

Structural Variant Calling

Michael Schatz

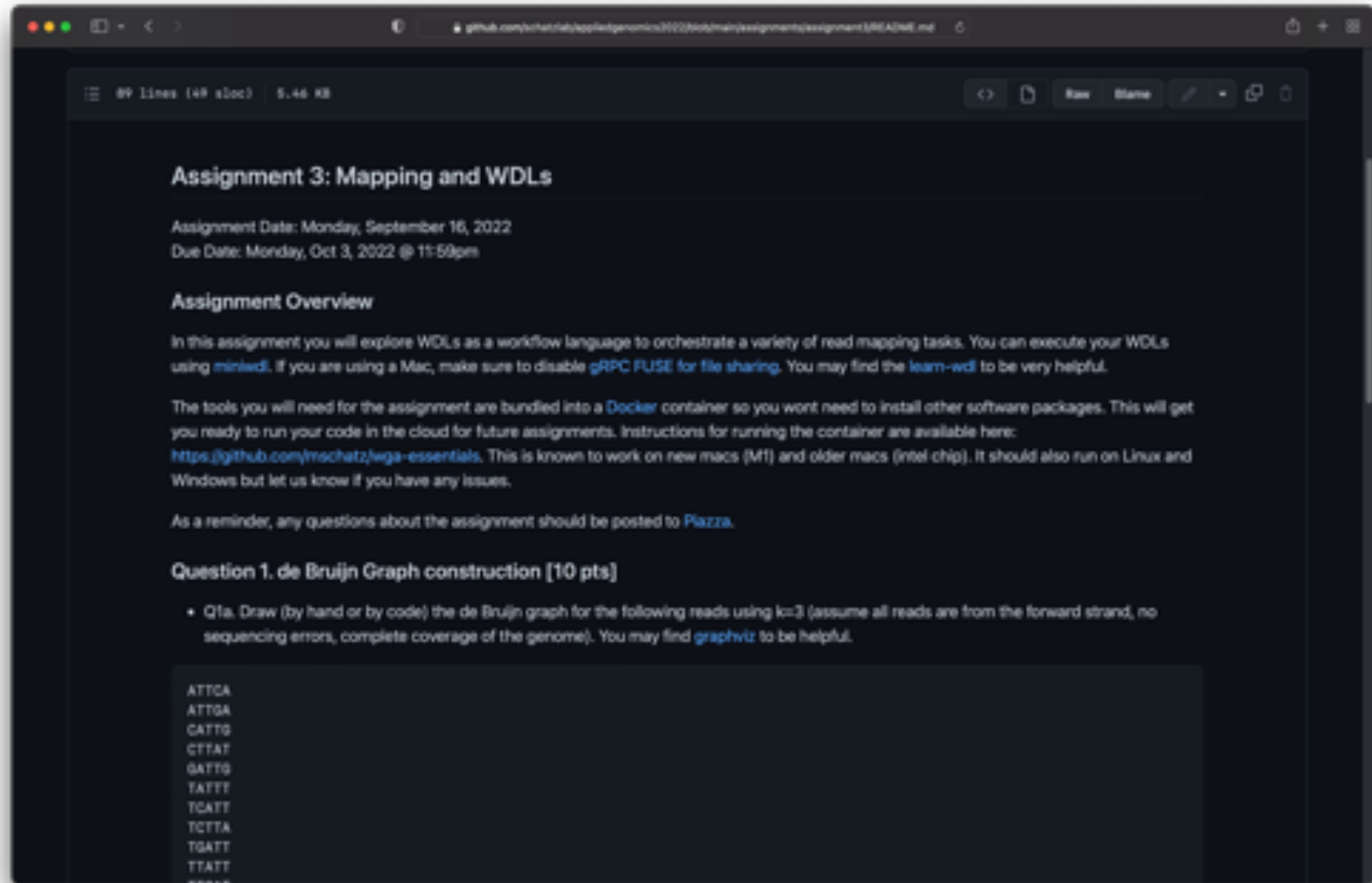
October 4, 2022

Lecture 11: Applied Comparative Genomics



Assignment 3: Mapping and WDL

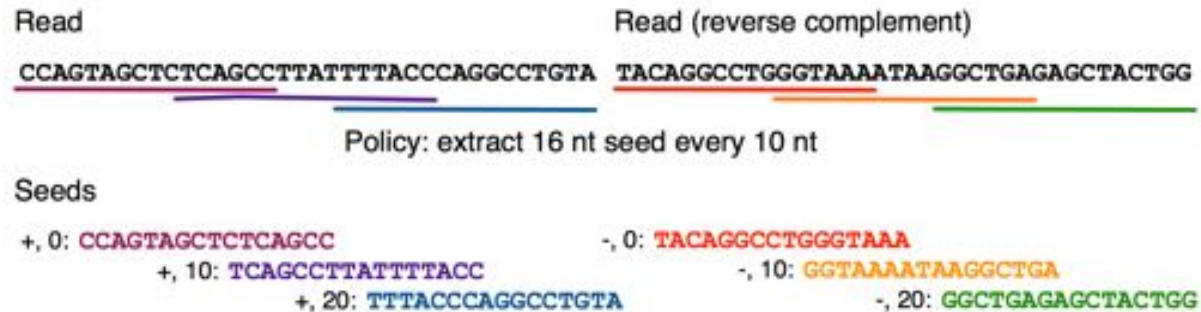
Due Monday Oct 3 by 11:59pm



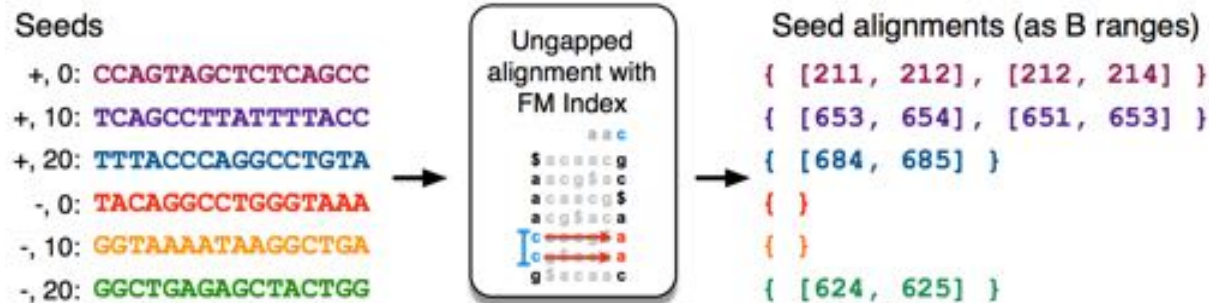
<https://github.com/schatzlab/appliedgenomics2022/tree/main/assignments/assignment3>
Check Piazza for questions!

Algorithm Overview

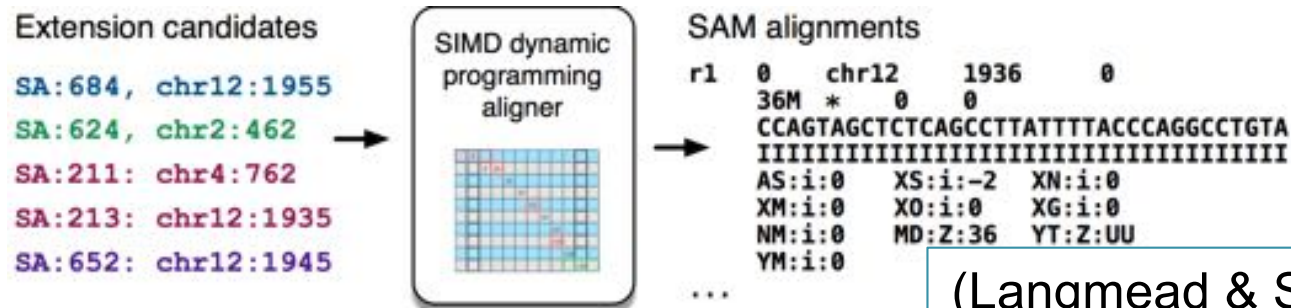
1. Split read into segments



2. Lookup each segment and prioritize



3. Evaluate end-to-end match



(Langmead & Salzberg, 2012)

Similarity metrics

- Hamming distance

- Count the number of substitutions to transform one string into another

MIKESCHATZ

| | x | | x x x x |

MICESHATZZ

5

- Edit distance

- The minimum number of substitutions, insertions, or deletions to transform one string into another

MIKESCHAT-Z

| | x | | x | | | x |

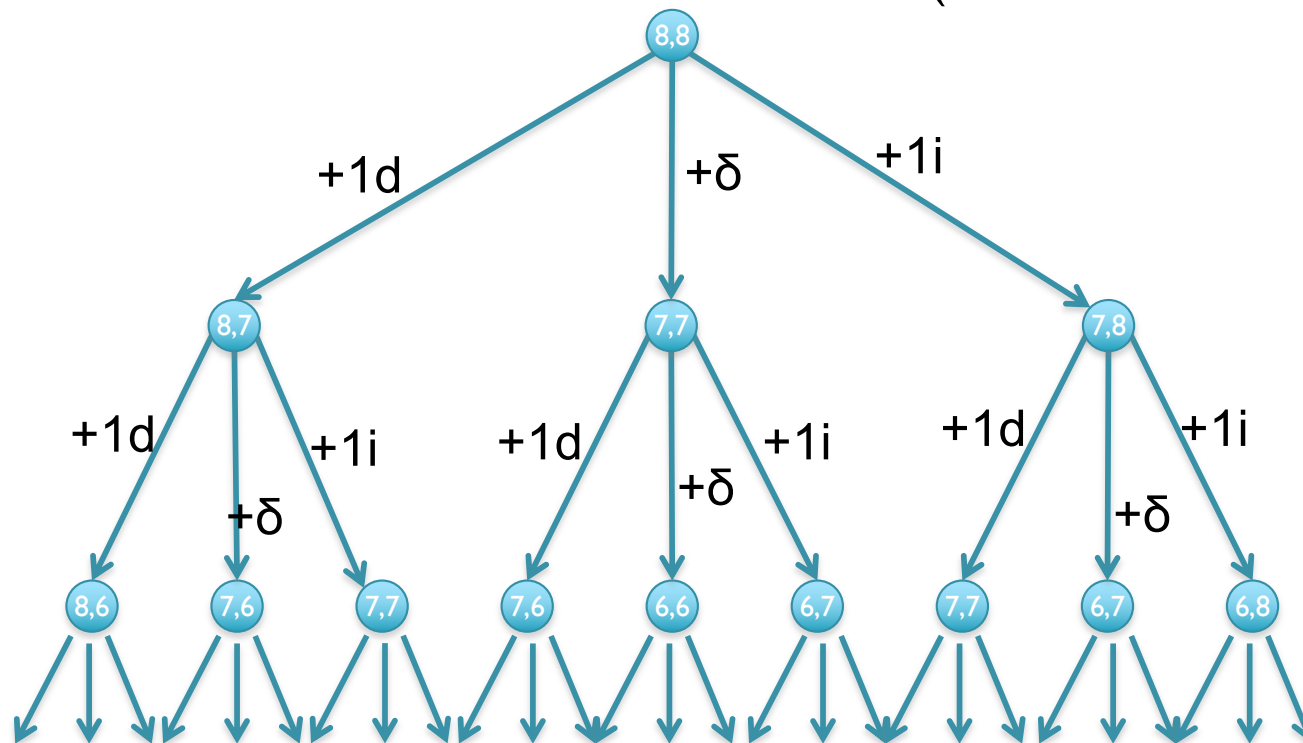
MICES-HATZZ

3

Recursive solution

- Computation of D is a recursive process.
 - At each step, we only allow matches, substitutions, and indels
 - $D(i,j)$ in terms of $D(i',j')$ for $i' \leq i$ and $j' \leq j$.

$$D(\text{AGCACACA}, \text{ACACACTA}) = \min\{D(\text{AGCACACA}, \text{ACACACT}) + 1, \\ D(\text{AGCACAC}, \text{ACACACTA}) + 1, \\ D(\text{AGCACAC}, \text{ACACACT}) + \delta(\text{A}, \text{A})\}$$



[What is the running time?]

Dynamic Programming Matrix

		A	C	A	C	A	C	T	A
	<u>0</u>	1	2	3	4	5	6	7	8
A	1	<u>0</u>	1	2	3	4	5	6	7
G	2	<u>1</u>	1	2	3	4	5	6	7
C	3	2	<u>1</u>	2	2	3	4	5	6
A	4	3	2	<u>1</u>	2	2	3	4	5
C	5	4	3	2	<u>1</u>	2	2	3	4
A	6	5	4	3	2	<u>1</u>	2	3	3
C	7	6	5	4	3	2	<u>1</u>	<u>2</u>	3
A	8	7	6	5	4	3	2	2	<u>2</u>

$$D[\text{AGCACACA}, \text{ACACACTA}] = 2$$

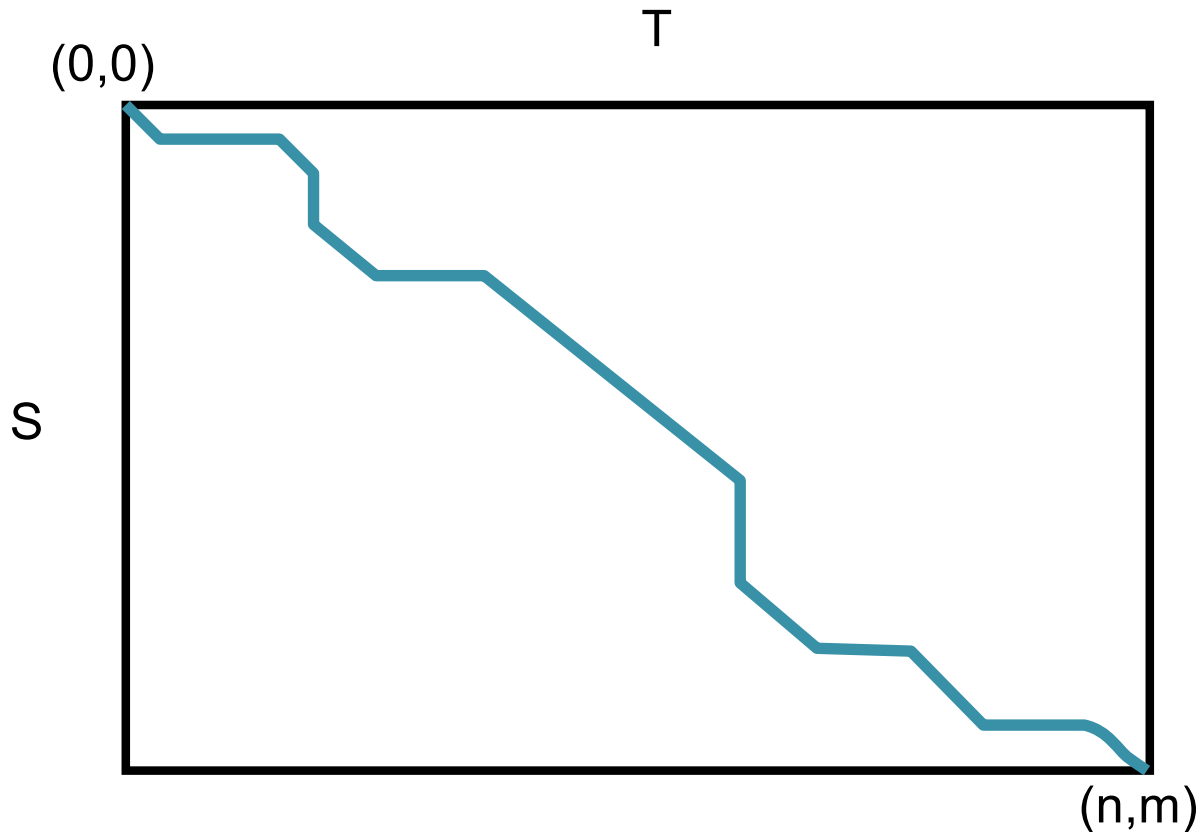
AGCACAC-A

| * | | | | * |

A-CACACTA

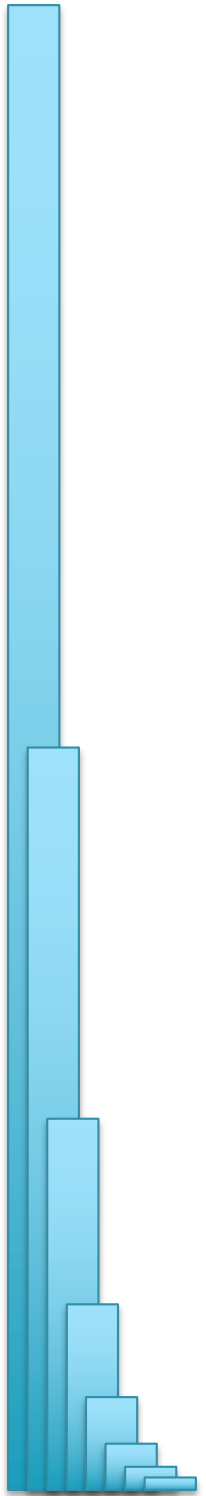
[Can we do it any better?]

Global Alignment Schematic

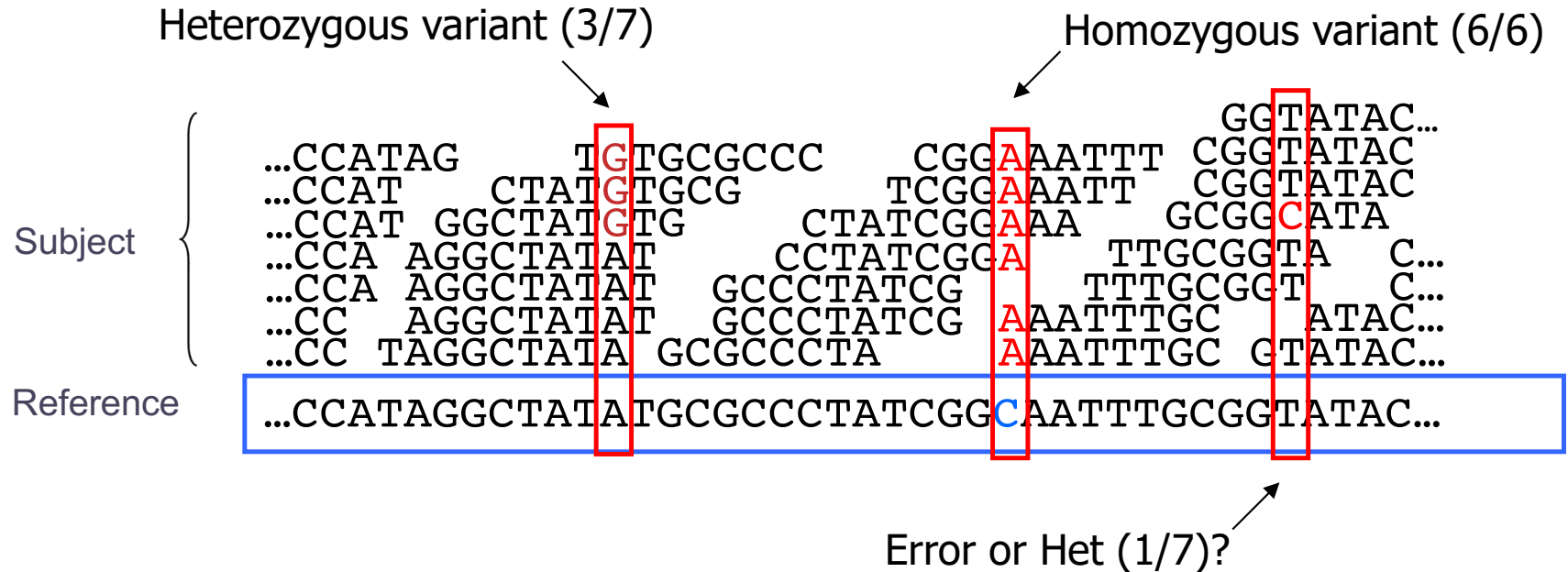


- A high quality alignment will stay close to the diagonal
 - If we are only interested in high quality alignments, we can skip filling in cells that can't possibly lead to a high quality alignment
 - Find the global alignment with at most edit distance d : $O(2dn)$

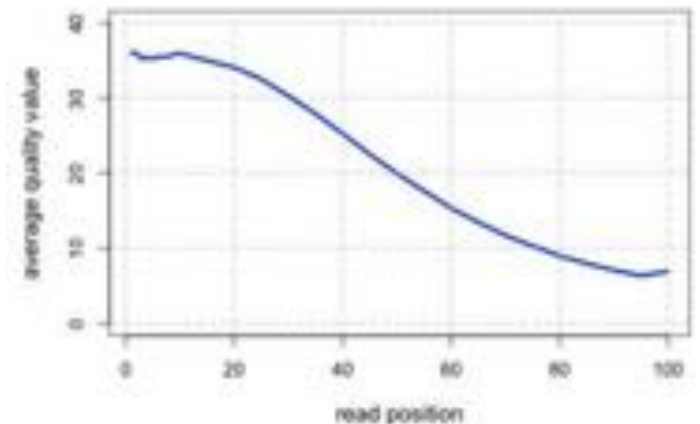
Variant Calling



Genotyping Theory



- If there were no sequencing errors, identifying SNPs would be very easy: any time a read disagrees with the reference, it must be a variant!
- Sequencing instruments make mistakes
 - Quality of read decreases over the read length
- A single read differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times



The Binomial Distribution: Adventures in Coin Flipping

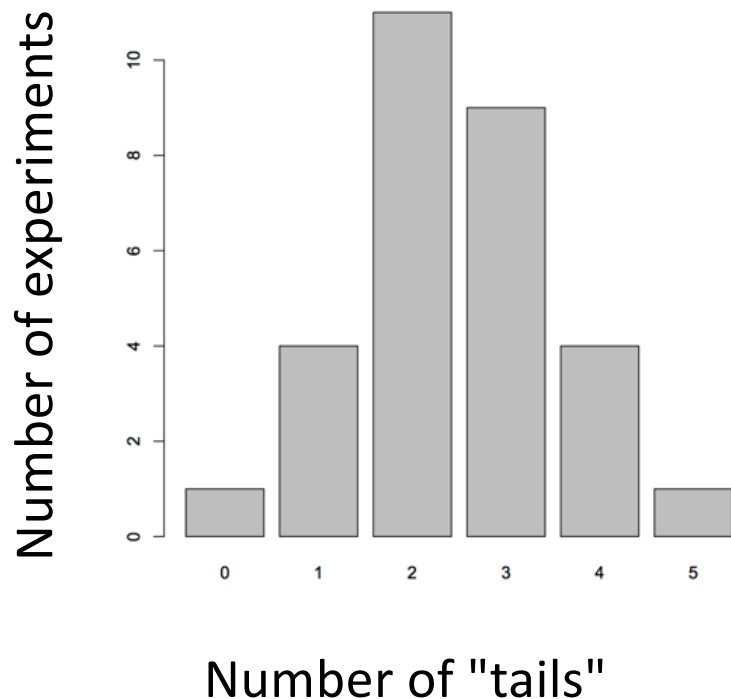


$P(\text{heads}) = 0.5$



$P(\text{tails}) = 0.5$

What is the distribution of tails
(alternate alleles) do we expect to see
after 5 tosses (sequence reads)?



R code:

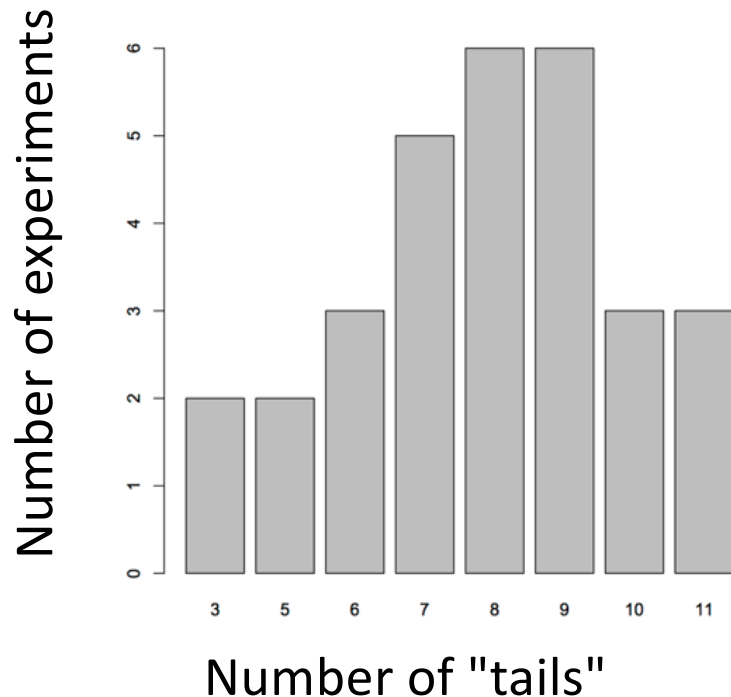
```
barplot(table(rbinom(30, 5, 0.5)))
```

30 experiments (students tossing coins)

5 tosses each

Probability of Tails

What is the distribution of tails
(alternate alleles) do we expect to see
after **15** tosses (sequence reads)?



R code:

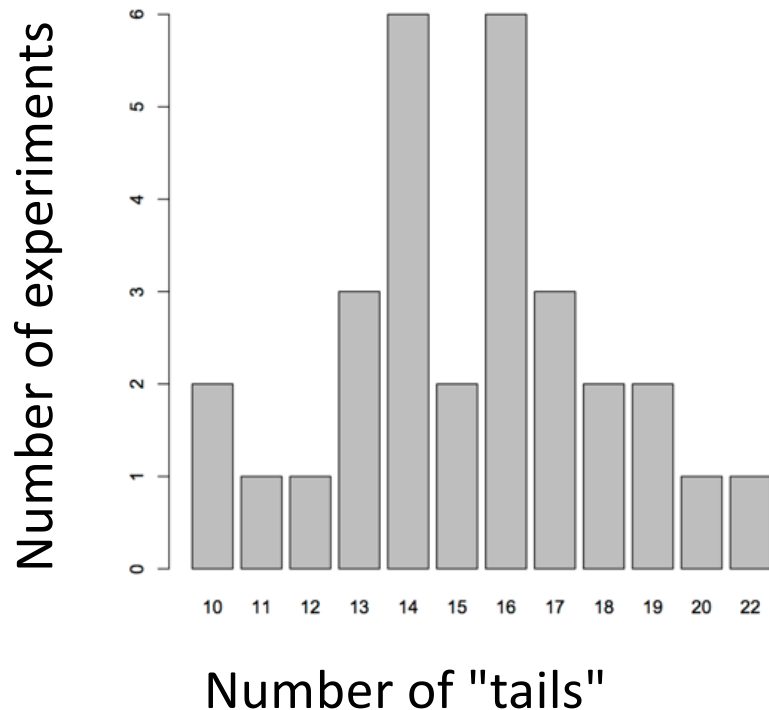
```
barplot(table(rbinom(30, 15, 0.5)))
```

30 experiments (students tossing coins)

15 tosses each

Probability of Tails

What is the distribution of tails
(alternate alleles) do we expect to see
after 30 tosses (sequence reads)?



R code:

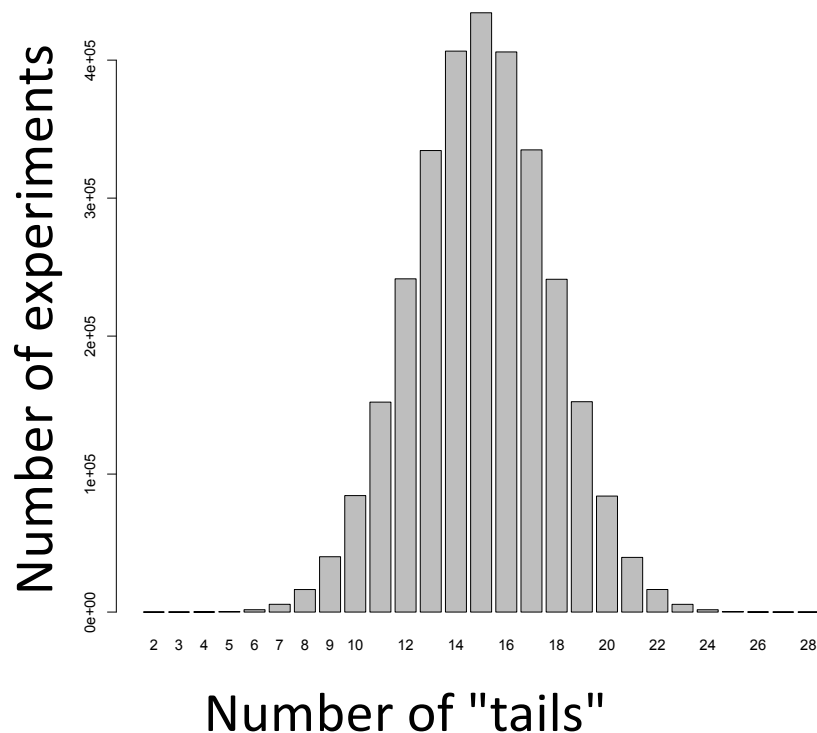
```
barplot(table(rbinom(30, 30, 0.5)))
```

30 experiments (students tossing coins)

30 tosses each

Probability of Tails

What is the distribution of tails
(alternate alleles) do we expect to see
after 30 tosses (sequence reads)?



R code:

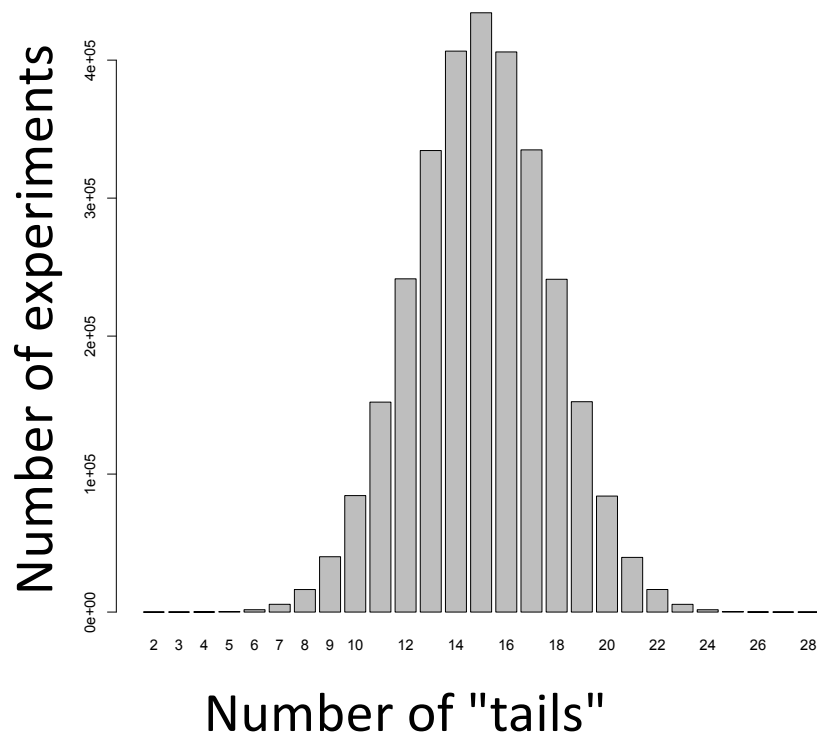
```
barplot(table(rbinom(3e6, 30, 0.5)))
```

3M experiments (students tossing coins)

30 tosses each

Probability of Tails

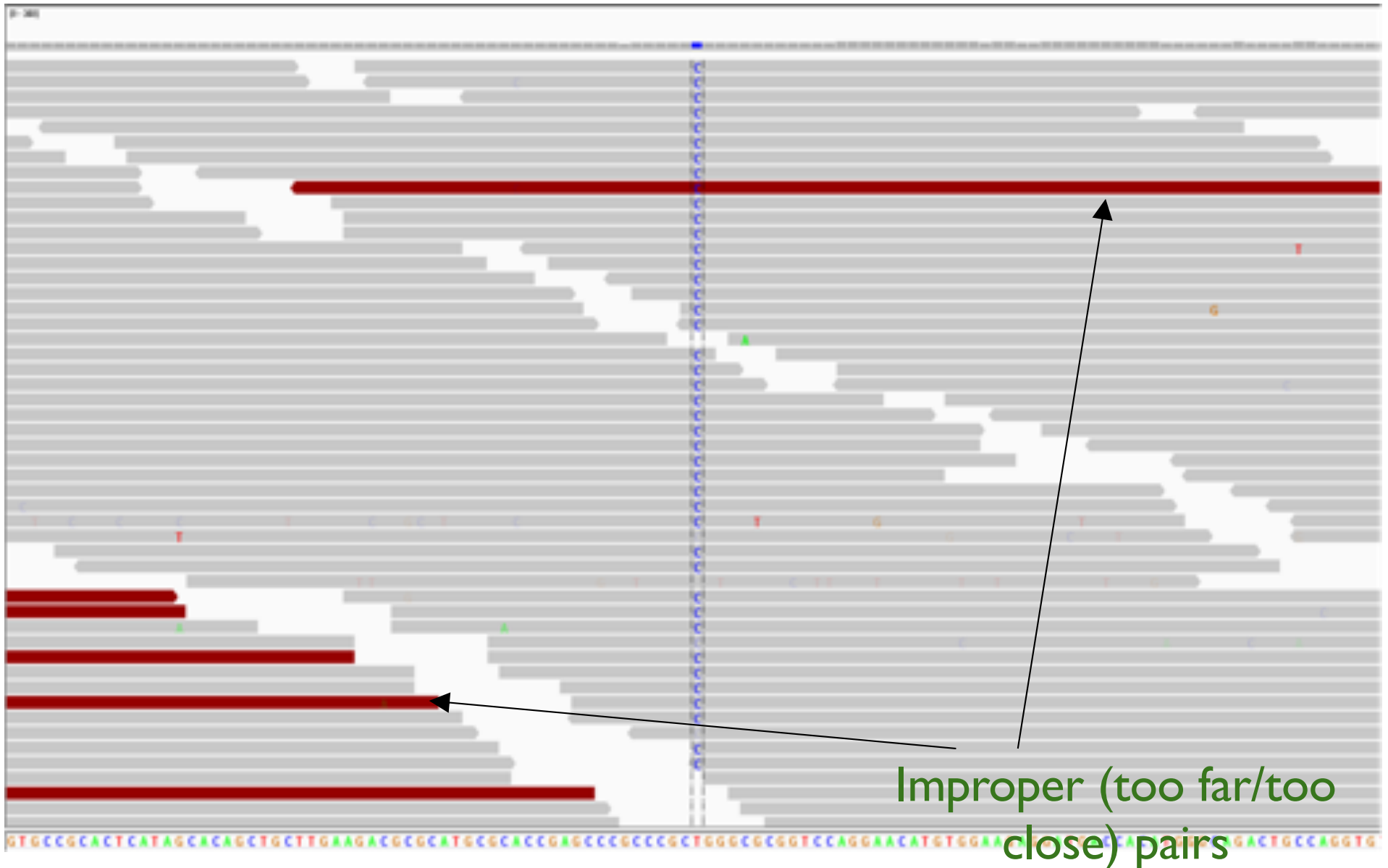
So, with 30 tosses (reads), we are much more likely to see an even mix of alternate and reference alleles at a heterozygous locus in a genome



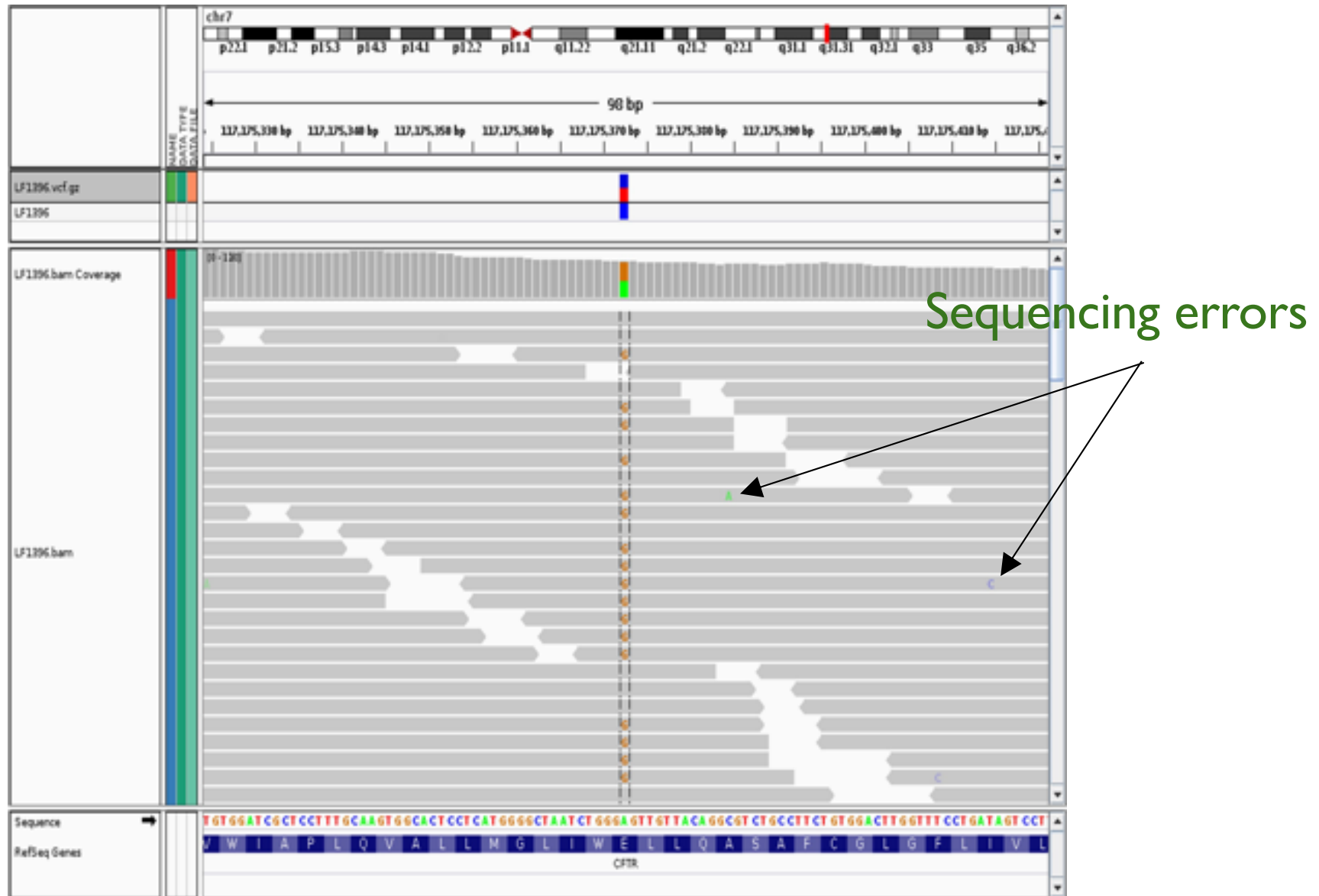
This is why at least a "30X" (30 fold sequence coverage) genome is recommended: it confers sufficient power to distinguish heterozygous alleles and from mere sequencing errors

$$P(3/30 \text{ het}) <?> P(3/30 \text{ err})$$

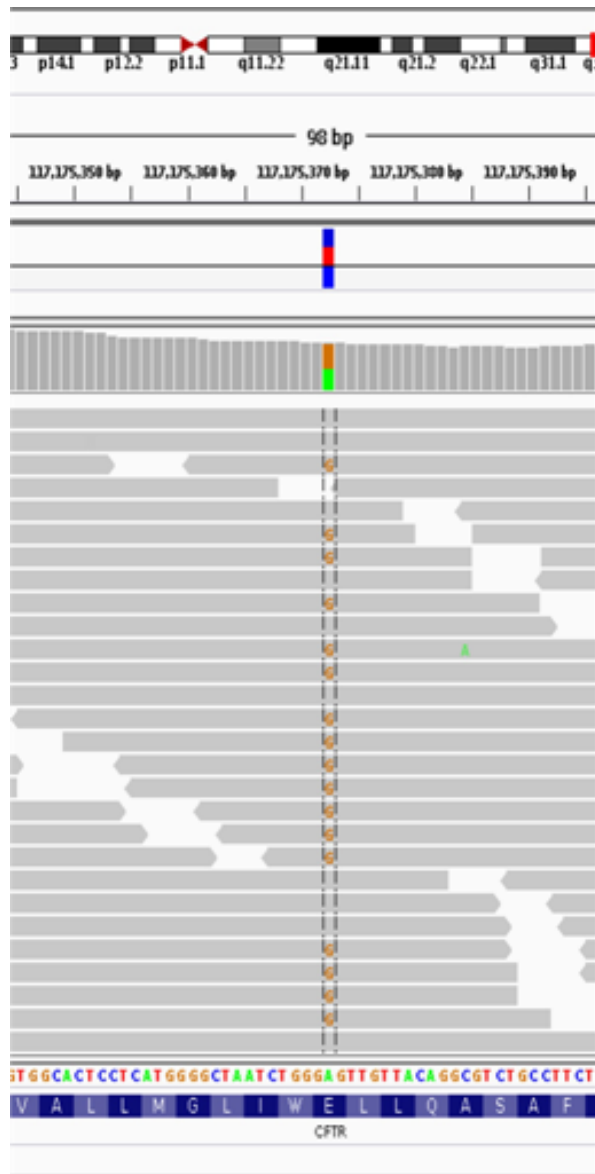
Homozygous for the "C" allele



Sequencing errors fall out as noise (most of the time)



What information is needed to decide if a variant exists?



- Depth of coverage at the locus
- Bases observed at the locus
- The base qualities of each allele
- The strand composition
- Mapping qualities
- Proper pairs?
- Expected polymorphism rate

PolyBayes: The first statistically rigorous variant detection tool.

letter

© 1999 Nature America Inc. • <http://genetics.nature.com>

A general approach to single-nucleotide polymorphism discovery

Gabor T. Marth¹, Ian Korf¹, Mark D. Yandell¹, Raymond T. Yeh¹, Zhijie Gu², Hamideh Zakeri², Nathan O. Stitzel¹, LaDeana Hillier¹, Pui-Yan Kwok² & Warren R. Gish¹

Its main innovation was the use of Bayes's theorem

Netscape: PolyBayes Web site


File Edit View Go Communicator Help

Back Forward Reload Home Search Netscape Print Security Shop Stop

Bookmarks Location: <http://genome.wustl.edu/gsc/Informatics/polybayes/> What's Related

WebMail Radio People Yellow Pages Download Calendar Channels

Site map

 **PolyBayes**

[Home](#) [About](#) [Features](#) [Analysis](#) [Publication](#) [FAQ](#) [Authors](#) [Site map](#) [Documentation](#) [Download](#) [Links](#) [Contact](#)

14	-	30
15	-	30
16	-	30
17	-	30
18	-	30
19	A	40
20	G	38

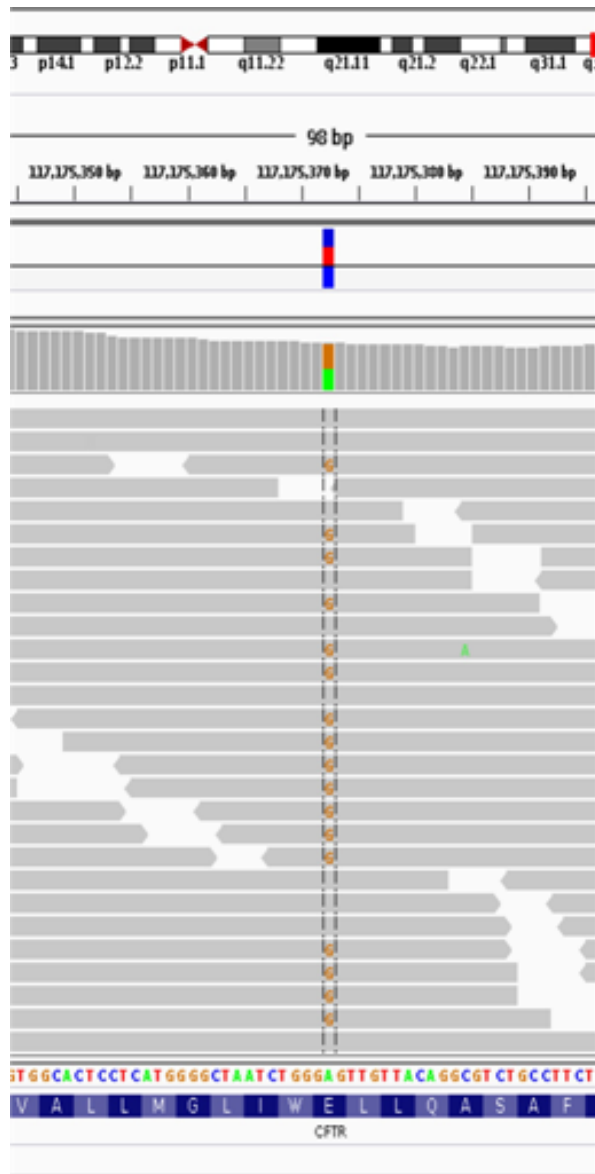
Results

Description	Symbol	Value
Probability of SNP	P(SNP)	0.853076599574195
Most likely variation	VAR	A/G
Probability of variation	P(VAR)	0.853003076184499
Alignment depth	D	2

Comments to: Gabor Marth, marth@genetics.wustl.edu, Washington University Genome Sequencing Center
Last modified: Mon Feb 12 17:06:10 2001

100%

Bayesian SNP calling



$$P(\text{SNP} | \text{Data}) = \frac{P(\text{Data} | \text{SNP}) * P(\text{SNP})}{P(\text{Data})}$$

PolyBayes: The first statistically rigorous variant detection tool.

letter

© 1999 Nature America Inc. • <http://genetics.nature.com>

A general approach to single-nucleotide polymorphism discovery

Gabor T. Marth¹, Ian Korf¹, Mark D. Yandell¹, Raymond T. Yeh¹, Zhijie Gu², Hamideh Zakeri², Nathan O. Stitzel¹, LaDeana Hillier¹, Pui-Yan Kwok² & Warren R. Gish¹

Bayesian
posterior
probability

Base call +
Base quality

Expected (prior)
polymorphism rate

$$P(SNP) = \sum_{\text{all variable } S} \frac{\frac{P(S_1 | R_1)}{P_{Prior}(S_1)} \cdots \frac{P(S_N | R_N)}{P_{Prior}(S_N)} \cdot P_{Prior}(S_1, \dots, S_N)}{\sum_{S_{i_1} \in [A, C, G, T]} \cdots \sum_{S_{i_N} \in [A, C, G, T]} \frac{P(S_{i_1} | R_1)}{P_{Prior}(S_{i_1})} \cdots \frac{P(S_{i_N} | R_N)}{P_{Prior}(S_{i_N})} \cdot P_{Prior}(S_{i_1}, \dots, S_{i_N})}$$

Probability of observed base composition
(should model sequencing error rate)

PolyBayes: The first statistically rigorous variant detection tool.

letter

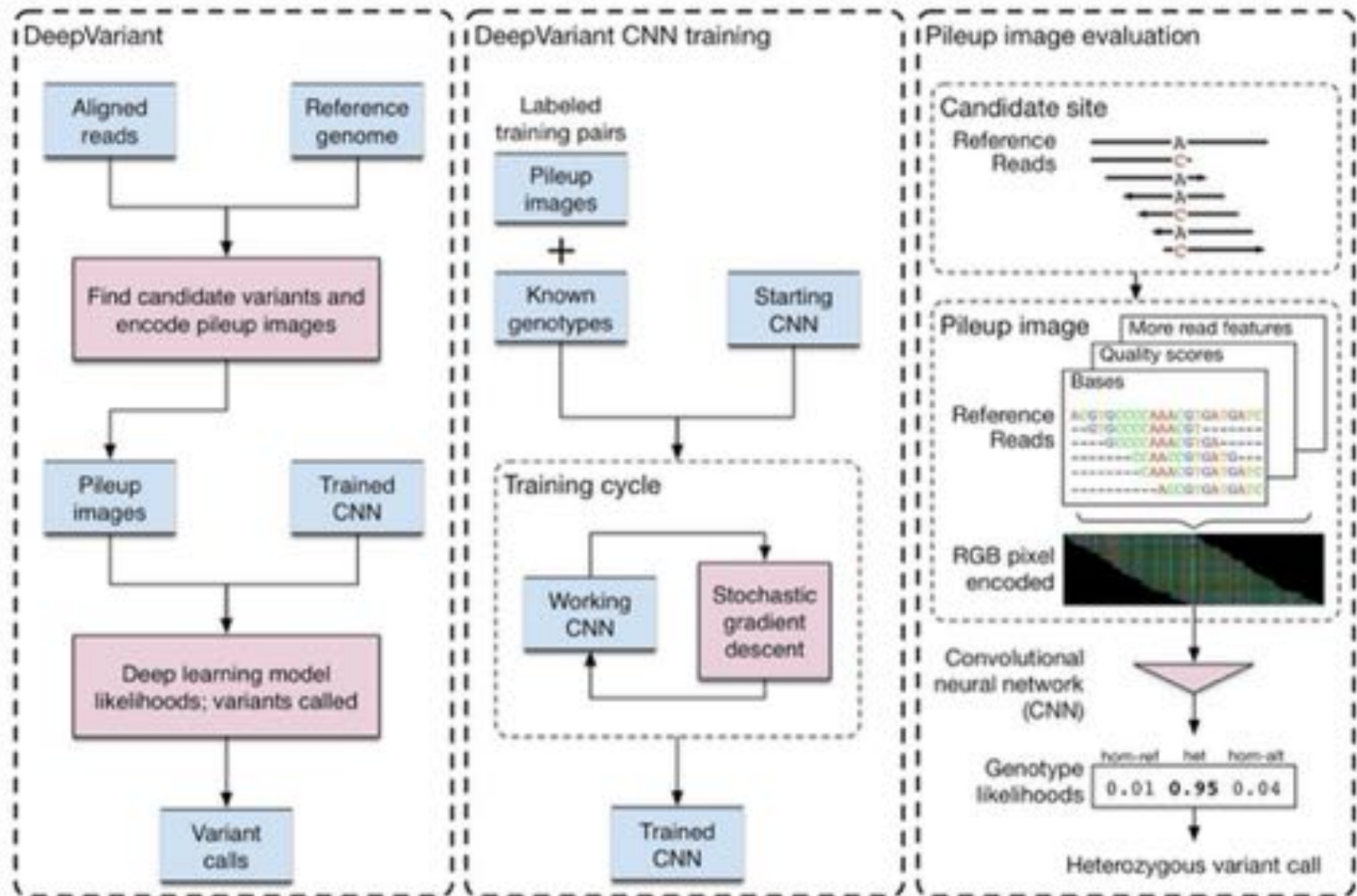
© 1999 Nature America Inc. • <http://genetics.nature.com>

A general approach to single-nucleotide polymorphism discovery

Gabor T. Marth¹, Ian Korf¹, Mark D. Yandell¹, Raymond T. Yeh¹, Zhijie Gu², Hamideh Zakeri²,
Nathan O. Stitzel¹, LaDeana Hillier¹, Pui-Yan Kwok² & Warren R. Gish¹

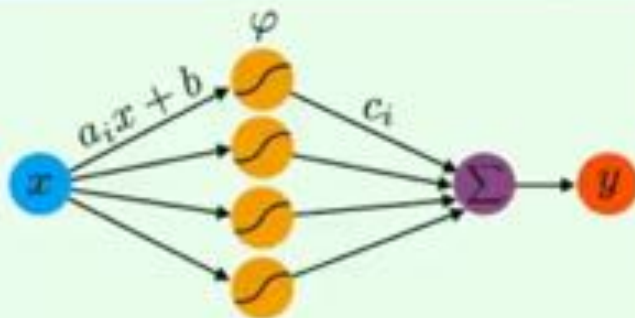
This Bayesian statistical framework has been adopted by many modern statistical SNP/INDEL callers such as FreeBayes, GATK, and samtools

DeepVariant



Creating a universal SNP and small indel variant caller with deep neural networks

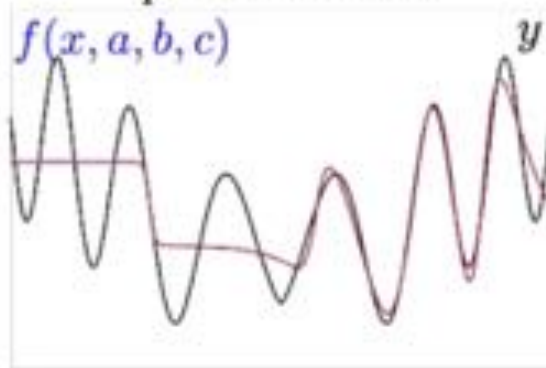
Poplin et al. (2018) Nature Biotechnology. <https://www.nature.com/articles/nbt.4235>



1 hidden layer perceptron:

$$y \approx f(x, a, b, c) \stackrel{\text{def.}}{=} \sum_{i=1}^p c_i \varphi(a_i x + b_i)$$

$p = 6$ neurons



$p = 20$ neurons



Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of a real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

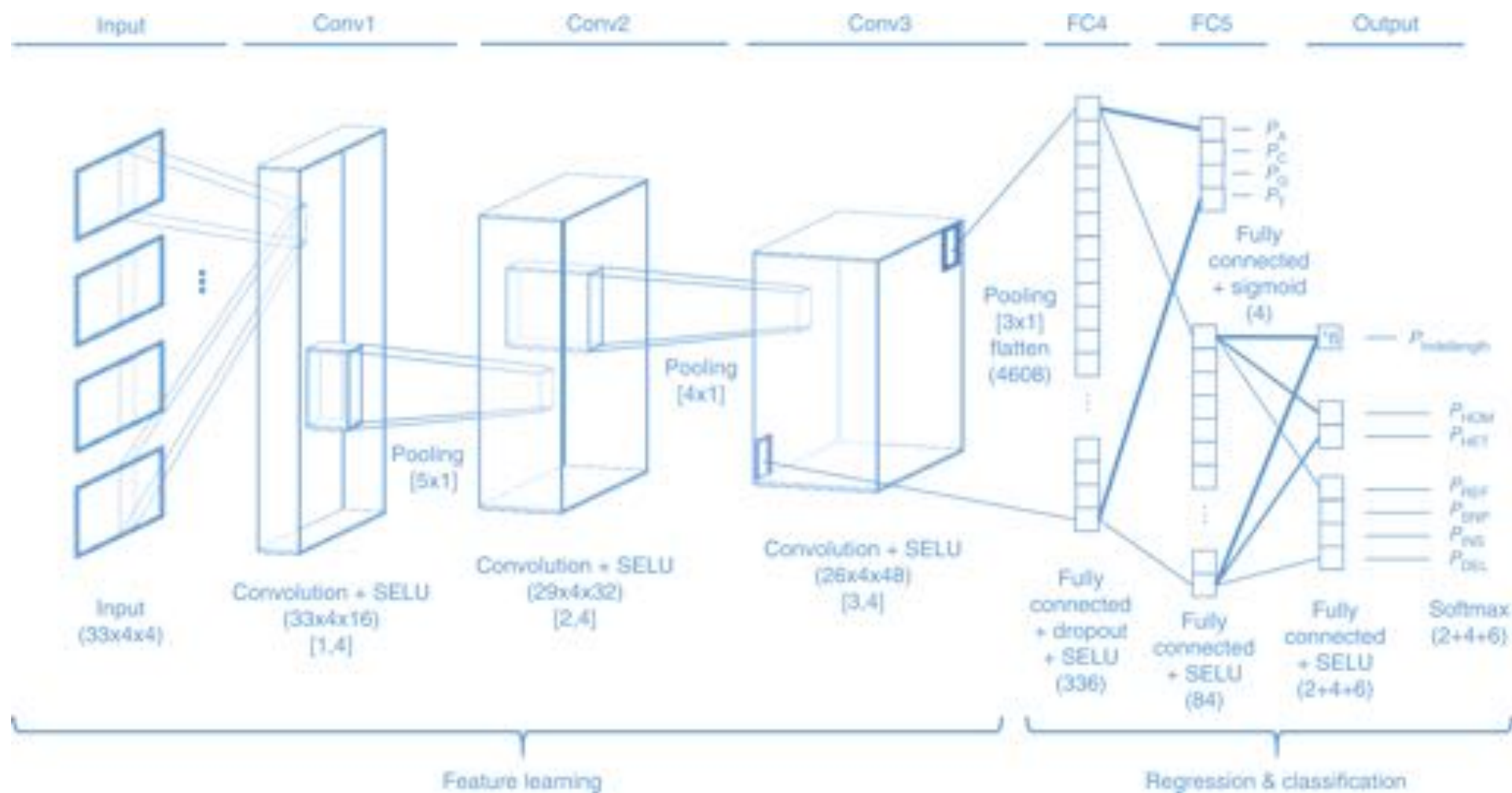
Key words. Neural networks, Approximation, Completeness.



George Cybenko

Approximation by superpositions of a sigmoidal function

Cybenko, G. (1989) Mathematics of Control Signal Systems doi: 10.1007/BF02551274



A multi-task convolutional deep neural network for variant calling in single molecule sequencing

Luo et al. (2019) Nature Communication. <https://doi.org/10.1038/s41467-019-09025-z>

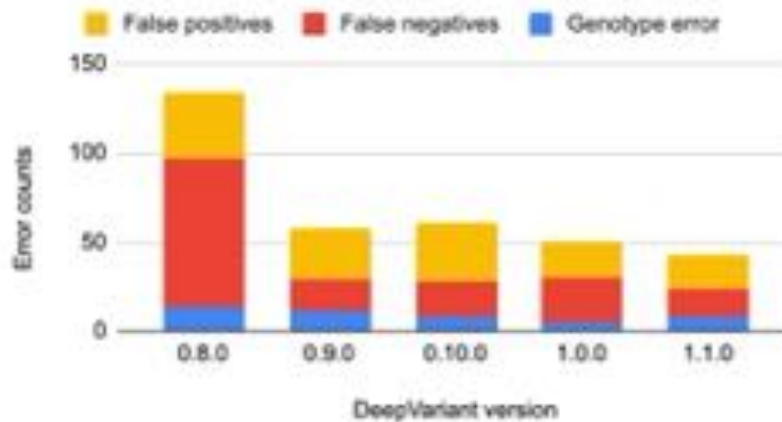
WGS SNP error counts (HG003)



WGS INDEL error counts (HG003)



PacBio SNP error counts (HG003)



PacBio INDEL error counts (HG003)

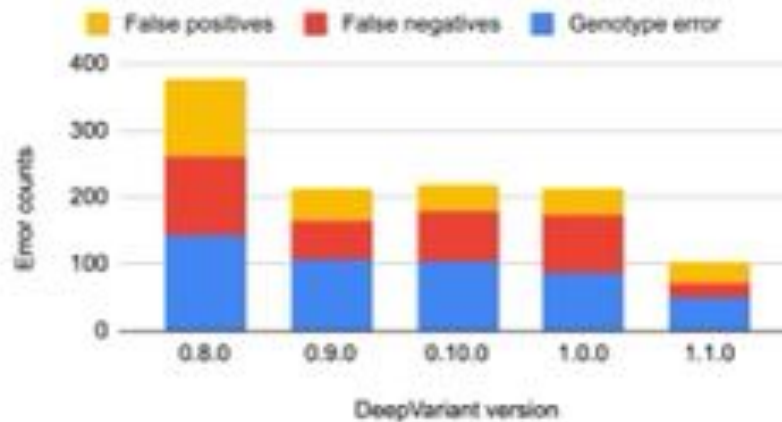


Figure 2: Error counts over the years for HG003. For Illumina WGS, we use a HiSeqX PCR-free dataset at 30x coverage. For PacBio HiFi, we use the same BAM as the one in our [case study](#).

DeepVariant over the years

<https://google.github.io/deepvariant/posts/2021-06-08-deepvariant-over-the-years/>

VCF Format

Example

VCF header

```
##fileformat=VCFv4.8
##fileDate=20180707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
```

Mandatory header lines

Optional header lines (meta-data about the annotations in the VCF body)

Body

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	SAMPLE1	SAMPLE2
1	1	.	ACG	A,AT	.	PASS	.	GT:DP	1/2:13	0/0:29
1	2	rs1	C	T,CT	.	PASS	H2;AA=T	GT:GQ	0/1:100	2/2:70
1	5	.	A	G	.	PASS	.	GT:GQ	1/0:77	1/1:93
1	100	.	T		.	PASS	SVTYPE=DEL;END=300	GT:GQ:DP	1/1:12:3	0/0:20

Deletion

SNP

Large SV

Insertion

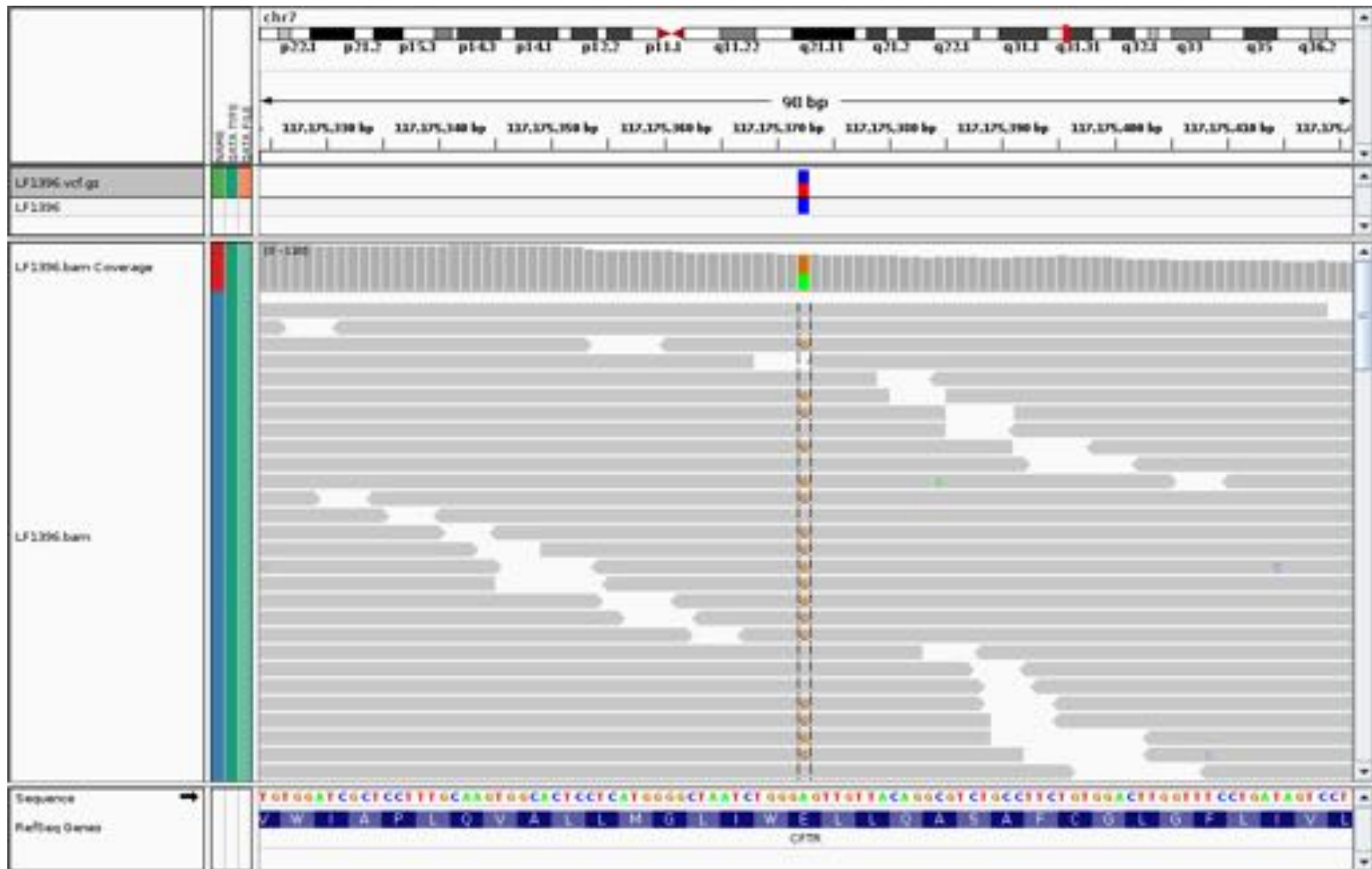
Other event

Reference alleles (GT=0)

Alternate alleles (GT>0 is an index to the ALT column)

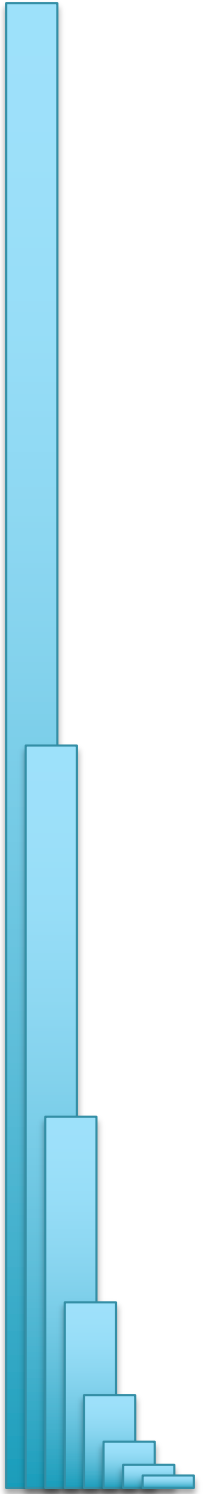
Phased data (G and C above are on the same chromosome)

VCF Format

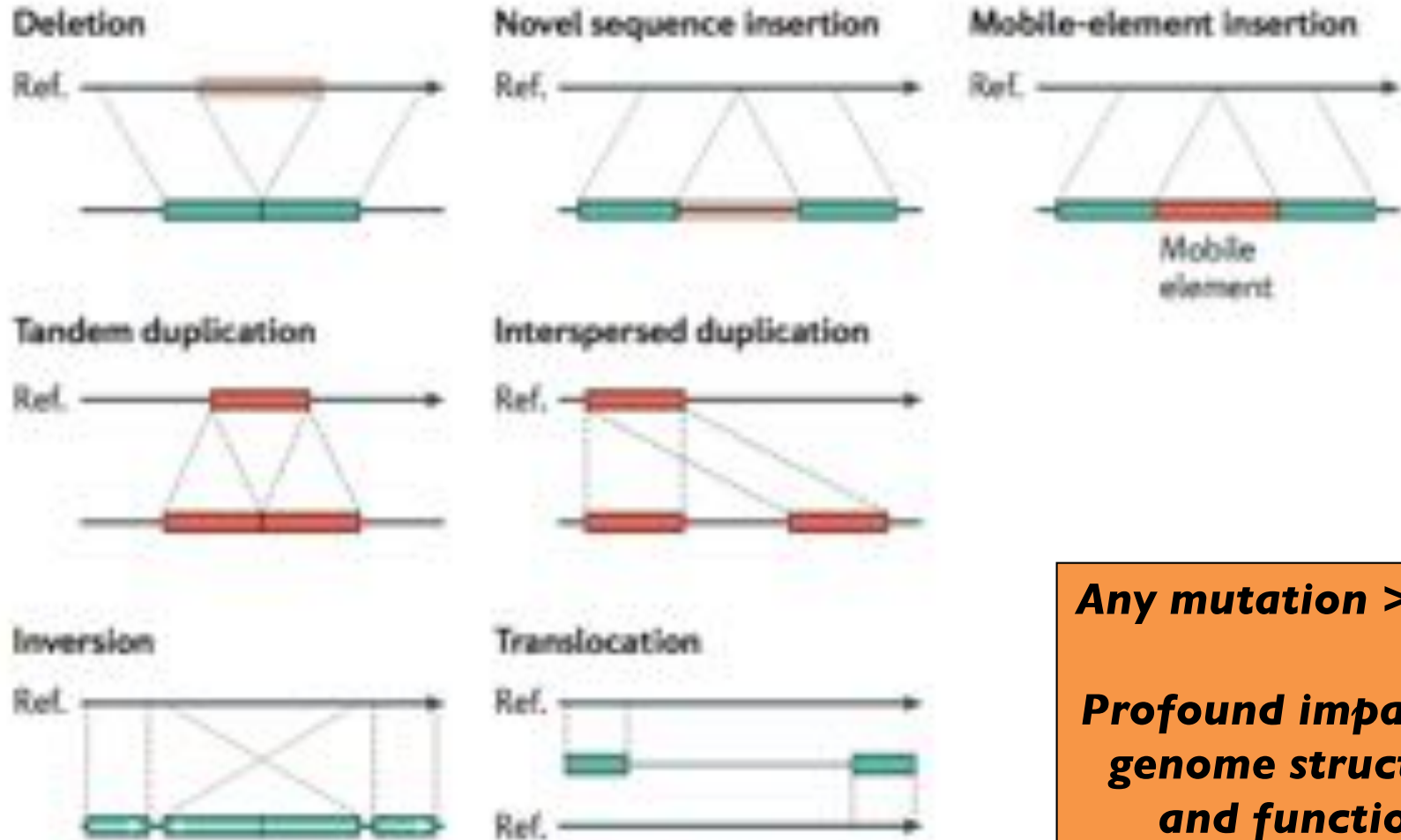


#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	LF1396
chr7	117175373	.	A	G	90	PASS	AF=0.5	GT	0/1

What about indels & structural variants



Structural Variations



Any mutation >50bp

**Profound impact on
genome structure
and function**

Genome structural variation discovery and genotyping

Alkan, C, Coe, BP, Eichler, EE (2011) *Nature Reviews Genetics*. May;12(5):363-76. doi: 10.1038/nrg2958.

Early 2000s dogma: SNPs account for most human genetic variation



Discovery of abundant copy-number variation

Science, July 2004

Large-Scale Copy Number Polymorphism in the Human Genome

Jonathan Sebat,¹ B. Lakshmi,¹ Jennifer Troge,¹ Joan Alexander,¹ Janet Young,² Pär Lundin,³ Susanne Månér,³ Hillary Massa,² Megan Walker,² Maoyen Chi,¹ Nicholas Navin,¹ Robert Lucito,¹ John Healy,¹ James Hicks,¹ Kenny Ye,⁴ Andrew Reiner,¹ T. Conrad Gilliam,⁵ Barbara Trask,² Nick Patterson,⁶ Anders Zetterberg,³ Michael Wigler^{1*}

76 CNVs in 20 individuals
70 genes

Nature Genetics, Aug. 2004

Detection of large-scale variation in the human genome

A John Iafrate^{1,2}, Lars Feuk³, Miguel N Rivera^{1,2}, Marc L Listewnik¹, Patricia K Donahoe^{2,4}, Ying Qi³, Stephen W Scherer^{3,5} & Charles Lee^{1,2,5}

255 CNVs in 55 individuals
127 genes

- 331 CNVs, only 11 in common
- Half observed in only 1 individual
- Impact "plenty" of genes
- Correlated with segmental duplications in the reference genome

Why is structural variation relevant / important?

- ▶ They are common and affect a large fraction of the genome
 - ▶ In total, SVs impact more base pairs than all single-nucleotide differences.
- ▶ They are a major driver of genome evolution
 - ▶ Speciation can be driven by rapid changes in genome architecture
 - ▶ Genome instability and aneuploidy: hallmarks of solid tumor genomes