# Variant Calling

Michael Schatz
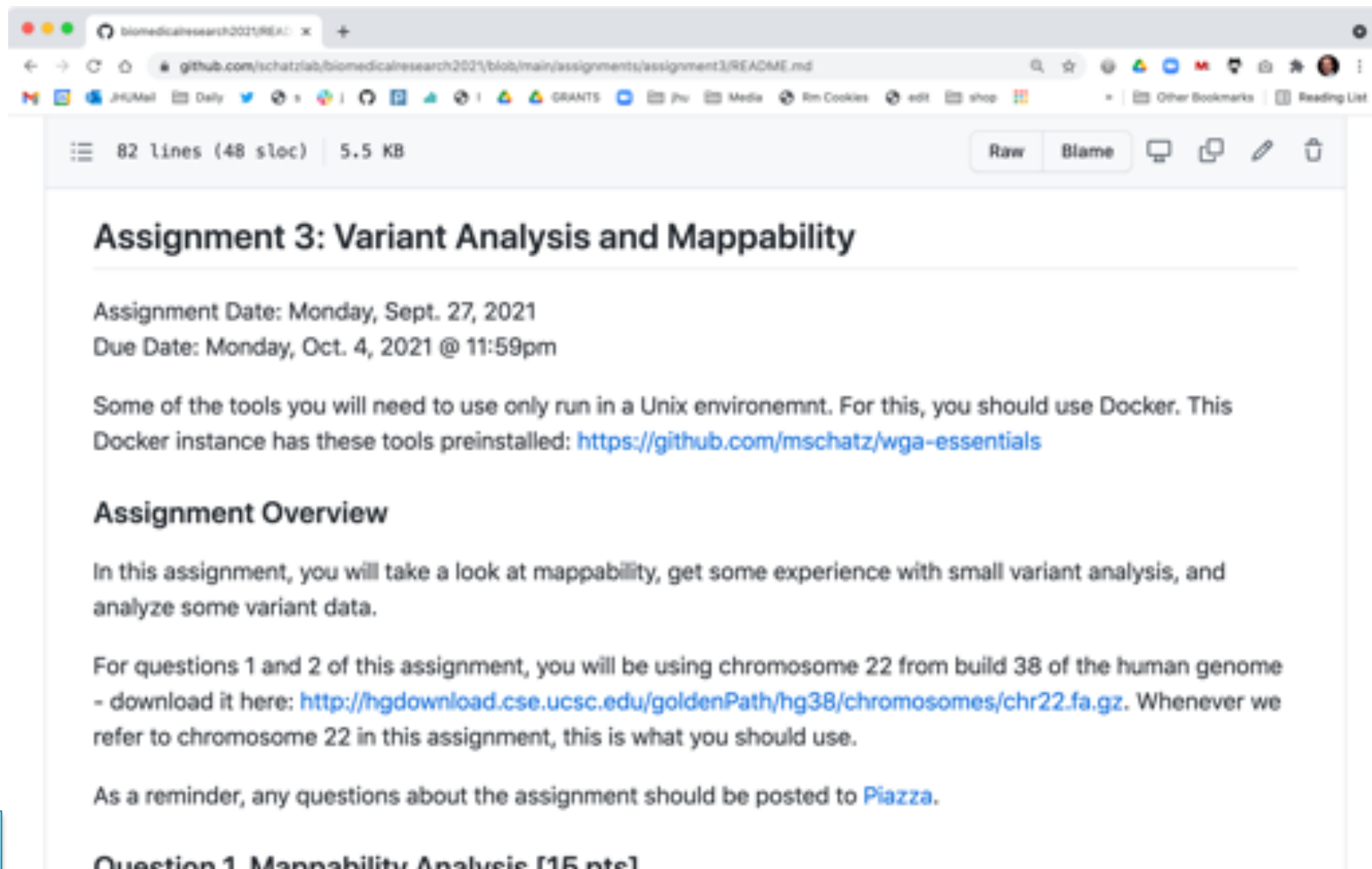
Sept 29, 2021
Lecture 9: Computational Biomedical Research

# Assignment 3: Variant Analysis & Mappability
# Due Oct 4 @ 11:59pm



https://github.com/schatzlab/biomedicalresearch2021

# Monday's class

Browser window showing the National Human Genome Research Institute website:

**URL:** genome.gov/event-calendar/Bold-Predictions-for-Human-Genomics-by-2030

**Page title:** Bold Predictions for Human G...

COVID-19   CDC health information   NIH research   Español

**NIH** National Human Genome Research Institute

Begin your search here

About Genomics    Research Funding    Research at NHGRI    Health    Careers & Training    News & Events    About NHGRI

## Upcoming

Session 8 - October 4, 2021, 3 p.m. to 4:30 p.m.

Bold Prediction #8: A person's complete genome sequence along with informative annotations can be securely and readily accessible on their smartphone.

Speakers:

Michael Schatz, Ph.D.
Johns Hopkins University & Cold Spring Harbor Laboratory

Gillian Hooker, Ph.D., ScM, LCGC
Concert Genetics

Moderator:

Sarah Bates, M.S.
NHGRI

Session 9 - November 1, 2021, 3 p.m. to 4:30 p.m.

Bold Prediction #9: Individuals from ancestrally diverse backgrounds will benefit equitably from advances in human genomics.

Speakers:

Registration: bit.ly/2XXhLYJ

# Read Mapping

# Personal Genomics

How does your genome compare to the reference?



Heart Disease

Cancer

Presidential smile

# Similarity metrics

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

$$\begin{array}{c}\texttt{MIKESCHATZ}\\\texttt{||x||xxxx|}\\\texttt{MICESHATZZ}\\5\end{array}$$

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

$$\begin{array}{c}\texttt{MIKESCHAT-Z}\\\texttt{||x||x|||x|}\\\texttt{MICES-HATZZ}\\3\end{array}$$

# Reverse Engineering Edit Distance

$$D(AGCACACA, ACACACTA) = ?$$

Imagine we already have the optimal alignment of the strings, the last column can only be 1 of 3 options:

```
...M              ...I              ...D
...A              ...-              ...A
...A              ...A              ...-
```

The optimal alignment of last two columns is then 1 of 9 possibilities

```
...MM  ...IM  ...DM        ...MI  ...II  ...DI        ...MD  ...ID  ...DD
...CA  ...-A  ...CA        ...A-  ...--  ...A-        ...CA  ...-A  ...CA
...TA  ...TA  ...-A        ...TA  ...TA  ...-A        ...A-  ...A-  ...--
```

The optimal alignment of the last three columns is then 1 of 27 possibilities...

```
    ...M...           ...I...           ...D...

    ...X...           ...-...           ...X...

    ...Y...           ...Y...           ...-...
```

Eventually spell out every possible sequence of {I,M,D}

# 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' ≤ i and j' ≤ 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

- We could code this as a recursive function call...

   ...with an exponential number of function evaluations

- There are only $(n+1) \times (m+1)$ pairs $i$ and $j$
  - We are evaluating $D(i,j)$ multiple times

- Compute $D(i,j)$ bottom up.
  - Start with smallest $(i,j) = (1,1)$.
  - Store the intermediate results in a table.
    - Compute $D(i,j)$ *after* $D(i-1,j)$, $D(i,j-1)$, and $D(i-1,j-1)$

# Recurrence Relation for D

Find the edit distance (minimum number of operations to convert one string into another) in $O(mn)$ time

- Base conditions:
  - $D(i,0) = i$, for all $i = 0,...,n$
  - $D(0,j) = j$, for all $j = 0,...,m$

- For $i > 0$, $j > 0$:

$$D(i,j) = \min \begin{cases} D(i-1,j) + 1, & \text{// align 0 chars from S, 1 from T} \\ D(i,j-1) + 1, & \text{// align 1 chars from S, 0 from T} \\ D(i-1,j-1) + \delta(S(i),T(j)) & \text{// align 1+1 chars} \end{cases}$$

[Why do we want the min?]

# Dynamic Programming Matrix

|   |   | A | C | A | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| **A** | 1 |   |   |   |   |   |   |   |   |
| **G** | 2 |   |   |   |   |   |   |   |   |
| **C** | 3 |   |   |   |   |   |   |   |   |
| **A** | 4 |   |   |   |   |   |   |   |   |
| **C** | 5 |   |   |   |   |   |   |   |   |
| **A** | 6 |   |   |   |   |   |   |   |   |
| **C** | 7 |   |   |   |   |   |   |   |   |
| **A** | 8 |   |   |   |   |   |   |   |   |

[What does the initialization mean?]

# Dynamic Programming Matrix

|   |   | A | C | A | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 1 | 0 |   |   |   |   |   |   |   |
| G | 2 |   |   |   |   |   |   |   |   |
| C | 3 |   |   |   |   |   |   |   |   |
| A | 4 |   |   |   |   |   |   |   |   |
| C | 5 |   |   |   |   |   |   |   |   |
| A | 6 |   |   |   |   |   |   |   |   |
| C | 7 |   |   |   |   |   |   |   |   |
| A | 8 |   |   |   |   |   |   |   |   |

$$D[A,A] = min\{D[A,]+1, D[,A]+1, D[,]+\delta(A,A)\}$$

# Dynamic Programming Matrix

|   |   | A | C | A | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 1 | 0 | 1 |   |   |   |   |   |   |
| G | 2 |   |   |   |   |   |   |   |   |
| C | 3 |   |   |   |   |   |   |   |   |
| A | 4 |   |   |   |   |   |   |   |   |
| C | 5 |   |   |   |   |   |   |   |   |
| A | 6 |   |   |   |   |   |   |   |   |
| C | 7 |   |   |   |   |   |   |   |   |
| A | 8 |   |   |   |   |   |   |   |   |

$D[A,AC] = \min\{D[A,A]+1, D[,AC]+1, D[,A]+\delta(A,C)\}$

# Dynamic Programming Matrix

|   |   | A | C | A | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 1 | 0 | 1 | 2 |   |   |   |   |   |
| G | 2 |   |   |   |   |   |   |   |   |
| C | 3 |   |   |   |   |   |   |   |   |
| A | 4 |   |   |   |   |   |   |   |   |
| C | 5 |   |   |   |   |   |   |   |   |
| A | 6 |   |   |   |   |   |   |   |   |
| C | 7 |   |   |   |   |   |   |   |   |
| A | 8 |   |   |   |   |   |   |   |   |

$$D[A,ACA] = \min\{D[A,AC]+1, D[,ACA]+1, D[,AC]+\delta(A,A)\}$$

# Dynamic Programming Matrix

|   |   | A | C | A | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|
|   |   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| G | 2 |   |   |   |   |   |   |   |   |
| C | 3 |   |   |   |   |   |   |   |   |
| A | 4 |   |   |   |   |   |   |   |   |
| C | 5 |   |   |   |   |   |   |   |   |
| A | 6 |   |   |   |   |   |   |   |   |
| C | 7 |   |   |   |   |   |   |   |   |
| A | 8 |   |   |   |   |   |   |   |   |

D[A,ACACACTA] = 7

```
-------A
*******|
ACACACTA
```

[What about the other A?]

# Dynamic Programming Matrix

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

D[AG,ACACACTA] = 7

```
----AG--
****|***
ACACACTA
```

# Dynamic Programming Matrix

|   |   | A | C | A | C | A | C | T | A |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| G | 2 | 1 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| C | 3 | 2 | 1 | 2 | 2 | 3 | 4 | 5 | 6 |
| A | 4 | 3 | 2 | 1 | 2 | 2 | 3 | 4 | 5 |
| C | 5 | 4 | 3 | 2 | 1 | 2 | 2 | 3 | 4 |
| A | 6 | 5 | 4 | 3 | 2 | 1 | 2 | 3 | 3 |
| C | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 2 | 3 |
| A | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 2 | 2 |

D[AGCACACA,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: $O(2dn)$

Nathan Edwards

# Local vs. Global Alignment

- The <u>Global Alignment Problem</u> tries to find the best end-to-end alignment between the two strings

  - Only applicable for very closely related sequences

- The <u>Local Alignment Problem</u> tries to find pairs of **substrings** with highest similarity.

  - Especially important if one string is substantially longer than the other

  - Especially important if there is only a distant evolutionary relationship

# Global vs Local Alignment Schematic



Nathan Edwards

# Local vs. Global Alignment (cont'd)

- ## Global Alignment

```
--T—-CC-C-AGT—-TATGT-CAGGGGACACG—A-GCATGCAGA-GAC
  |   ||  |    ||    |  |   |  |||        ||  |   |   |    |  ||||    |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG—T-CAGAT--C
```

- ## Local Alignment—better alignment to find conserved segment

```
          tccCAGTTATGTCAGggggacacgagcatgcagagac
             ||||||||||||
aattgccgccgtcgttttcagCAGTTATGTCAGatc
```

# Part 2: Variant Calling

# Genotyping Theory

Heterozygous variant (3/7)    Homozygous variant (6/6)

```
                                                   GGTATAC...
Subject  ...CCATAG     TGTGCGCCC     CGGAAATTT  CGGTATAC
         ...CCAT    CTATGTGCG         TCGGAAATT  CGGTATAC
         ...CCAT GGCTATGTG       CTATCGGAAA     GCGGCATA
         ...CCA AGGCTATAT        CCTATCGGA     TTGCGGTA    C...
         ...CCA AGGCTATAT    GCCCTATCG      TTTGCGGT     C...
         ...CC   AGGCTATAT    GCCCTATCG   AAATTTGC    ATAC...
         ...CC TAGGCTATA GCGCCCTA      AAATTTGC GTATAC...

Reference ...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...
```

Error or Het (1/7)?

- 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(heads) = 0.5

P(tails) = 0.5

Aaron Quinlan

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



Number of experiments
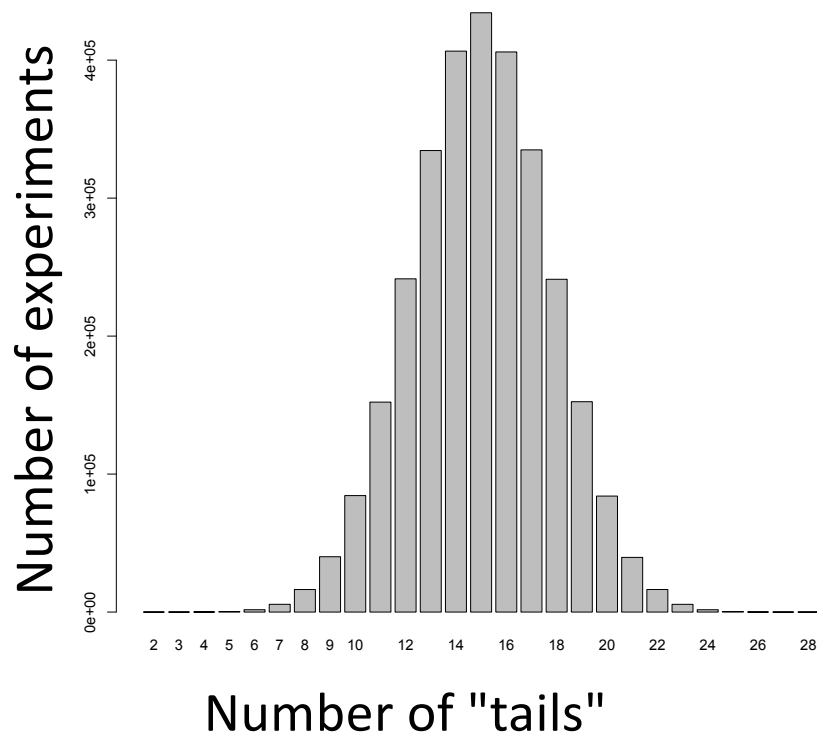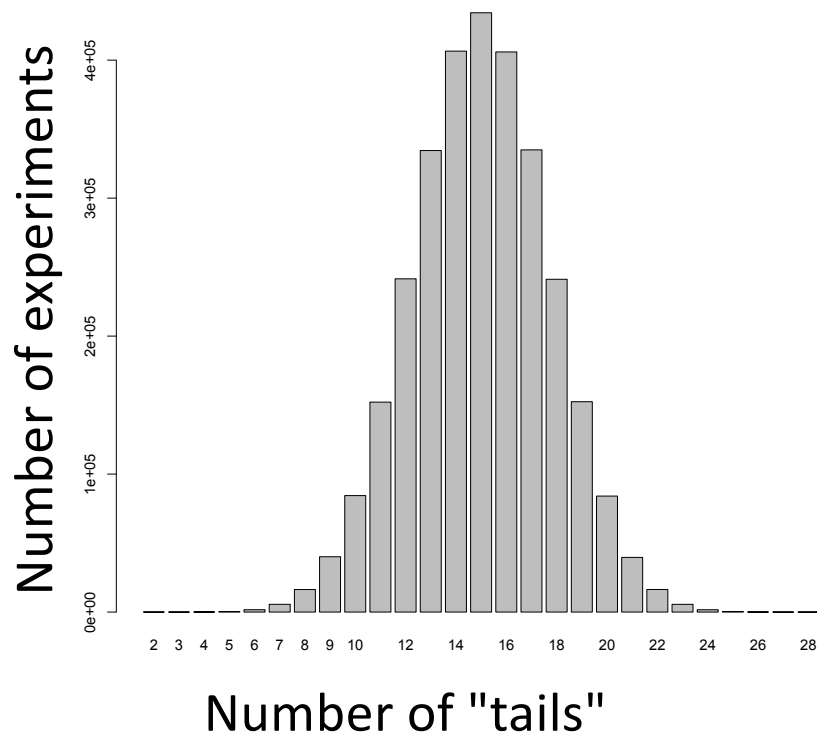
Number of "tails"

R code:

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

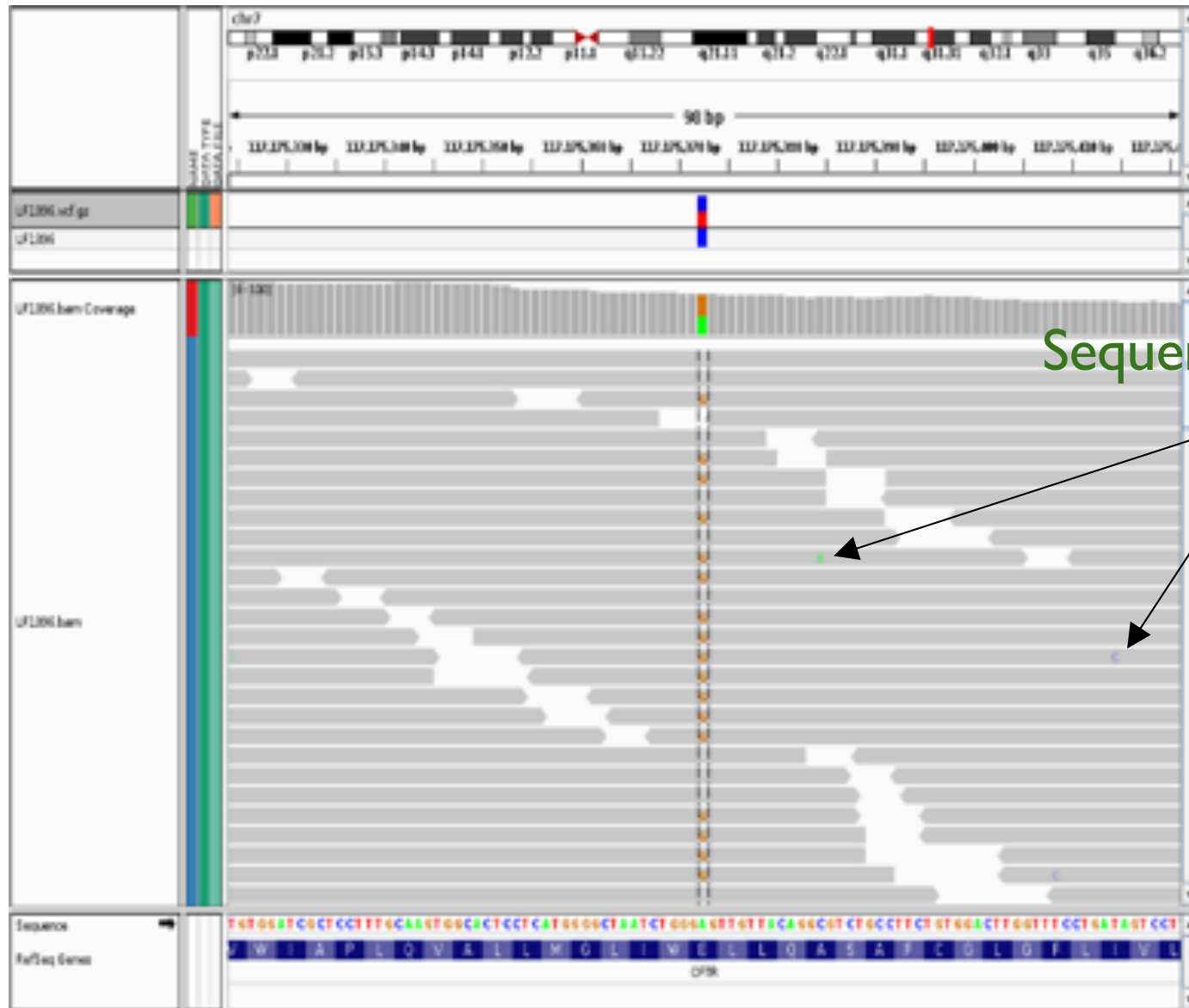30 experiments (students tossing coins)
5 tosses each
Probability of Tails

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



Number of experiments (y-axis)

Number of "tails" (x-axis)

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)?



Number of experiments (y-axis)

Number of "tails" (x-axis)

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 <u>at least</u> 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 het) <?> P(3/30 err)

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



Sequencing errors

# 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.

## A general approach to single-nucleotide polymorphism discovery

Gabor T. Marth[1], Ian Korf[1], Mark D. Yandell[1], Raymond T. Yeh[1], Zhijie Gu[2], Hamideh Zakeri[2], Nathan O. Stitziel[1], LaDeana Hillier[1], Pui-Yan Kwok[2] & Warren R. Gish[1]

**Its main innovation was the use of Bayes's theorem**

# Bayesian SNP calling

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

# PolyBayes: The first statistically rigorous variant detection tool.

## A general approach to single-nucleotide polymorphism discovery

Gabor T. Marth[1], Ian Korf[1], Mark D. Yandell[1], Raymond T. Yeh[1], Zhijie Gu[2], Hamideh Zakeri[2], Nathan O. Stitziel[1], LaDeana Hillier[1], Pui-Yan Kwok[2] & Warren R. Gish[1]

Bayesian posterior probability

Base call + Base quality

Expected (prior) polymorphism rate

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

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

http://www.nature.com/ng/journal/v23/n4/full/ng1299_452.html

# 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. Marth[1], Ian Korf[1], Mark D. Yandell[1], Raymond T. Yeh[1], Zhijie Gu[2], Hamideh Zakeri[2], Nathan O. Stitziel[1], LaDeana Hillier[1], Pui-Yan Kwok[2] & Warren R. Gish[1]

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

**A universal SNP and small-indel variant caller using deep neural networks**
Poplin *et al.* (2018) *Nature Biotechnology.* doi: https://doi.org/10.1038/nbt.4235

# VCF Format

# VCF Format



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