

# Variant Calling

Michael Schatz

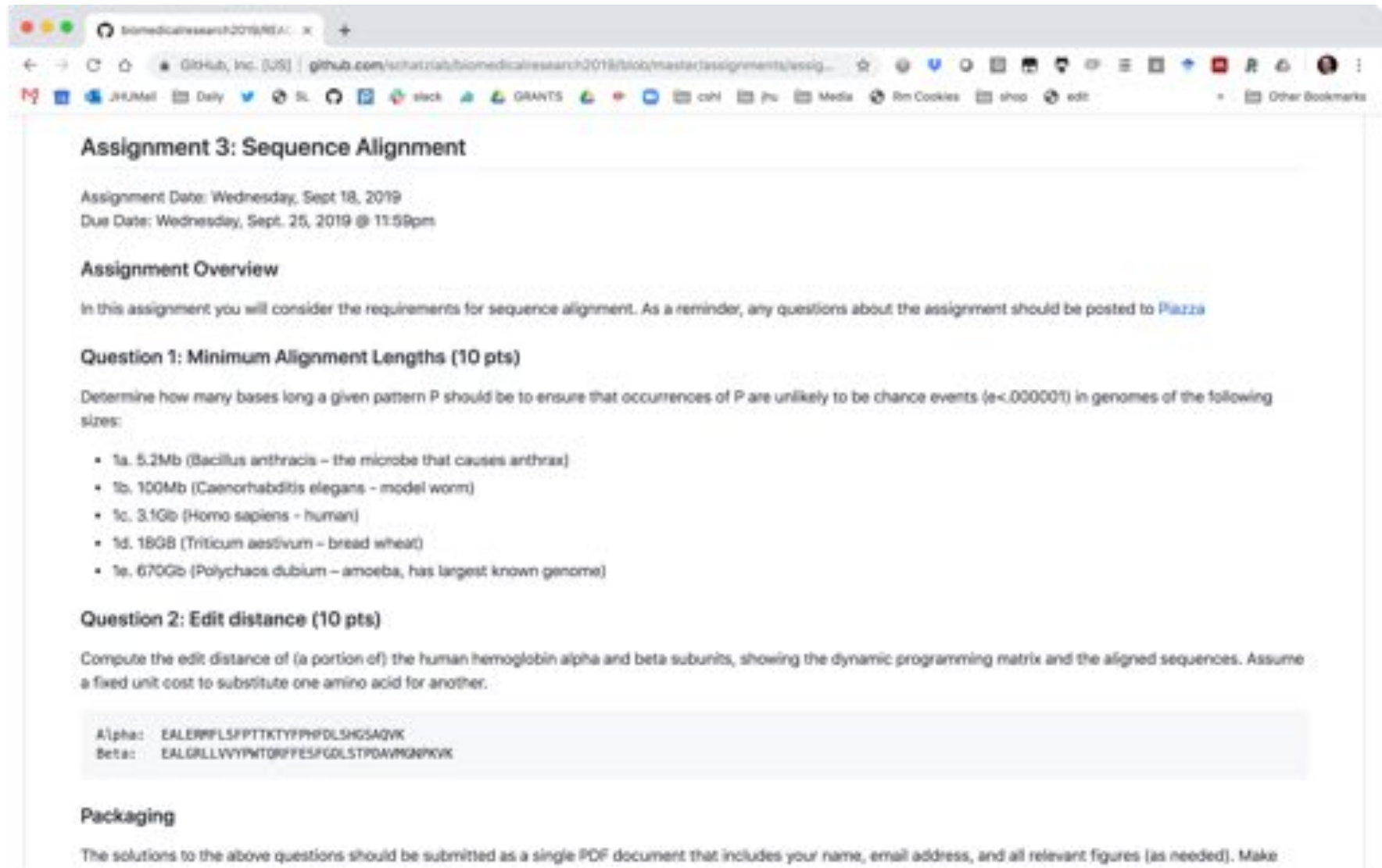
Sept 23, 2019

Lecture 7: Computational Biomedical Research



# Assignment 3: Sequence Alignment

## Due Monday Sept 30 @ 11:59pm



The screenshot shows a web browser displaying a GitHub repository page for 'Assignment 3: Sequence Alignment'. The page includes the assignment date (Wednesday, Sept 18, 2019), due date (Wednesday, Sept. 25, 2019 @ 11:59pm), and an overview of the assignment. It lists two questions: 'Question 1: Minimum Alignment Lengths (10 pts)' and 'Question 2: Edit distance (10 pts)'. Question 1 asks to determine the minimum length of a pattern P for various genomes. Question 2 asks to compute the edit distance between human hemoglobin alpha and beta subunits. The page also includes a 'Packaging' section with submission instructions.

**Assignment 3: Sequence Alignment**

Assignment Date: Wednesday, Sept 18, 2019  
Due Date: Wednesday, Sept. 25, 2019 @ 11:59pm

**Assignment Overview**

In this assignment you will consider the requirements for sequence alignment. As a reminder, any questions about the assignment should be posted to [Piazza](#)

**Question 1: Minimum Alignment Lengths (10 pts)**

Determine how many bases long a given pattern P should be to ensure that occurrences of P are unlikely to be chance events ( $p < .000001$ ) in genomes of the following sizes:

- 1a. 5.2Mb (Bacillus anthracis – the microbe that causes anthrax)
- 1b. 100Mb (Caenorhabditis elegans – model worm)
- 1c. 3.1Gb (Homo sapiens – human)
- 1d. 18Gb (Triticum aestivum – bread wheat)
- 1e. 670Gb (Polychaos dubium – amoeba, has largest known genome)

**Question 2: Edit distance (10 pts)**

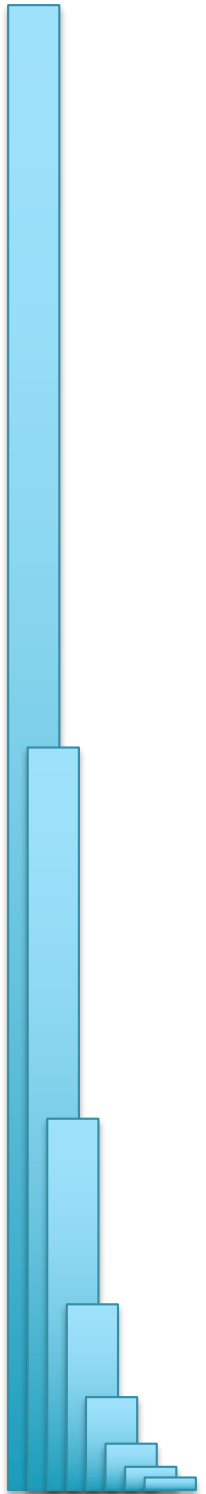
Compute the edit distance of (a portion of) the human hemoglobin alpha and beta subunits, showing the dynamic programming matrix and the aligned sequences. Assume a fixed unit cost to substitute one amino acid for another.

```
Alpha: EALERMFLSFPTTKTYFPHFDSHGSAGVK
Beta:  EALGRLLVYYPWTDRFFESFGDLSTPDVWGKPKVK
```

**Packaging**

The solutions to the above questions should be submitted as a single PDF document that includes your name, email address, and all relevant figures (as needed). Make

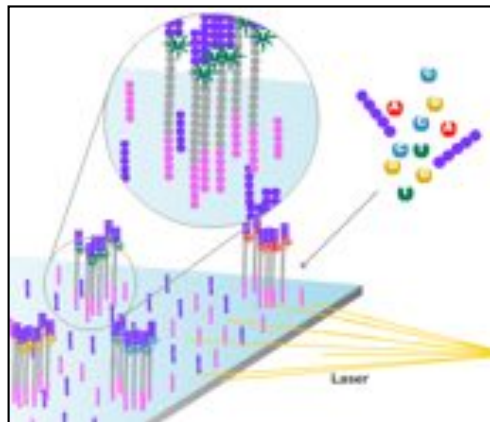
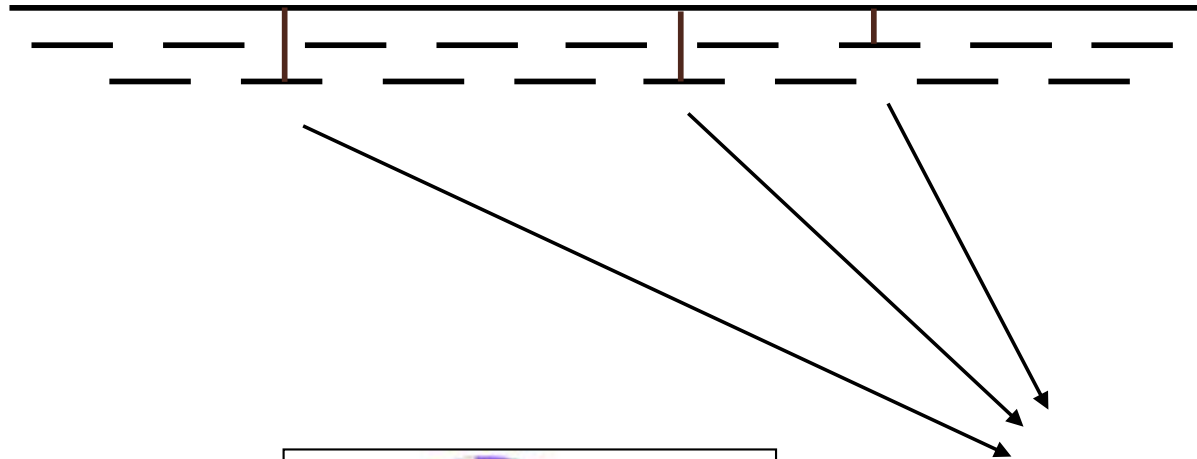
<https://github.com/schatzlab/biomedicalresearch2019>



## Part I: Recap

# Personal Genomics

How does your genome compare to the reference?



Heart Disease

Cancer

Presidential Smile

# Brute Force Analysis



- Brute Force:
  - At every possible offset in the genome:
    - Do all of the characters of the query match?
- Analysis
  - Simple, easy to understand
  - Genome length =  $n$  [3B]
  - Query length =  $m$  [7]
  - Comparisons:  $(n-m+1) * m$  [21B]
- Overall runtime:  $O(nm)$ 
  - [How long would it take if we double the genome size, read length?]
  - [How long would it take if we double both?]

# Searching the Index

- Strategy 2: Binary search
  - Compare to the middle, refine as higher or lower
- Searching for GATTACA
  - $Lo = 1; Hi = 15; Mid = (1+15)/2 = 8$
  - $Middle = Suffix[8] = CC$   
=> Higher:  $Lo = Mid + 1$
  - $Lo = 9; Hi = 15; Mid = (9+15)/2 = 12$
  - $Middle = Suffix[12] = TACC$   
=> Lower:  $Hi = Mid - 1$
  - $Lo = 9; Hi = 11; Mid = (9+11)/2 = 10$
  - $Middle = Suffix[10] = GATTACC$   
=> Lower:  $Hi = Mid - 1$
  - $Lo = 9; Hi = 9; Mid = (9+9)/2 = 9$
  - $Middle = Suffix[9] = GATTACA...$   
=> Match at position 2!

#	Sequence	Pos
1	ACAGATTACC...	6
2	ACC...	13
3	AGATTACC...	8
4	ATTACAGATTACC...	3
5	ATTACC...	10
6	C...	15
7	CAGATTACC...	7
8	CC...	14
9	GATTACAGATTACC...	2
10	GATTACC...	9
11	TACAGATTACC...	5
12	TACC...	12
13	TGATTACAGATTACC...	1
14	TTACAGATTACC...	4
15	TTACC...	11

Lo  
Hi  
→

# Binary Search Analysis

- Binary Search

Initialize search range to entire list

$mid = (hi+lo)/2$ ;  $middle = suffix[mid]$

if query matches middle: done

else if query < middle: pick low range

else if query > middle: pick hi range

Repeat until done or empty range

[WHEN?]

- Analysis

- More complicated method

- How many times do we repeat?

- How many times can it cut the range in half?

- Find smallest  $x$  such that:  $n/(2^x) \leq 1$ ;  $x = \lg_2(n)$

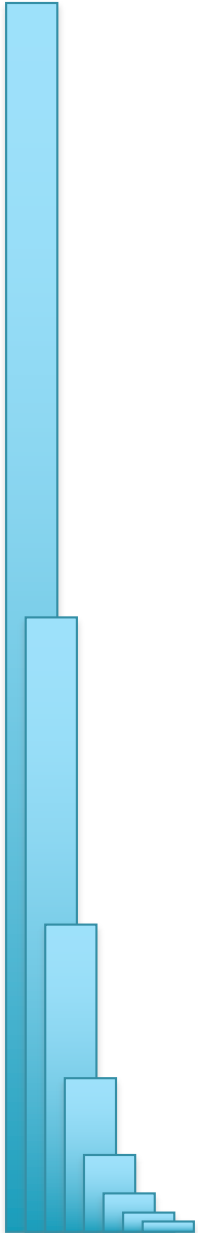
[32]

- Total Runtime:  $O(m \lg n)$

- More complicated, but **much** faster!

- Looking up a query loops 32 times instead of 3B

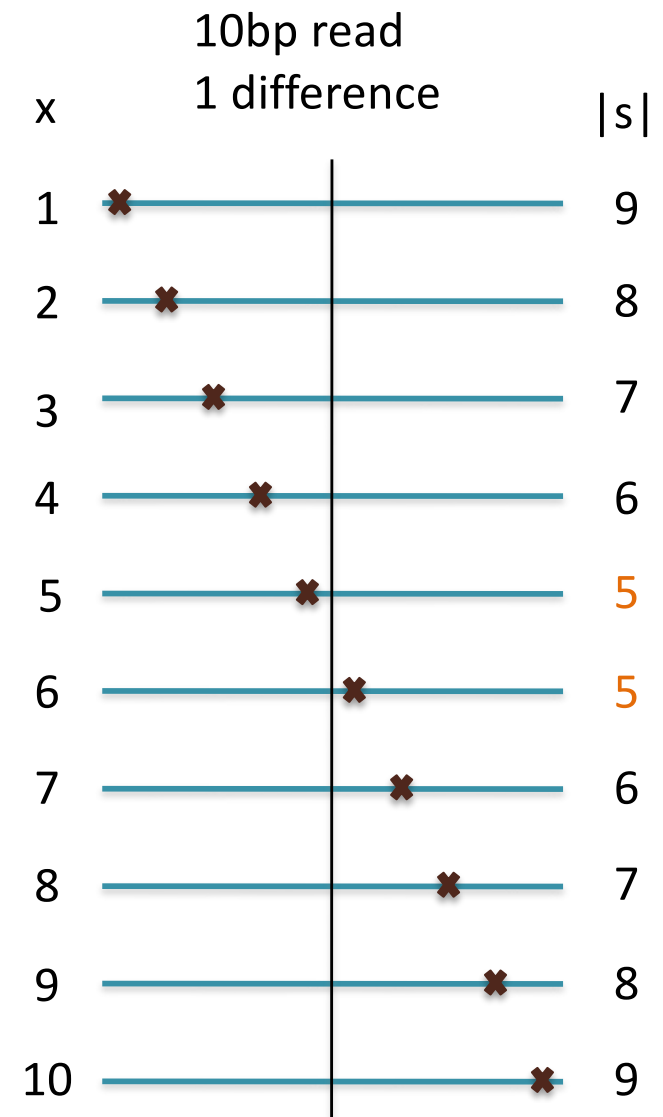
[How long does it take to search 6B or 24B nucleotides?]



# Seed-and-Extend Alignment

Theorem: An alignment of a sequence of length  $m$  with at most  $k$  differences **must** contain an exact match at least  $s = m / (k + 1)$  bp long  
(Baeza-Yates and Perleberg, 1996)

- Proof: Pigeonhole principle
  - 1 pigeon can't fill 2 holes
- Seed-and-extend search
  - Use an index to rapidly find short exact alignments to seed longer in-exact alignments
    - BLAST, MUMmer, Bowtie, BWA, SOAP, ...
  - Specificity of the depends on seed length
    - Guaranteed sensitivity for  $k$  differences
    - Also finds some (but not all) lower quality alignments <- heuristic





# 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

# Edit Distance Example

AGCACACA → ACACACTA in 4 steps

AGCACACA → (1. change G to C)

ACCACACA → (2. delete C)

ACACACA → (3. change A to T)

ACACACT → (4. insert A after T)

ACACACTA → done

[Is this the best we can do?]

# Edit Distance Example

AGCACACA → ACACACTA in 3 steps

AGCACACA → (1. change G to C)

ACCACACA → (2. delete C)

ACACACA → (3. insert T after 3<sup>rd</sup> C)

ACACACTA → done

[Is this the best we can do?]

# Reverse Engineering Edit Distance

$$D(\text{AGCACACA}, \text{ACACACTA}) = ?$$

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...

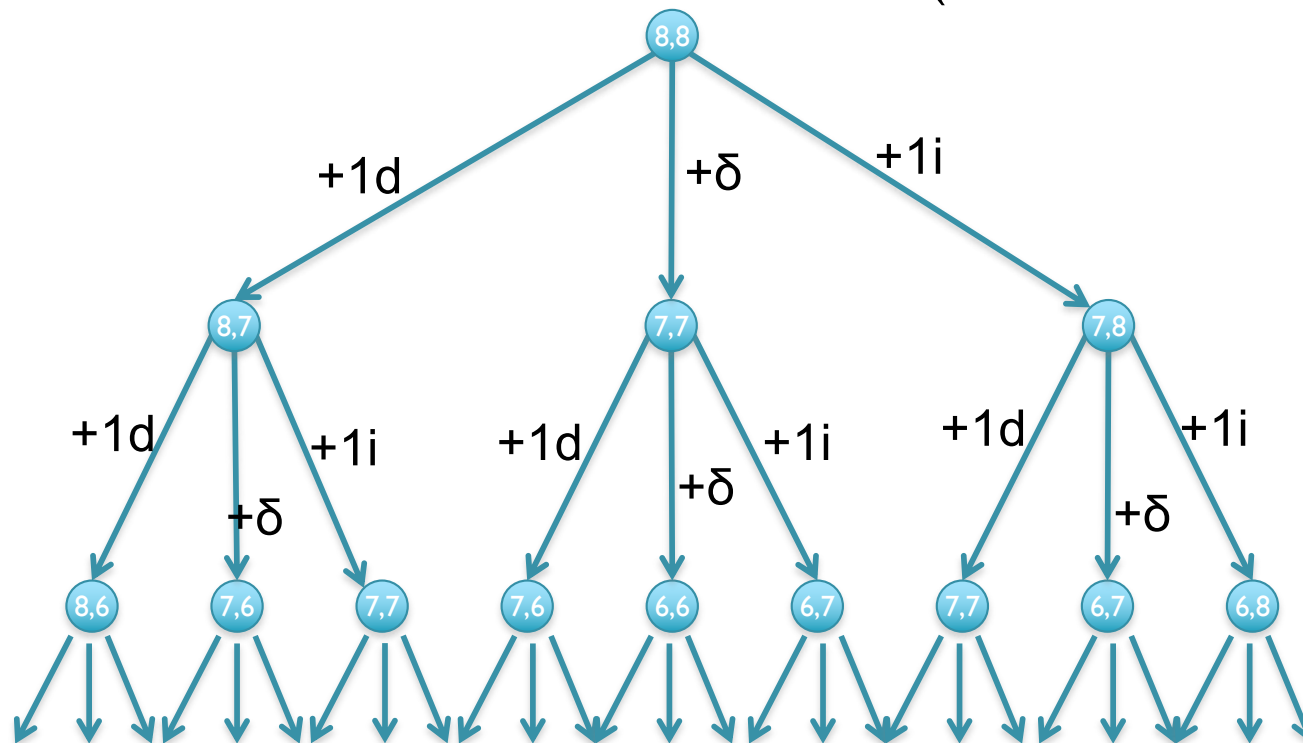
...M...	...I...	...D...
...X...	...-...	...X...
...Y...	...Y...	...-...

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  $i' \leq i$  and  $j' \leq j$ .

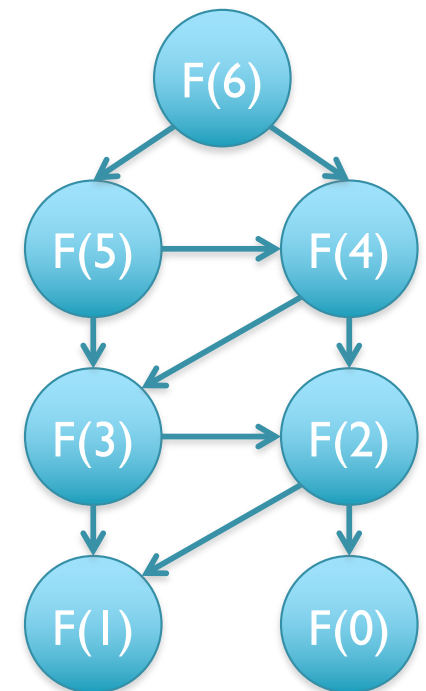
$$D(\text{AGCACACA}, \text{ACACACTA}) = \min \{ D(\text{AGCACACA}, \text{ACACACT}) + 1, \\ D(\text{AGCACAC}, \text{ACACACTA}) + 1, \\ D(\text{AGCACAC}, \text{ACACACT}) + \delta(\text{A}, \text{A}) \}$$



[What is the running time?]

# Dynamic Programming

- General approach for solving (some) complex problems
  - When applicable, the method takes far less time than naive methods.
    - Polynomial time ( $O(n)$  or  $O(n^2)$ ) instead of exponential time ( $O(2^n)$  or  $O(3^n)$ )
- Requirements:
  - **Overlapping subproblems**
  - **Optimal substructure**
- Applications:
  - Fibonacci
  - Longest Increasing Subsequence (Bonus Slides!)
  - Sequence alignment, Dynamic Time Warp, Viterbi
- Not applicable:
  - Traveling salesman problem, Clique finding, Subgraph isomorphism, ...
  - The cheapest flight from airport A to airport B involves a single connection through airport C, but the cheapest flight from airport A to airport C involves a connection through some other airport D.



# 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  $i = 0, \dots, n$
- $D(0,j) = j$ , for all  $j = 0, \dots, m$

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

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

[Why do we want the min?]



# Dynamic Programming Matrix

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

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

$$D[AGCACACA, ACACACTA] = 2$$

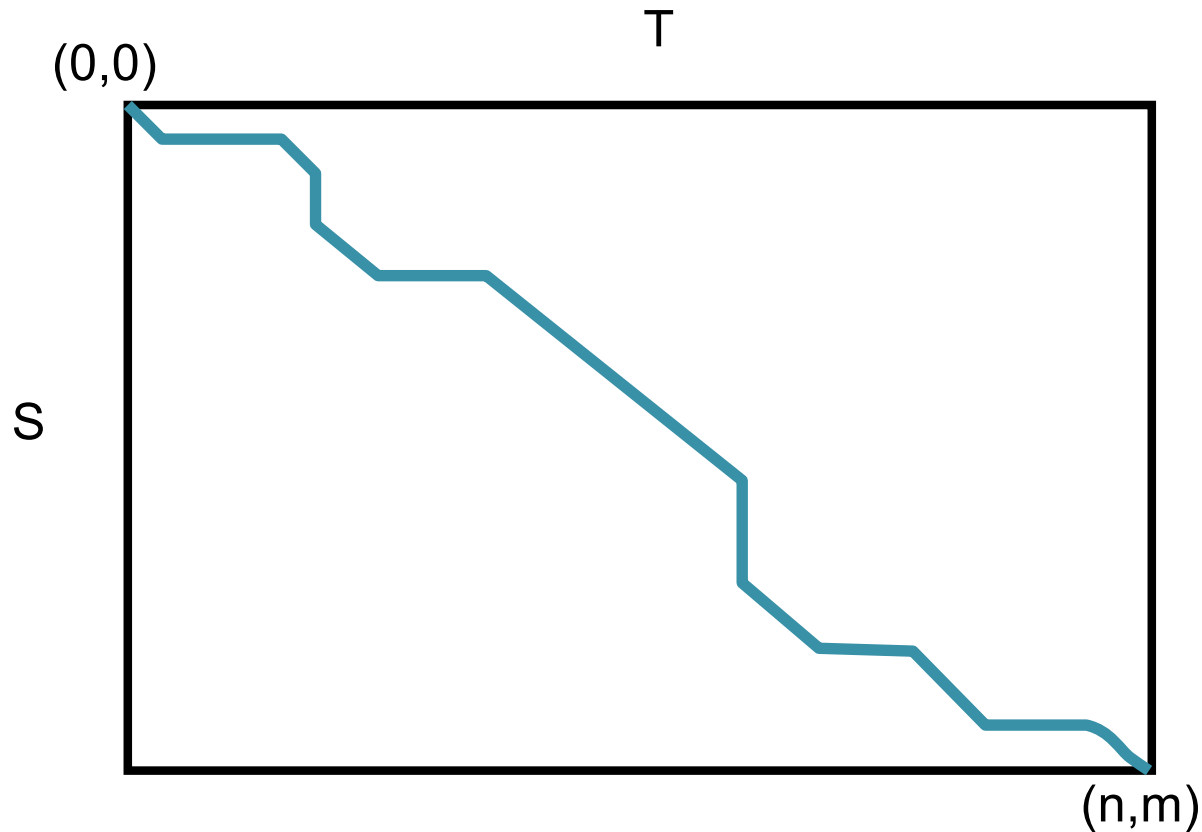
AGCACAC-A

| \* | | | | \* |

A-CACACTA

[Can we do it any better?]

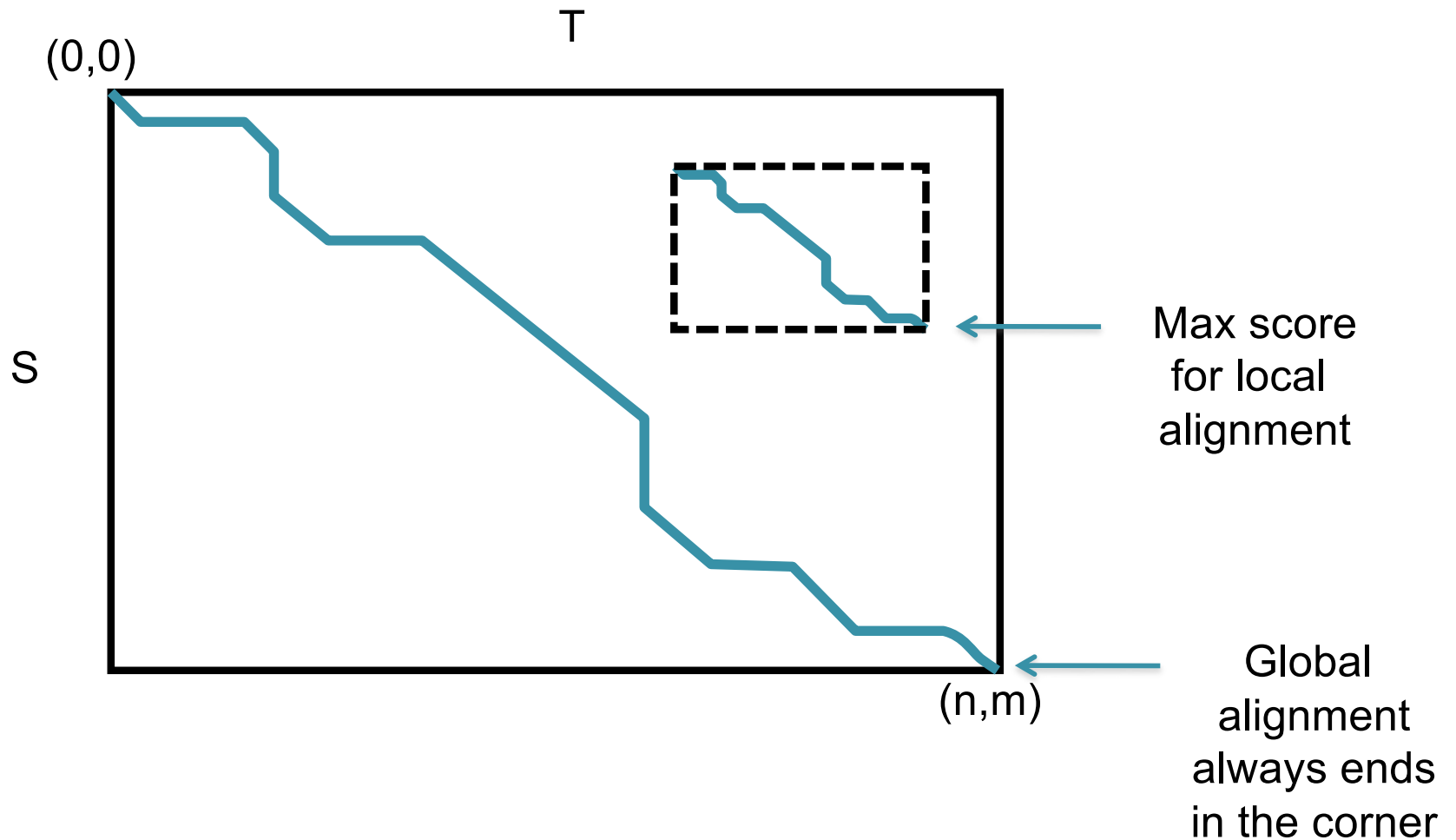
# Global Alignment Schematic



- 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  $d$ :  $O(2dn)$



# Global vs Local Alignment Schematic



# 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

```
          tccCAGTTATGTCAGgggacacgagcatgcagagac
            |||||
aattgccgccgctcgtttttcagCAGTTATGTCAGatc
```

# Alignment Ambiguity

Notice that the edit distance of GATTTACA and GATACA is 2,  
but there are multiple possible optimal alignments:

GATTTACA  
GAT--ACA

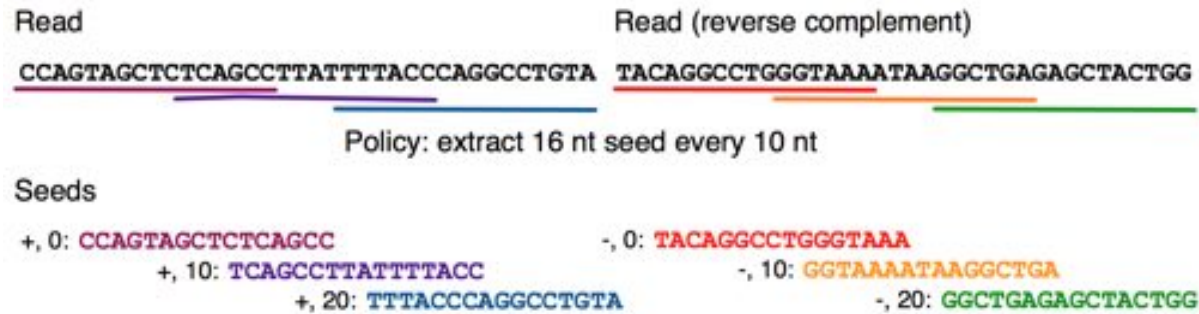
GATTTACA  
GA-T-ACA

GATTTACA  
GA--TACA

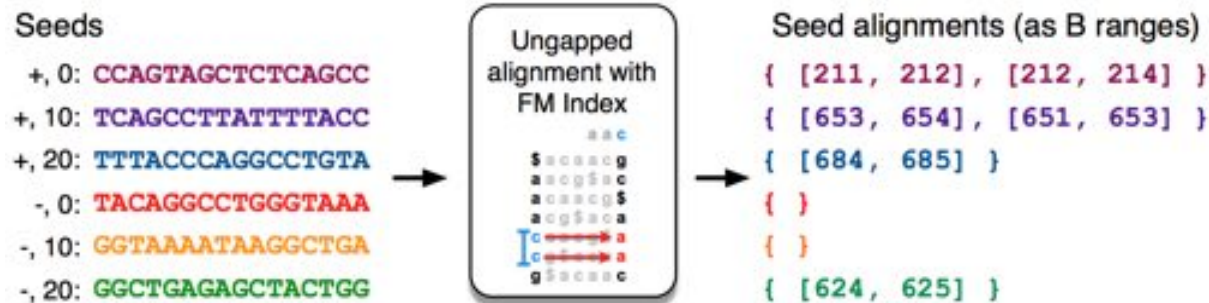
Which one is right?

# Algorithm Overview

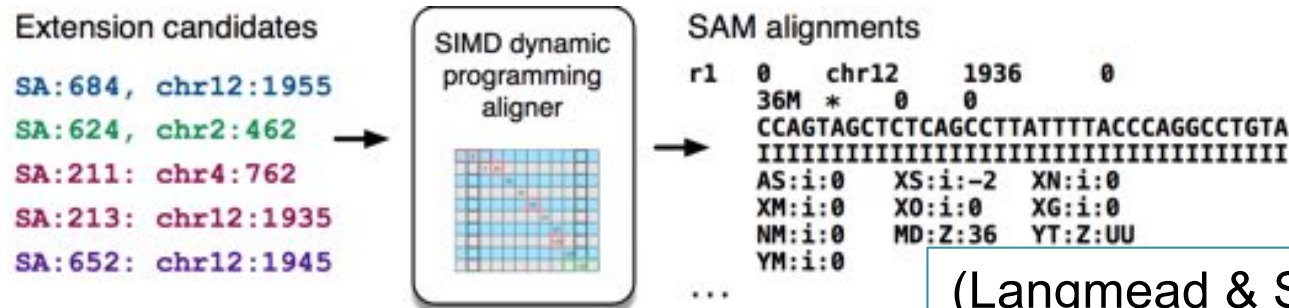
## 1. Split read into segments



## 2. Lookup each segment and prioritize

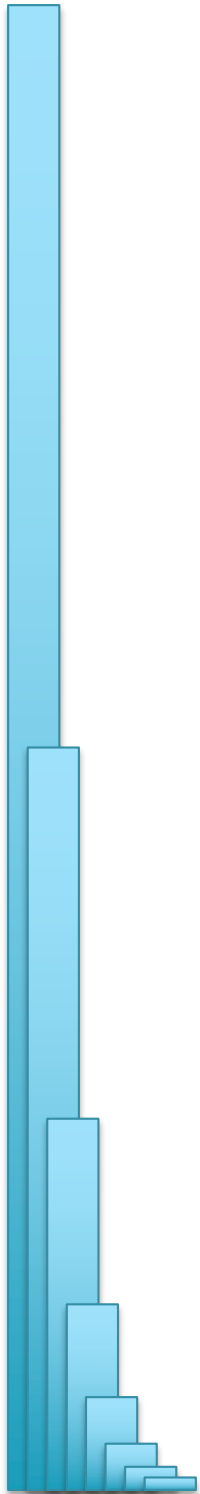


### 3. Evaluate end-to-end match



(Langmead & Salzberg, 2012)

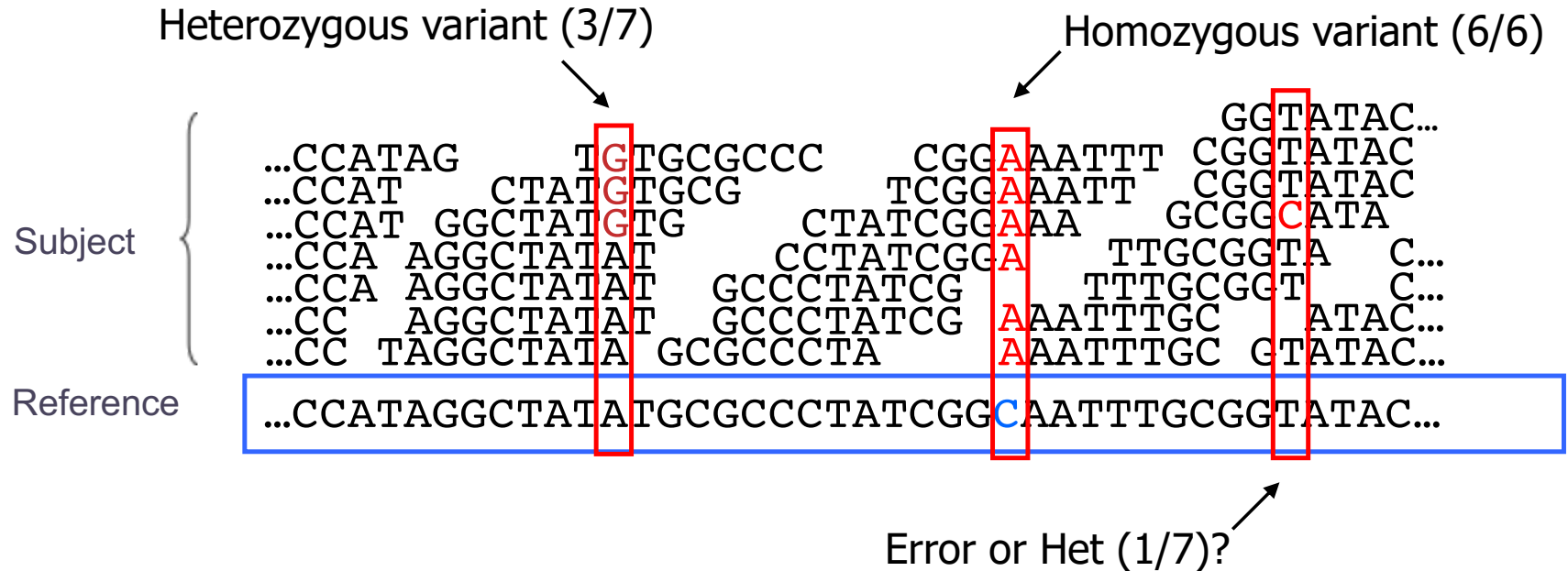
## Part 2: Variant Calling



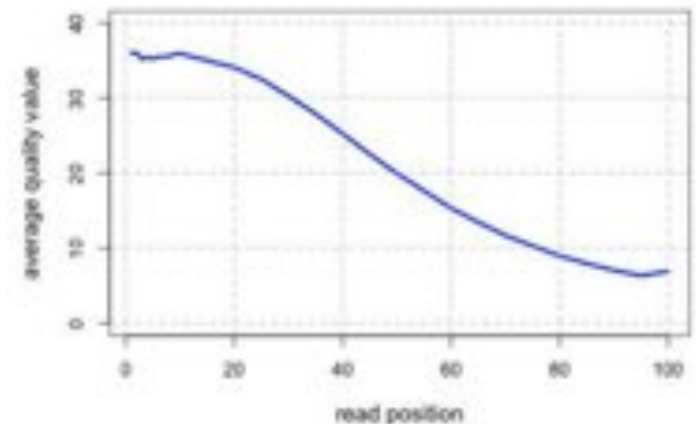
# Variant Calling Overview



# 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



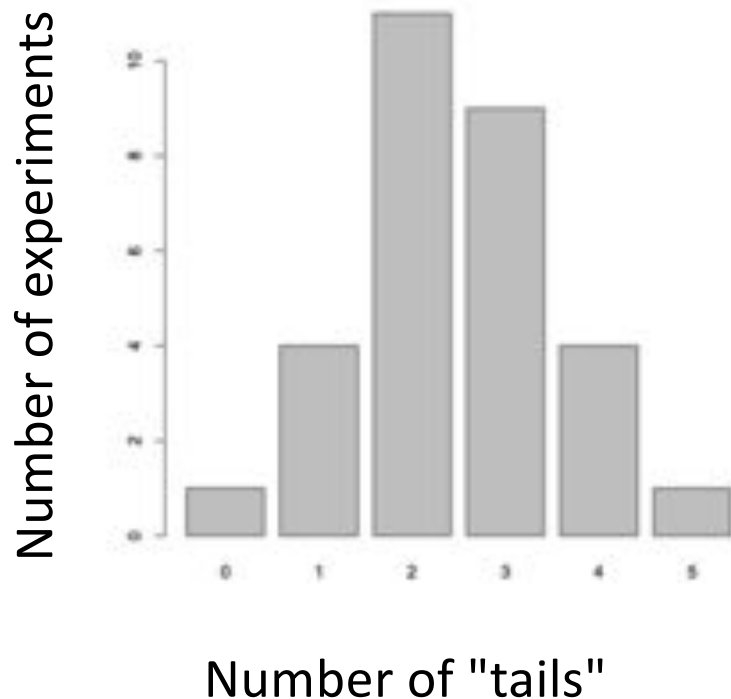
$P(\text{heads}) = 0.5$



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



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



R code:

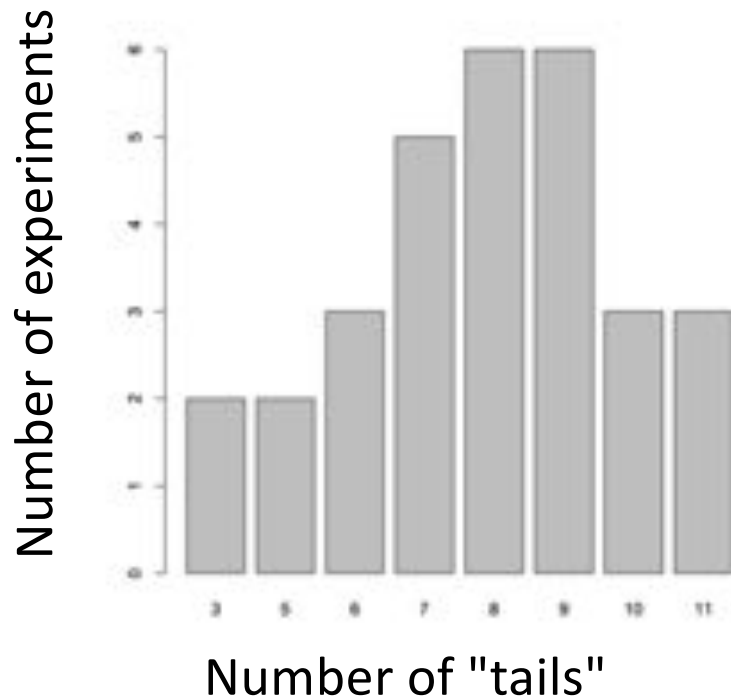
```
barplot(table(rbinom(30, 5, 0.5)))
```

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:

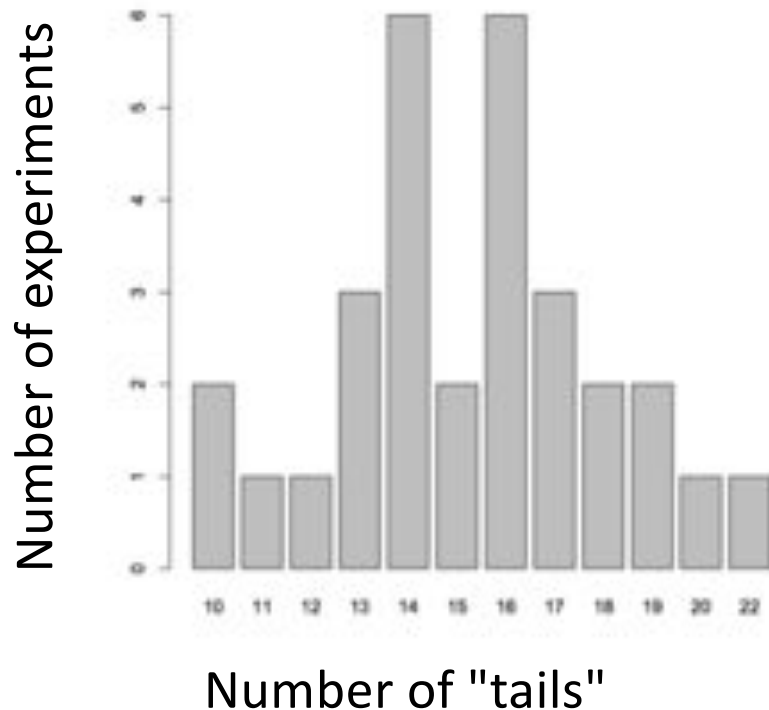
```
barplot(table(rbinom(30, 15, 0.5)))
```

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:

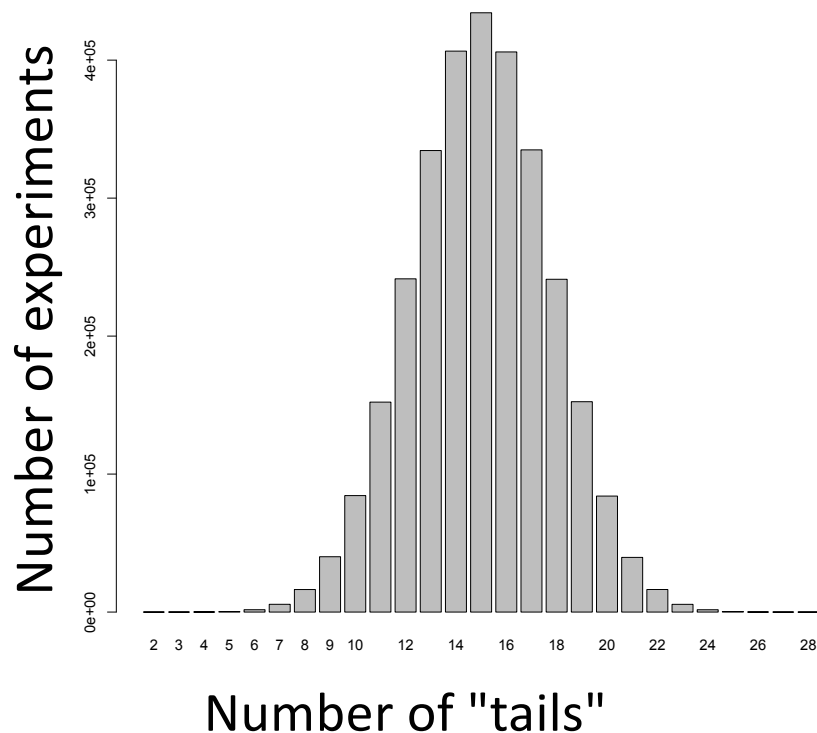
```
barplot(table(rbinom(30, 30, 0.5)))
```

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:

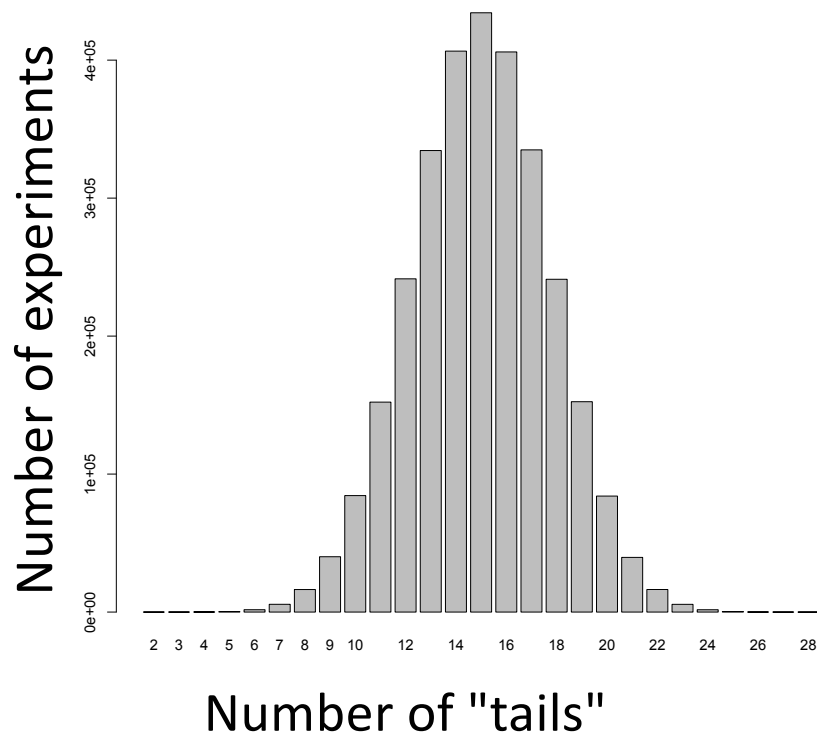
```
barplot(table(rbinom(3e6, 30, 0.5)))
```

3M experiments (students tossing coins)

30 tosses each

Probability of Tails

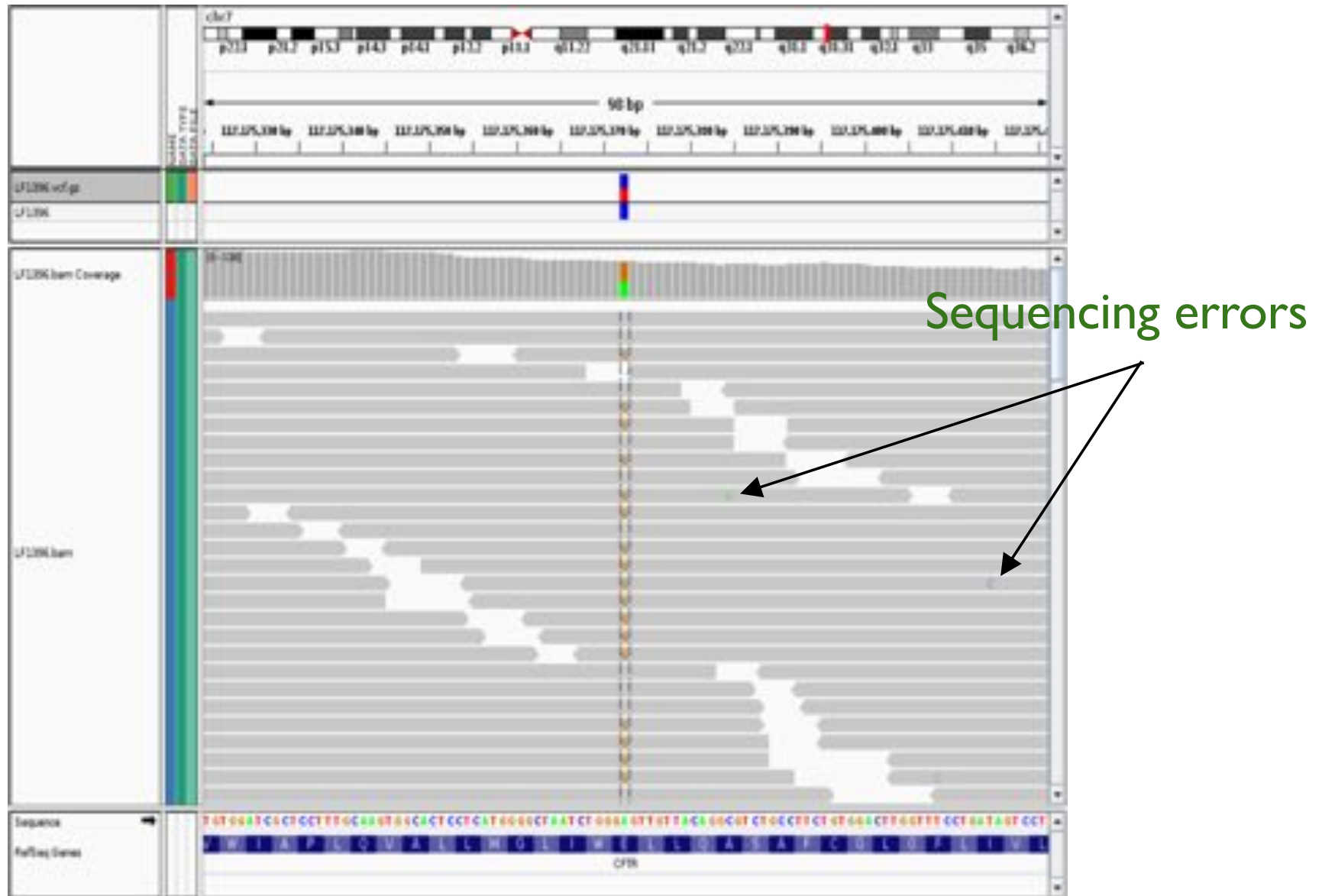
So, with 30 tosses (reads), we are much more likely to see an even mix of alternate and reference alleles at a heterozygous locus in a genome



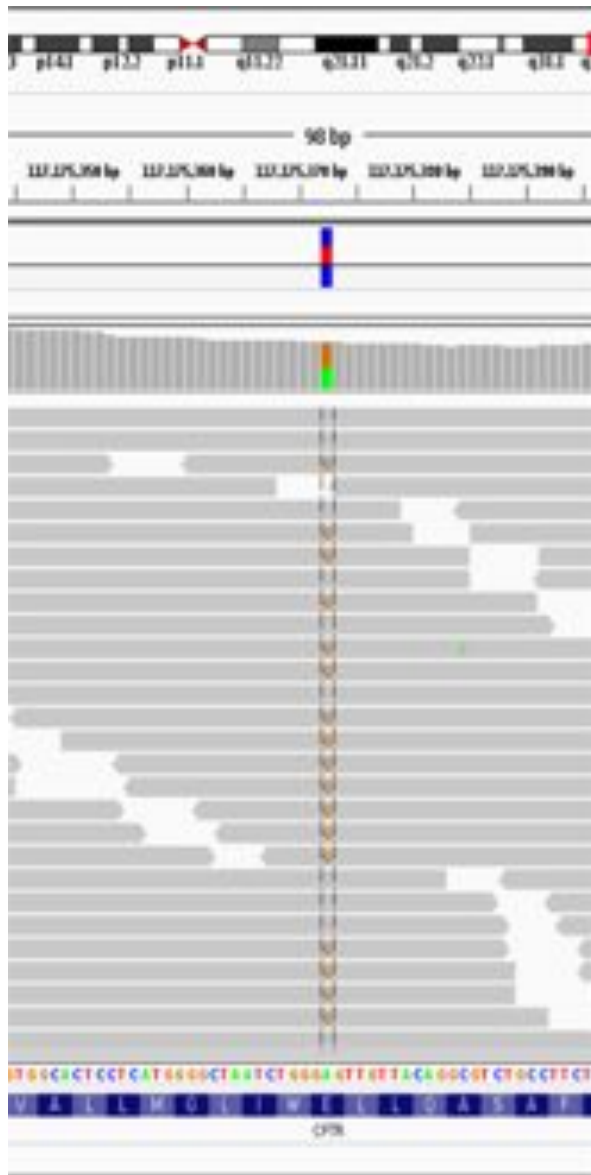
This is why at least a "30X" (30 fold sequence coverage) genome is recommended: it confers sufficient power to distinguish heterozygous alleles and from mere sequencing errors

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

# 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

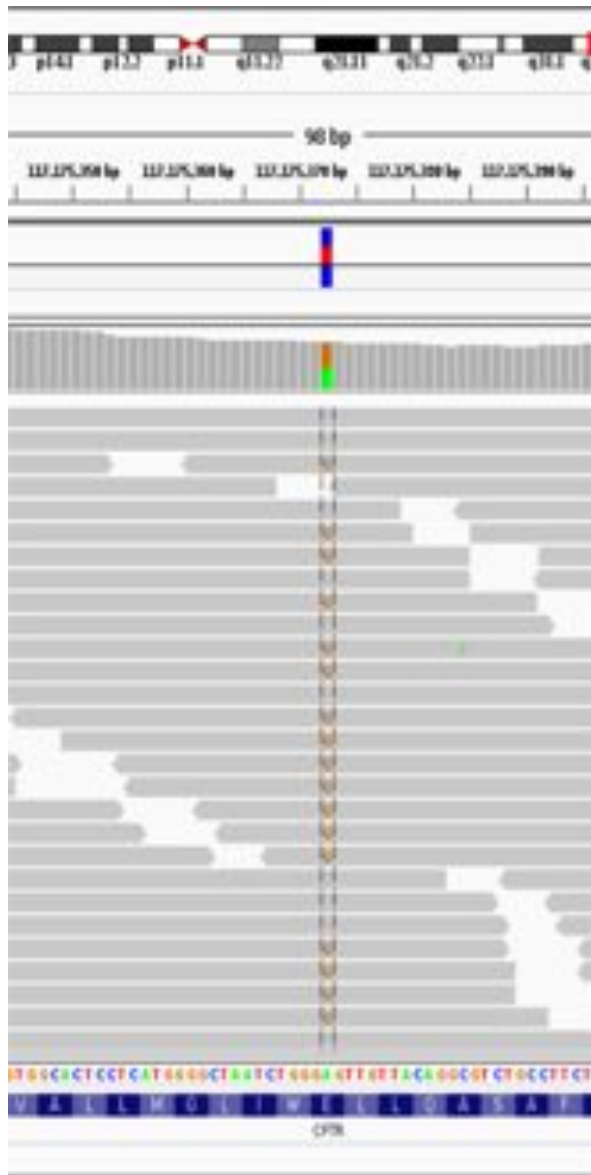
Gabor T. Marth<sup>1</sup>, Ian Korf<sup>1</sup>, Mark D. Yandell<sup>1</sup>, Raymond T. Yeh<sup>1</sup>, Zhijie Gu<sup>2</sup>, Hamideh Zakeri<sup>2</sup>, Nathan O. Stitzel<sup>1</sup>, LaDeana Hillier<sup>1</sup>, Pui-Yan Kwok<sup>2</sup> & Warren R. Gish<sup>1</sup>

Its main innovation was the use of Bayes's theorem





# 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. Marth<sup>1</sup>, Ian Koef<sup>1</sup>, Mark D. Yandell<sup>1</sup>, Raymond T. Yeh<sup>1</sup>, Zhijie Gu<sup>2</sup>, Hamideh Zakeri<sup>2</sup>, Nathan O. Stitzel<sup>1</sup>, LaDeana Hillier<sup>1</sup>, Pui-Yan Kwok<sup>2</sup> & Warren R. Gish<sup>1</sup>

Bayesian  
posterior  
probability

Base call +  
Base quality

Expected (prior)  
polymorphism rate

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

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

# 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<sup>1</sup>, Ian Korf<sup>1</sup>, Mark D. Yandell<sup>1</sup>, Raymond T. Yeh<sup>1</sup>, Zhijie Gu<sup>2</sup>, Hamideh Zakeri<sup>2</sup>,  
Nathan O. Stitzel<sup>1</sup>, LaDeana Hillier<sup>1</sup>, Pui-Yan Kwok<sup>2</sup> & Warren R. Gish<sup>1</sup>

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

# VCF Format

## Example

**VCF header**

```
##fileformat=VCFv4.0
##fileDate=20100707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=8,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	2/2:20
1	5	.	A	G	.	PASS	.	GT:GQ	1/0:77	1/1:93
1	100	.	T	<DEL>	.	PASS	SVTYPE=DEL;END=300	GT:GQ:DP	1/1:12:3	0/0:20

**Deletion**

**SNP**

**Large SV**

**Insertion**

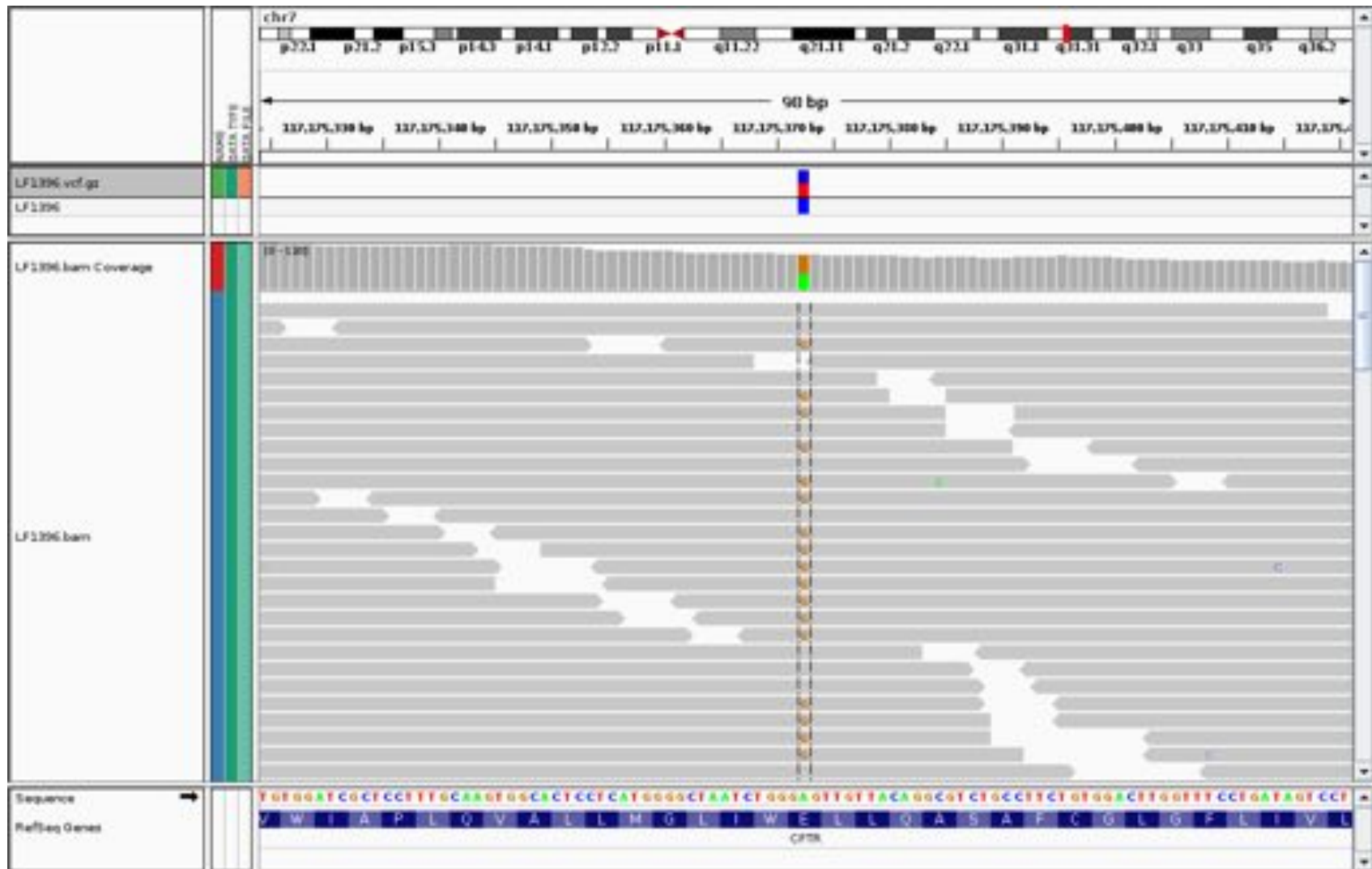
**Other event**

**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

# Next Steps

1. Reflect on the magic and power of DNA 😊
2. Check out the course webpage
3. Work on HW3

