

# Variant Calling

Michael Schatz

Sept 23, 2024

Lecture 8: Applied Comparative Genomics



# Assignment 3: Variant Finding

## Due Monday Sept 30 by 11:59pm

The screenshot shows a web browser window displaying the README.md file for Assignment 3 on a GitHub repository. The URL in the address bar is <https://github.com/schatzlab/appliedgenomics2024/blob/main/assignments/assignment3/README.md>. The browser interface includes a sidebar with a 'Files' section showing the directory structure of the assignment folder, and the main content area showing the assignment details.

**Assignment 3: BWT and Variant Calling**

Assignment Date: Monday, September 23, 2024  
Due Date: Monday, September 30, 2023 @ 11:59pm

**Assignment Overview**

In this assignment you will implement the BWT and explore the requirements for variant calling. The programming exercises can be computed in any programming language, although we recommend python (or C++, Java, or Rust). R is generally inefficient at string processing unless you take great care. See the resources at the bottom of the page for tips for the variant calling exercises.

As a reminder, any questions about the assignment should be posted to [Piazza](#).

**Question 1. BWT Encoding [20 pts]**

In the language of your choice, implement a BWT encoder and encode the string below. Faster (Linear time) methods exist for computing the BWT, although for this assignment you can use the simple method based on standard sorting techniques. Your solution does *not* need to be an optimal algorithm and can use  $O(n^2)$  space and  $O(n^2 \lg n)$  time.

Here is the recommended pseudo code (make sure to submit your code as well as the encoded string):

```
computeBWT(string s)
    ## add the magic end-of-string character
    s = s + "$"

    ## build up the BWT from the cyclic permutations
    ## note the ith cyclic permutation is just "s[i..n] + s[0..i]"
    rows = []
    for (i = 0; i < length(s); i++)
        rows.append(cyclic_permutation(s, i))

    ## just use the builtin sort command to sort the cyclic permutations
    sort(rows)
```

<https://github.com/schatzlab/appliedgenomics2024/tree/main/assignments/assignment3>

Check Piazza for questions!

# Read Mapping

# Exact Matching Review & Overview

Where is GATTACA in the human genome?

Brute Force (3 GB)	Suffix Array (>15 GB)	Hash Table (>15 GB)	BWT (3 GB)
BANANA BAN ANA NAN ANA	<p>6 \$ 5 A\$ 3 ANA\$ 1 ANANAS\$ 0 BANANA\$ 4 NA\$ 2 NANA\$</p>	<p>0 → BAN → 0 → NULL 0 → ANA → 1 → ANA → 3 → NULL 0 → NAN → 2 → NULL</p>	<p>BANANA\$ =&gt; \$BANANA<b>A</b> A\$BANA<b>N</b> ANA\$BA<b>N</b> ANANA\$b<b>B</b> BANANA\$<b>N</b> NA\$BAN<b>A</b> NANA\$b<b>A</b> =&gt; <b>ANNB\$AA</b></p>
$O(m * n)$	$O(m + \lg n)$	$O(1)$	$O(m)$
Slow & Easy	Full-text index	Fixed-length lookup	Full-text and concise

\*\*\* These are general techniques applicable to any text search problem \*\*\*

# Bowtie2 Algorithm Overview

## 1. Split read into segments

Read  
CCAGTAGCTCTCAGCCTTATTTACCCAGGCCTGTA      Read (reverse complement)  
    TACAGGCCTGGGTAAAATAAGGCTGAGAGCTACTGG

Policy: extract 16 nt seed every 10 nt

Seeds

+ , 0: CCAGTAGCTCTCAGCC	- , 0: TACAGGCCTGGGTAAA
+ , 10: TCAGCCTTATTTACC	- , 10: GGTAAAATAAGGCTGA
+ , 20: TTTACCCAGGCCTGTA	- , 20: GGCTGAGAGCTACTGG

## 2. Lookup each segment and prioritize

Seeds

+ , 0: CCAGTAGCTCTCAGCC	→	Ungapped alignment with FM Index	→	Seed alignments (as B ranges)
+ , 10: TCAGCCTTATTTACC				{ [211, 212], [212, 214] }
+ , 20: TTTACCCAGGCCTGTA				{ [653, 654], [651, 653] }
- , 0: TACAGGCCTGGGTAAA				{ [684, 685] }
- , 10: GGTAAAATAAGGCTGA				{ }
- , 20: GGCTGAGAGCTACTGG				{ }

## 3. Evaluate end-to-end match

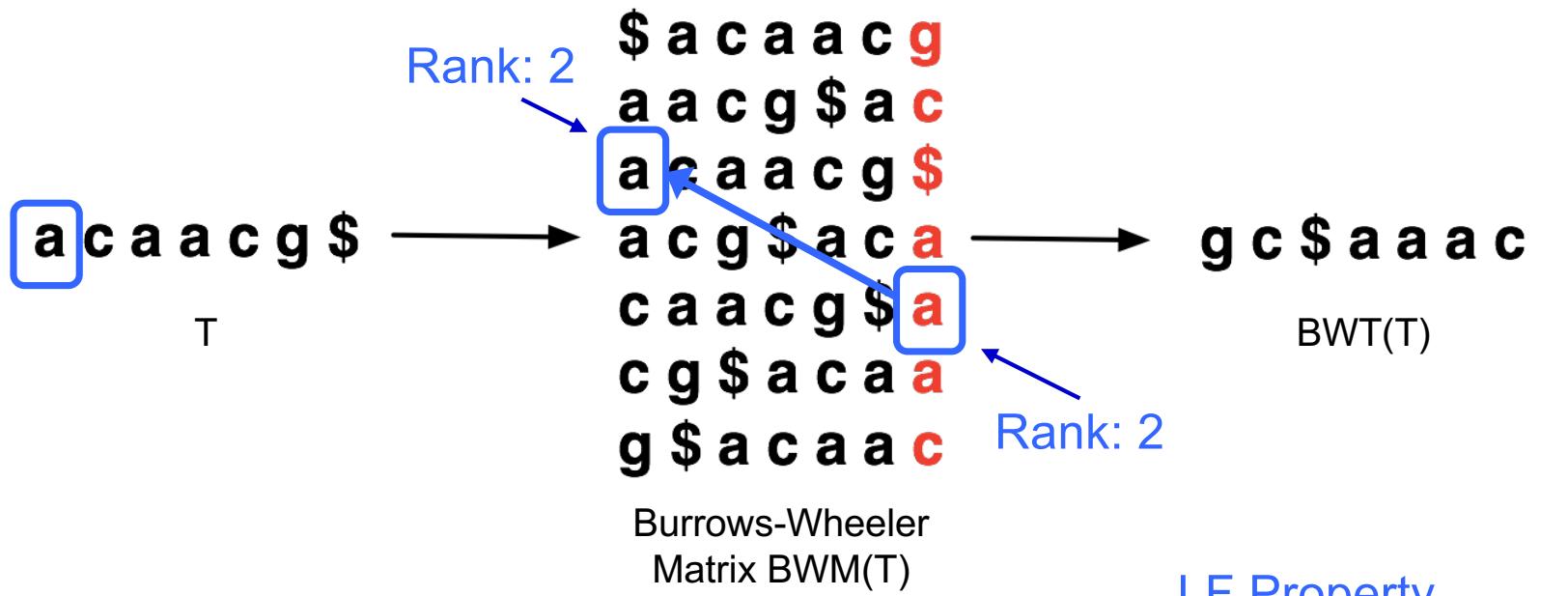
Extension candidates

SA:684, chr12:1955	→	SIMD dynamic programming aligner	→	SAM alignments
SA:624, chr2:462				r1 0 chr12 1936 0
SA:211: chr4:762				36M * 0 0
SA:213: chr12:1935				CCAGTAGCTCTCAGCCTTATTTACCCAGGCCTGTA
SA:652: chr12:1945				II

(Langmead & Salzberg, 2012)

# Burrows-Wheeler Transform

- Reversible permutation of the characters in a text



- $\text{BWT}(T)$  is the index for  $T$

LF Property  
implicitly encodes  
Suffix Array

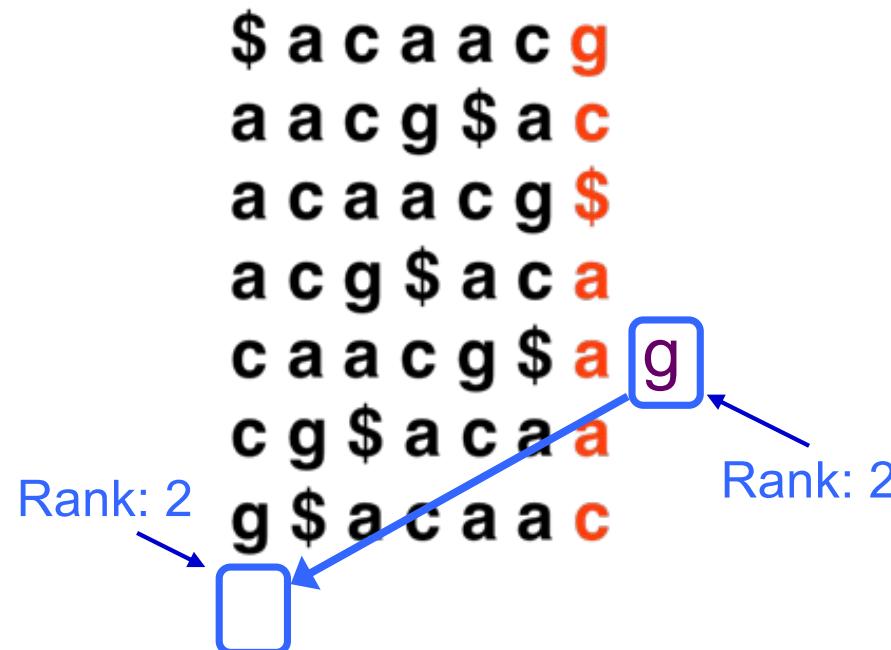
A block sorting lossless data compression algorithm.

Burrows M, Wheeler DJ (1994) Digital Equipment Corporation. Technical Report 124

# BWT Exact Matching

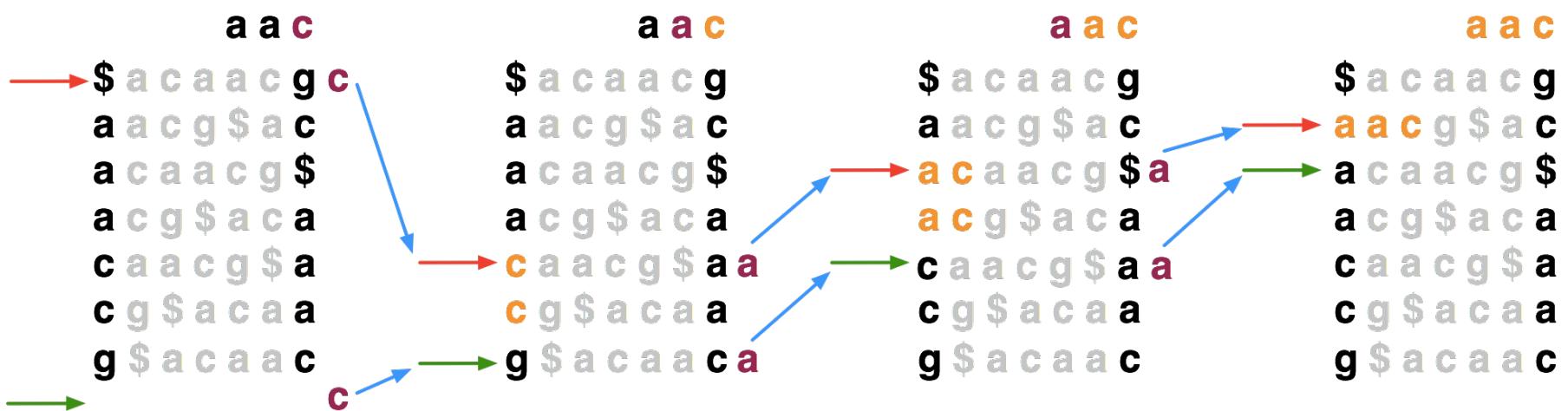
- $\text{LFc}(r, c)$  does the same thing as  $\text{LF}(r)$  but it ignores  $r$ 's actual final character and “pretends” it's  $c$ :

$$\text{LFc}(5, \text{g}) = 8$$



# BWT Exact Matching

- Start with a range, (**top**, **bot**) encompassing all rows and repeatedly apply **LFc**:  
**top** = **LFc**(**top**, **qc**); **bot** = **LFc**(**bot**, **qc**)  
**qc** = the next character to the left in the query



Ferragina P, Manzini G: Opportunistic data structures with applications. *FOCS. IEEE Computer Society*; 2000.

[Search for TTA this BWT string: ACTGA\$TTA ]

# FM Index

\$ a c a a c g  
a a c g \$ a c  
a c a a c g \$  
a c g \$ a c a ←  
c a a c g \$ a ←  
c g \$ a c a a  
g \$ a c a a c

**Where is this row in the original text?**

- Could repeatedly apply LF transformation, but could require many steps:  $O(N)$
- Instead, periodically record the position: sample the suffix array every 100th row

**Which 'a' is this?**

- Could count from the top of the matrix, but this will be very slow:  $O(n)$
- Instead, periodically (every 100<sup>th</sup> row) keep track of the number of occurrences of the characters that have seen so far: \$:1, a:2, c:1, g:1, t:0

FM-index: BWT + SA sample + count table  
3Gb + (30M ints) + (30M \* 5 ints)  
3Gb + 120Mb + 600Mb = 3.72Gb

## Opportunistic Data Structures with Applications

Ferragina and Manzini (2000). Proceedings of the 41st Annual Symposium on Foundations of Computer Science. p.390.

# Algorithm Overview

## 1. Split read into segments

Read  
  
Read (reverse complement)  
  
Policy: extract 16 nt seed every 10 nt

Seeds

+ , 0: CCAGTAGCTCTCAGCC	- , 0: TACAGGCCTGGGTAAA
+ , 10: TCAGCCTTATTTACC	- , 10: GGTAAAATAAGGCTGA
+ , 20: TTTACCCAGGCCTGTA	- , 20: GGCTGAGAGCTACTGG

## 2. Lookup each segment and prioritize

Seeds

+ , 0: CCAGTAGCTCTCAGCC	→	Ungapped alignment with FM Index
+ , 10: TCAGCCTTATTTACC	→	[Seed alignments (as B ranges)]
+ , 20: TTTACCCAGGCCTGTA	→	{ [211, 212], [212, 214] }
- , 0: TACAGGCCTGGGTAAA	→	{ [653, 654], [651, 653] }
- , 10: GGTAAAATAAGGCTGA	→	{ [684, 685] }
- , 20: GGCTGAGAGCTACTGG	→	{ }

Seed alignments (as B ranges)

{ [211, 212], [212, 214] }
{ [653, 654], [651, 653] }
{ [684, 685] }
{ }
{ }
{ [624, 625] }

## 3. Evaluate end-to-end match

Extension candidates

SA:684, chr12:1955	→	SIMD dynamic programming aligner
SA:624, chr2:462	→	SIMD dynamic programming aligner
SA:211: chr4:762	→	SIMD dynamic programming aligner
SA:213: chr12:1935	→	SIMD dynamic programming aligner
SA:652: chr12:1945	→	SIMD dynamic programming aligner

SAM alignments

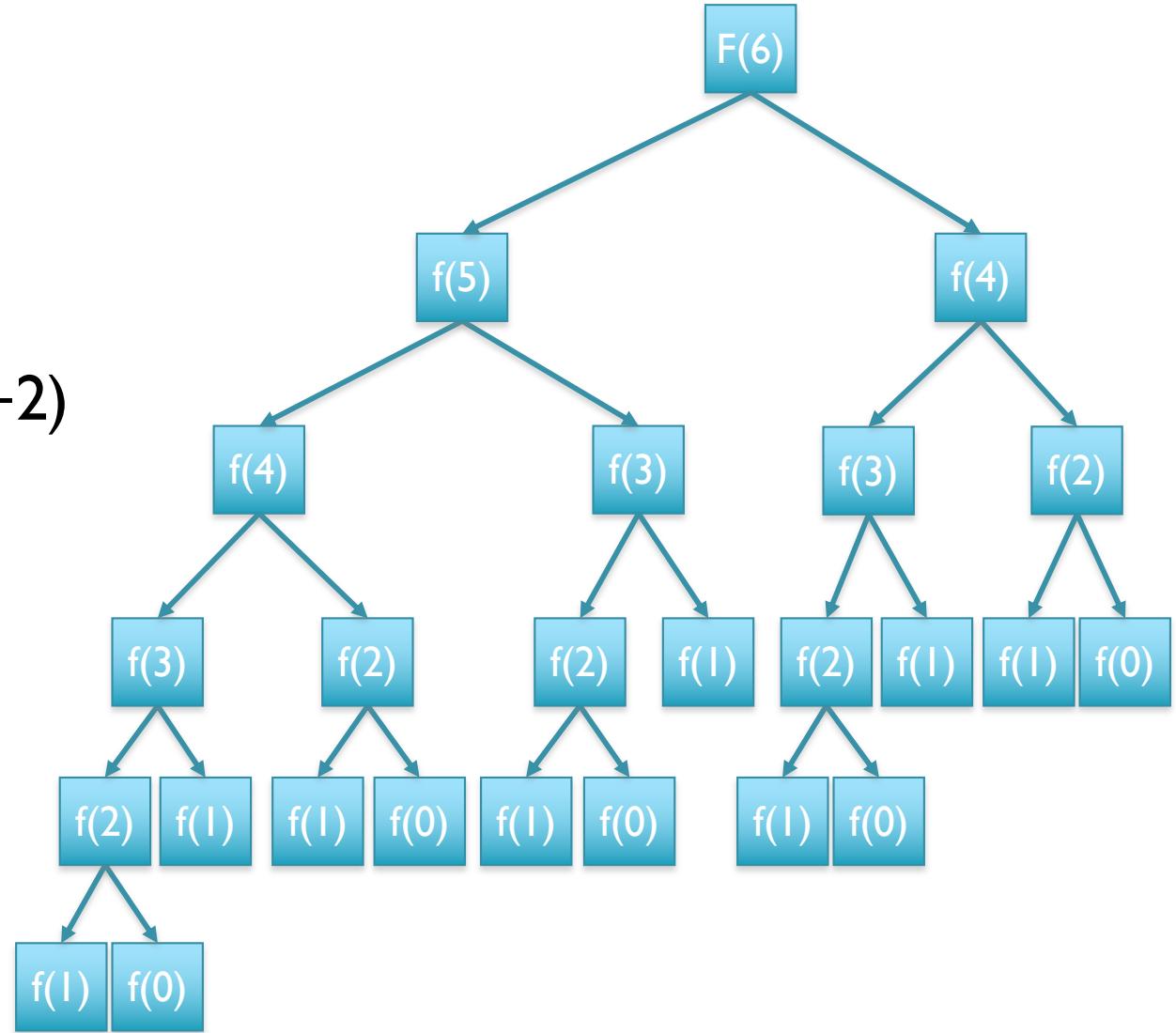
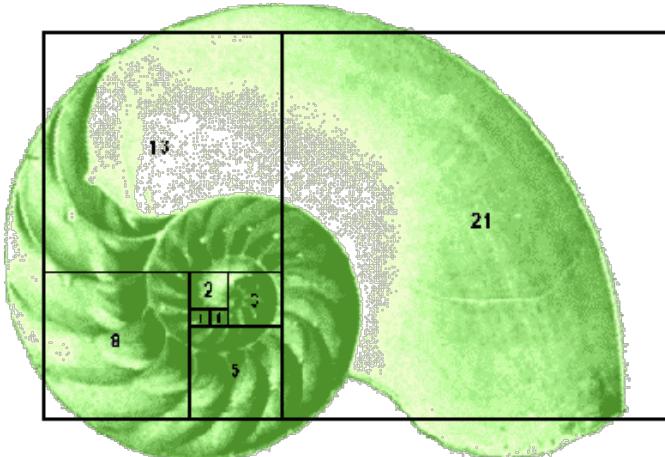
r1 0 chr12 1936 0
36M * 0 0
CCAGTAGCTCTCAGCCTTATTTACCCAGGCCTGTA
II
AS:i:0 XS:i:-2 XN:i:0
XM:i:0 XO:i:0 XG:i:0
NM:i:0 MD:Z:36 YT:Z:UU
...

(Langmead & Salzberg, 2012)

# Dynamic Programming

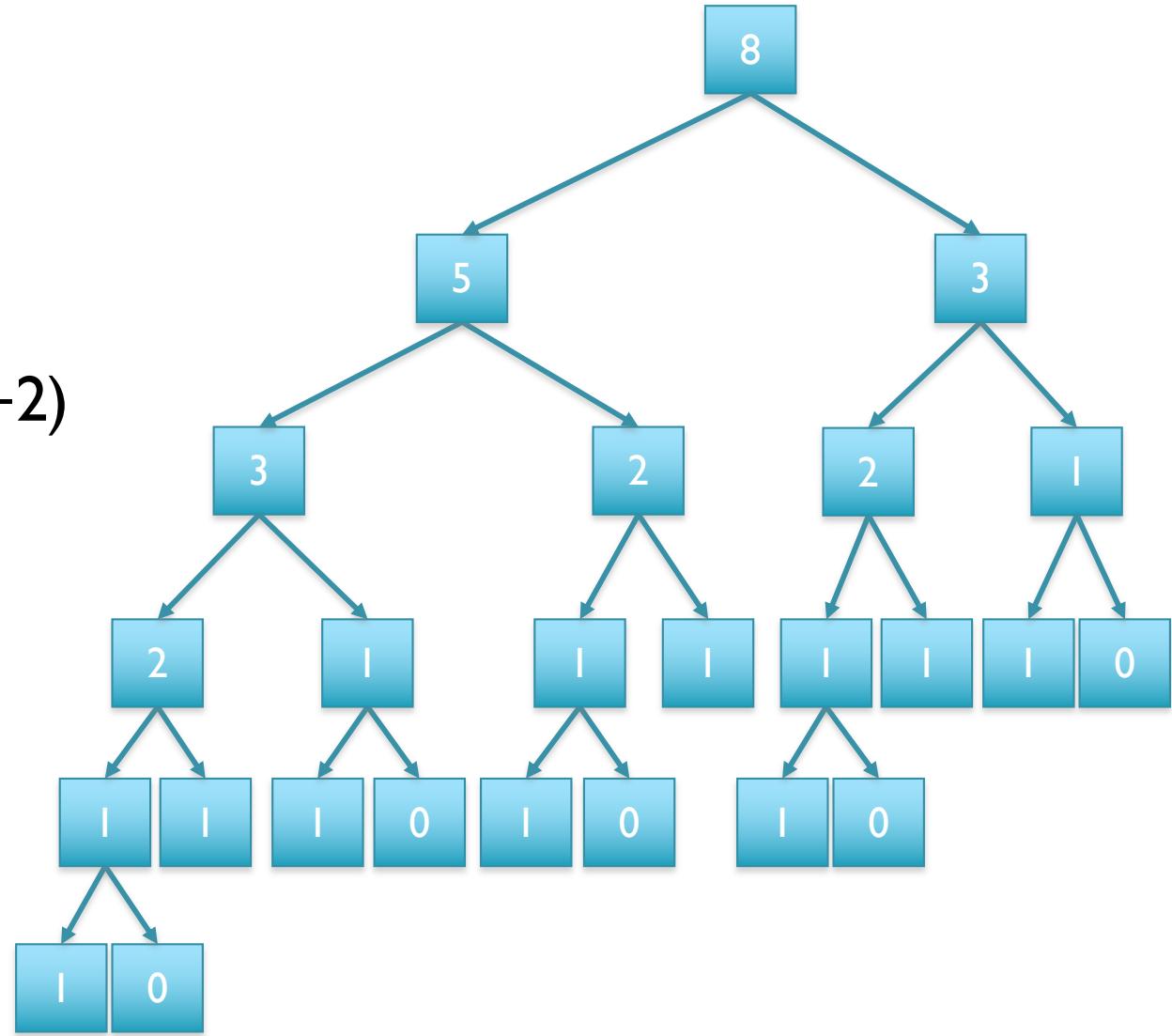
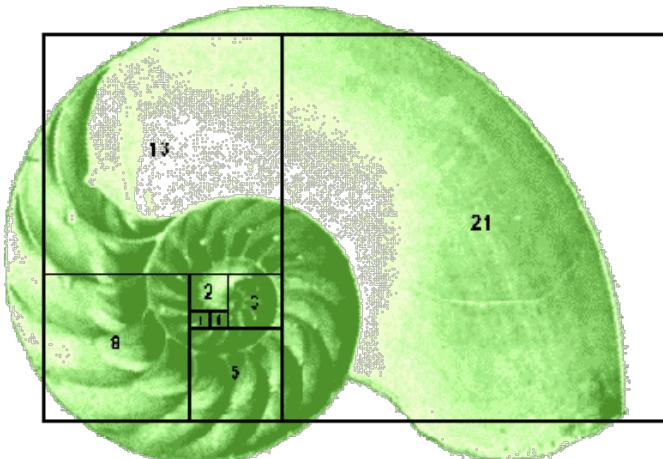
# Fibonacci Sequence

```
def fib(n):  
    if n == 0 or n == 1:  
        return n  
  
    else:  
        return fib(n-1) + fib(n-2)
```



# Fibonacci Sequence

```
def fib(n):
    if n == 0 or n == 1:
        return n
    else:
        return fib(n-1) + fib(n-2)
```



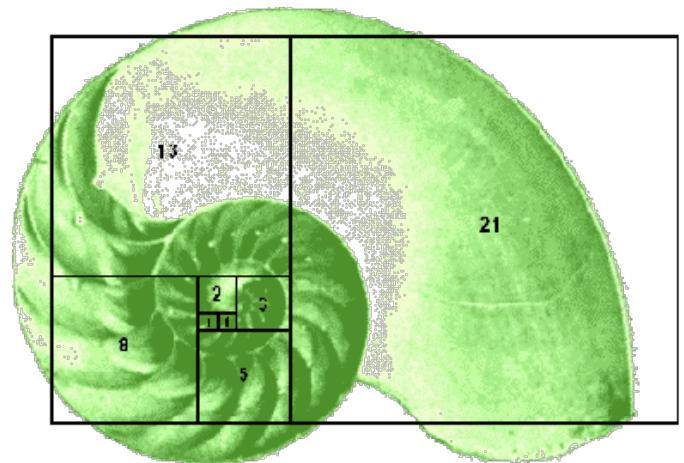
## What is the running time?

# Bottom-up Fibonacci Sequence

```
def fib(n):  
    table = [0] * (n+1)  
    table[0] = 0  
    table[1] = 1  
    for i in range(2,n+1):  
        table[i] = table[i-2] + table[i-1]  
    return table[n]
```

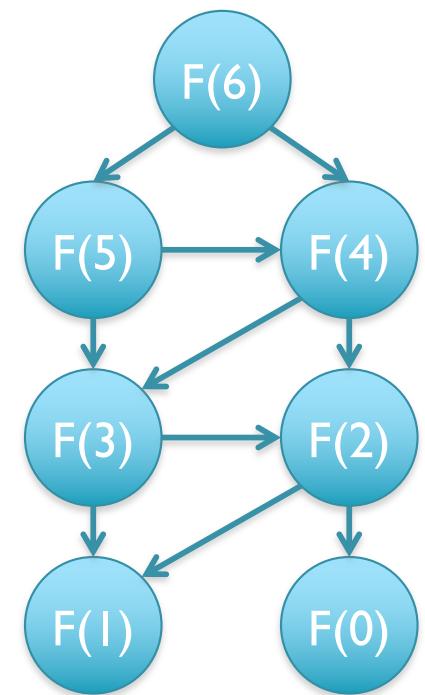
0	1	2	3	4	5	6
0	1	1	2	3	5	8

[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.



# Similarity metrics

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

MIKESCHATZ

| | X | | XXXX |

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

# 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:

... <b>M</b>	... <b>I</b>	... <b>D</b>
...A	...-	...A
...A	...A	...-

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

... <b>MM</b>	... <b>IM</b>	... <b>DM</b>	... <b>MI</b>	... <b>II</b>	... <b>DI</b>	... <b>MD</b>	... <b>ID</b>	... <b>DD</b>
...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...

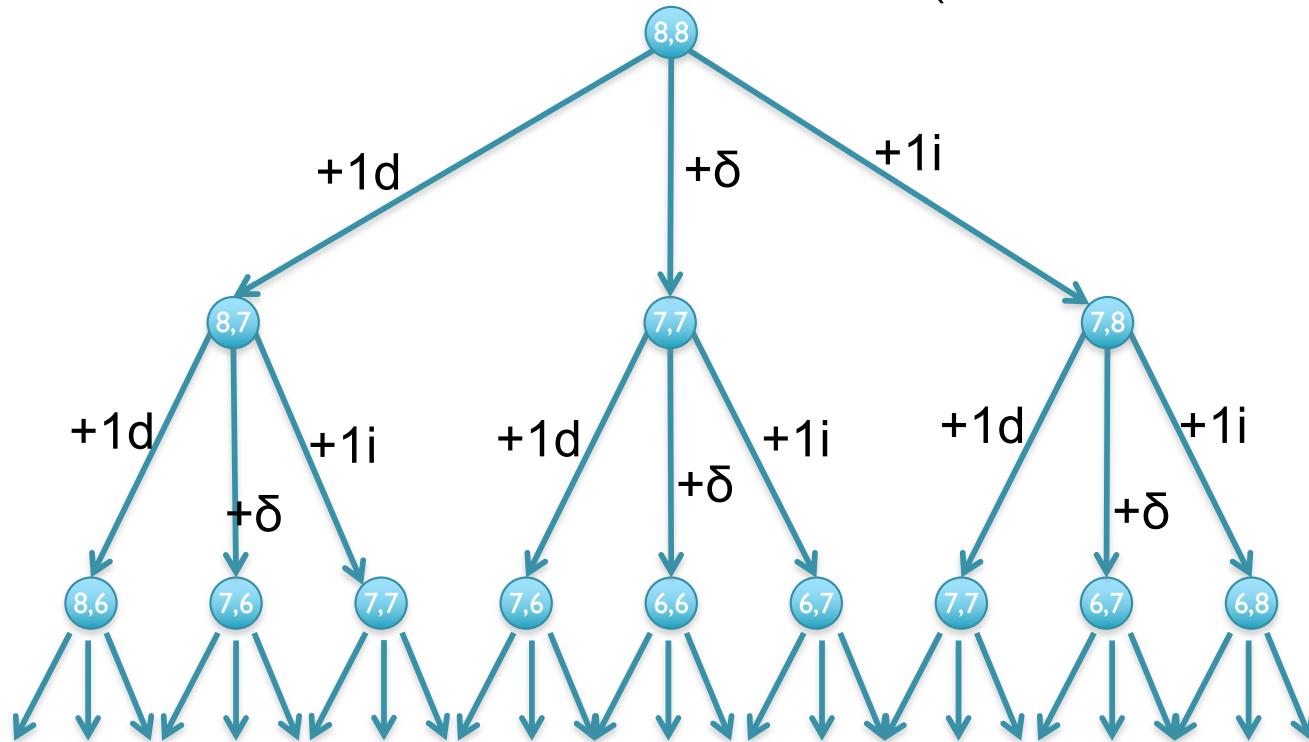
... <b>M</b> ...	... <b>I</b> ...	... <b>D</b> ...
...X...	...-...	...X...
...Y...	...Y...	...-...

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  $j' \leq j$ .

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



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

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

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

[Why do we want the min?]

# Dynamic Programming Matrix

		A	C	A	C	A	C	T	A
	0	I	2	3	4	5	6	7	8
A	I								
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	I	2	3	4	5	6	7	8
A	I	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	I	2	3	4	5	6	7	8
A	I	0	I						
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	I	2	3	4	5	6	7	8
A	I	0	I	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	I	0	I	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

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

$$D[AG, ACACACTA] = 7$$

----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?]

# Alignment Ambiguity

Which of these is the correct alignment?

GATTTACA  
GATT-ACA

GATTTACA  
GAT-TACA

GATTTACA  
GA-TTACA

Which of these is the correct alignment?

GATATATAACA  
GATAT--ACA

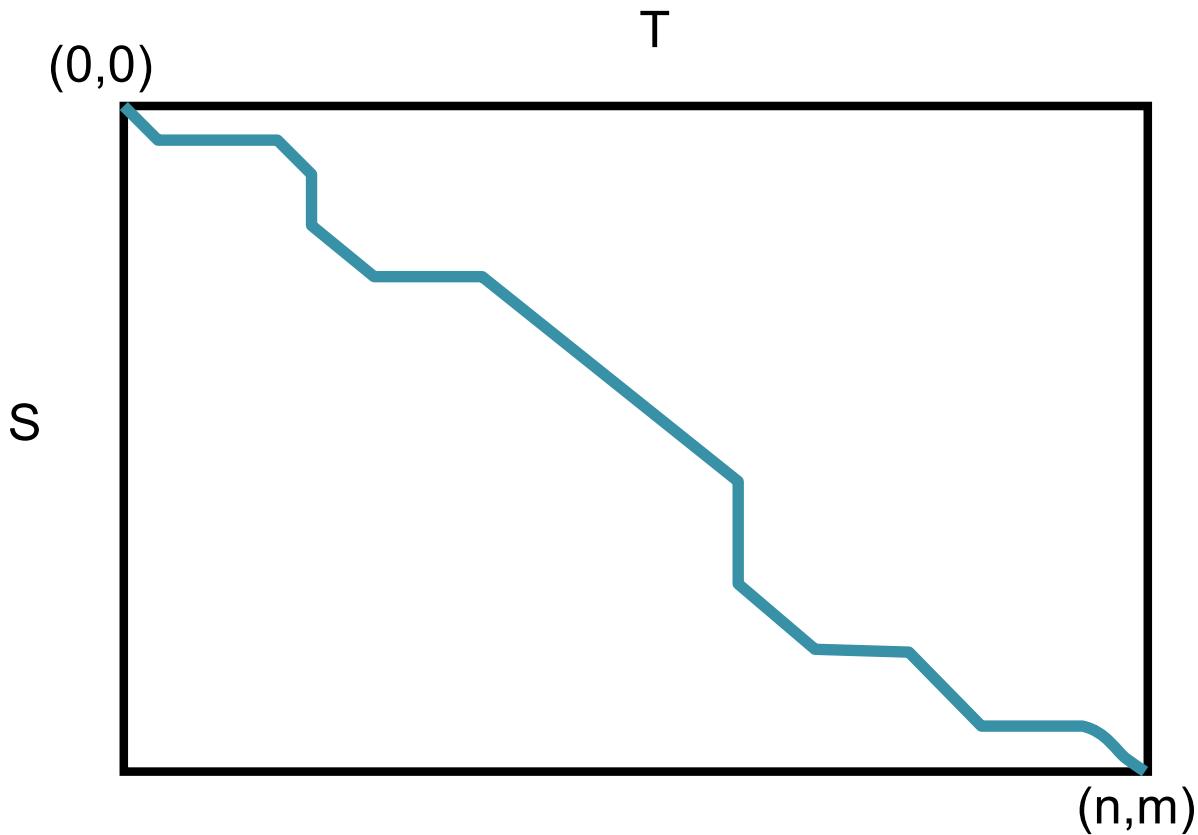
GATATATAACA  
GAT--ATACA

GATATATAACA  
GAT-TA-ACA

They are all equally good!

- Generally prefer to merge gaps together when possible (affine gaps)
- Use variant “normalization” to shift gaps into a canonical position
  - Commonly shift gaps as far to the left as possible

# Global Alignment Schematic

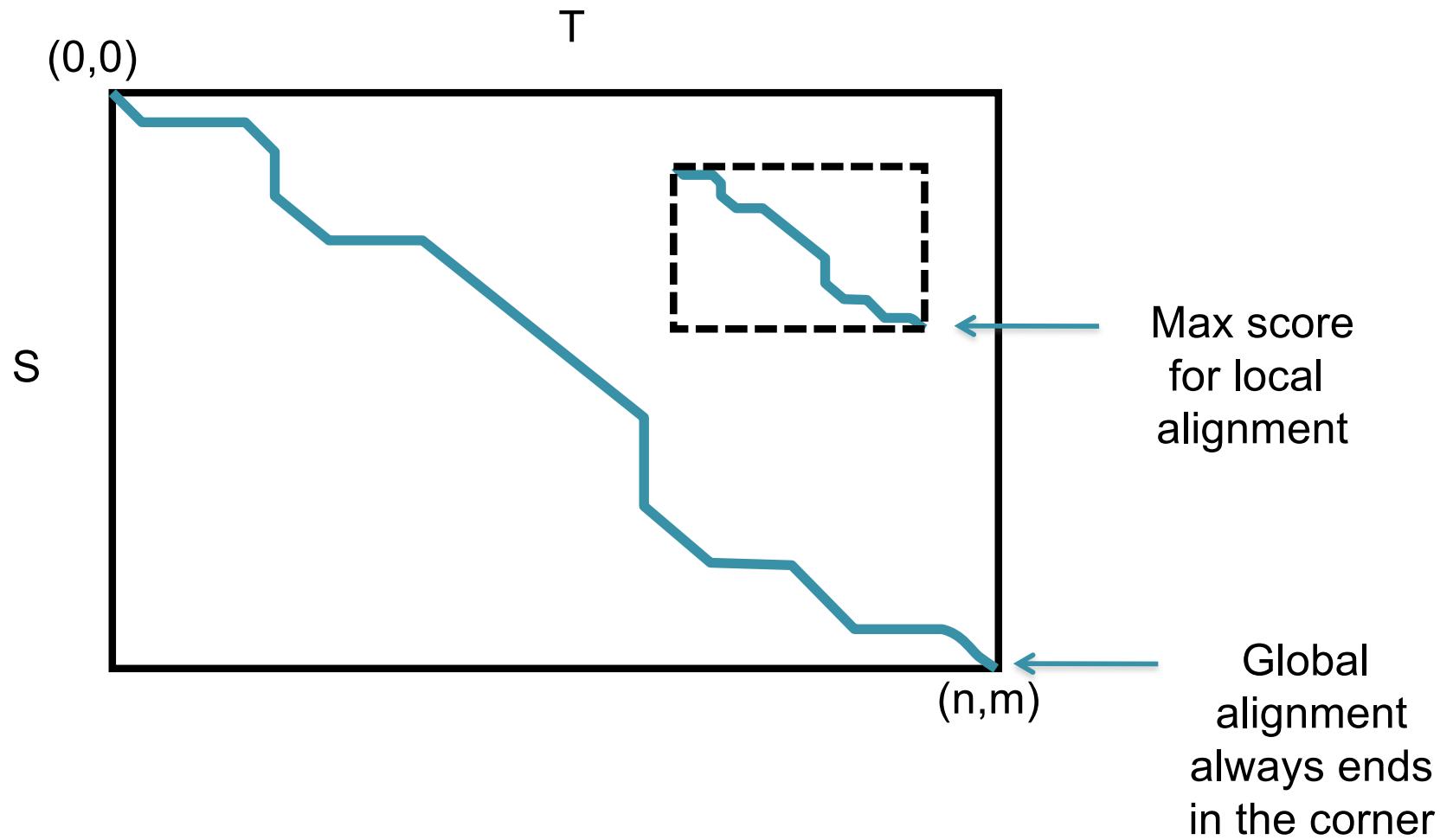


- A high quality alignment will stay close to the diagonal
  - If we are only interested in high quality alignments, we can skip filling in cells that can't possibly lead to a high quality alignment
  - Find the global alignment with at most edit distance  $d$ :  $O(2dn)$

# Local vs. Global Alignment

- The Global Alignment Problem tries to find the best end-to-end alignment between the two strings
  - Only applicable for very closely related sequences
- The Local Alignment Problem 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



# 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

```
tccCAGTTATGTCAGggacacgagcatgcagagac  
|||||||||||||  
aattgccgcgtcgatcagCAGTTATGTCAGatc
```

# SAM / BAM Files

## 1.4 The alignment section: mandatory fields

In the SAM format, each alignment line typically represents the linear alignment of a segment. Each line consists of 11 or more TAB-separated fields. The first eleven fields are always present and in the order shown below; if the information represented by any of these fields is unavailable, that field's value will be a placeholder, either '0' or '\*' as determined by the field's type. The following table gives an overview of these mandatory fields in the SAM format:

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,254}	Query template NAME
2	FLAG	Int	[0, 2 <sup>16</sup> - 1]	bitwise FLAG
3	RNAME	String	\*  [:rname:^*]=[:rname:]*	Reference sequence NAME <sup>11</sup>
4	POS	Int	[0, 2 <sup>31</sup> - 1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0, 2 <sup>8</sup> - 1]	MAPping Quality
6	CIGAR	String	\*  ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	\* =  [:rname:^*]=[:rname:]*	Reference name of the mate/next read
8	PNEXT	Int	[0, 2 <sup>31</sup> - 1]	Position of the mate/next read
9	TLEN	Int	[-2 <sup>31</sup> + 1, 2 <sup>31</sup> - 1]	observed Template LENgth
10	SEQ	String	\*  [A-Za-z.=.]+	segment SEQuence
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33

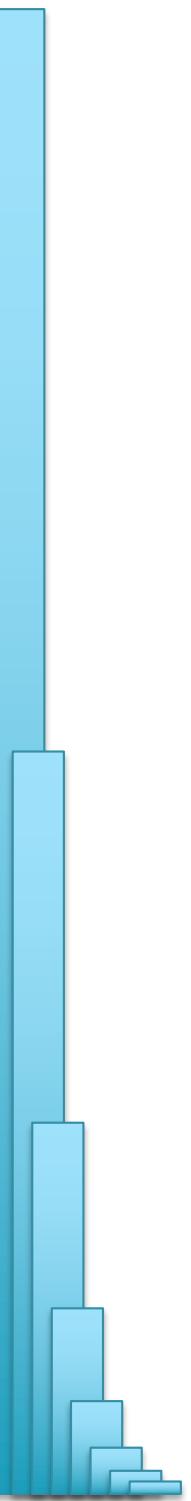
<sup>11</sup>Reference sequence names may contain any printable ASCII characters with the exception of certain punctuation characters, and may not start with '\*' or '='. See Section 1.2.1 for details and an explanation of the [:rname:] notation.

# SAM / BAM Files

```
Coor      12345678901234 5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

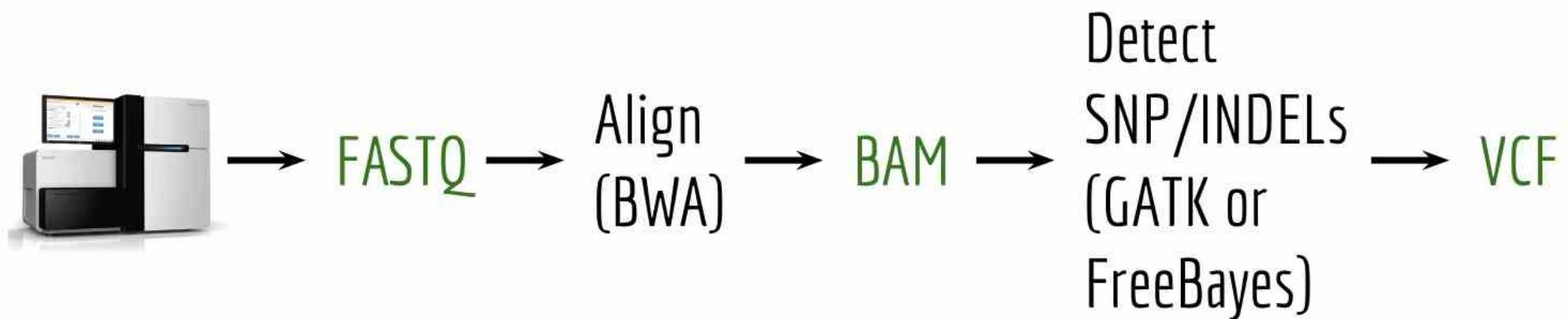
+r001/1    TTAGATAAAGGATA*CTG
+r002    aaaAGATAA*GGATA
+r003    gcctaAGCTAA
+r004    ATAGCT.....TCAGC
-r003    ttagctTAGGC
-r001/2    CAGCGGCAT
```

```
@HD VN:1.6 SO:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

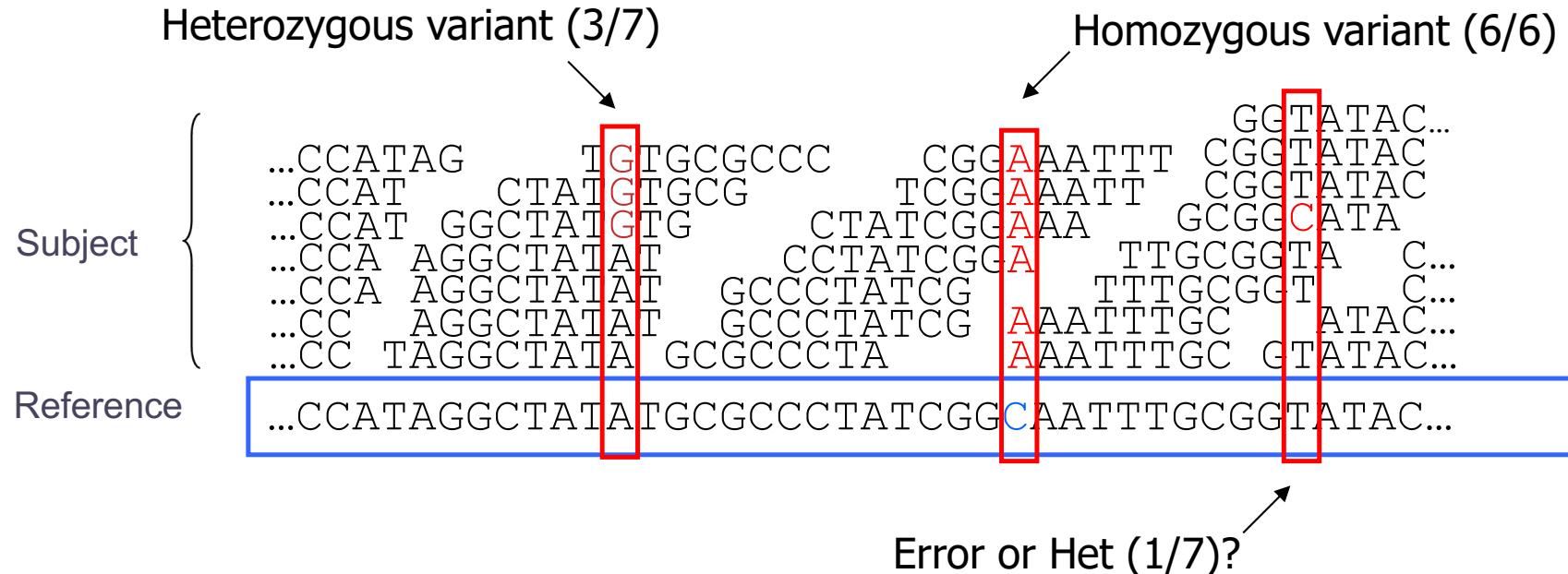


## Part 2: Variant Calling

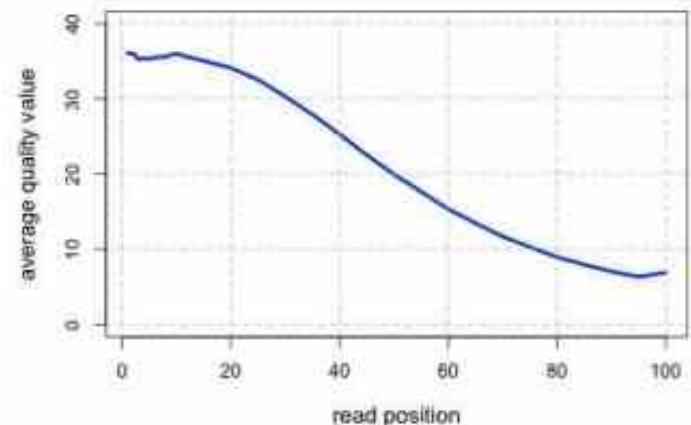
# Variant Calling Overview



# Genotyping Theory



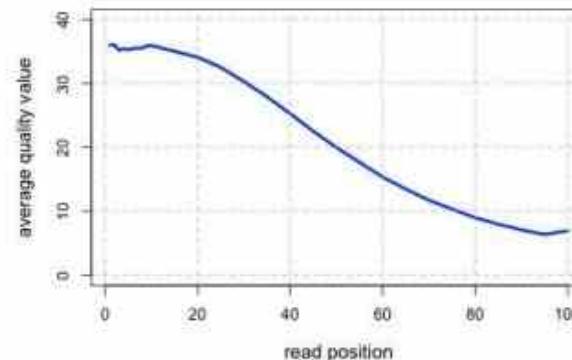
- If there were no sequencing errors, identifying SNPs would be very easy: any time a read disagrees with the reference, it must be a variant!
- Sequencing instruments make mistakes
  - Quality of read decreases over the read length
- A single read differing from the reference is probably just an error,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



# Illumina Quality

QV	p <sub>error</sub>
40	1/10000
30	1/1000
20	1/100
10	1/10

$$Q_{\text{sanger}} = -10 \log_{10} p$$



S - Sanger Phred+33, raw reads typically (0, 40)  
 X - Solexa Solexa+64, raw reads typically (-5, 40)  
 I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)  
 J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)  
     with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (**bold**)  
     (Note: See discussion above).  
 L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

# The Binomial Distribution: Adventures in Coin Flipping

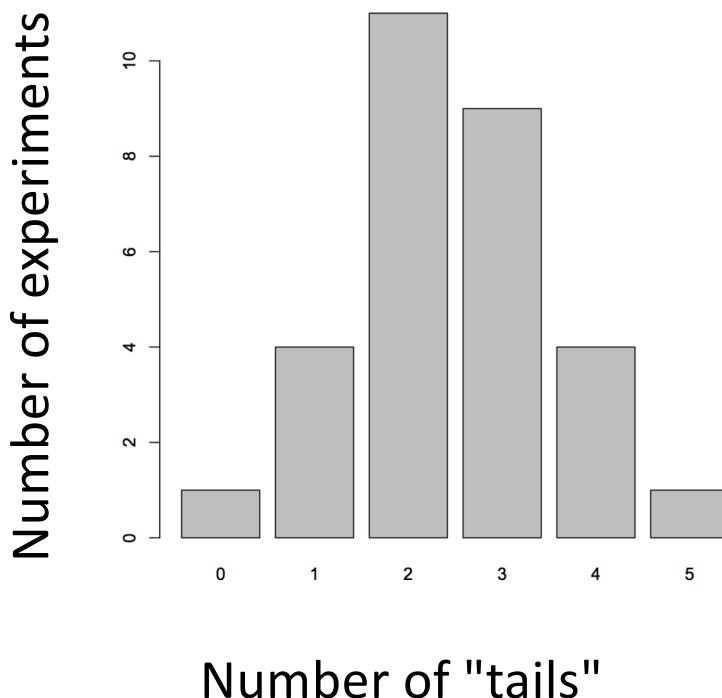


$P(\text{heads}) = 0.5$



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

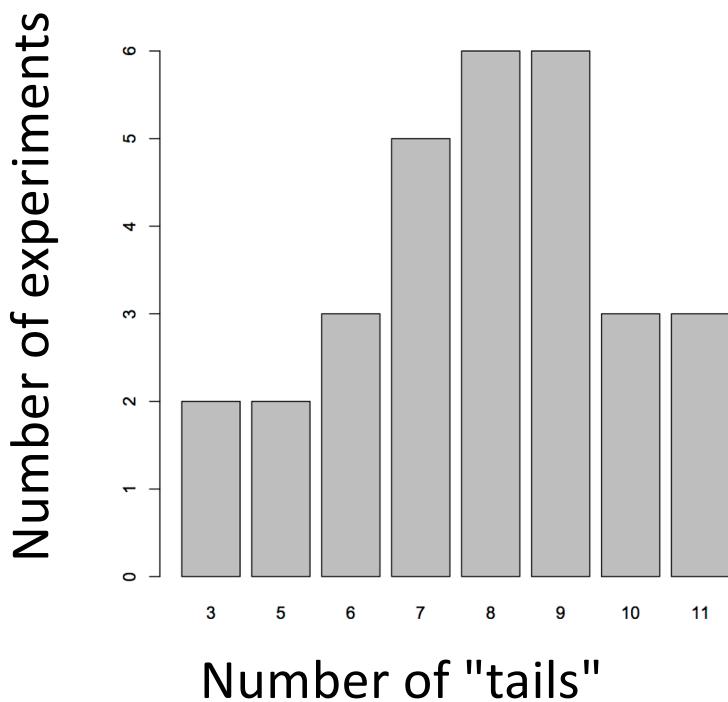
# What is the distribution of tails (alternate alleles) do we expect to see after 5 tosses (sequence reads)?



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:

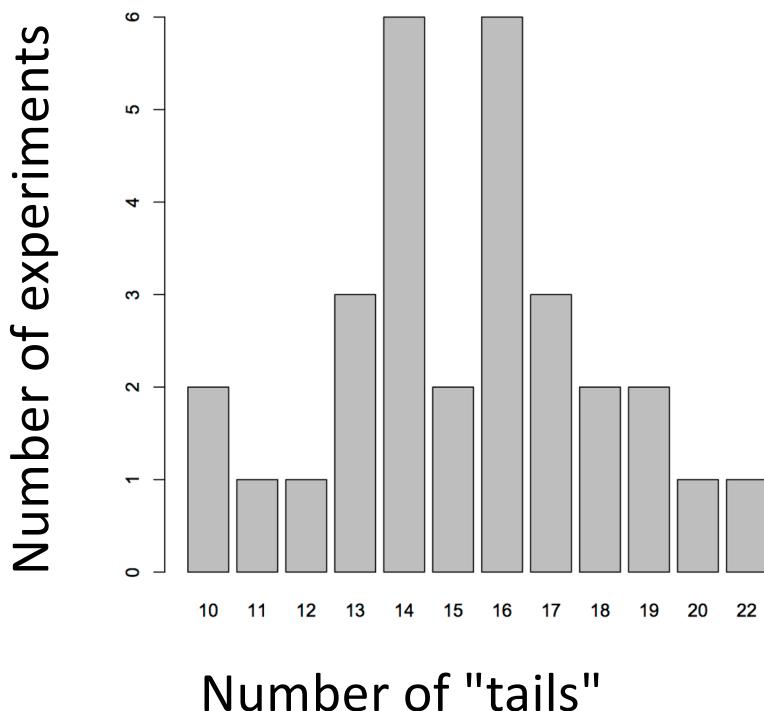
```
barplot(table(rbinom(30, 15, 0.5)))
```

30 experiments (students tossing coins)

15 tosses each

Probability of Tails

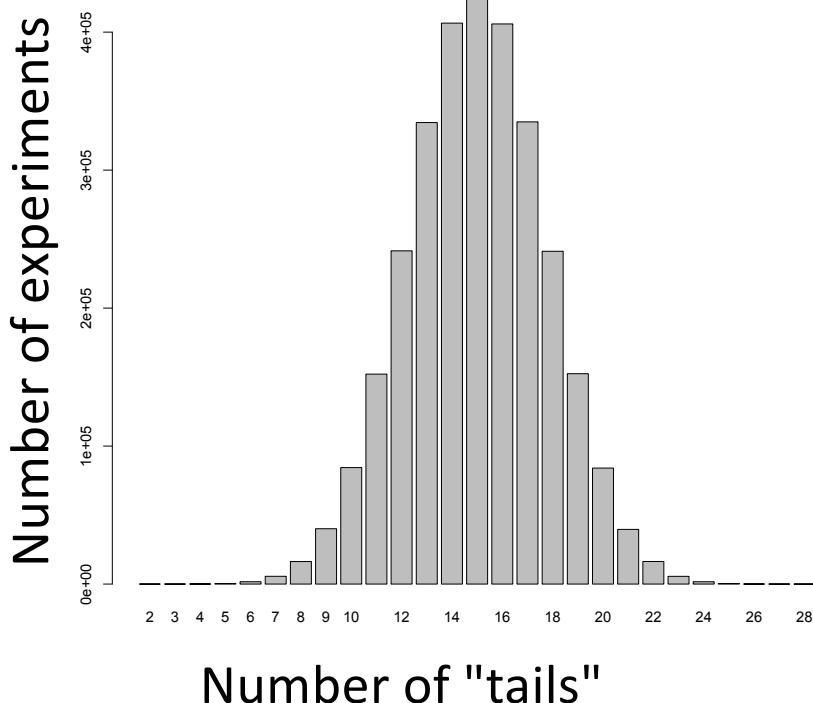
# What is the distribution of tails (alternate alleles) do we expect to see after 30 tosses (sequence reads)?



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:

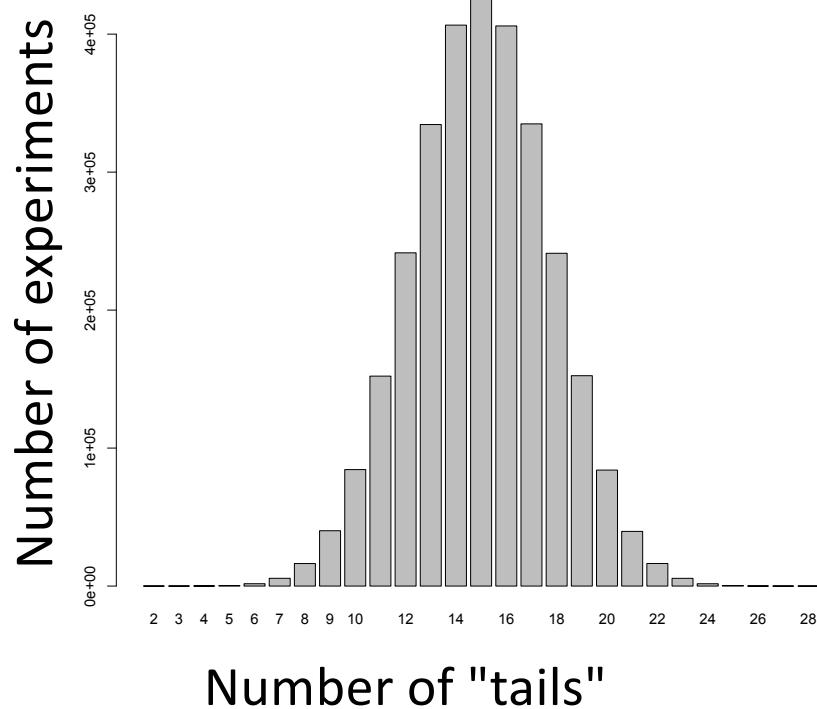
```
barplot(table(rbinom(3e6, 30, 0.5)))
```

3M experiments (students tossing coins)

30 tosses each

Probability of Tails

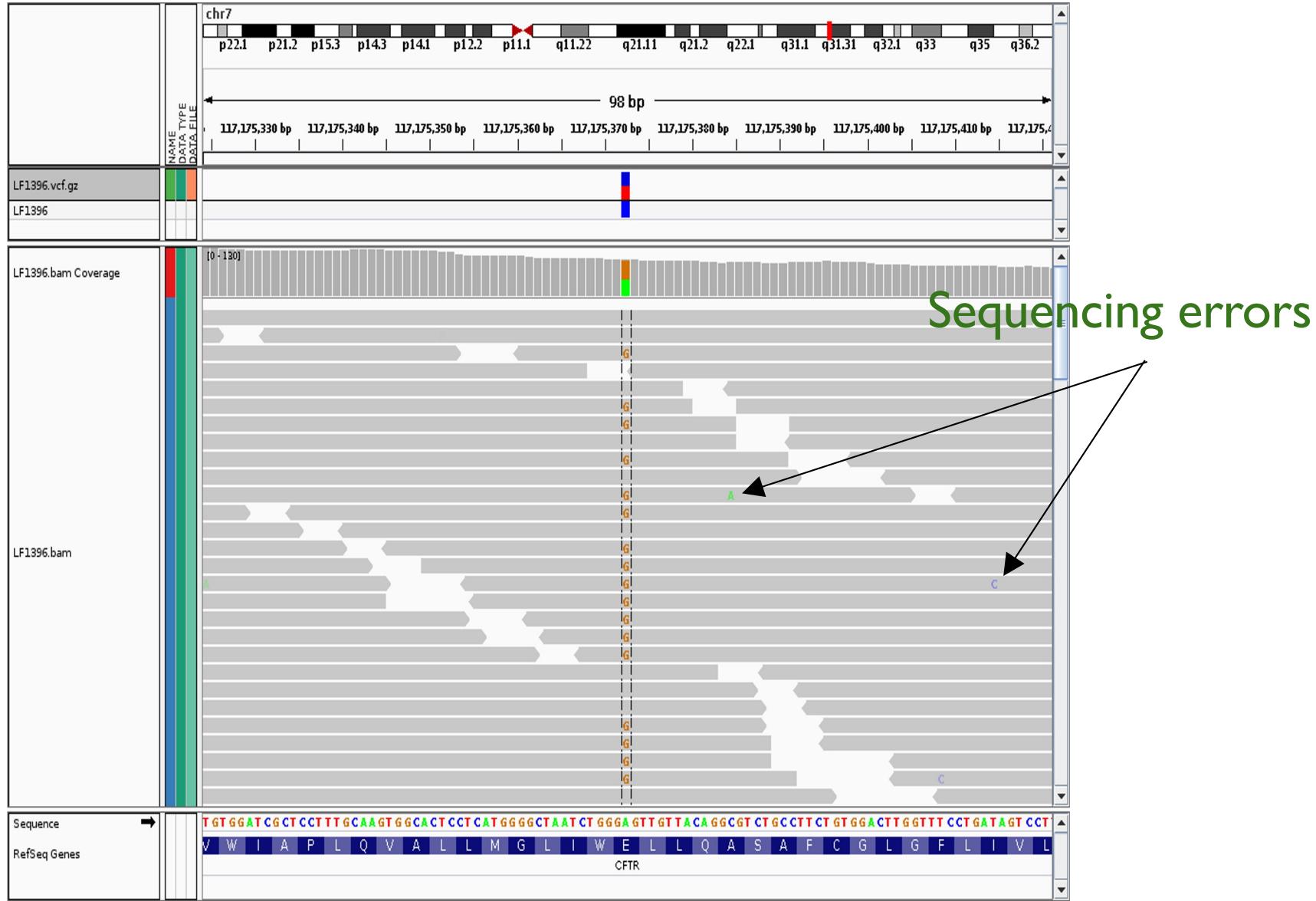
So, with 30 tosses (reads), we are much more likely to see an even mix of alternate and reference alleles at a heterozygous locus in a genome



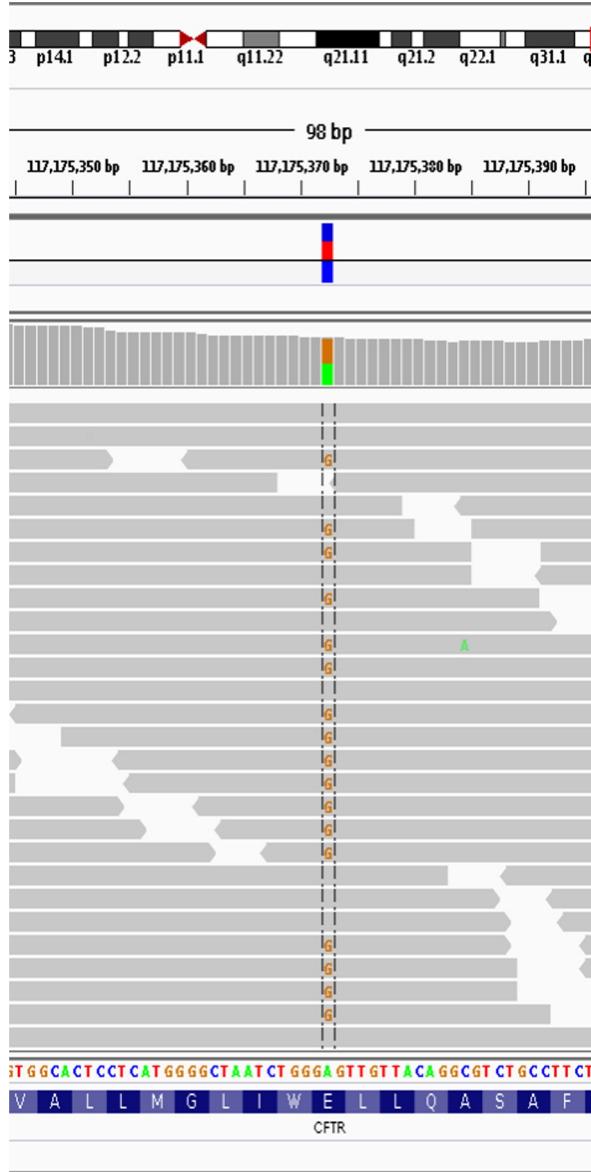
This is why at least a "30X" (30 fold sequence coverage) genome is recommended: it confers sufficient power to distinguish heterozygous alleles and from mere sequencing errors

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

# 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

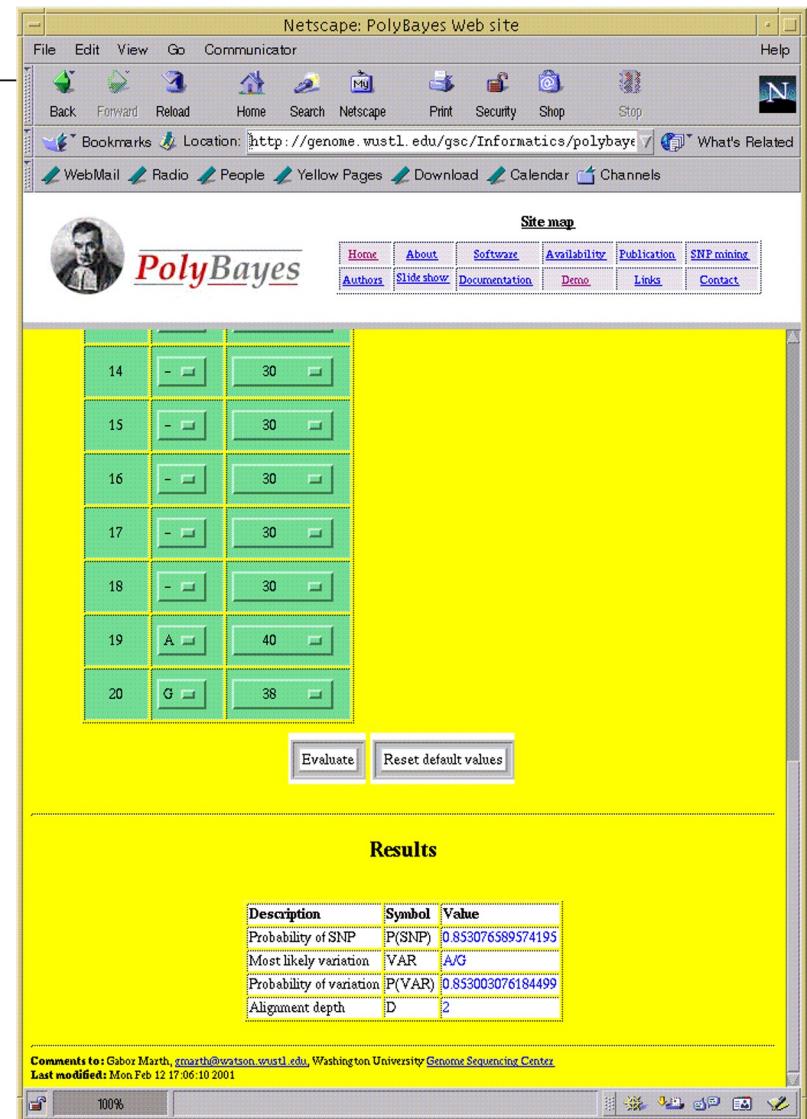


© 1999 Nature America Inc. • <http://genetics.nature.com>

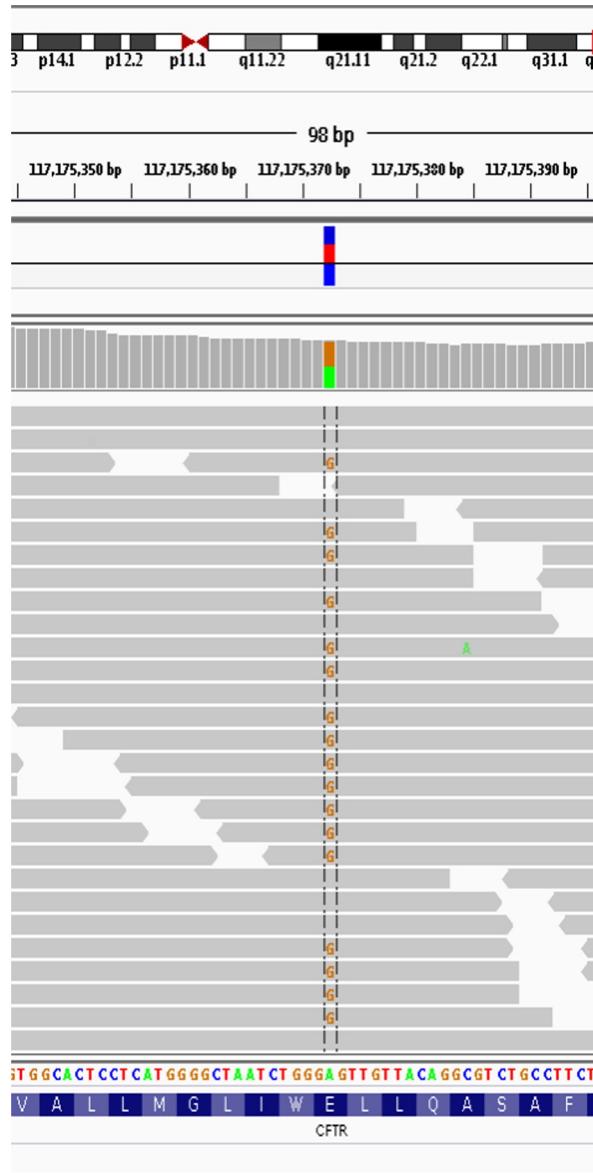
## A general approach to single-nucleotide polymorphism discovery

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 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>

Bayesian posterior probability

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

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

Base call +  
Base quality

Expected (prior)  
polymorphism rate

# PolyBayes: The first statistically rigorous variant detection tool.

*letter*



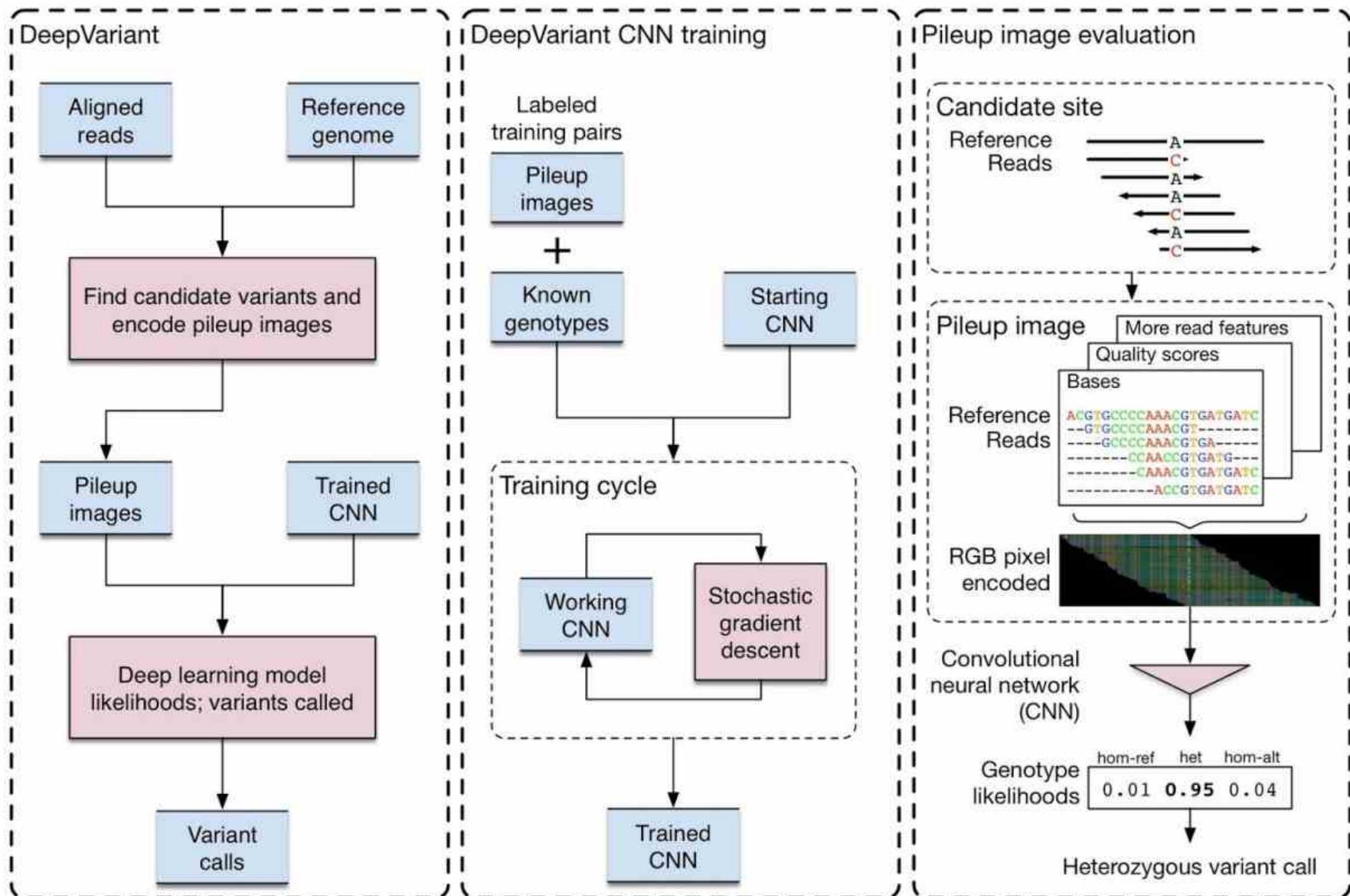
© 1999 Nature America Inc. • <http://genetics.nature.com>

## 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>

This Bayesian statistical framework has been adopted by most other modern 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>

# 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=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">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE1 SAMPLE2
1 1 . ACG A,AT .
1 2 rsl C T,CT .
1 5 . A G .
1 100 T <DEL> .
.
```

Mandatory header lines

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

Body {

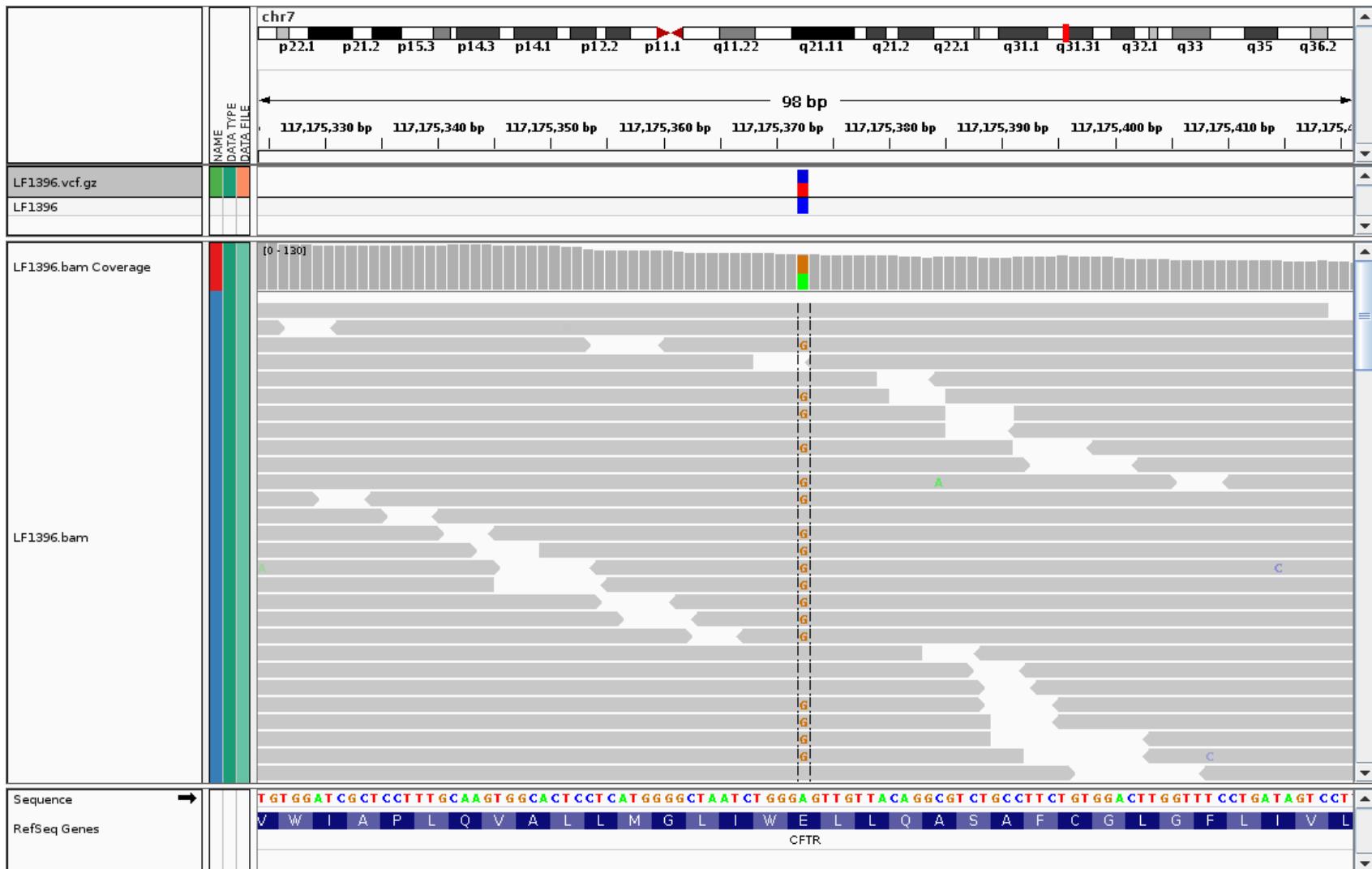
Deletion      SNP      Large SV      Insertion      Other event

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