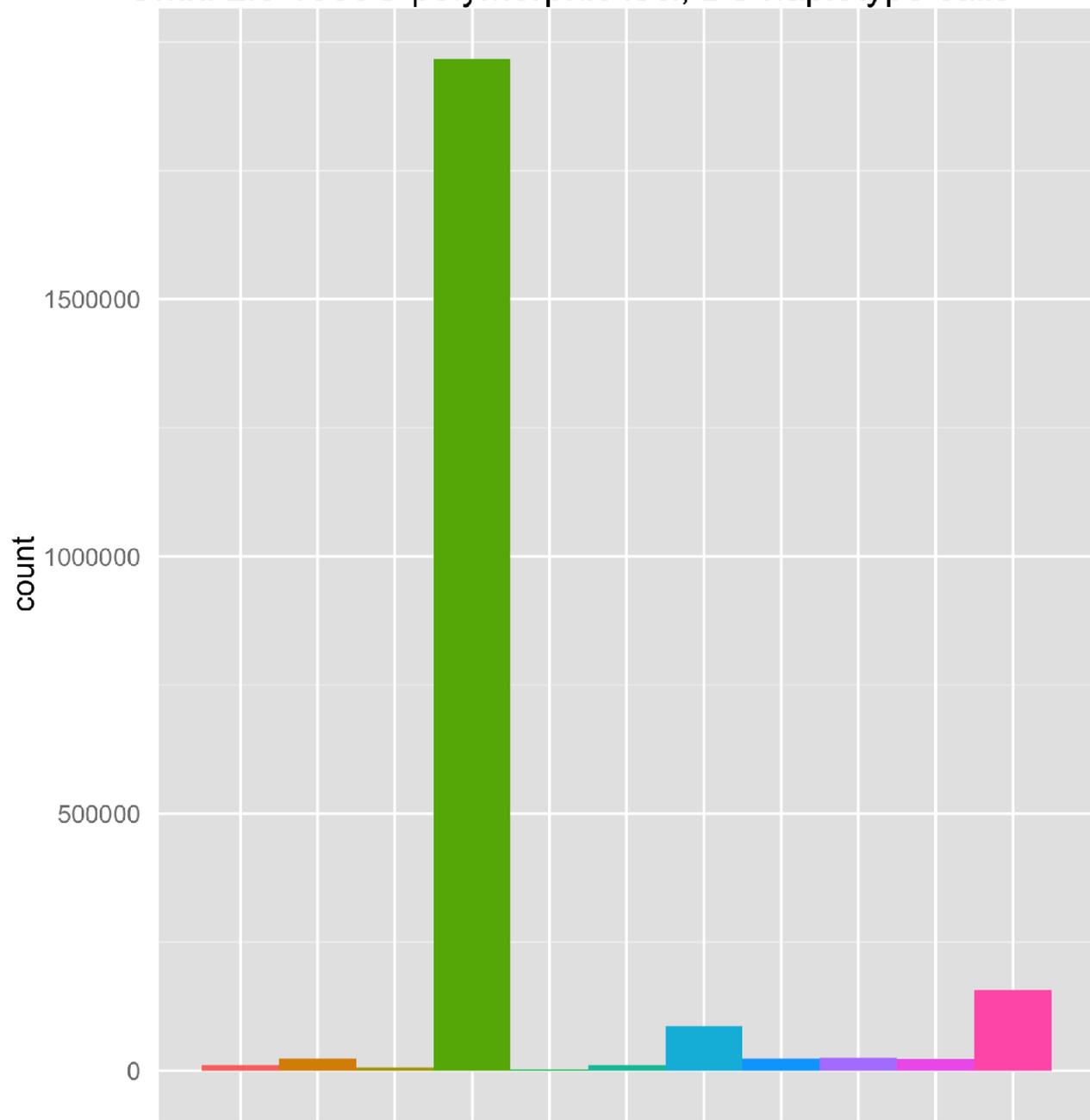
Omni 2.5 1000G monomorphic loci, BC haplotype calls



CLASS

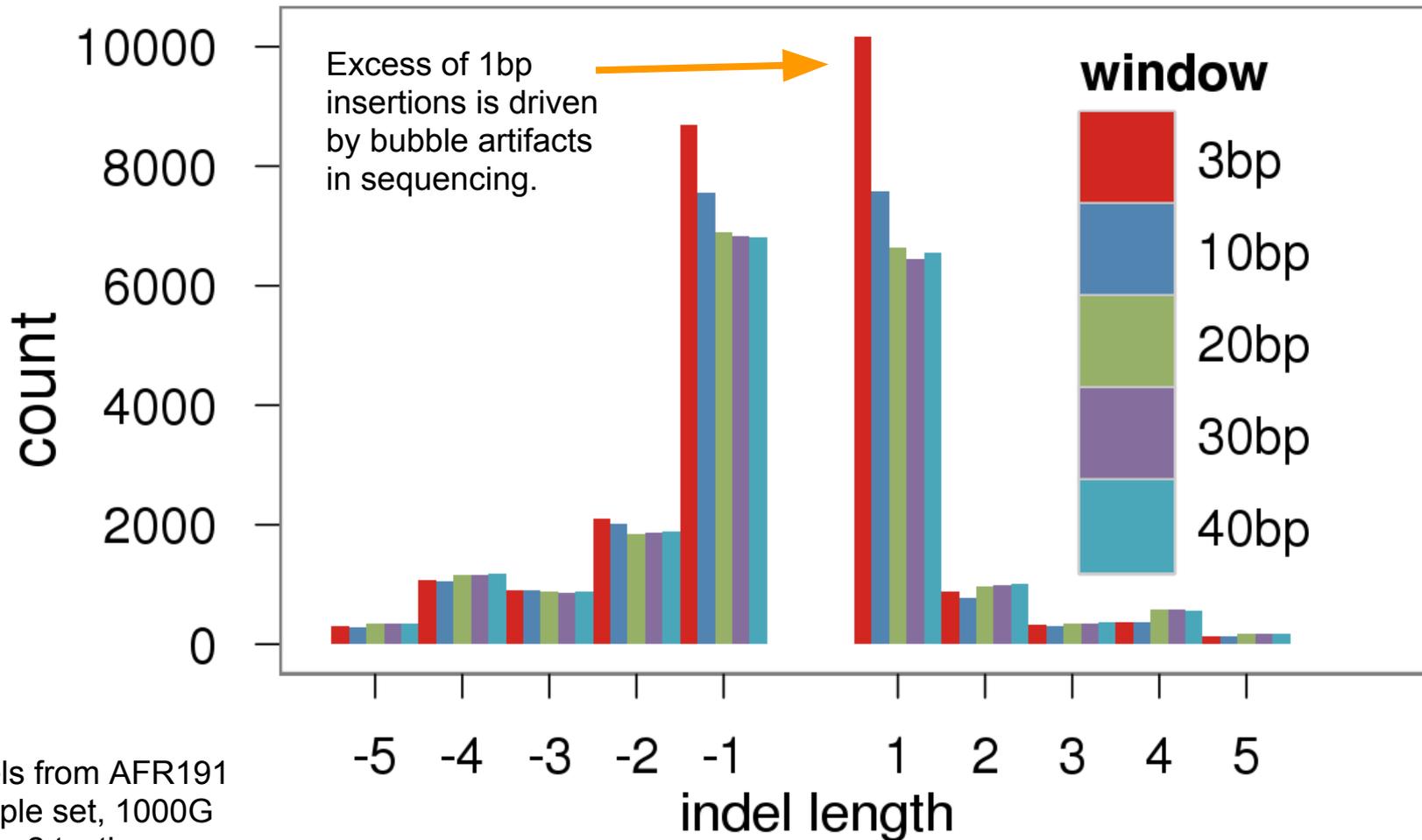
- biallelic complex
- biallelic INDEL
- biallelic MNP
- biallelic SNP
- multiallelic complex
- multiallelic INDEL
- multiallelic INDEL, SNP, and MNP
- multiallelic mixed
- multiallelic SNP
- multiallelic SNP and MNP
- multiallelic SNP, MNP, and complex

Omni 2.5 1000G polymorphic loci, BC haplotype calls



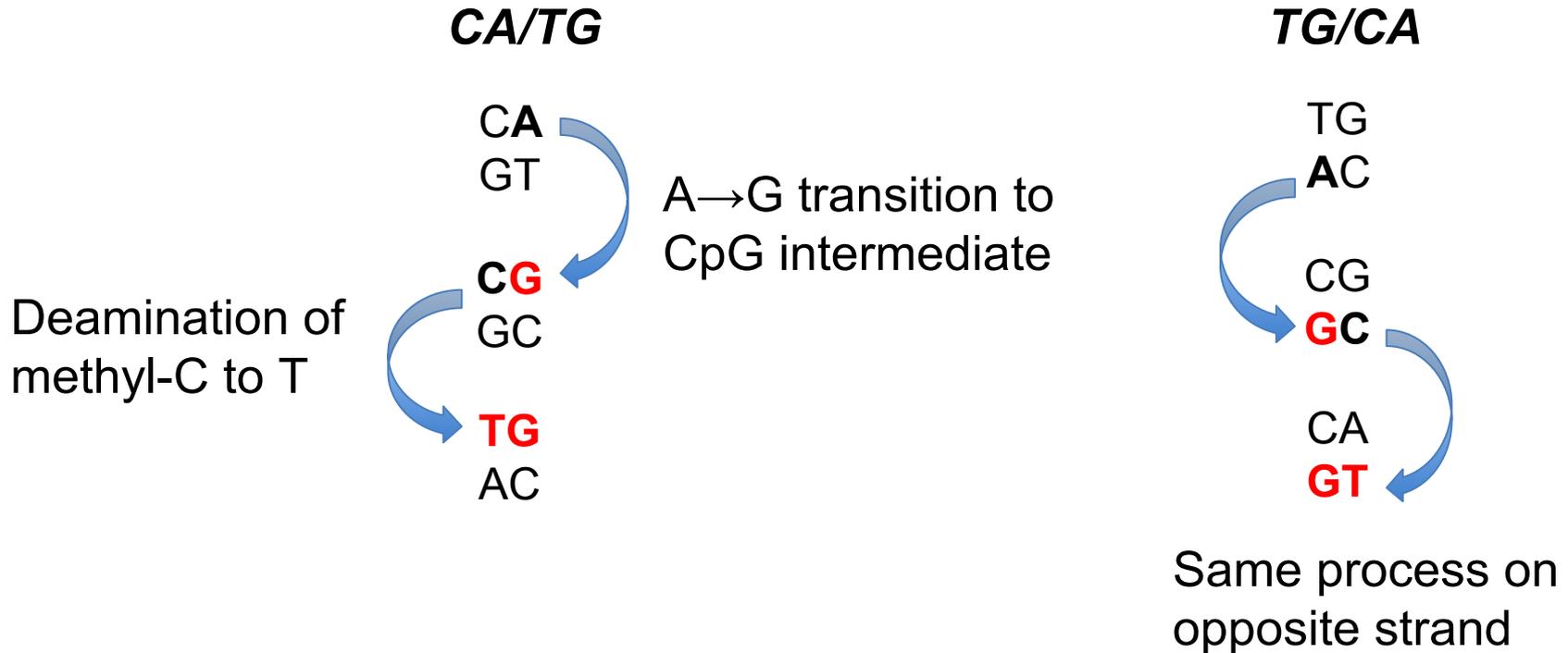
- CLASS**
- biallelic complex
 - biallelic INDEL
 - biallelic MNP
 - biallelic SNP
 - multiallelic complex
 - multiallelic INDEL
 - multiallelic INDEL, SNP, and MNP
 - multiallelic mixed
 - multiallelic SNP
 - multiallelic SNP and MNP
 - multiallelic SNP, MNP, and complex

Measuring haplotypes improves specificity



Indels from AFR191 sample set, 1000G phase2 testing.

Direct detection of haplotypes can remove directional bias associated with alignment-based detection



Overview

1. Genesis of variation (SNPs and indels)
2. Causes of sequencing error
3. Sequence alignment
4. Alignment-based variant detection
5. Assembly-based variant detection
6. Haplotype-based variant detection
7. **Graph-based methods**

Graph-based methods

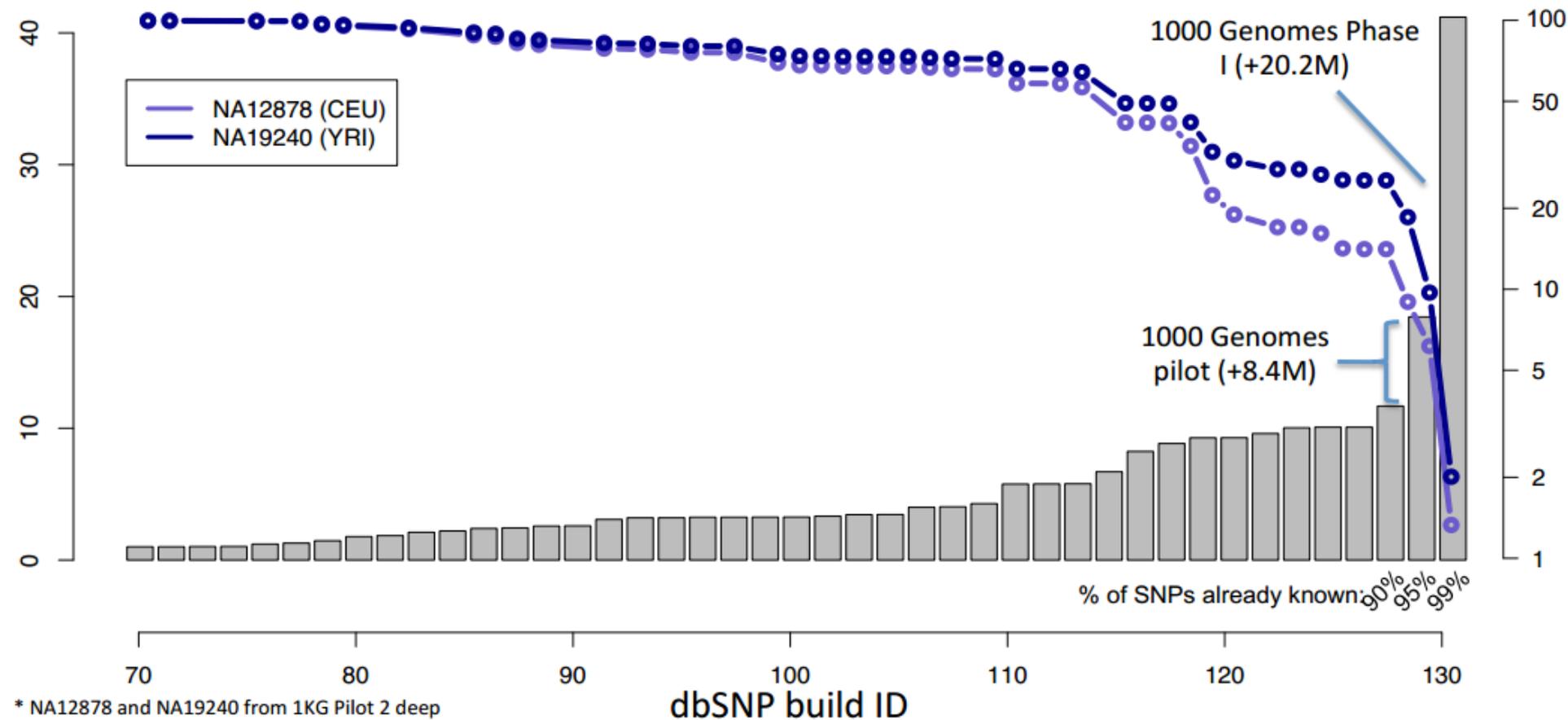
1. Variation + sequence = variation graphs
2. Local realignment to a variation graph
3. Constructing a whole genome variation graph
4. Aligning to a whole genome variation graph

**(1) Variation +
sequence
= variation graph**

Most SNPs (>99%) are now known

Number of SNPs in dbSNP (Millions)

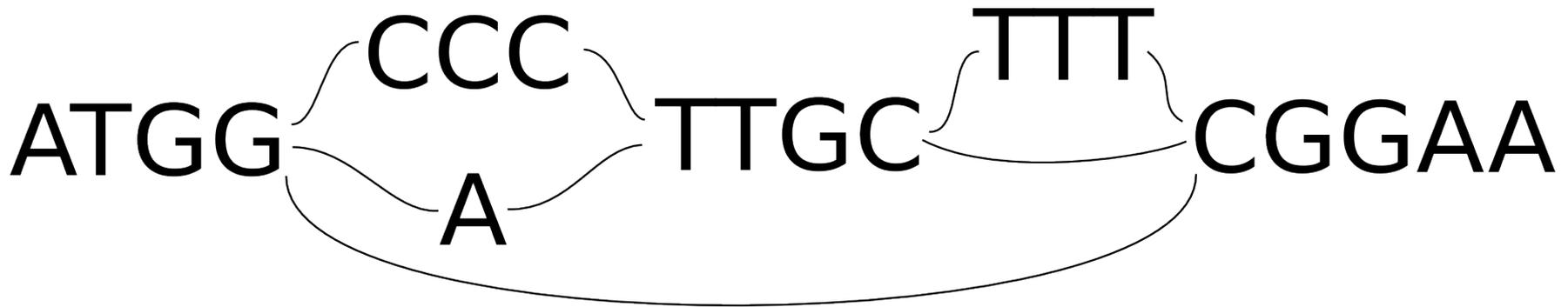
% of novel SNPs discovered in whole-genome sequencing



* NA12878 and NA19240 from 1KG Pilot 2 deep whole genome sequencing

Variation graphs

A variation graph represents many genomes in the same context.



Nodes contain sequence and directed edges represent potential links between successive sequences.

A multiple sequence alignment is a variation graph

traditional MSA

(a) . . P K M I V R P Q K N E T V .
T H . K M L V R . . . N E T I M

consensus sequence

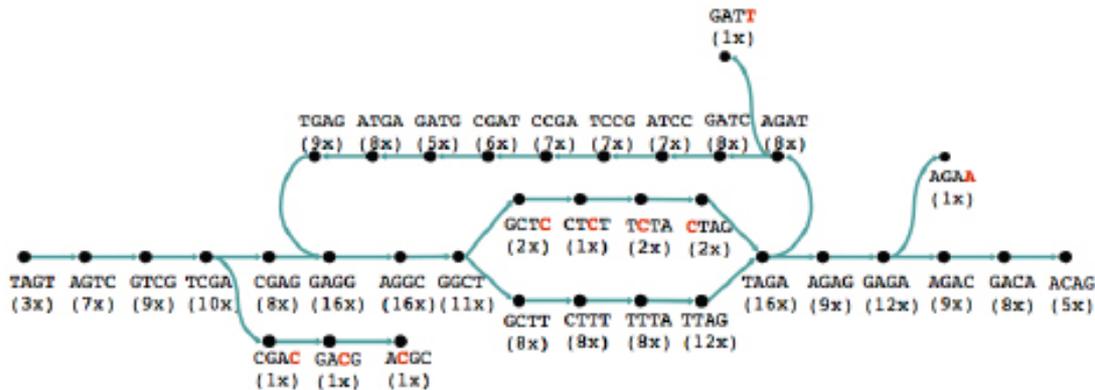
(b) (P) → (K) → (M) → (I) → (V) → (R) → (P) → (Q) → (K) → (N) → (E) → (T) → (V)

positionally-matching regions aligned

(c)

(d)

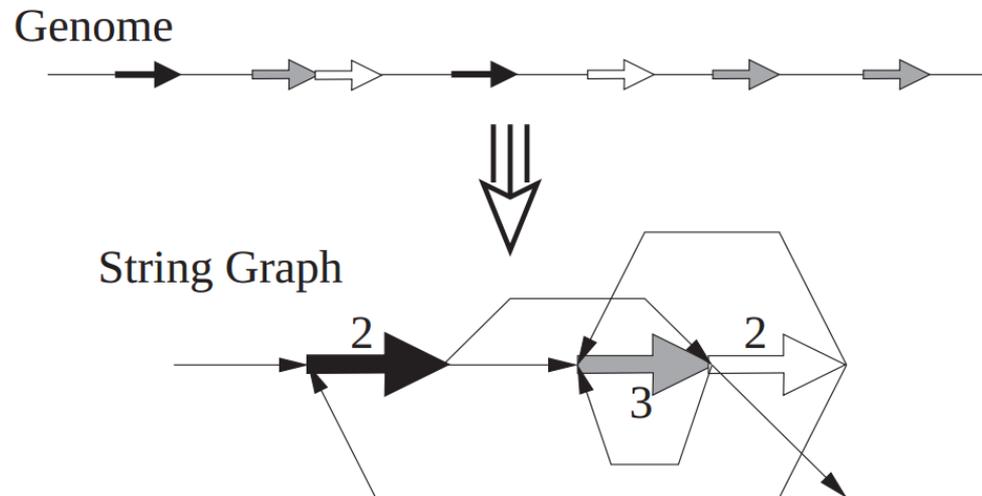
Assembly graphs are variation graphs



<http://plus.maths.org/content/os/issue55/features/sequencing/index>,
credit Daniel Zerbino

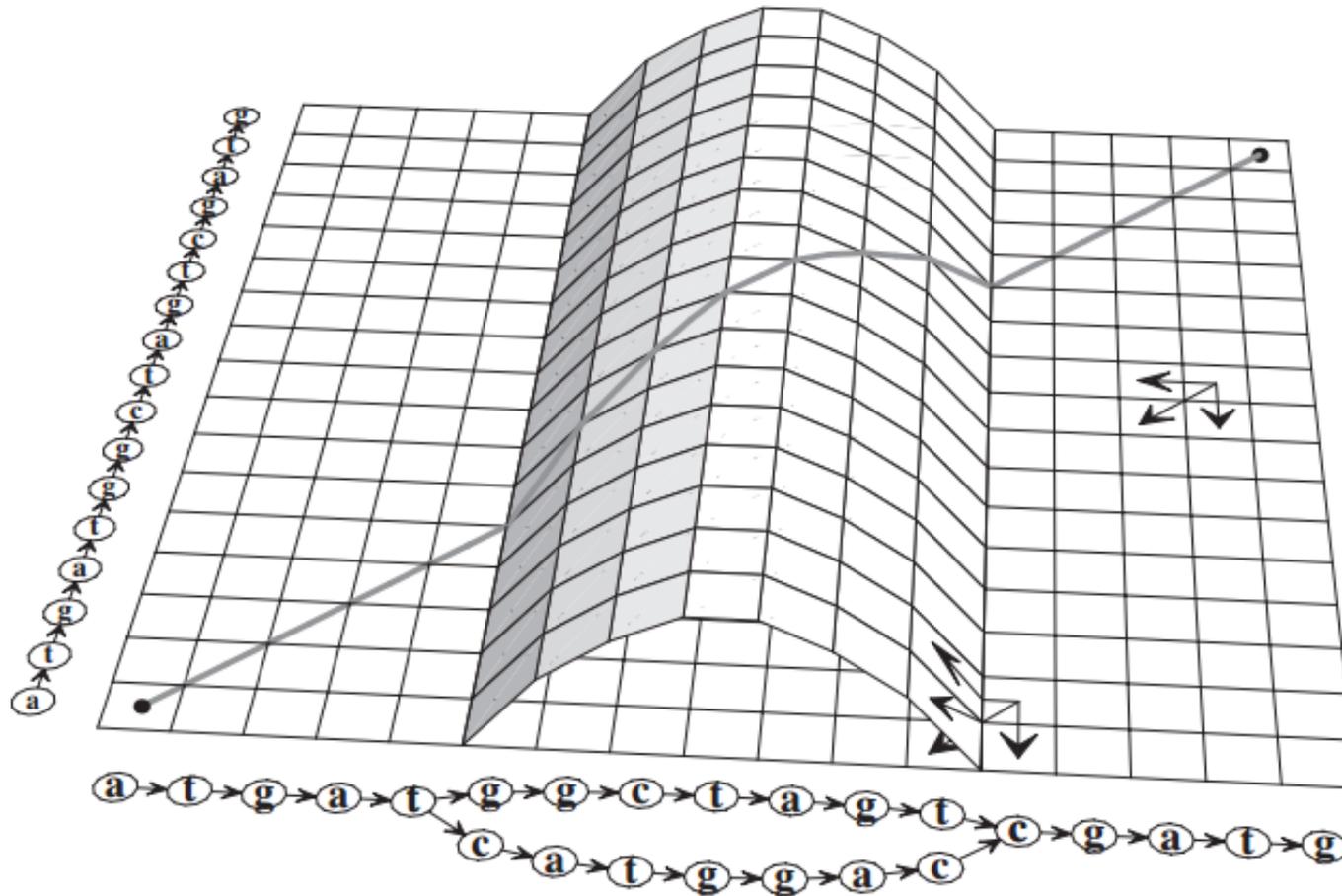
A string graph follows the same format as the variation graph.

A de Bruijn graph can be converted into a variation graph (as previous) by setting node sequences to the first letter in the kmer and compressing non-branching runs into single nodes.

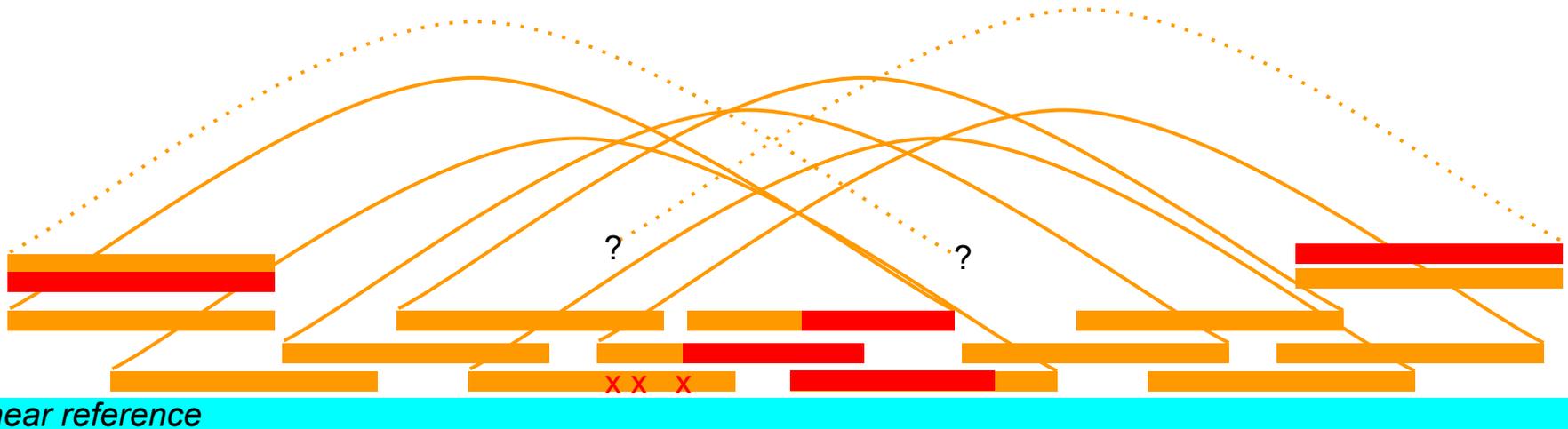


(2) Local realignment to a variation graph

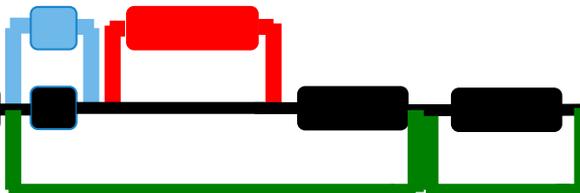
Local alignment against a graph



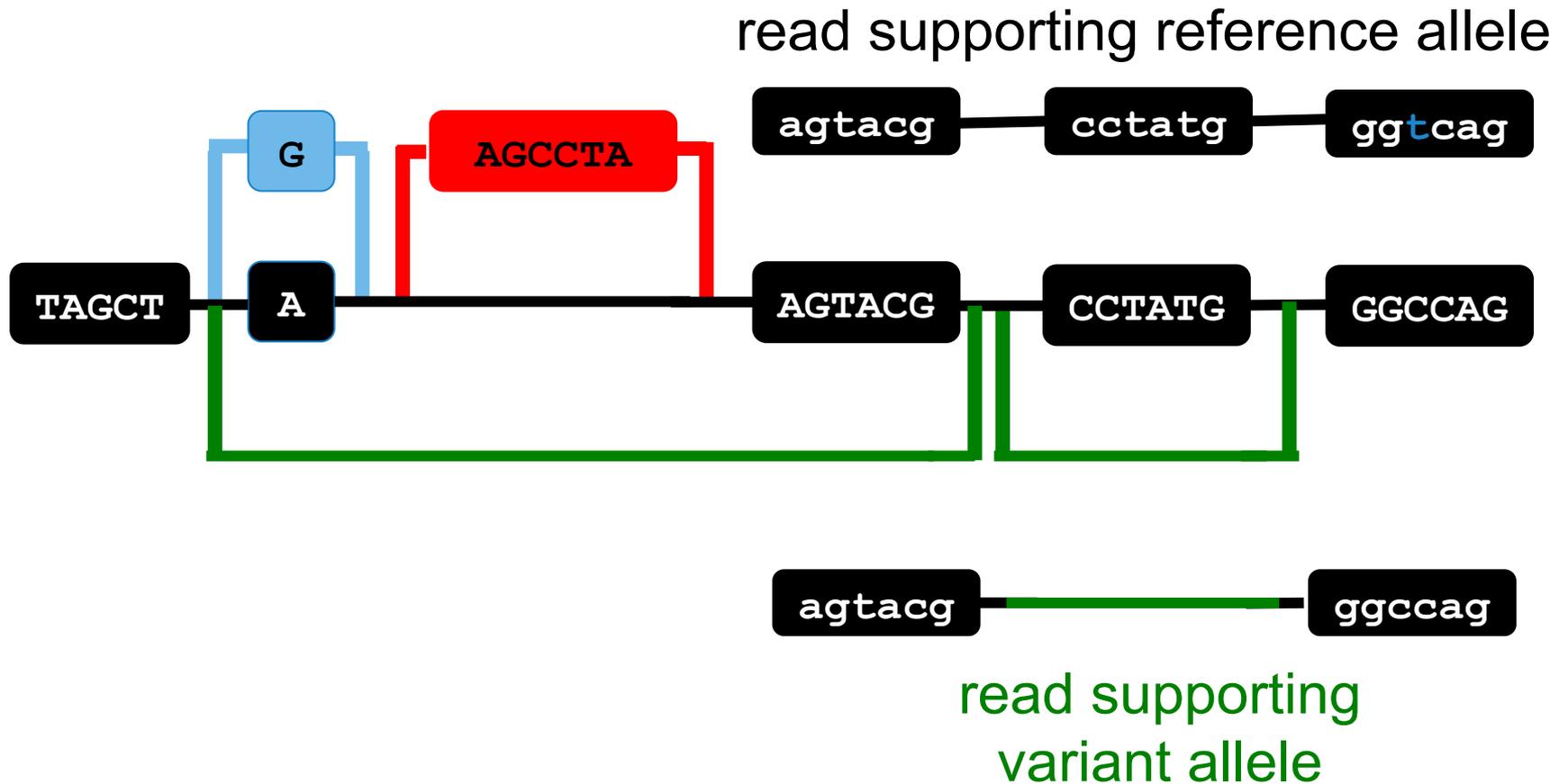
Seeding graph-based alignments



Test imperfectly-mapped reads against graph.



Detecting variation on the graph



TATTTGGGTTACAGTTTTTTGACTATTACATGTAAAGCCAAAAAACTGTAGGATAAATTCTC

ACCCTTGAAGAA

TAGGATAAATG

TATTTGGGTTACAGTTTTTTGACTATTACATGTAAAGCCAAAAAACTGTAGGATAAATTCTC



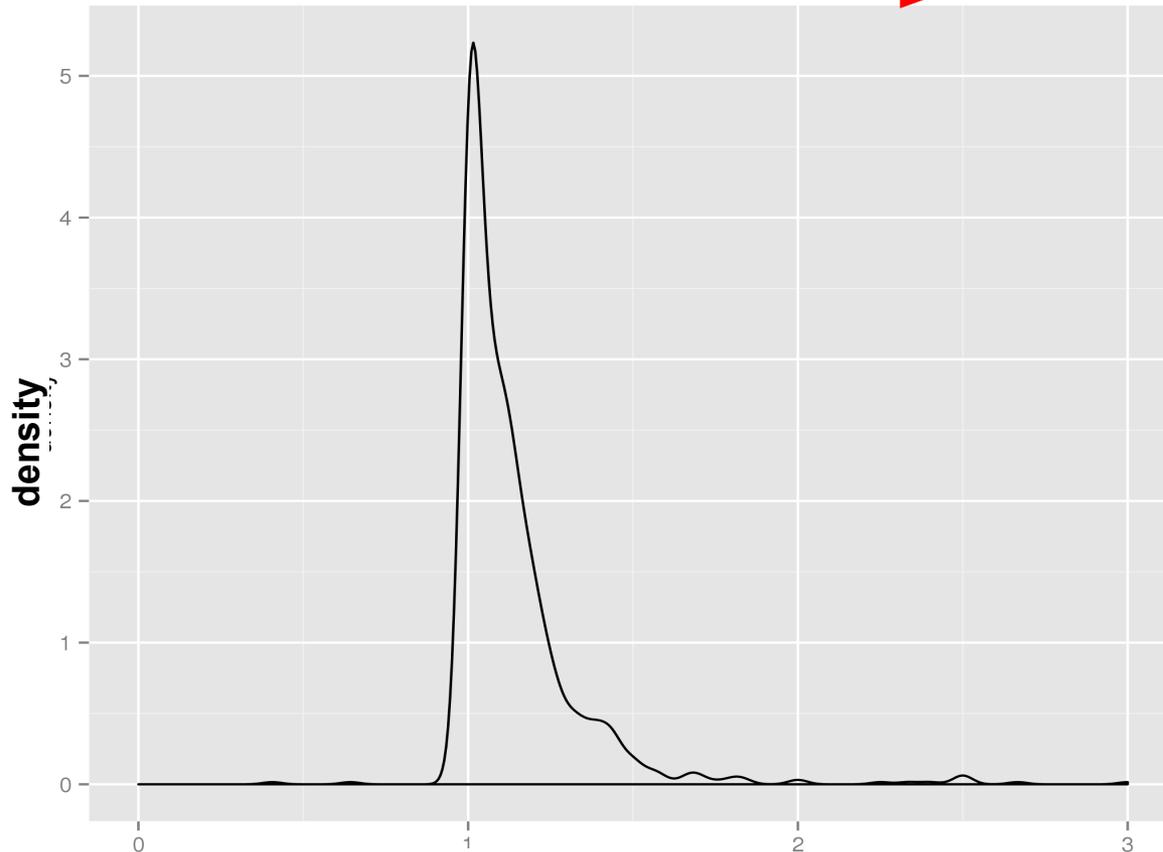
raw alignments

realigned to variation graph

AATATTTGGGTTACAGTTTTTTGACTATTACATGTAAAGCCAAAAAACTGTAGGATAAATTCTCTAGC

realignment to variation graph reduces reference bias

Improvement in observation support



Standard alignment is frustrated even by small variants

(tested against 1mb segment on chr20 using 1000Gp3 union allele list)

Ratio between observations with and without realignment to graph of union variants

(3) Constructing a whole genome variation graph

Objective

Construct a whole genome variation graph to serve as our reference for sequence analysis.

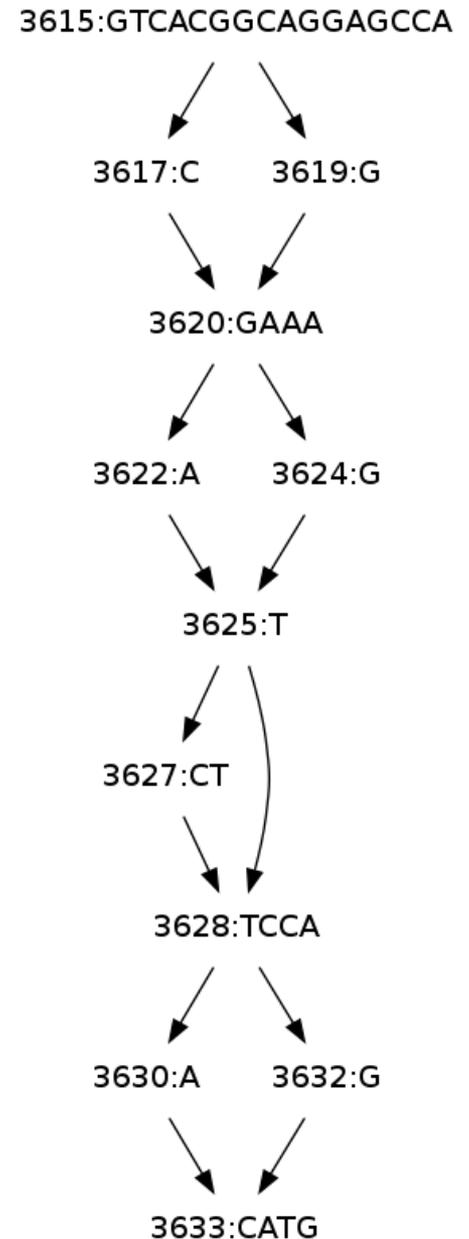
We'll need:

- (1) a data model
- (2) a serialization format
- (3) an API (construct, align, view, ...)

Conceptual model

We have sequences (nodes), and allowed linkages between them (edges).

Paths (walks through the graph) are also useful.



Data model

Use protocol buffers to define objects

Graph, **Node**, **Edge**, **Path**, Alignment, Edit

```
message Graph {  
  repeated Node node = 1;  
  repeated Edge edge = 2;  
  repeated Path path = 3;  
}
```

```
message Path {  
  optional int32 target_position = 1;  
  optional string name = 2;  
  repeated Mapping mapping = 3;  
}
```

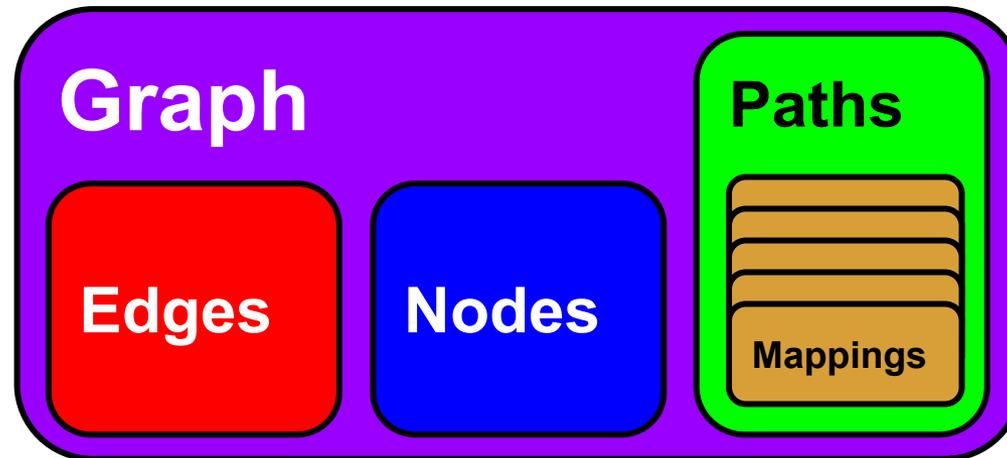
```
message Edge {  
  required int64 from = 1;  
  required int64 to = 2;  
  optional bytes data = 3;  
}
```

```
message Node {  
  optional string sequence = 1;  
  optional string name = 2;  
  required int64 id = 3;  
  optional double quality = 4;  
  optional int32 coverage = 5;  
  optional bytes data = 6;  
}
```

```
message Mapping {  
  required int64 node_id =  
  1;  
  repeated Edit edit = 2;  
}
```

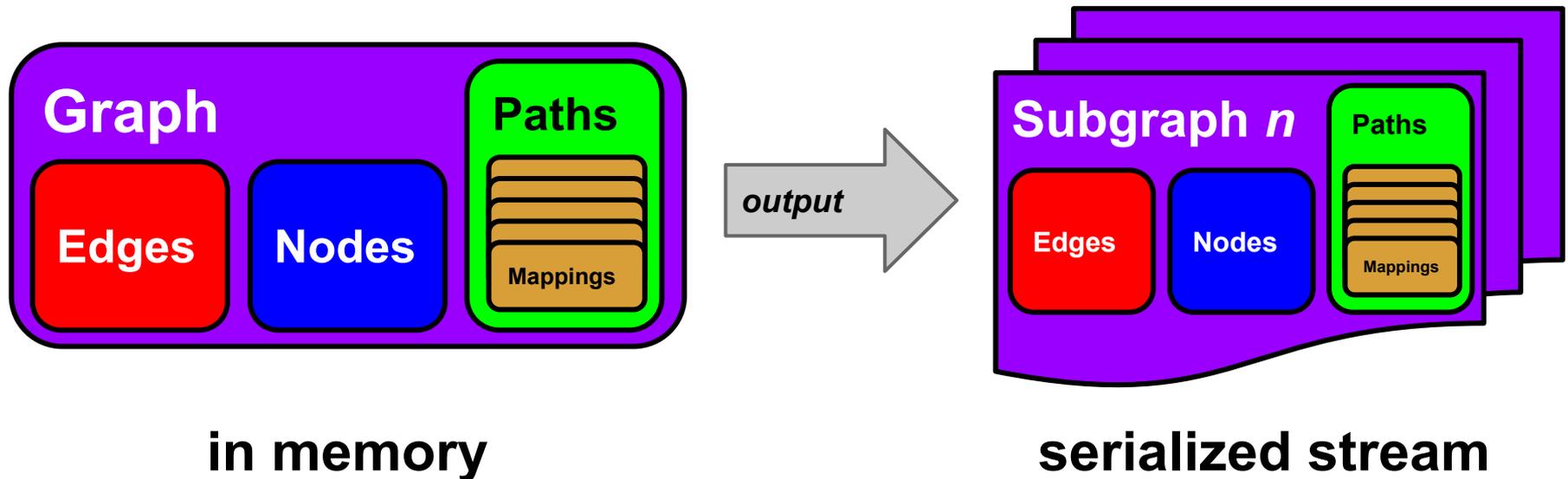
Format

Basic entity is a Graph, which is composed of nodes, edges, and paths.



Serialization

To serialize the graph, we generate a stream of sub-graphs that can be reassembled into the whole.



POS	ID	REF	ALT
...			

Construction

I:TGGGAGAGAACTGGAACAAGAACCCAGTGCTCTTTCTGCTCTA

For each variant

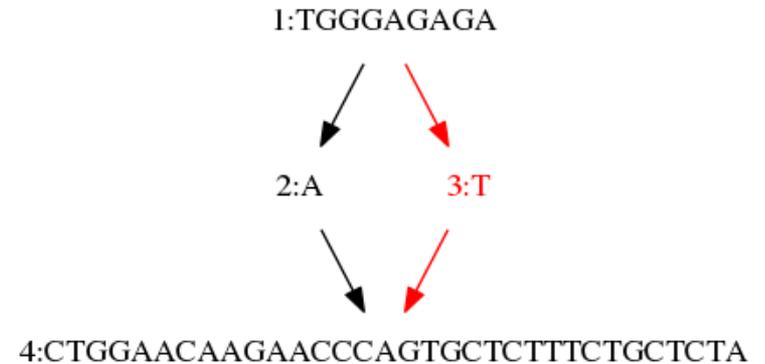
1. cut the reference path around the variant
2. add the novel (ALT) sequence to the graph

Construction

POS	ID	REF	ALT
10	.	A	T
...			

For each variant

1. cut the reference path around the variant
2. add the novel (ALT) sequence to the graph

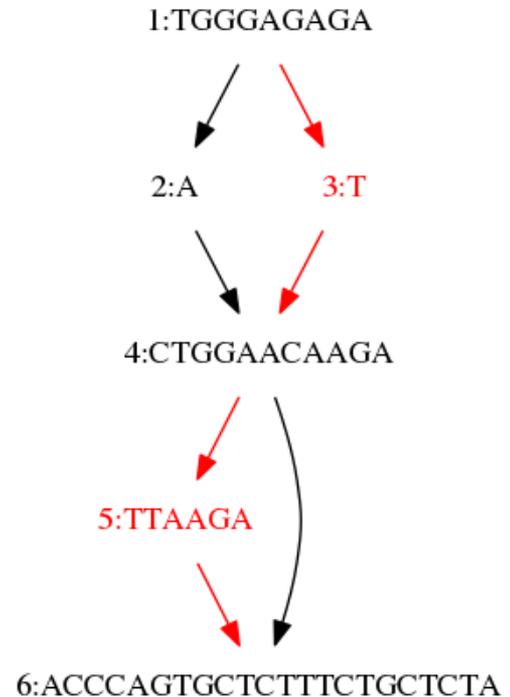


Construction

POS	ID	REF	ALT
10	.	A	T
21	.	A	ATTAAGA
...			

For each variant

1. cut the reference path around the variant
2. add the novel (ALT) sequence to the graph

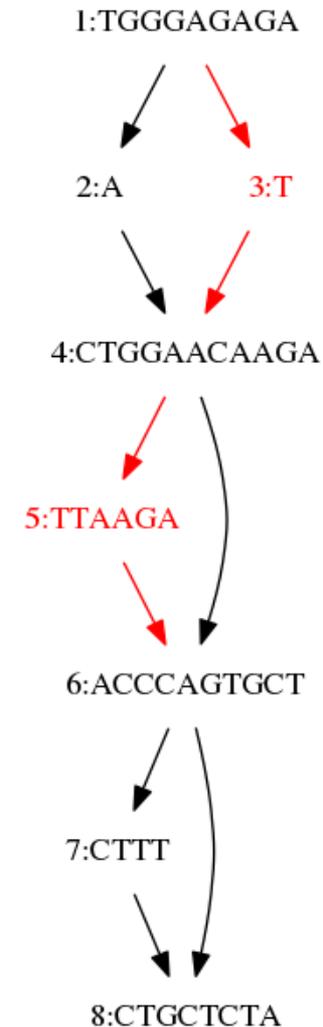


Construction

POS	ID	REF	ALT
10	.	A	T
21	.	A	ATTAAGA
31	.	TCTTT	T

For each variant

1. cut the reference path around the variant
2. add the novel (ALT) sequence to the graph



Indexing

Instead of indexing the serialized graph representation (e.g. BAM, VCF), write the graph into a disk-backed key/value store:

```
{"key":"+g+26+n", "value":{"id": 26, "sequence": "TATTTGAAGT"}}
```

```
{"key":"+g+26+p+1+103", "value":{"node_id": 26, "edit": []}}
```

```
{"key":"+g+26+t+24", "value":{"from": 24, "to": 26}}
```

```
{"key":"+g+26+t+25", "value":{"from": 25, "to": 26}}
```

```
{"key":"+g+27+f+29", "value":{"from": 27, "to": 29}}
```

...

```
{"key":"+k+TCATACTACTG+69", "value":-2}
```

```
{"key":"+k+TCATACTACTG+70", "value":-4}
```

```
{"key":"+k+TCATACTACTG+72", "value":-5}
```

```
{"key":"+k+TCATATGTCCA+167", "value":29}
```

```
{"key":"+k+TCATATGTCCA+169", "value":-1}
```

```
{"key":"+k+TCATATGTCCA+171", "value":-2}
```

...

This allows our graphs and kmer indexes to have > main memory size. (Uses rocksdb.)

Command-line API

usage: `vg <command> [options]`

commands:

<code>-- construct</code>	graph construction
<code>-- view</code>	format conversions for graphs and alignments
<code>-- index</code>	index features of the graph in a disk-backed key/value store
<code>-- find</code>	use an index to find nodes, edges, kmers, or positions
<code>-- paths</code>	traverse paths in the graph
<code>-- align</code>	local alignment
<code>-- map</code>	global alignment
<code>-- stats</code>	metrics describing graph properties
<code>-- join</code>	combine graphs via a new head
<code>-- ids</code>	manipulate node ids
<code>-- concat</code>	concatenate graphs tail-to-head
<code>-- kmers</code>	enumerate kmers of the graph
<code>-- sim</code>	simulate reads from the graph
<code>-- mod</code>	filter and transform the graph
<code>-- subject</code>	map alignments onto specific paths

(4) Aligning to a whole genome variation graph

Constructing a whole human genome variation graph

Constructed 1000G phase3 + GRCh37 variation graph using `vg construct`.

5h20m on 32-core system @Sanger

3.07G on disk

3.181 Gbp of sequence in graph

286 million nodes (11 bp/node)

376 million edges (8.5 bp/edge)

Indexing 1000G+GRCh37

for each node:

for each k -path extending no more than n edges or k bp away from the node:

→ index the k -mers of the k -path

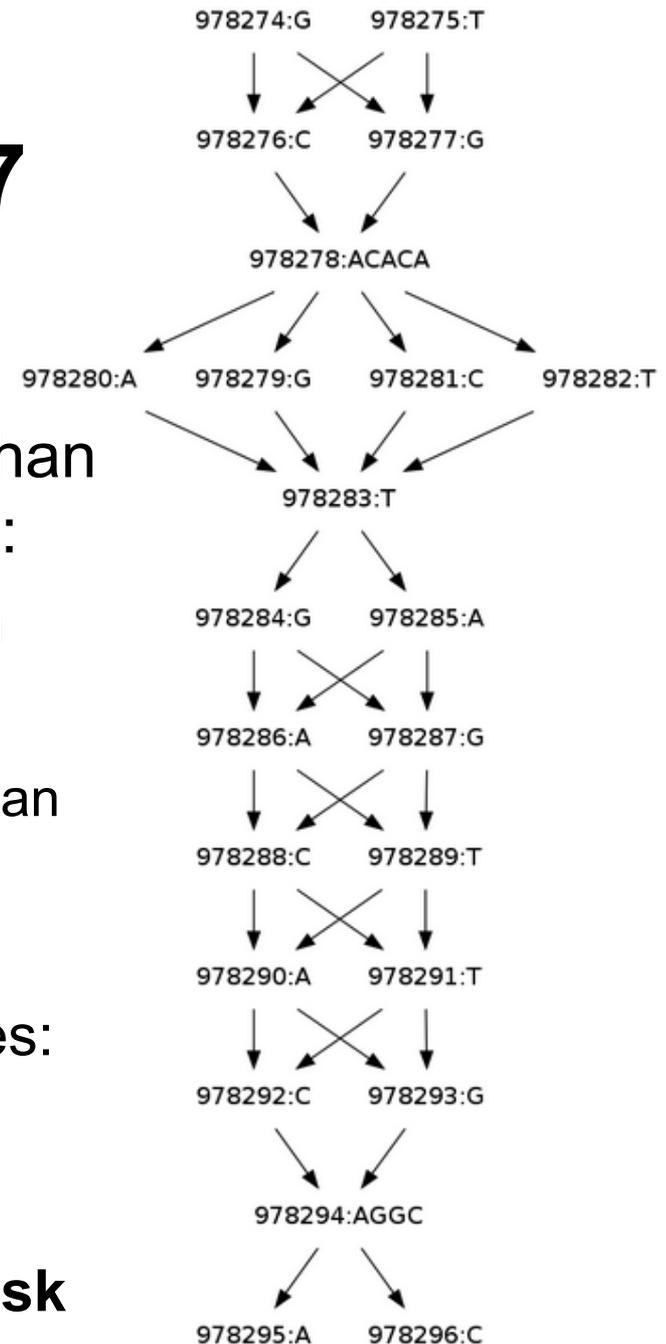
(Edge bounding limits the k -mer space, which can be very important in real data. See →)

for 27-mers crossing no more than 11 edges:

52h on a 32-core host @Sanger

Index is **175G** and on Lustre

Can align **70 read/s/CPU** with index on disk



Aligning to 1000G+GRCh37

For each k-mer in the read

If the k-mer is informative

→ sort hits by node id

For each cluster of hits : $\max(\text{id}) - \min(\text{id}) < M$

Get local region of graph from index

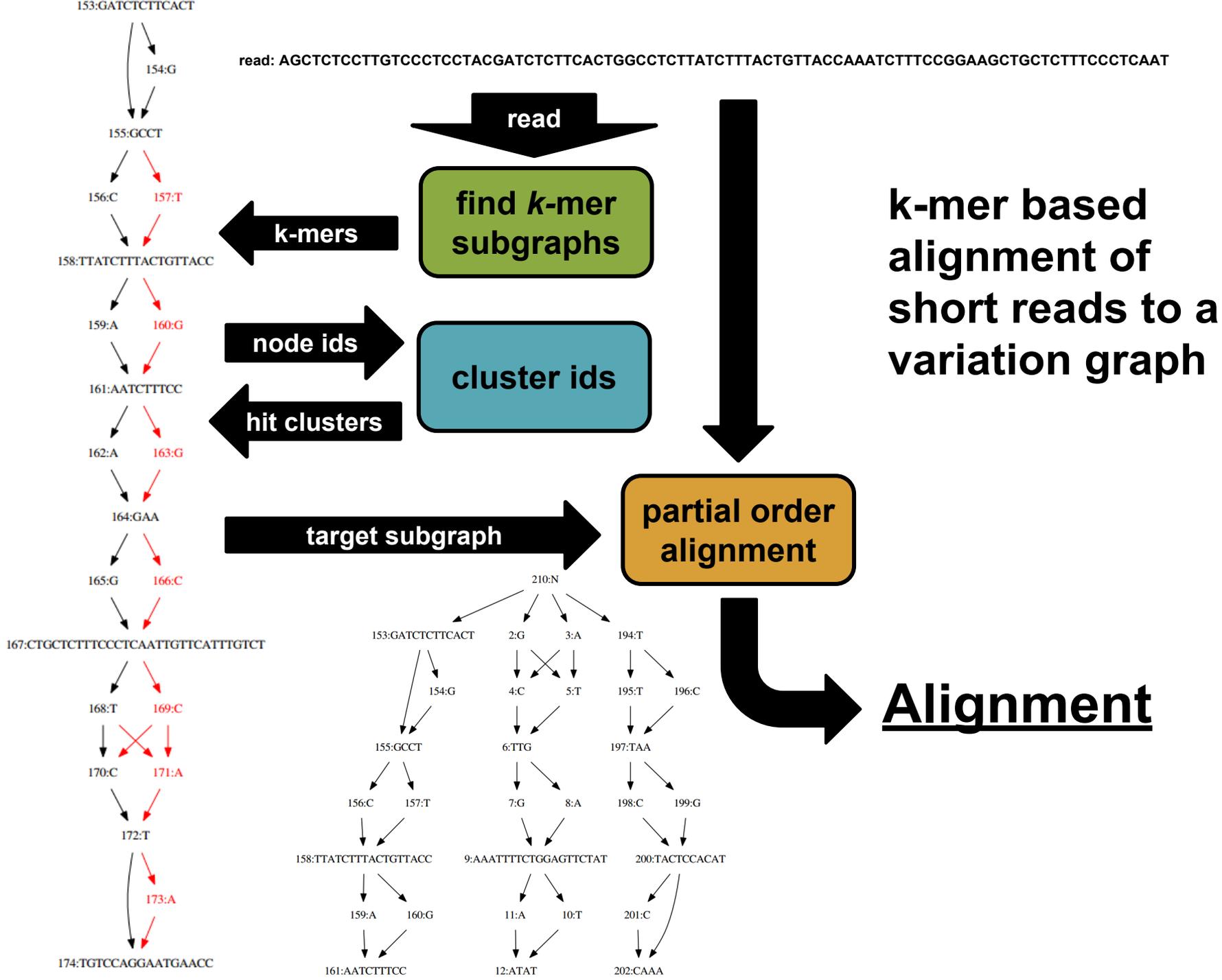
Join the subgraphs together

Align to the joined graph using POA → alignment

If alignment failed, decrease k-mer length or stride

8000 read/s on 32-core machine @Sanger (cached)

~ 40 hours for a deep 150x2 Illumina X10 run



Thanks!

Gabor Marth

Deniz Kural

Wan-Ping Lee

Alistair Ward

Richard Durbin

Josh Randall

Benedict Patton

Zamin Iqbal

Jerome Kelleher

Jared Simpson

David Haussler

Daniel Zerbino

Jouni Siren



EMBL-EBI

