# Variant calling and filtering for INDELs

Erik Garrison
SeqShop @ University of Michigan

# Overview

1. **Genesis of insertion/deletion (indel) polymorphism**
2. Standard approaches to detecting indels
3. Assembly-based indel detection
4. Haplotype-based indel detection
5. Primary filtering: Bayesian variant calling
6. Post-call filtering: SVM
7. Graph-based resequencing approaches

# An INDEL

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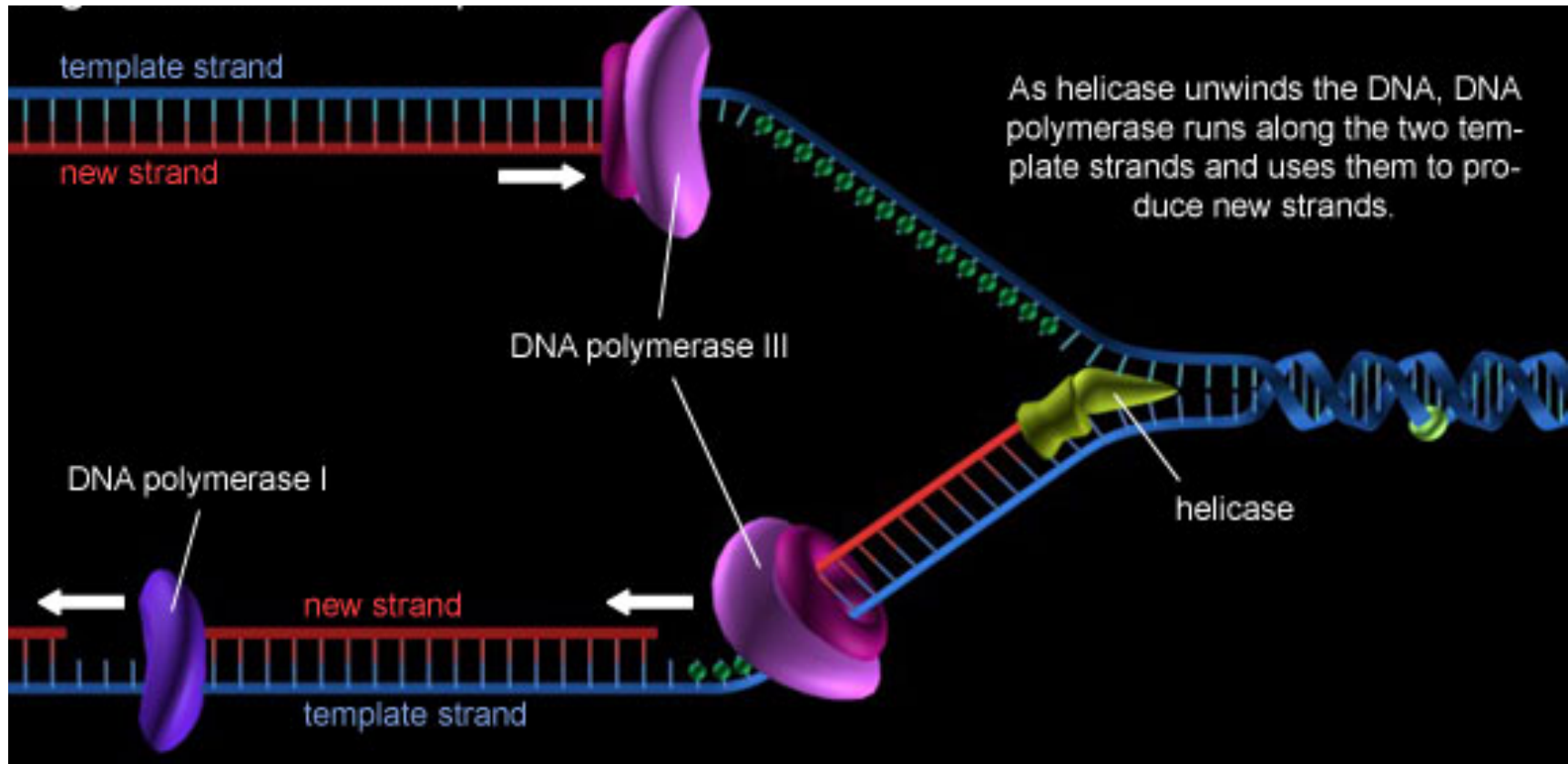
**AATTAGCCATTA**

**AATTA--CATTA**

# INDEL genesis

A number of processes are known to generate insertions and deletions in the process of DNA replication:
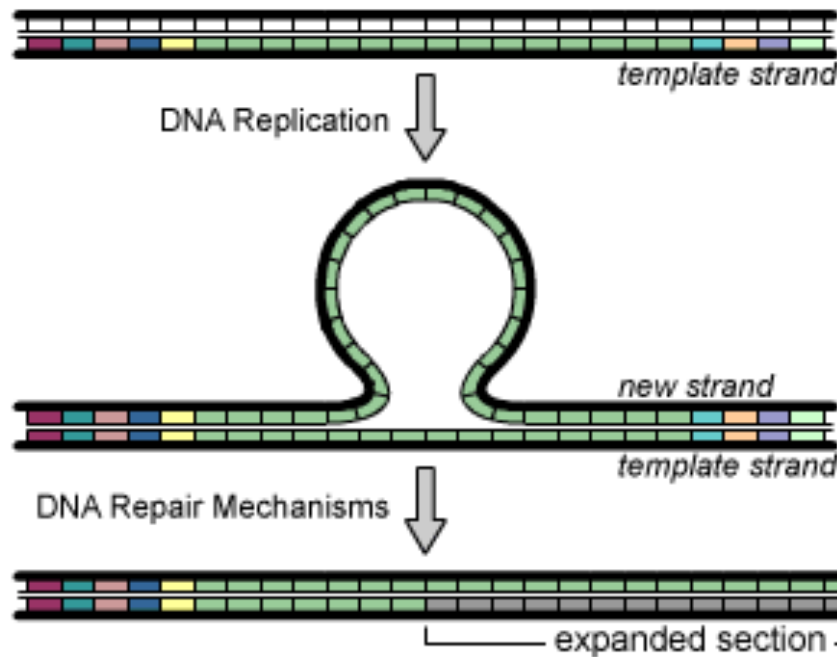
- Replication slippage
- Double-stranded break repair
- Structural variation (e.g. mobile element insertions, CNVs)

# DNA replication

# Polymerase *slippage*



A) Slippage Event

template strand

DNA Replication

new strand

template strand

DNA Repair Mechanisms

new strand
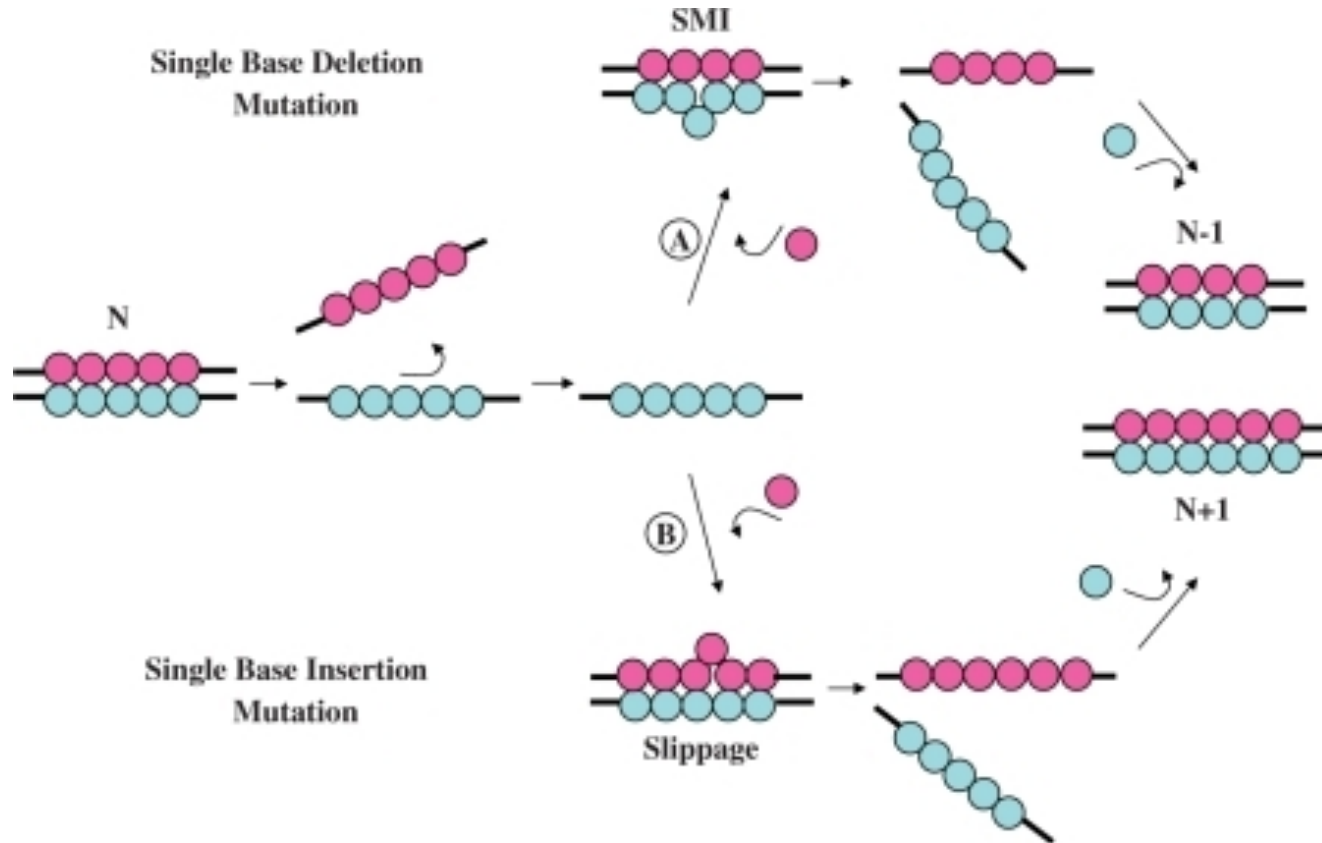
expanded section

template strand

B) No Slippage

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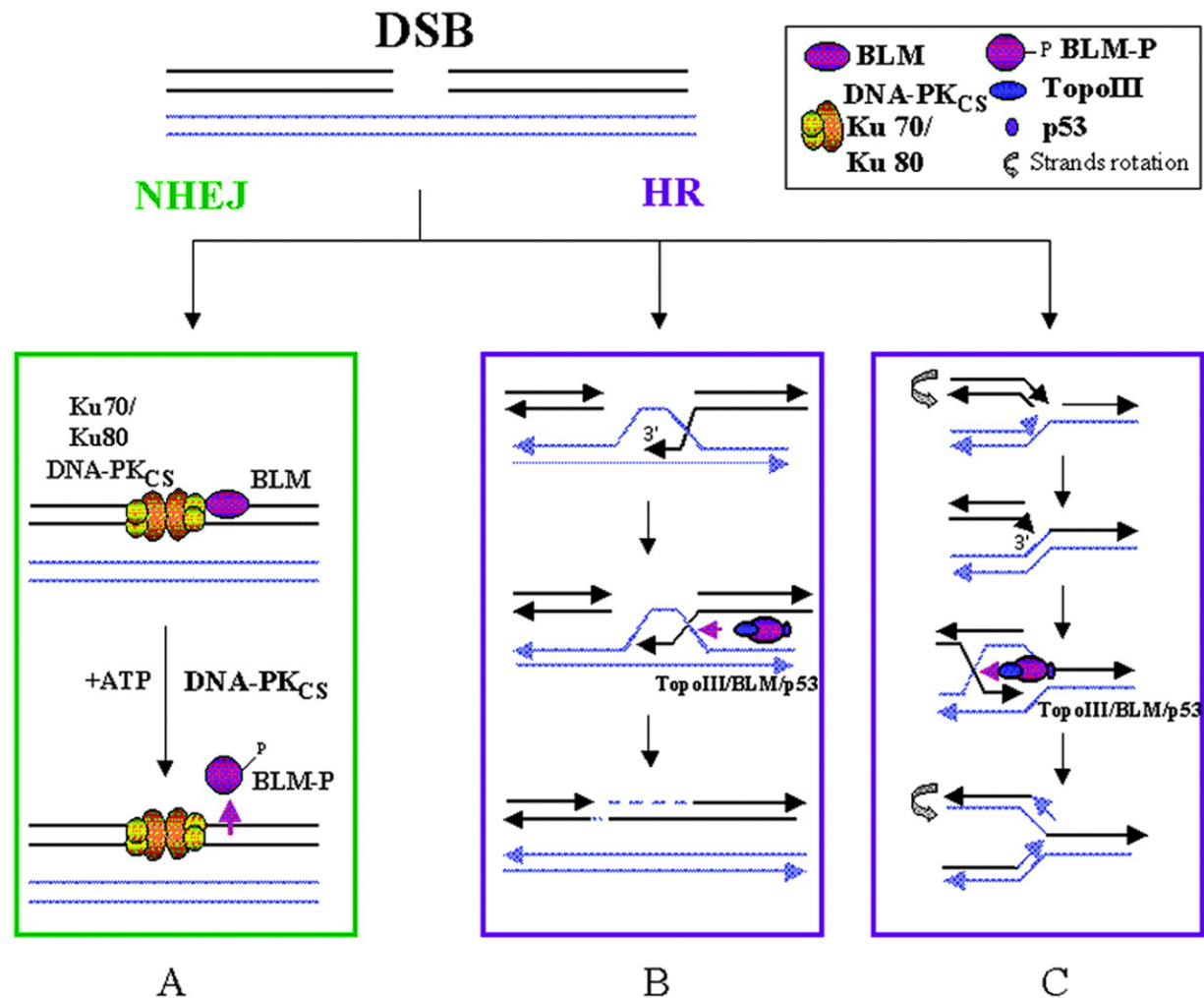
(B) Polymerase slippage, as theorized, cannot occur in DNA without repeating patterns of bases.

# Insertions and deletions via slippage



Energetic signatures of single base bulges: thermodynamic consequences and biological implications. Minetti CA, Remeta DP, Dickstein R, Breslauer KJ - Nucleic Acids Res. (2009)
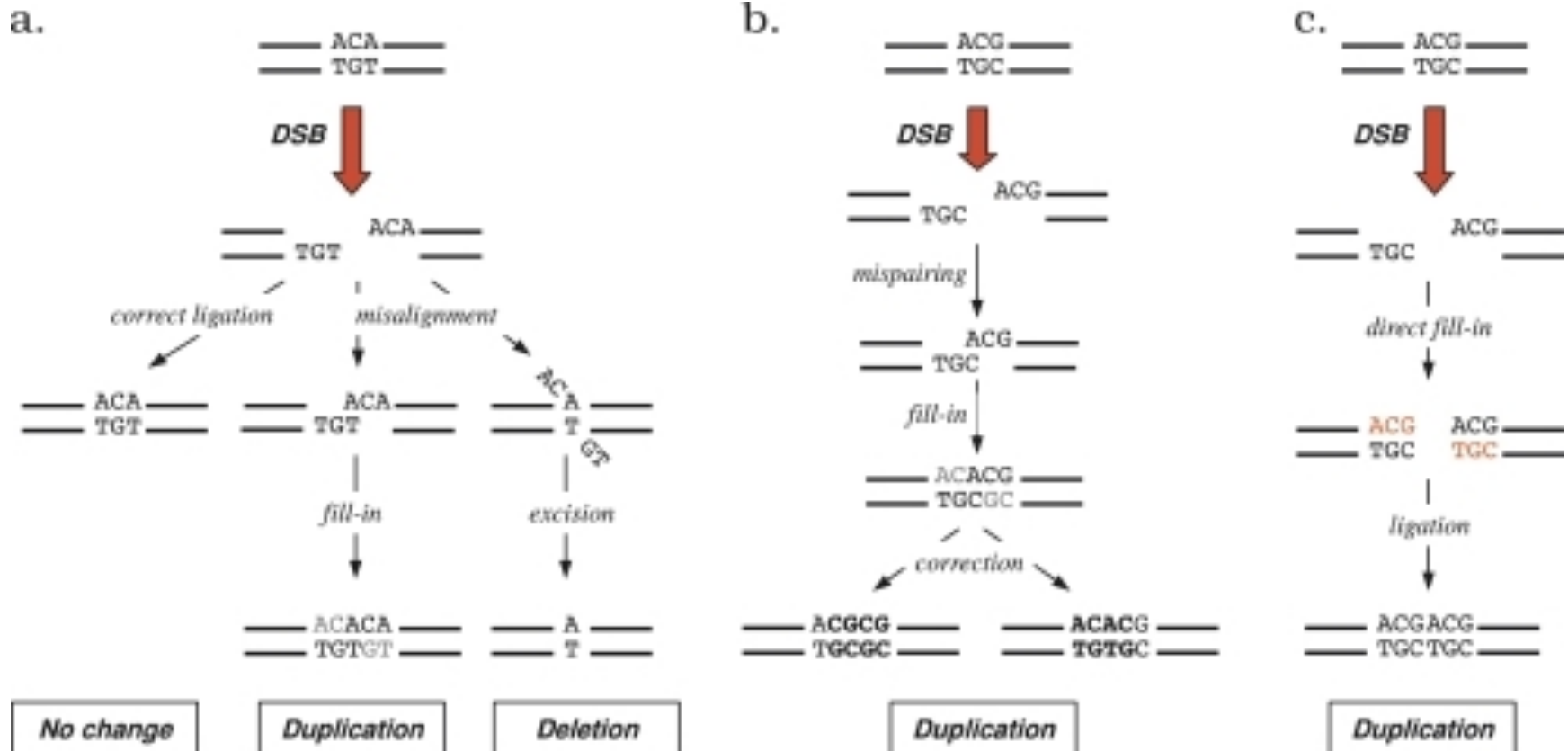
# Double-stranded break repair



Possible anti-recombinogenic role of Bloom's syndrome helicase in double-strand break processing.  doi:  10.1093/nar/gkg834
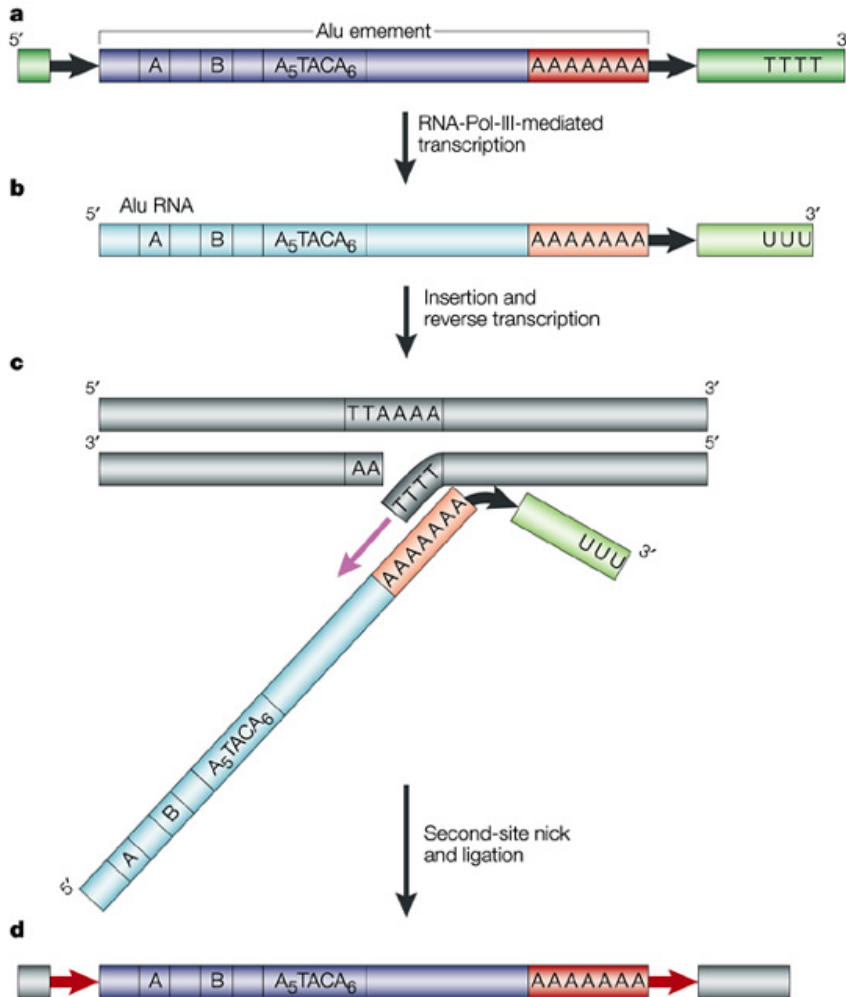
# NHEJ-derived indels



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)



a
5′ Alu element
A | B | $A_5TACA_6$ | AAAAAAA | TTTT 3′

RNA-Pol-III-mediated transcription

b
5′ Alu RNA
A | B | $A_5TACA_6$ | AAAAAAA | UUU 3′

Insertion and reverse transcription

c
5′ TTAAAA 3′
3′ AA 5′
TTTT
AAAAAAA
UUU 3′
$A_5TACA_6$
B
A
5′

Second-site nick and ligation

d
A | B | $A_5TACA_6$ | AAAAAAA

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 insertion/deletion (indel) polymorphism
2. **Standard approaches to detecting indels**
3. Assembly-based indel detection
4. Haplotype-based indel detection
5. Primary filtering: Bayesian variant calling
6. Post-call filtering: SVM
7. Graph-based resequencing approaches
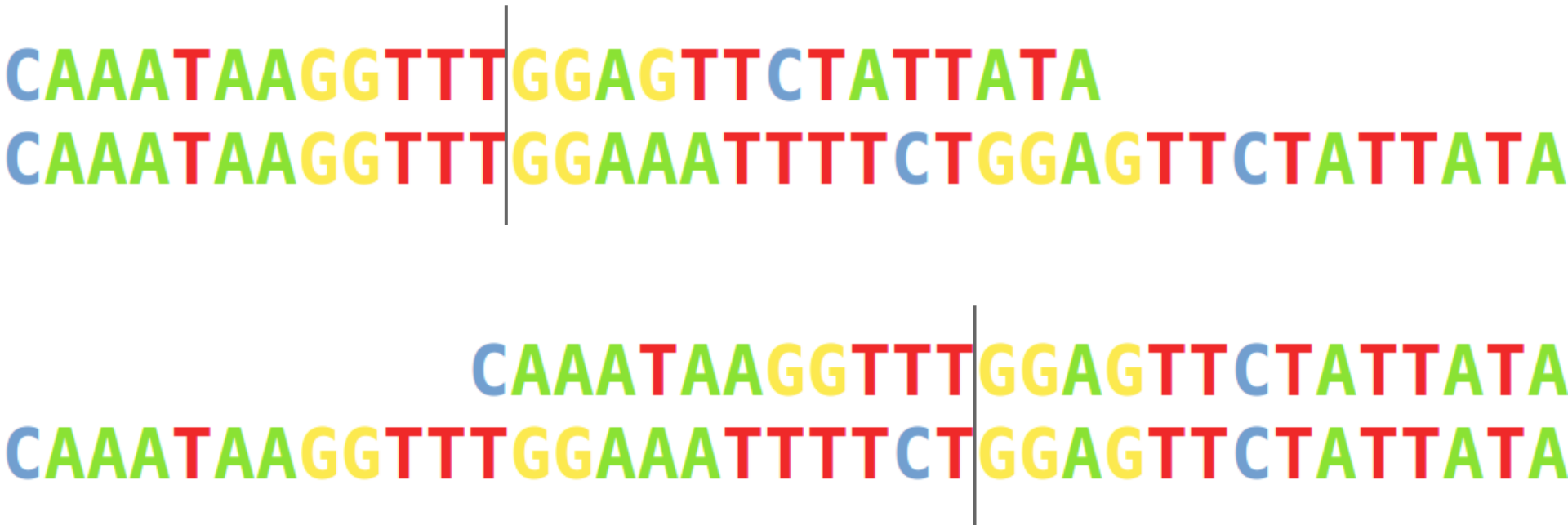
# Calling INDEL variation

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?

CAAATAAGGTTTGGAGTTCTATTATA
CAAATAAGGTTTGGAAATTTTCTGGAGTTCTATTATA

# Indel finder

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

# Indel finder

We can see this more easily like this:

**CAAATAAGGTTTGGAGTTCTATTATA**
**CAAATAAGGTTT***GGAAATTTTCT***GGAGTTCTATTATA**

CAAATAAGGTTT|GGAGTTCTATTATA
CAAATAAGGTTT|GGAAATTTTCTGGAGTTCTATTATA

CAAATAAGGTTT|GGAGTTCTATTATA
CAAATAAGGTTTGGAAATTTTCT|GGAGTTCTATTATA

# 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.

CAAATAAGGAAATTTTCTGGAGTTCTATTATA
CAAATAAGGTTTGCTATCTAGGTTATTATA

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----TCTGGAGTTCTATTATA
CAAATAAGG---TTTGCTATCT--AGGT-TATTATA

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

CAAATAAGGAAATTT--TCTGGAGTTCTATTATA
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?)

*First, let's review standard calling approaches.*

# Standard variant calling approach



Genome (FASTA)

Reads (FASTQ)

alignment and variant calling

Variation (VCF)

# Alignments to candidates



Reference

Reads

Variant observations

# The data exposed to the caller

Haplotype information is lost.

# 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.

CGTATGATCTAGCGCGCTAGCTAGCTAGC
CGTATGATCTA - - GCGCTAGCTAGCTAGC  ← Left aligned

CGTATGATCTAGCGCGCTAGCTAGCTAGC
CGTATGATCTAGC - - GCTAGCTAGCTAGC

CGTATGATCTAGCGCGCTAGCTAGCTAGC
CGTATGATCTAGCGC - -TAGCTAGCTAGC

example: 1000G PhaseI low coverage
chr15:81551110, ref:CTCTC alt:ATATA

ref: TGTCACTCGCTCTCTCTCTCTCTCTCTATATATATATATTTGTGCAT
alt: TGTCACTCGCTCTCTCTCTCTATATATATATATATATATTTGTGCAT

Interpreted as 3 SNPs

ref: TGTCACTCGCTCTCTCTCTCTCTCTCT------ATATATATATTTGTGCAT
alt: TGTCACTCGCTCTCTCTCTCT------ATATATATATATATATTTGTGCAT

Interpreted as microsatellite expansion/contraction

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

ref: TATAGAGAGAGAGAGAGAGCGAGAGAGAGAGAGAGAGGGGAGAGACGGAGTT
alt: TATAGAGAGAGAGAGAGAGCGAGAGAGAGAGAGAGAGAGGGGAGAGACGGAGTT

ref: TATAGAGAGAGAGAGAGAGC--GAGAGAGAGAGAGAGAGGGGAGAGACGGAGTT
alt: TATAGAGAGAGAGAGAG--CGAGAGAGAGAGAGAGAGAGGGGAGAGACGGAGTT

# Overview

1. Genesis of insertion/deletion (indel) polymorphism
2. Standard approaches to detecting indels
3. **Assembly-based indel detection**
4. Haplotype-based indel detection
5. Primary filtering: Bayesian variant calling
6. Post-call filtering: SVM
7. Graph-based resequencing approaches

## *Problem:* inconsistent indel 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)

# Assembly



Original sequence
GTAGTATAGTCAGTATCA

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

k-mers (2-mers)
GT  TA  AG  AT  TC  CA

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

de Bruijn graph

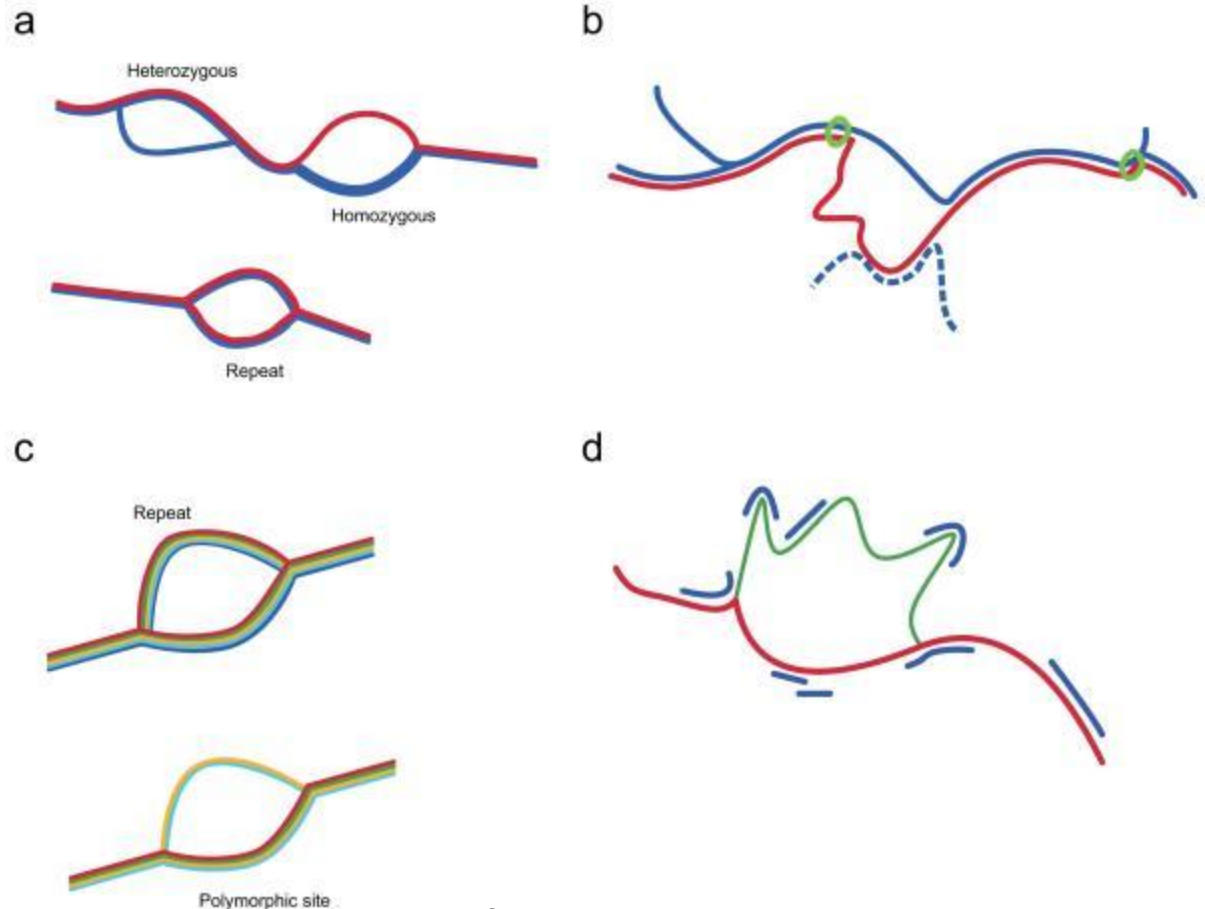http://www.nature.com/nrmicro/journal/v7/n4/full/nrmicro2088.html

# 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.
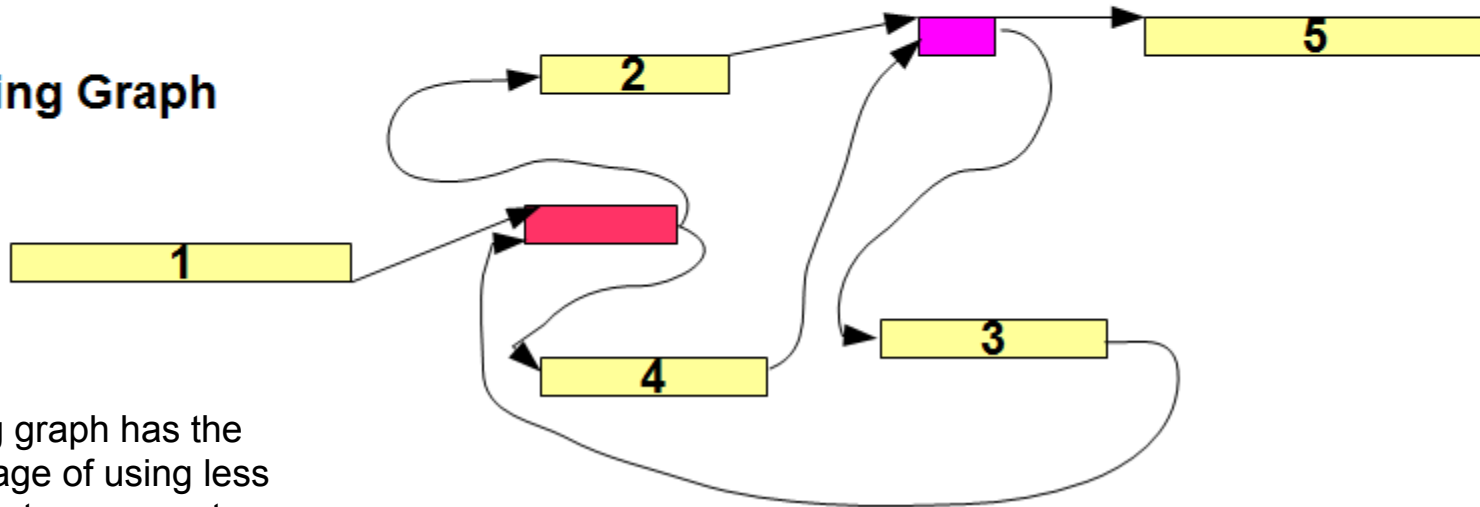
Many samples can be called at the same time.



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

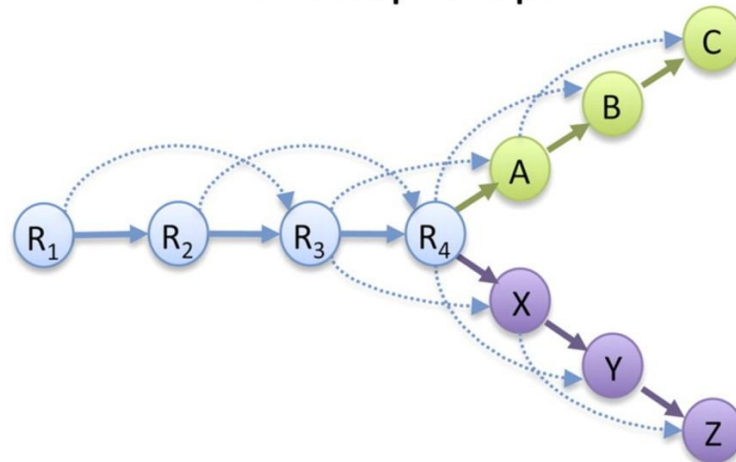# String graphs (SGA)

**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.

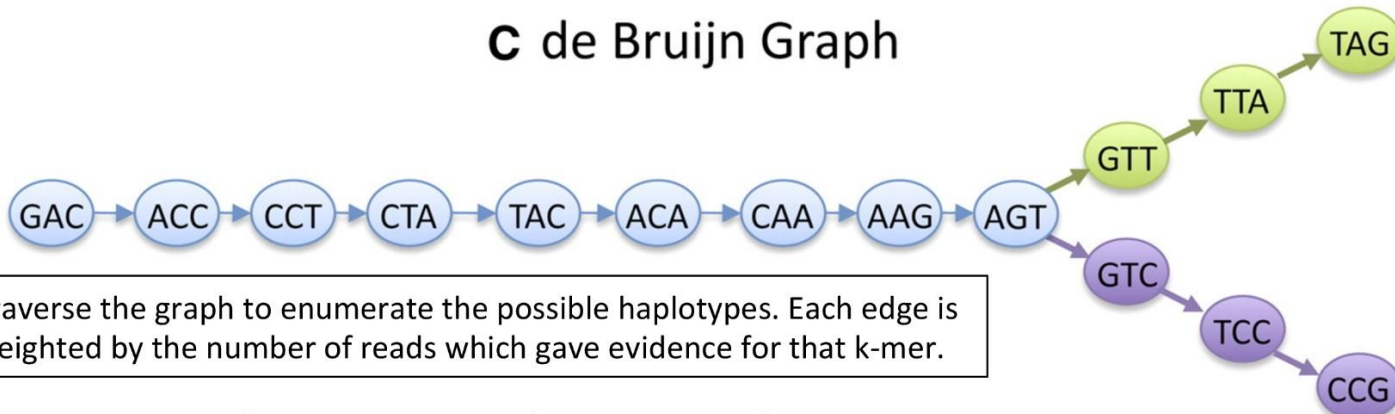# Discovering alleles using graphs (GATK HaplotypeCaller)



**A** Read Layout

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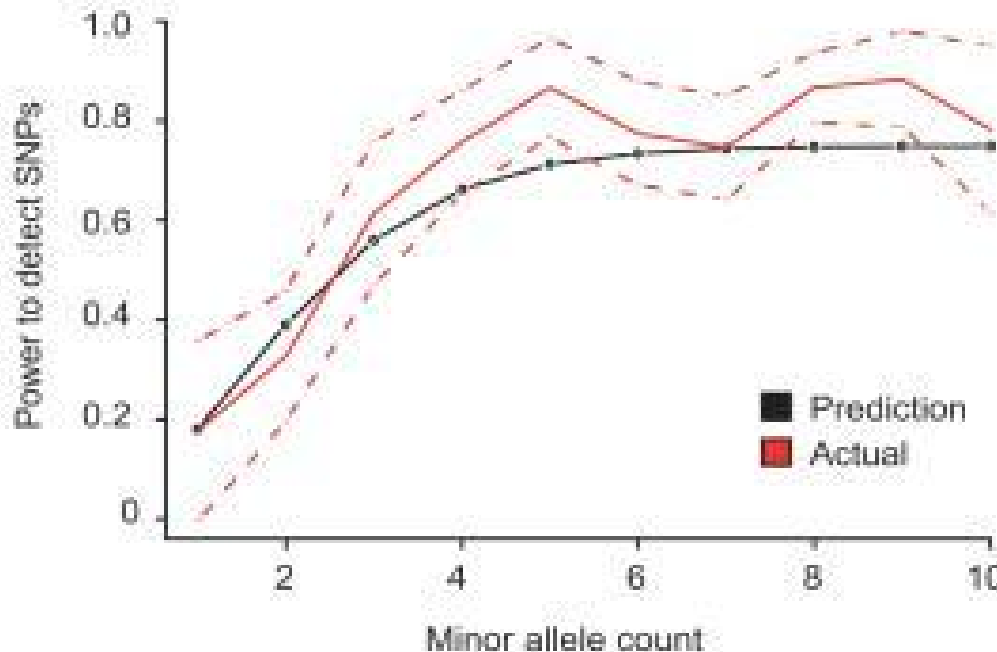Y: ACAAGTCC
Z: CAAGTCCG

**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.

*Assembly of large genomes using second-generation sequencing. Schatz. Genome Research. 2010.*

13

# **Why don't we just assemble?**

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



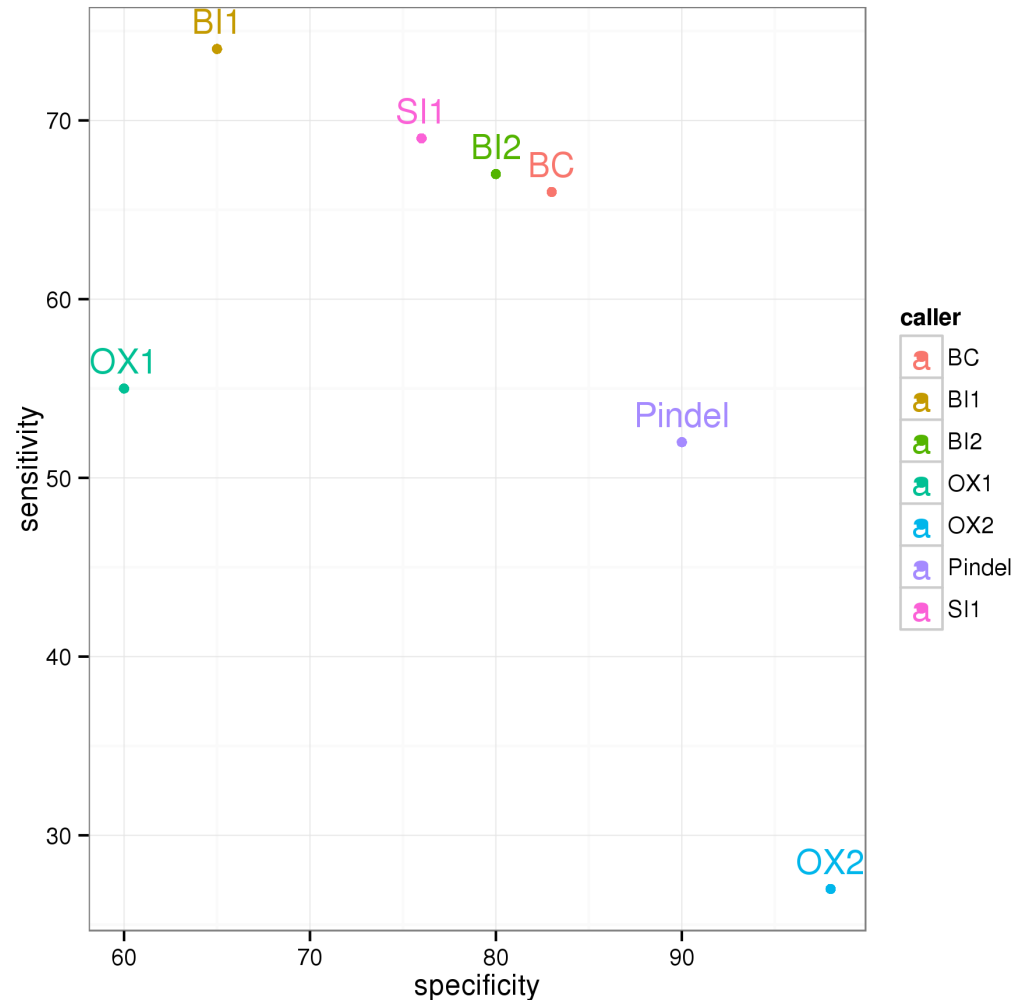The requirement of exact kmer matches means that errors disrupt coverage of alleles.

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)
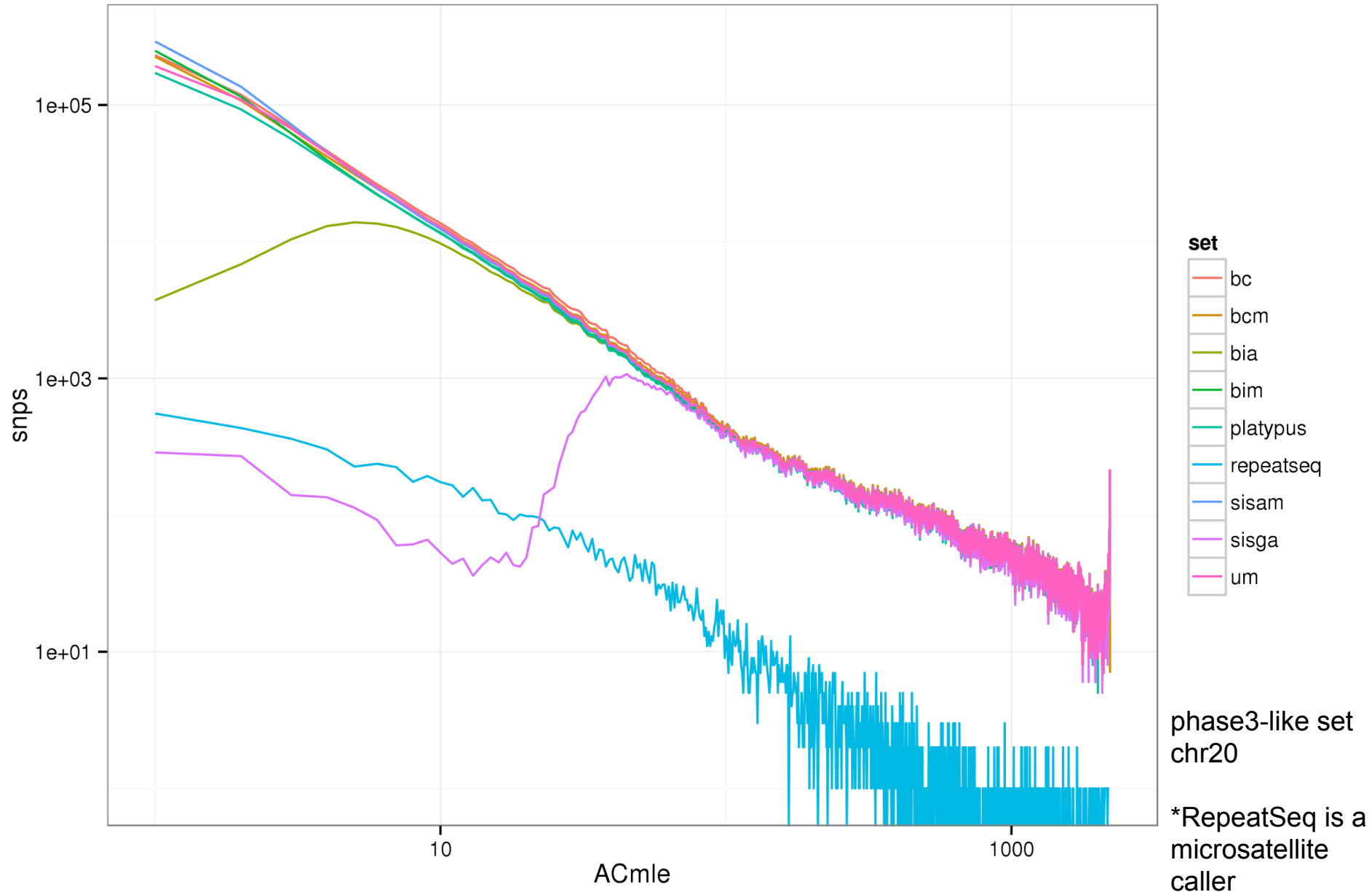
# Indel validation, 191 AFR samples

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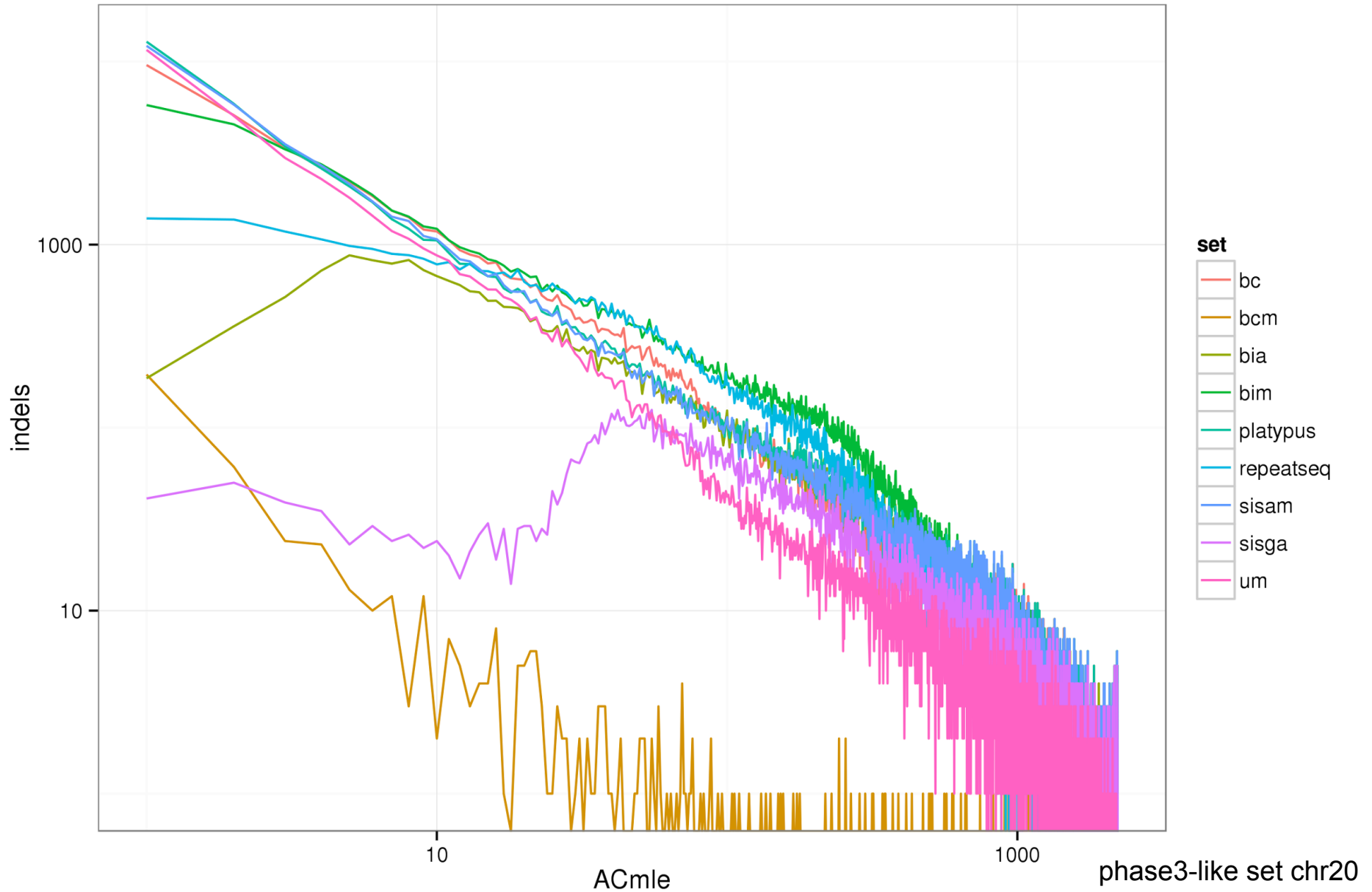
*The local assembly-based method Platypus (OX1) had a genotyping bug which caused poor performance.

# Site-frequency spectrum, SNPs



phase3-like set
chr20

*RepeatSeq is a microsatellite caller

# Site-frequency spectrum, indels



phase3-like set chr20

# Overview

1.  Genesis of insertion/deletion (indel) polymorphism
2.  Standard approaches to detecting indels
3.  Assembly-based indel detection
4.  **Haplotype-based indel detection**
5.  Primary filtering: Bayesian variant calling
6.  Post-call filtering: SVM
7.  Graph-based resequencing approaches

# Finding haplotype polymorphisms

Two reads

AGAACCCAGTGCTCTTTCTGCT

AGAACCCAGTGGTCTTTCTGCT

a SNP

AGAACCCAGTG C/G TCTTTCTGCT

Their alignment

AGAACCCAGTGCTCTATCTGCT

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

AGAACCCAGTG CTCTA / GTCTT TCTGCT

A variant locus implied by alignments

# Direct detection of haplotypes

Reference

Reads

Direct detection of haplotypes from reads resolves differentially-represented alleles (as the sequence is compared, not the alignment).

Allele detection is still alignment-driven.

Detection window

# 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, ½ of variants are within ~22bp 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
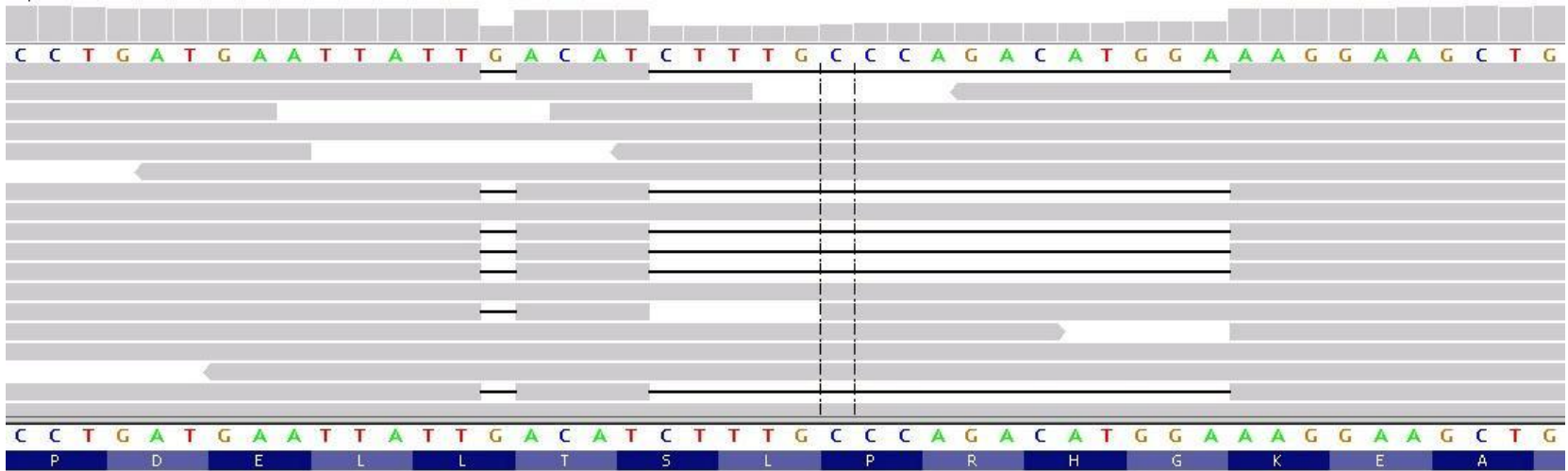
apparent:
AGG **T**AG CTG
Arg Ter ---

Apparent nonsense variant, one YRI homozygote

actual:
AGG **TT**G CTG
Arg Leu Leu

Actually a block substitution that results in a missense substitution

(Daniel MacArthur)

# Importance of haplotype effects: frame-restoring indels



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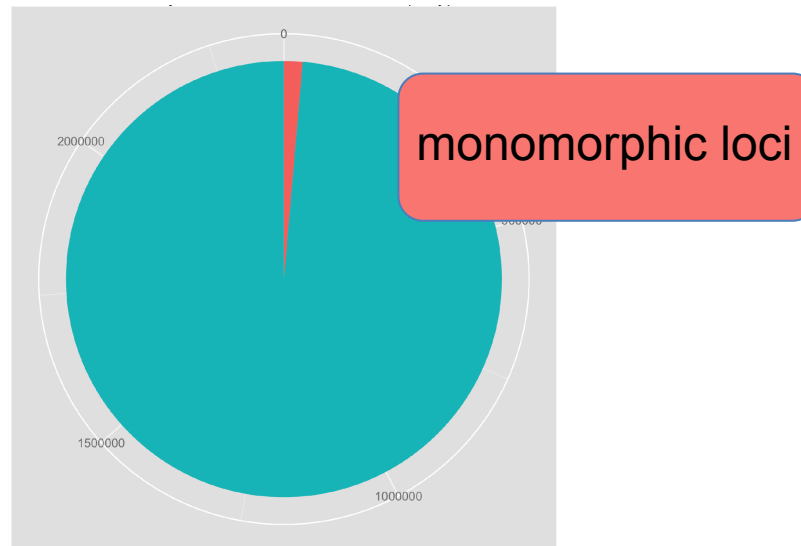
## chr6:117113761, GPRC6A **(~10% AF in 1000G)**

```
ref: ATTGTAATTCTCA--TA--TT--TGCCTTTGAAAGC
alt: ATTGTAATTCTCAGGTAATTTCCTGCCTTTGAAAGC
```

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

```
ref: CCACCGCGGCCCGCGCCTG-C-TCCAGGATGTCC
Alt: CCACCGCGG--CGCGCCTGTCTTCCAGGAGGTCC
```

# 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).
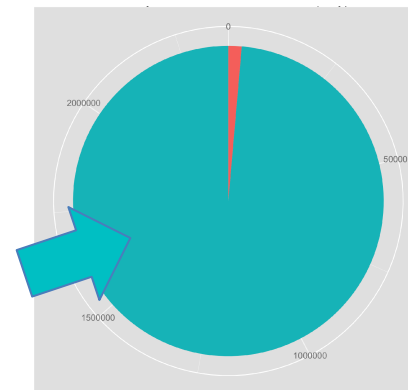


monomorphic loci

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
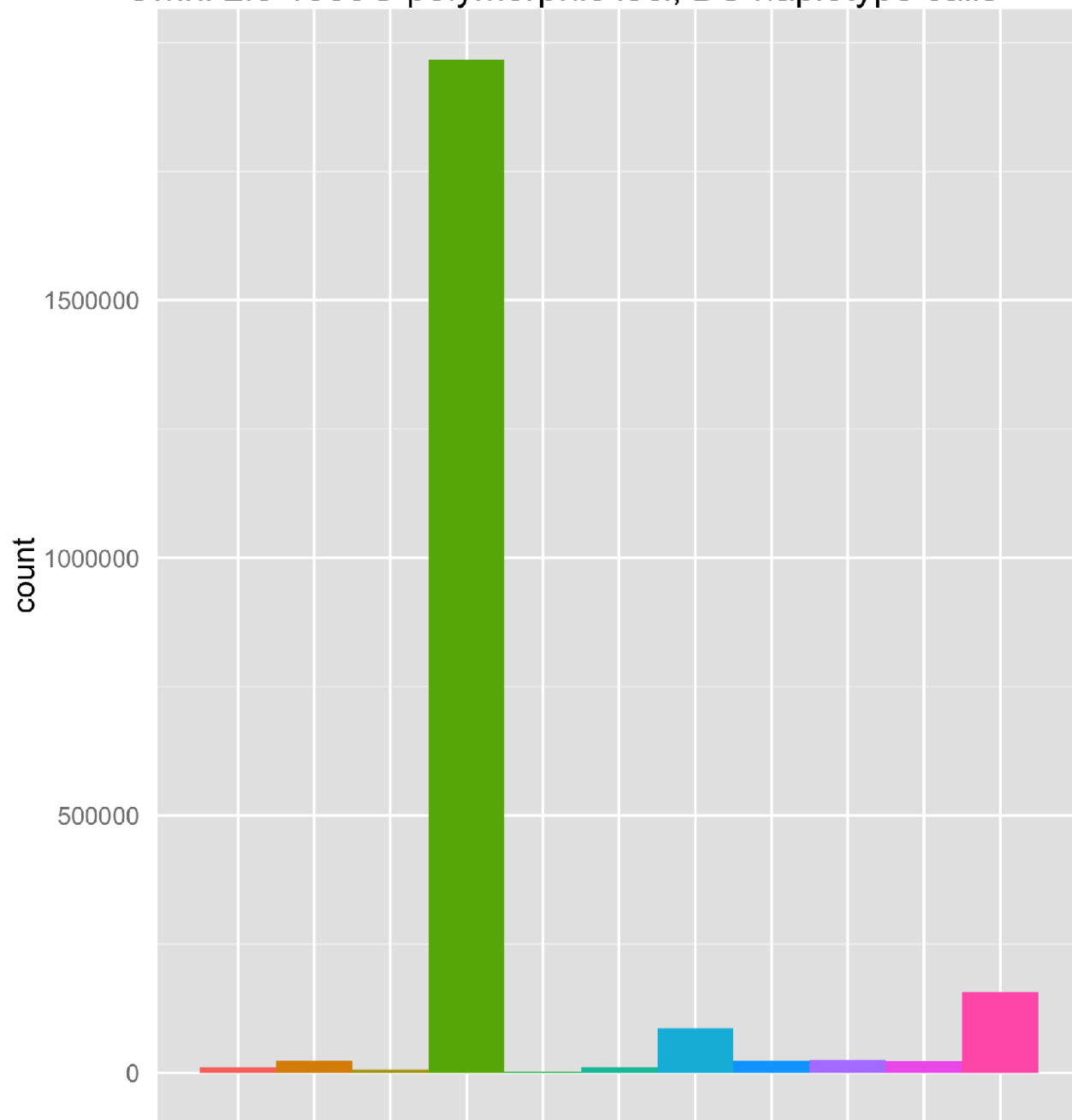- multiallelic SNP and MNP
- multiallelic SNP, MNP, and complex
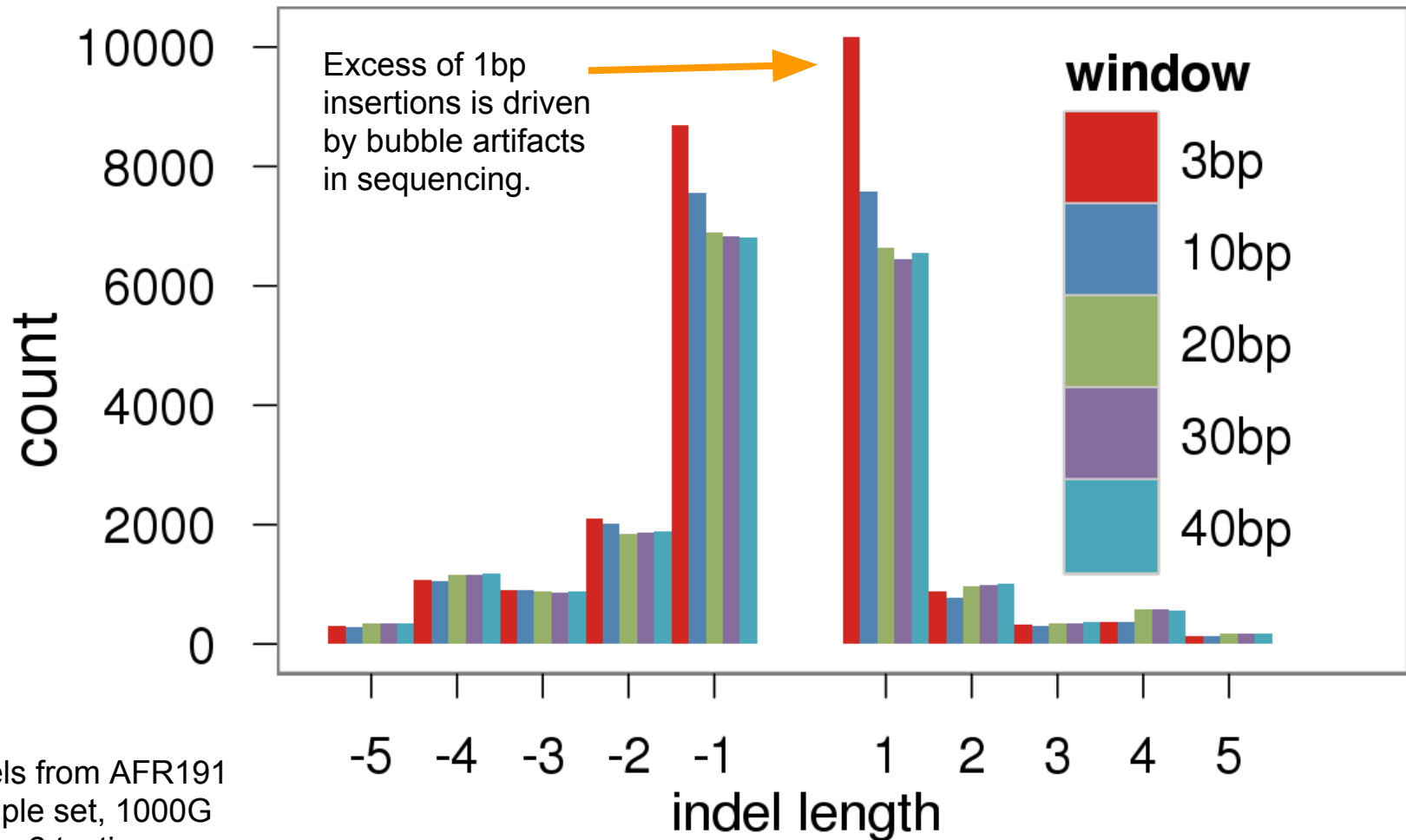
Omni 2.5 1000G polymorphic loci, BC haplotype calls

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

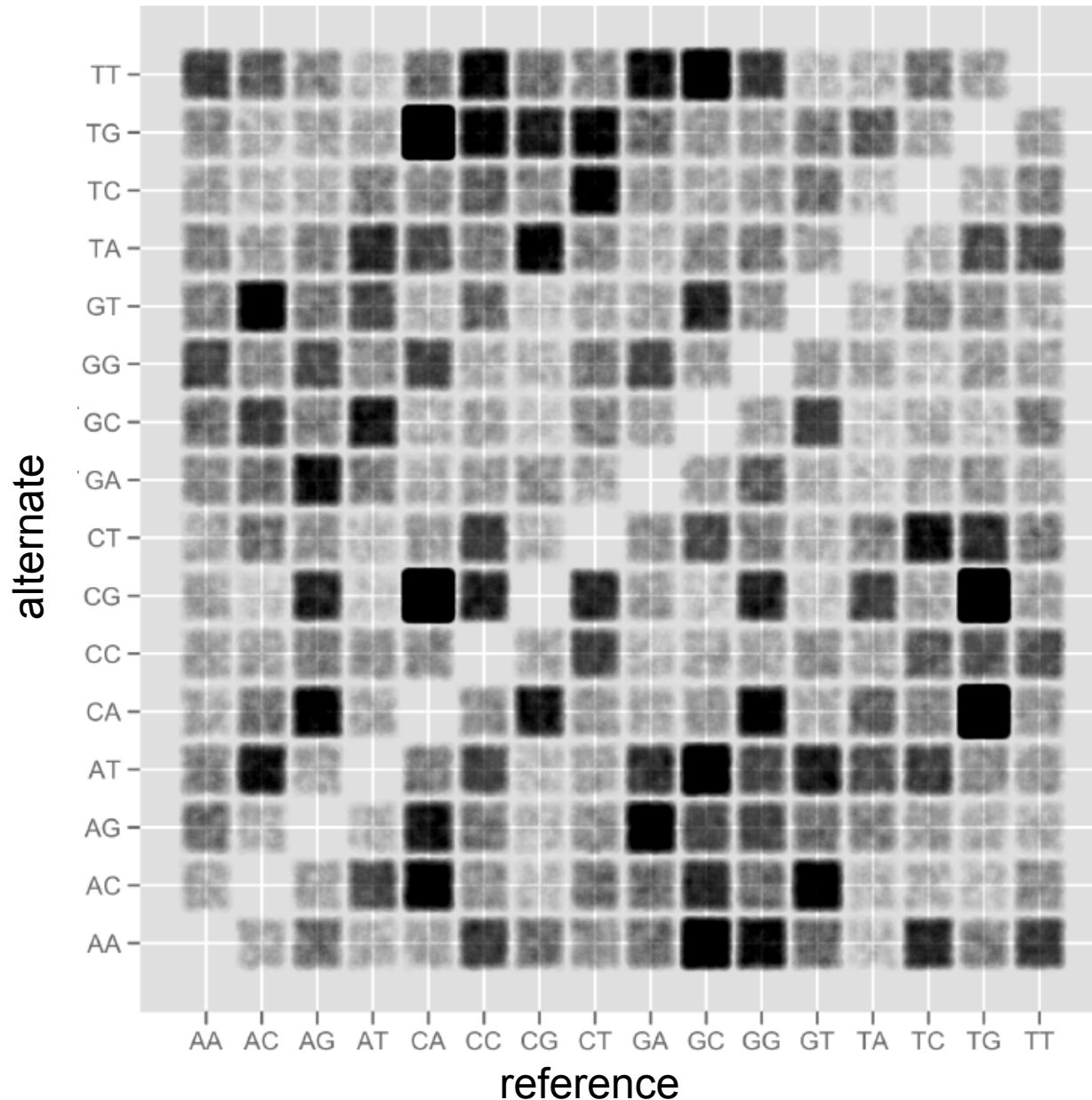# Measuring haplotypes improves specificity



Excess of 1bp insertions is driven by bubble artifacts in sequencing.

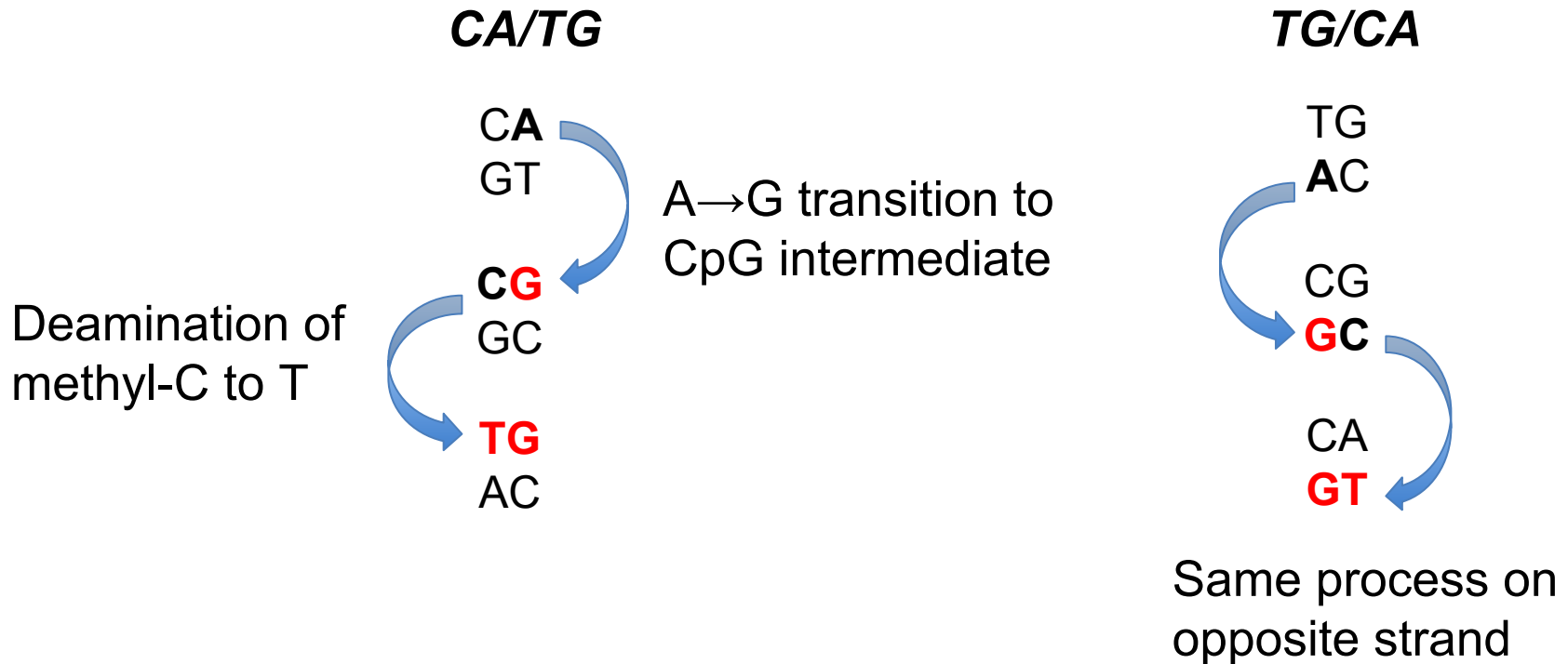window: 3bp, 10bp, 20bp, 30bp, 40bp

Indels from AFR191 sample set, 1000G phase2 testing.

# 2bp MNPs and dinucleotide intermediates

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

*CA/TG*

C**A**
GT

A→G transition to CpG intermediate

**C**<span style="color:red">**G**</span>
GC

Deamination of methyl-C to T

<span style="color:red">**TG**</span>
AC

*TG/CA*

TG
**A**C

CG
<span style="color:red">**G**</span>**C**

CA
<span style="color:red">**GT**</span>

Same process on opposite strand

# Overview

1. Genesis of insertion/deletion (indel) polymorphism
2. Standard approaches to detecting indels
3. Assembly-based indel detection
4. Haplotype-based indel detection
5. **Primary filtering: Bayesian variant calling**
6. Post-call filtering: SVM
7. Graph-based resequencing approaches

# Filtering INDELs

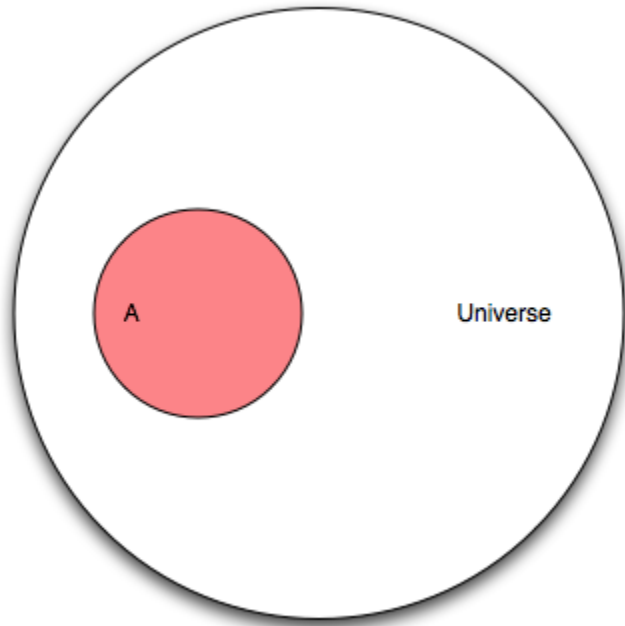As with SNPs, sequencing error rates are high.

So, we need to filter.

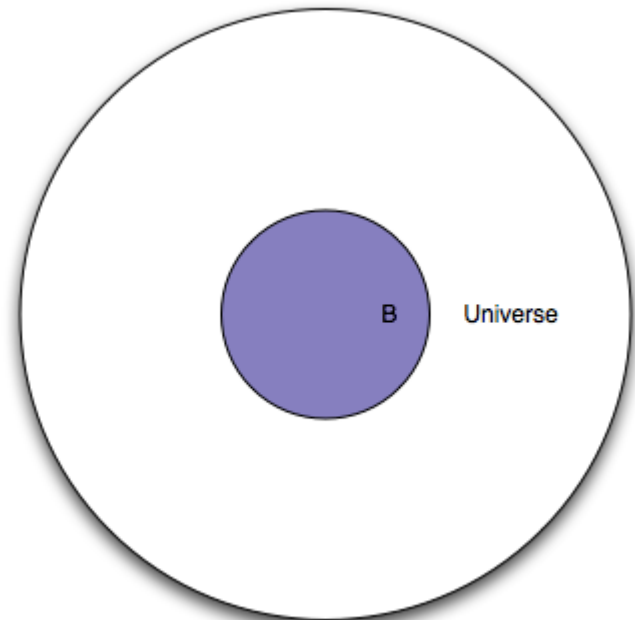The standard filter of NGS is the Bayesian variant caller.

Combines population-based priors and data from many samples to make high-quality calls.

# Bayesian (visual) intuition

We have a universe of individuals.

A = samples with a
variant at some locus

B = putative observations
of variant at some locus

Figures from http://oscarbonilla.com/2009/05/visualizing-bayes-theorem/
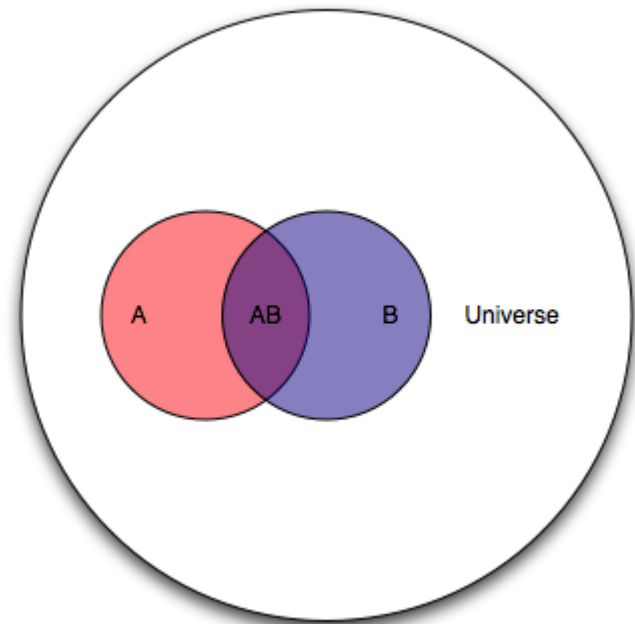
# probability(A|B)

We want to estimate the probability that we have a real polymorphism "A" given "|" that we observed variants in our alignments "B".
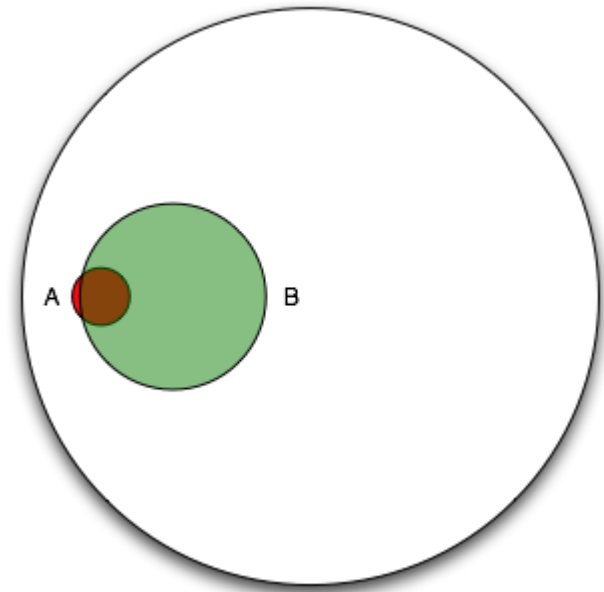
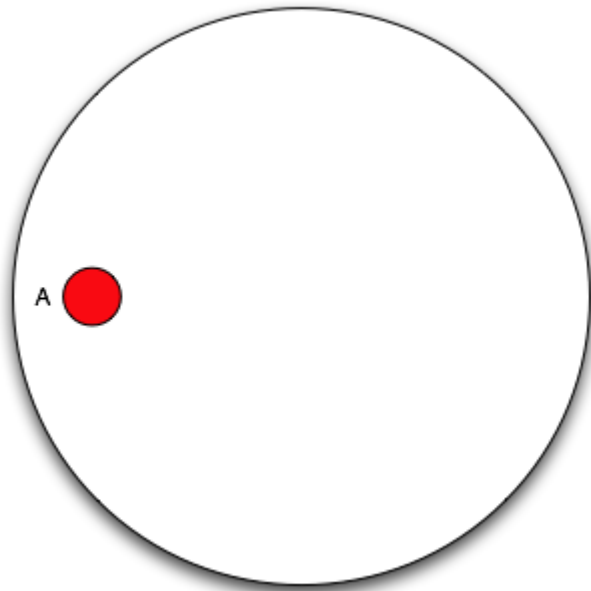$$P(A|B) = \frac{|AB|}{|B|}$$

$$P(A|B) = \frac{P(AB)}{P(B)}$$

$$P(B|A) = \frac{P(AB)}{P(A)}$$
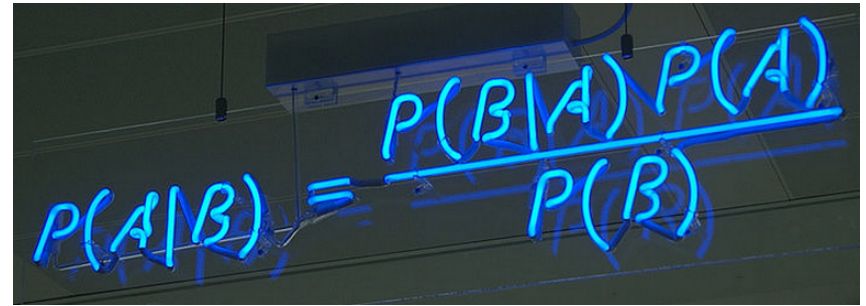
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

# In our case it's a bit more like this...



Observations (B) provide pretty good sensitivity, but poor specificity.

# The model



- Bayesian model estimates the probability of polymorphism at a locus given input data and the population mutation rate (~pairwise heterozygosity) and assumption of "neutrality" (random mating).
- Following Bayes theorem, the probability of a specific set of genotypes over some number of samples is:
  - **P(G|R) = ( P(R|G) P(G) ) / P(R)**
- Which in FreeBayes we extend to:
  - **P(G,S|R) = ( P(R|G,S) P(G)P(S) ) / P(R)**
  - **G** = genotypes, **R** = reads, **S** = locus is well-characterized/mapped
  - **P(R|G,S)** is our data likelihood, **P(G)** is our prior estimate of the genotypes, **P(S)** is our prior estimate of the mappability of the locus, **P(R)** is a normalizer.
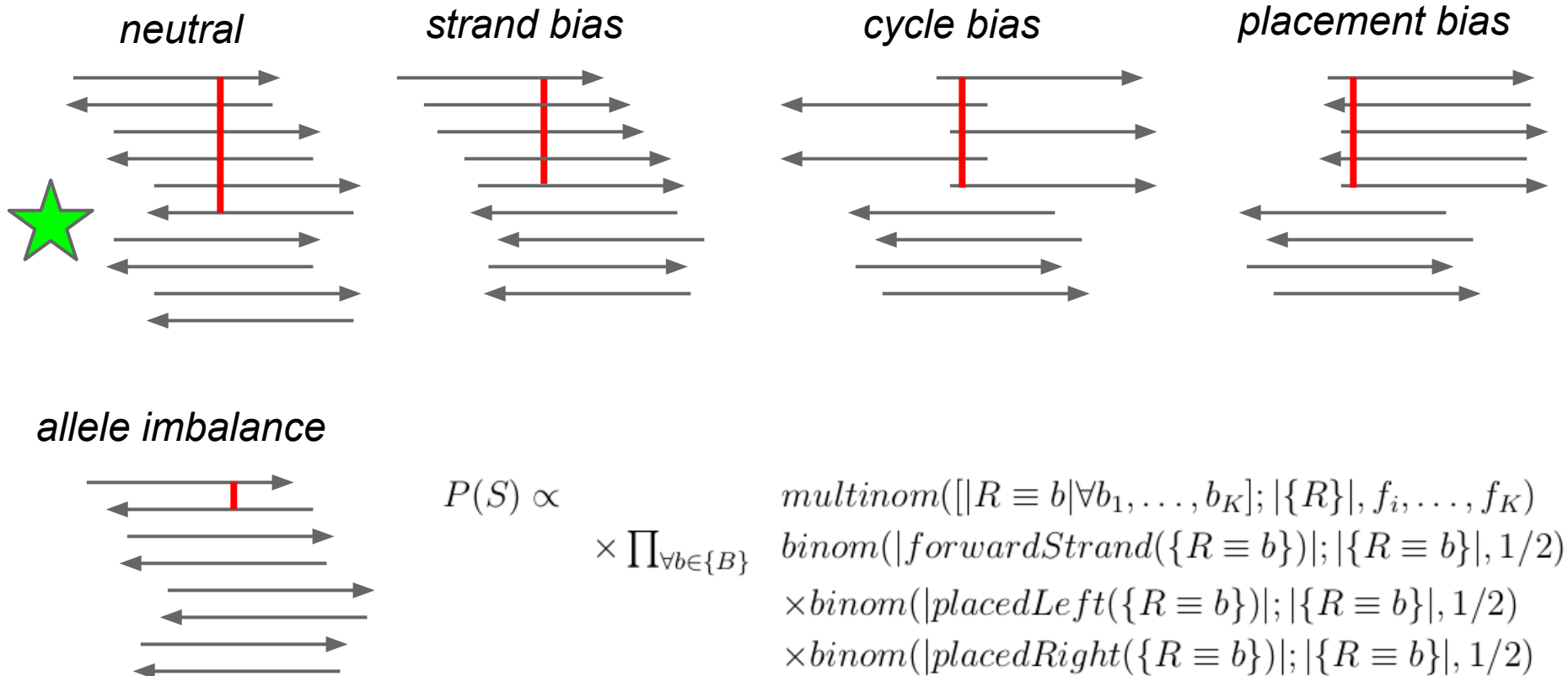
# Handling non-biallelic/diploid cases

We compose our data likelihoods, **P(Reads|Genotype)** using a discrete multinomial sampling probability:

$$P(reads|genoytpe) = \binom{|reads|}{|reads = A|, |reads = B| \dots}$$

$$X \prod_{\forall alleles \in genotype} freq(allele \in genotype)$$
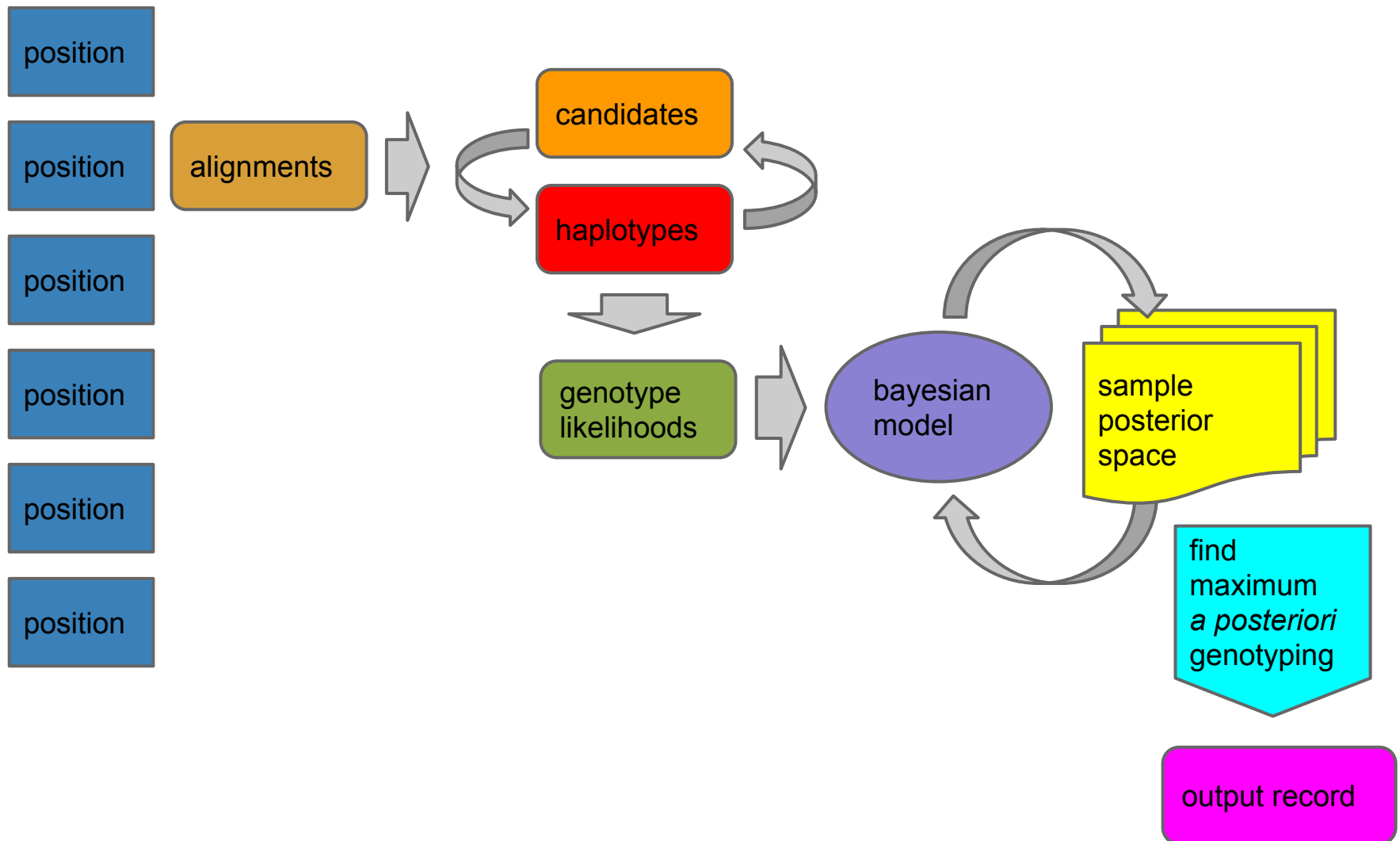
$$X \prod_{\forall reads} P(correct(read))$$

Our priors, **P(Genoypes)**, follow the Ewens Sampling Formula and the discrete sampling probability for genotypes.

# Are our locus and alleles sequenceable?

In WGS, biases in the way we observe an allele (placement, position, strand, cycle, or balance in heterozygotes) are often correlated with error. We include this in our posterior **P(G,S|R)**, and to do so we need an estimator of **P(S)**.

*neutral*          *strand bias*          *cycle bias*          *placement bias*

*allele imbalance*

$$P(S) \propto$$

$$\times \prod_{\forall b \in \{B\}}$$

$$multinom([|R \equiv b|\forall b_1, \dots, b_K]; |\{R\}|, f_i, \dots, f_K)$$

$$binom(|forwardStrand(\{R \equiv b\})|; |\{R \equiv b\}|, 1/2)$$

$$\times binom(|placedLeft(\{R \equiv b\})|; |\{R \equiv b\}|, 1/2)$$

$$\times binom(|placedRight(\{R \equiv b\})|; |\{R \equiv b\}|, 1/2)$$

# The detection process

position

position

alignments

position

position

position

position

candidates

haplotypes

genotype likelihoods

bayesian model

sample posterior space

find maximum *a posteriori* genotyping

output record

# Overview

1. Genesis of insertion/deletion (indel) polymorphism
2. Standard approaches to detecting indels
3. Assembly-based indel detection
4. Haplotype-based indel detection
5. Primary filtering: Bayesian variant calling
6. **Post-call filtering: SVM**
7. Graph-based resequencing approaches
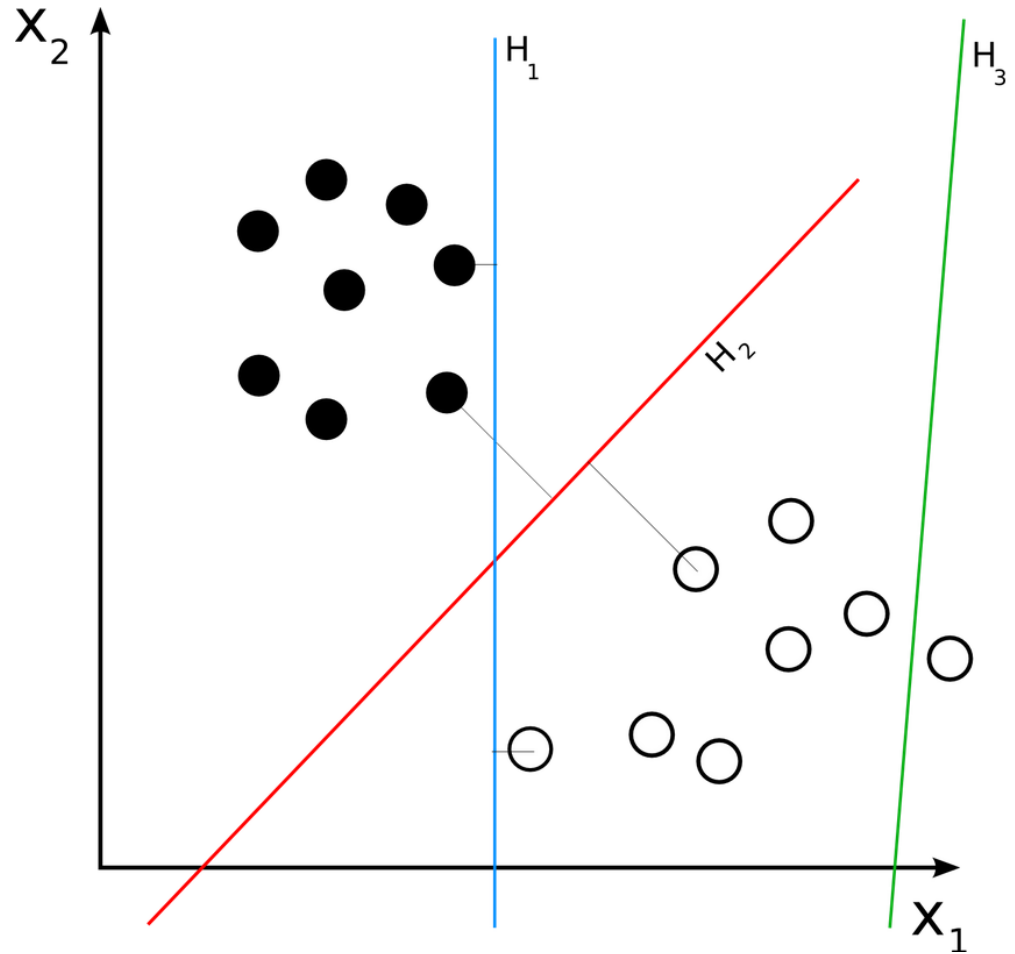
# SVM filtering

INDEL detection is hard.

*A priori* models can't capture all types of error.

It's especially difficult when we try to make a consensus set from lots of input variant callers.

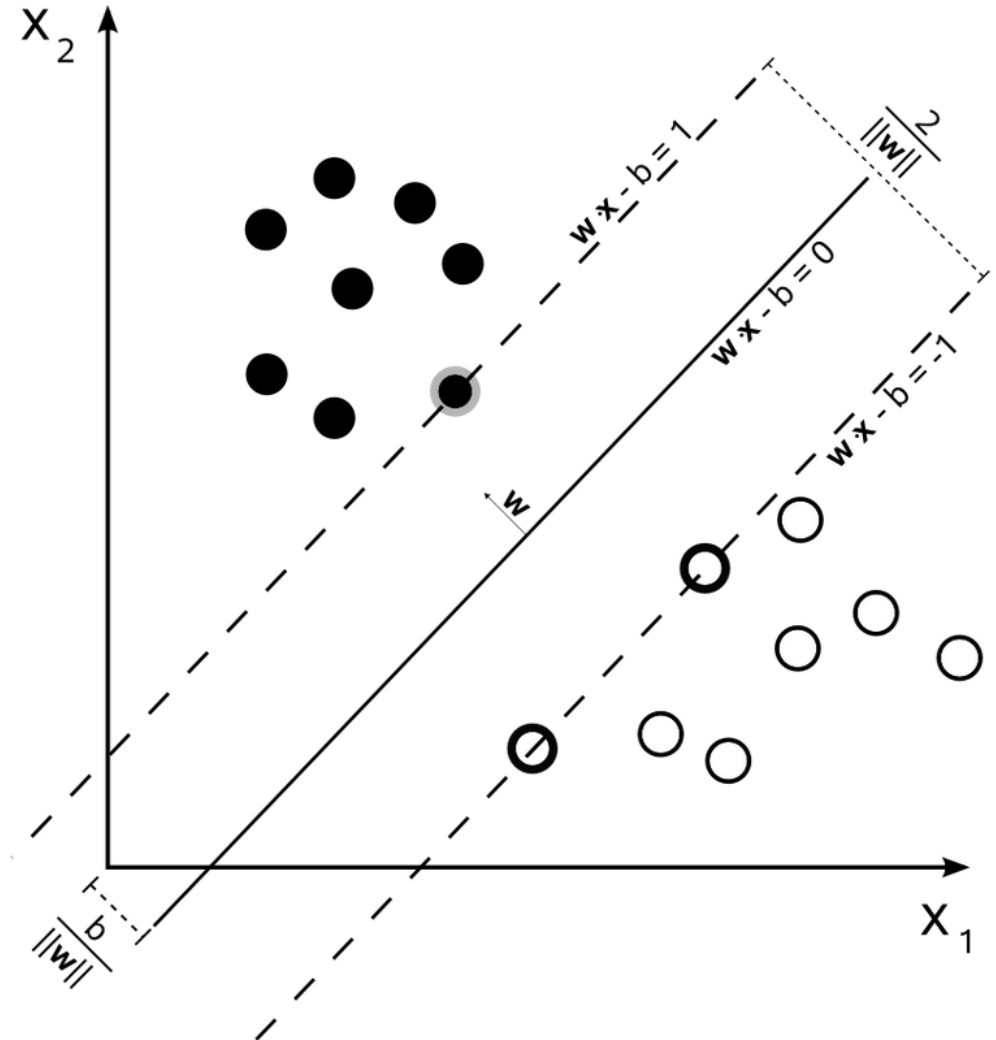We can use classifiers like Support Vector Machines (SVM) to further improve results.

# SVM classifier

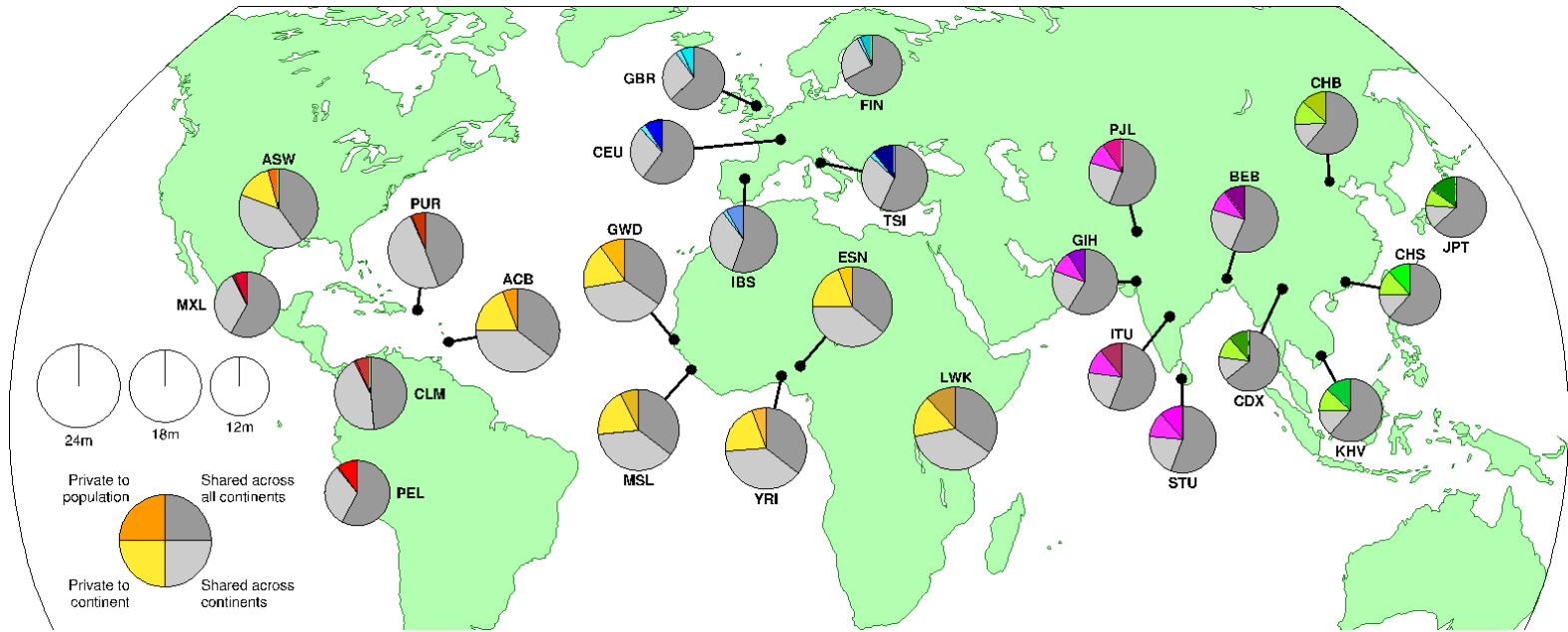Find a hyperplane (here a line in 2D) which separates observations.

# SVM classifier

The best separating hyperplane is determined by maximum margin between groups we want to classify.
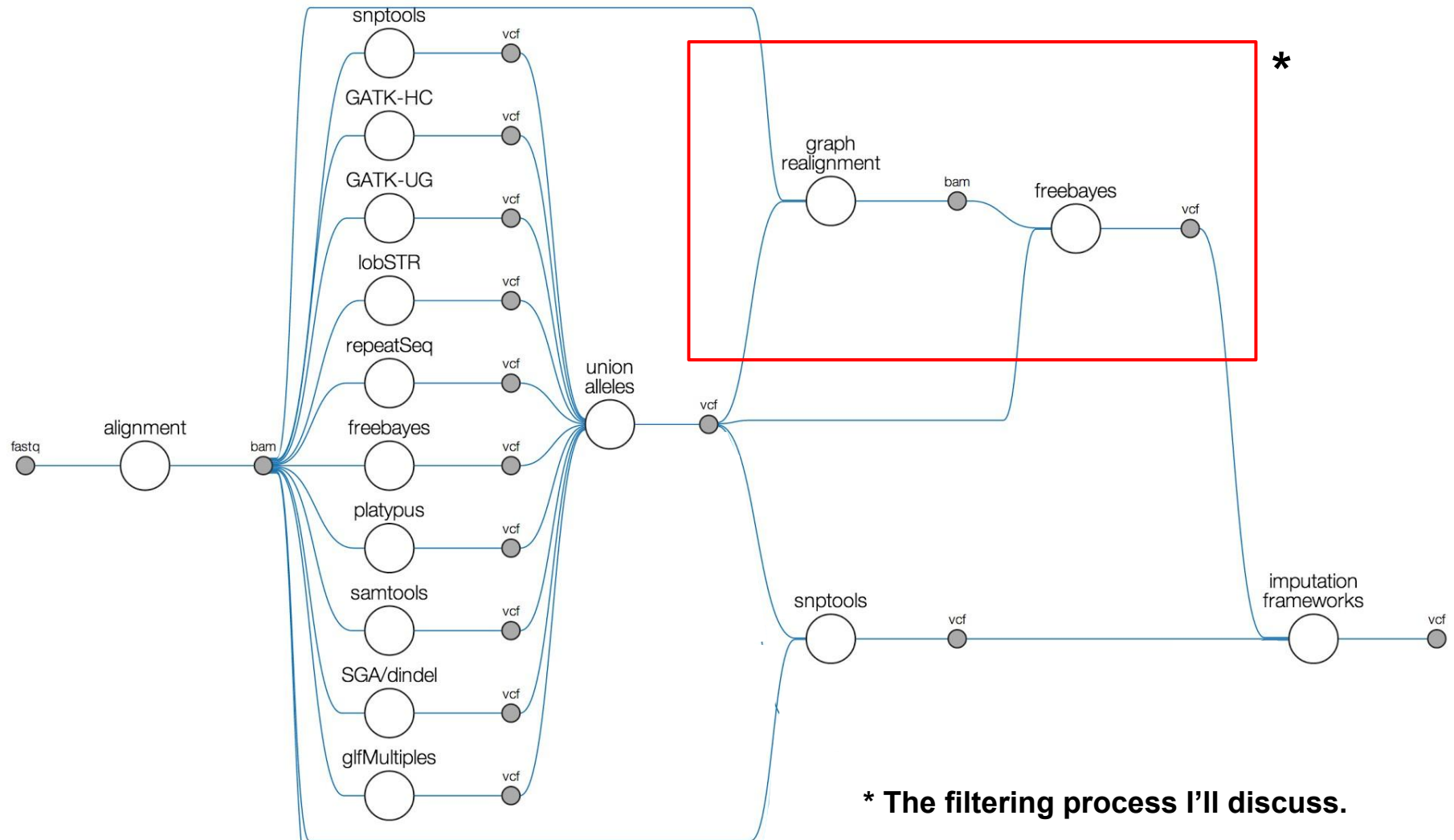
# SVM filtering in the 1000 Genomes



25 human populations X ~100 samples each.

# 1000G variant integration process



* The filtering process I'll discuss.

# SVM approach for INDEL filtering

Extract features that tend to vary with respect to call quality:

- call QUALity
- read depth
- sum of base qualities
- inbreeding coefficient (heterozygosity)
- entropy of sequence at locus
- mapping quality
- allele frequency in population
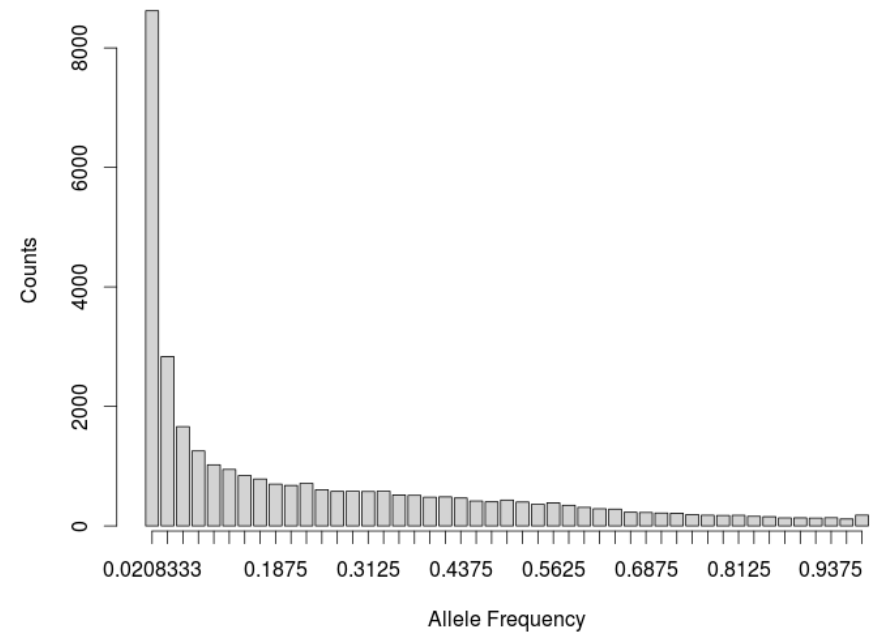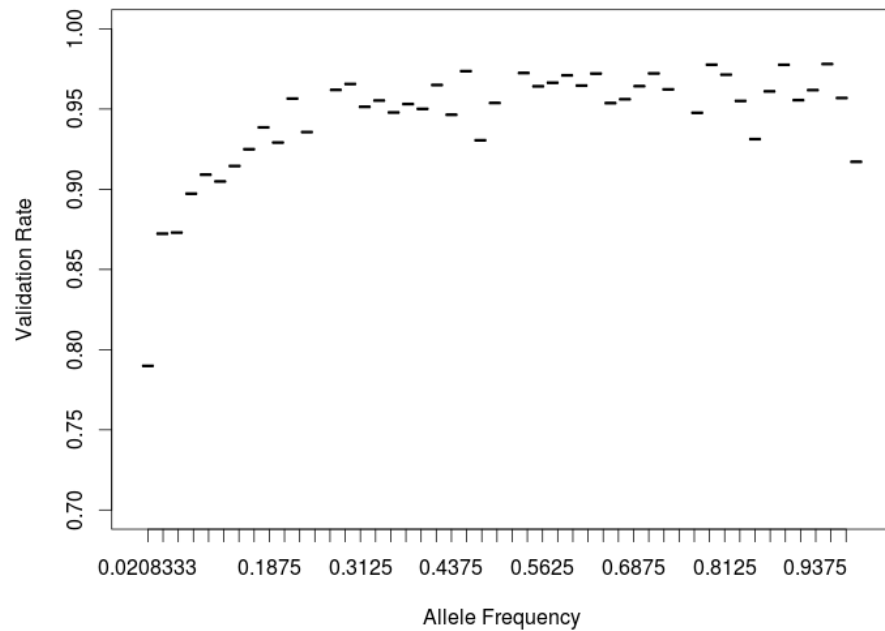- read pairing rate
- etc.

# SVM approach for INDEL filtering

Now, use overlaps in validation samples or sites to determine likely errors and true calls.

Use this list + annotations of the calls to train an SVM model.

Apply the model to all the calls, filter, and measure validation rate of the whole set.
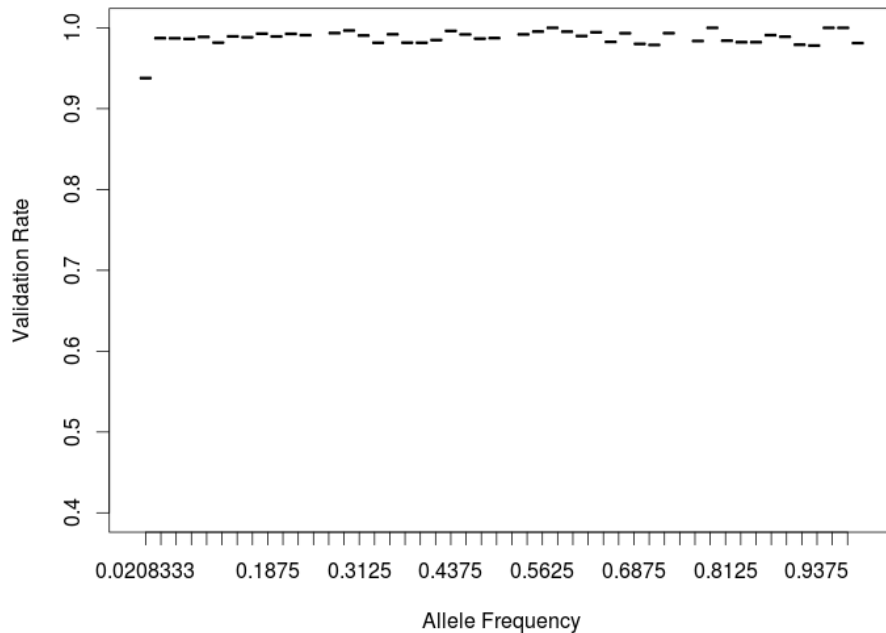
# Application of SVM to 1000G INDELs



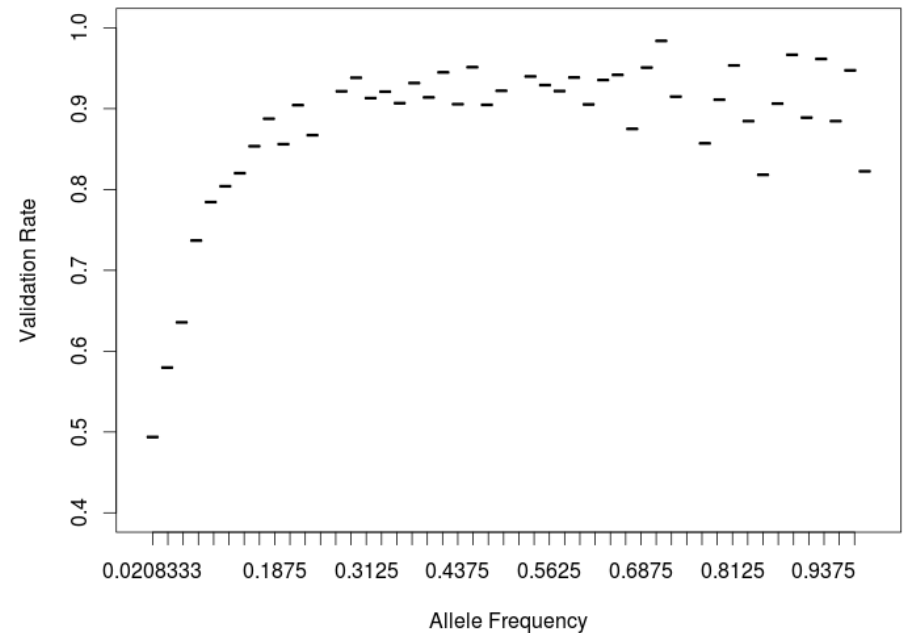Raw validation rates of indels in 1000G phase 3, "MVNCall" set.

*Tony Marcketta and Adam Auton*

# Application of SVM to 1000G INDELs
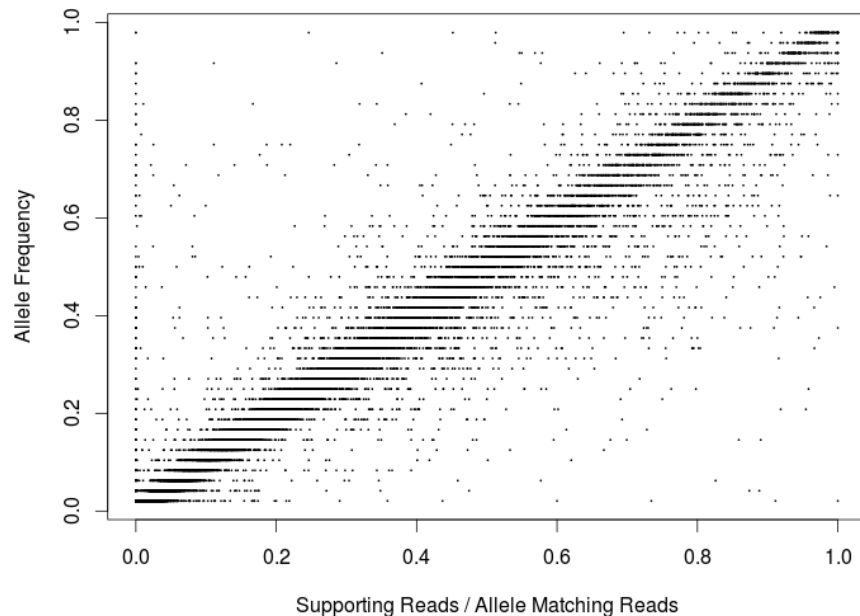


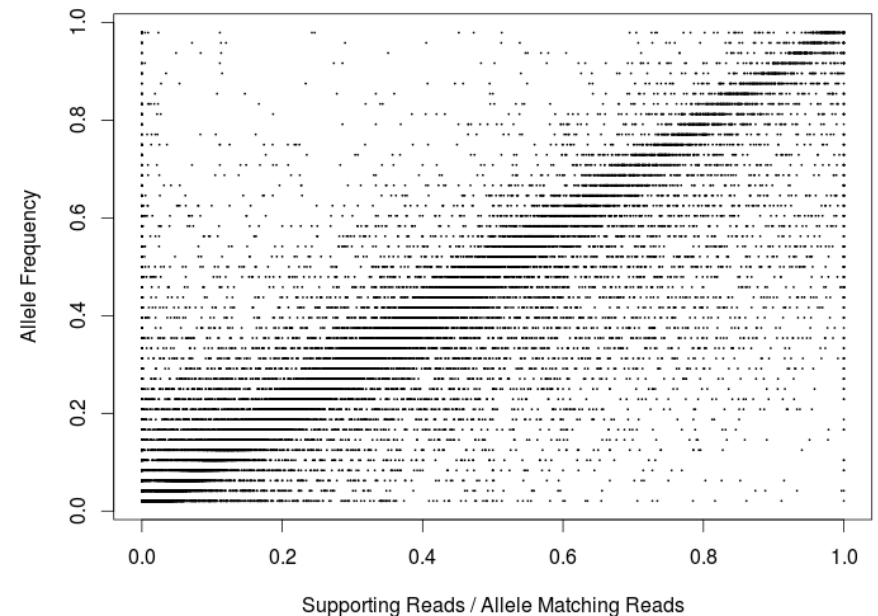Filtering results, using SVM-based method.

*Anthony Marcketta and Adam Auton*

# Application of SVM to 1000G INDELs

**Passing SVM**　　　　　　　　　　　　　　**Failing SVM**



Correlation between allele frequency and observation counts.
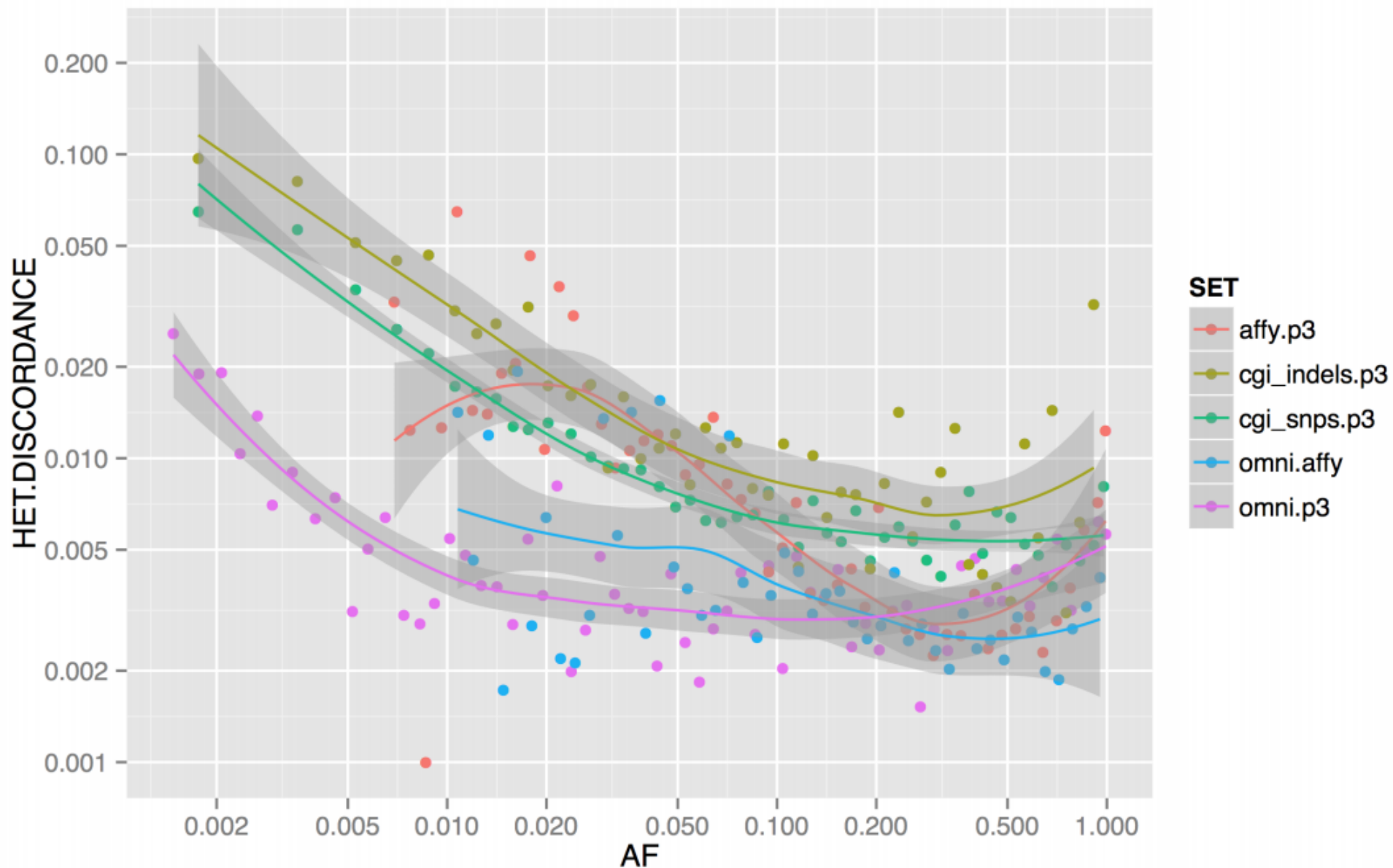
*Anthony Marcketta and Adam Auton*

# Indel results from 1000G

| Gold | Eval | R/R | R/A | A/A | All | NonRef |
|------|------|------|------|------|------|--------|
| CGI SNPs | Phase3 | 0.9998 | 0.9930 | 0.9983 | 0.9994 | 0.9920 |
| CGI Indels | Phase3 | 0.9990 | 0.9889 | 0.9923 | 0.9982 | 0.9805 |

Comparing the phase3 results to the genotypes for indels in the subset of samples for which we also had high-quality, high-coverage genomes from Complete Genomics.

**Genotype Accuracy by Allele Frequency**

# Overview

1. Genesis of insertion/deletion (indel) polymorphism
2. Standard approaches to detecting indels
3. Assembly-based indel detection
4. Haplotype-based indel detection
5. Primary filtering: Bayesian variant calling
6. Post-call filtering: SVM
7. **Graph-based resequencing approaches**

# We know the variants,
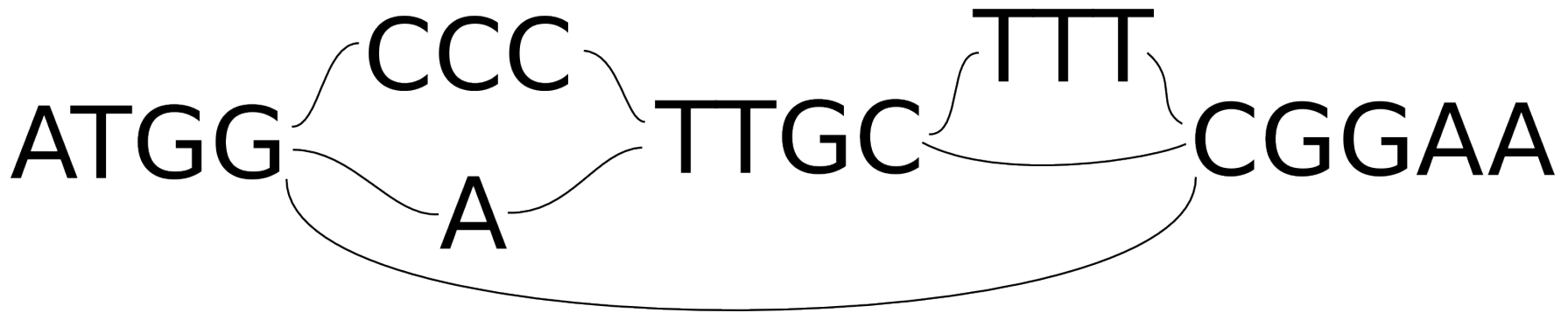## so why not use them in our analysis?

We resequence new genomes and compare them to a single reference haplotype.

To determine anything more than short variants, we must do everything *de novo.*

*If we could merge sequence and variation, we could detect known alleles of arbitrary scale and divergence with minimal cost.*
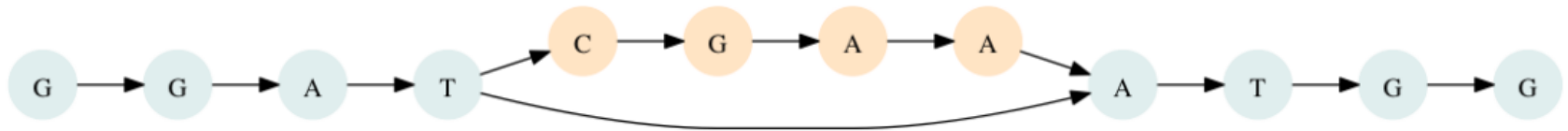
# Pan-genomes as graphs

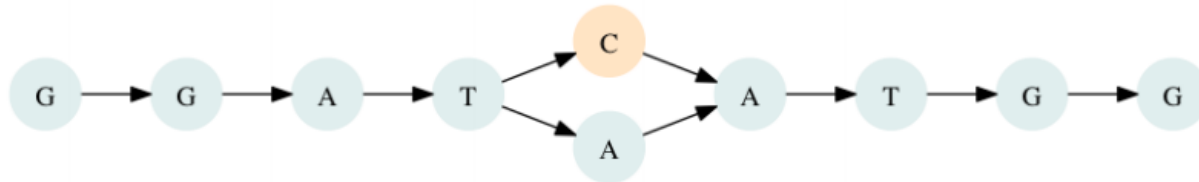We can combine sequence and variation using a *variant graphs*, or *graph reference*.


ATGG — CCC / A — TTGC — TTT — CGGAA

*This representation is directed (5' to 3'), and acyclic.

# Building the variant graph

# Local alignment against the graph

# Local alignment against the graph

# "Striped" string/DAG alignment

We improved performance of our aligner >10-fold by generalizing Farrar's striped Smith-Waterman algorithm to DAGs. *GSSW*

Data dependencies across DAG are limited to H and E vectors.

*Implemented using SSE2 instruction set.

$$E_{i,j} = \max \left\{ \begin{array}{c} E_{i,j-1} - G_{\text{ext}} \\ H_{i,j-1} - G_{\text{init}} \end{array} \right\}$$

$$F_{i,j} = \max \left\{ \begin{array}{c} F_{i-1,j} - G_{\text{ext}} \\ H_{i-1,j} - G_{\text{init}} \end{array} \right\}$$

$$H_{i,j} = \max \left\{ \begin{array}{c} 0 \\ E_{i,j} \\ F_{i,j} \\ H_{i-1,j-1} - W(q_i, d_j) \end{array} \right\}$$

copy H, E vectors

max of H, E vectors

Farrar, Bioinformatics (2006);  Rognes, BMC Bioinformatics (2011);  Zhao, PLoS One (2014)

# Seeding graph-based alignments



linear reference

**Test imperfectly-mapped reads against graph.**

graph reference

# Detecting variation on the graph



read supporting reference allele

read supporting variant allele

# Graph-based alignments with *glia*



BAM → *glia* → BAM → *freebayes*

1000G released alignments (bwa)

**"flattened" into reference space, with pseudo-reads of large insertions.

VCF

union alleles with exact breakpoints, used to build local graph

VCF

genotyped input alleles

# Application to 1000G variant integration



*Brian D'Astous*

# Unifying calls from many methods



*Tests from 1000G "phase3-like" chr20.

# *glia* reduces reference bias

# *glia* reduces reference bias



Improvement in observation support

**density**

Ratio between observations before and after realignment to graph of union variants

**Standard alignment is frustrated even by small variants!**

# Improving genotype likelihoods

Genotype Likelihood = P(data|genotype)

| SET | GRP | N | RR | RA | AA | ALT | ALL |
|---|---|---|---|---|---|---|---|
| SVM indels | UM | 6743 | 0.285 | 1.008 | 2.947 | 1.698 | 0.561 |
| SVM indels | BC* | 6743 | **0.034** | **0.673** | **0.245** | **0.521** | **0.129** |
| | | | | | | | |
| SNPs | BCM | 404270 | 0.029 | 1.373 | 0.445 | 1.093 | 0.111 |

\* includes glia realignment

Imputation of variant calls on chr20 via SHAPEIT 2.  Imputed results are tested against Complete Genomics samples in 1000 Genomes.

## We do as well for high-quality indels as SNPs!

*Olivier Delaneau, Androniki Menelaou, Jonathan Marchini*

# Mobile element detection



Using *glia+freebayes* to re-genotype an AluY insertion at 20:2252139 in the YRI population.  Insertion structure is estimated from split-read mappings.

# Alu genotyping efficiency

We re-call a substantial fraction of known (validated) Alus in 1000G low-coverage bwa alignments.

| set | re-genotyped Alus | % |
|---|---|---|
| Pilot 2 data (source) | 282 | 99.6 |
| PCR-free NA12878 | 281 | 99.3 |
| 5x NA12878 (low-coverage) | 173 | 61.1 |

Stewart et. al 2011.  A Comprehensive Map of Mobile Element Insertion Polymorphisms in Humans.  *PLoS Genetics.*

# Genotyping large deletions



Input SVs were generated by DELLY on deep, PCR-free samples used for validation in the 1000 Genomes Project.

When using this set as our reference, we can regenotype around 70% of such events in low-coverage samples.

# Performance using 1000G phase 3 SNPs and indels >1% frequency



**set**

- 10x.freebayes
- 10x.freebayes+glia.1kg
- 20x.freebayes
- 20x.freebayes+glia.1kg
- 30x.freebayes
- 30x.freebayes+glia.1kg
- 50x.freebayes
- 50x.freebayes+glia.1kg
- 5x.freebayes
- 5x.freebayes+glia.1kg

**Deep-coverage 100bp Illumina data on NA12878** was downsampled to 5, 10, 20, 30, and 50-fold. Calling by both freebayes and freebayes+glia (realigning to 1000G variants >1% MAF), and comparing the results to the **Genome In a Bottle truth set** demonstrates marked improvement in sensitivity, particularly at low-coverage.

| depth | snp AUC diff | indel AUC diff |
|-------|--------------|----------------|
| 5     | 6.02%        | 1.87%          |
| 10    | 1.07%        | 0.78%          |
| 20    | 0.26%        | 0.37%          |
| 30    | 0.08%        | 0.40%          |
| 50    | 0.02%        | 1.2%           |

# Questions?

…