



The Elizabeth H.  
and James S. McDonnell III

# McDONNELL GENOME INSTITUTE

at Washington University

## Alignment concepts

Felicia Gomez, PhD

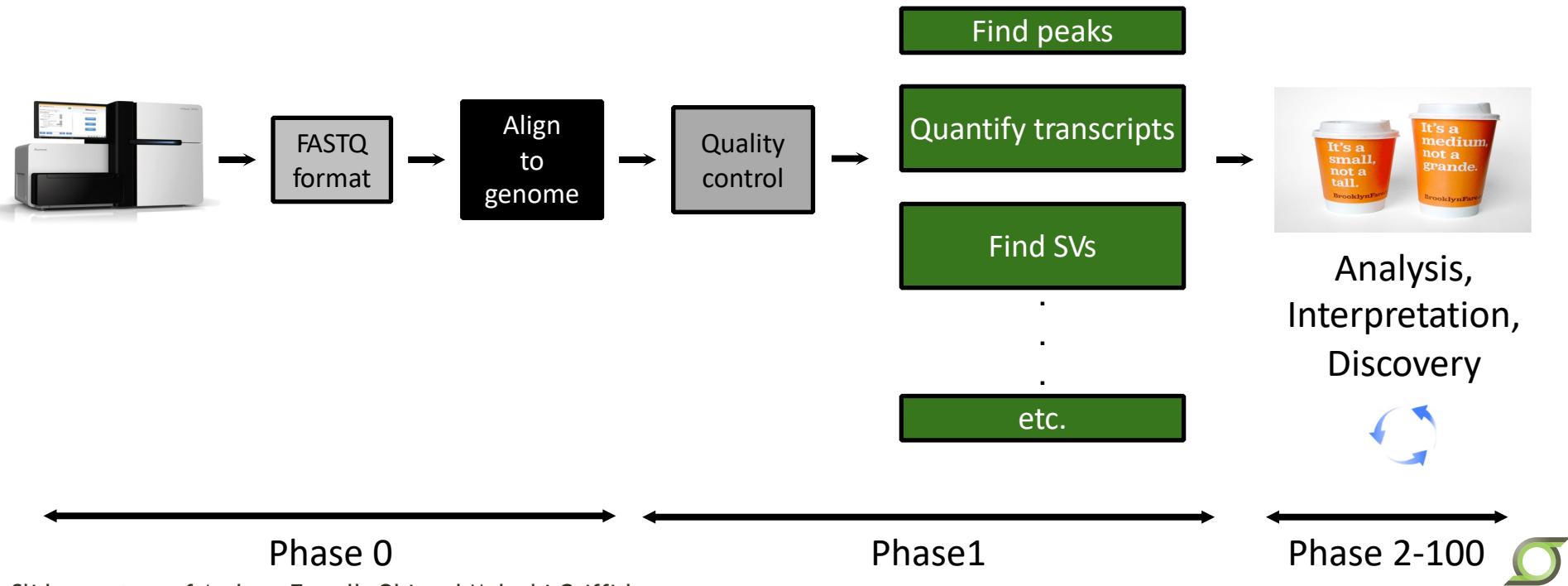
Assistant Professor

Washington University in Saint Louis

11/10/2023 - Advanced Sequencing Technologies and  
Applications

Cold Spring Harbor

# Alignment is central to most genomic research

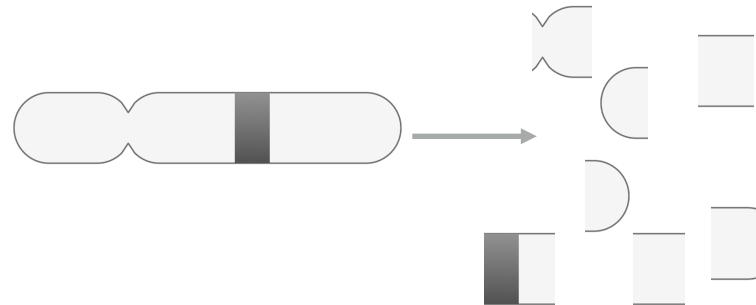


# Generate sequence reads

Fragment a genome → DNA library

PCR amplification

Sequence reads (ends of DNA fragment for mate pairs)

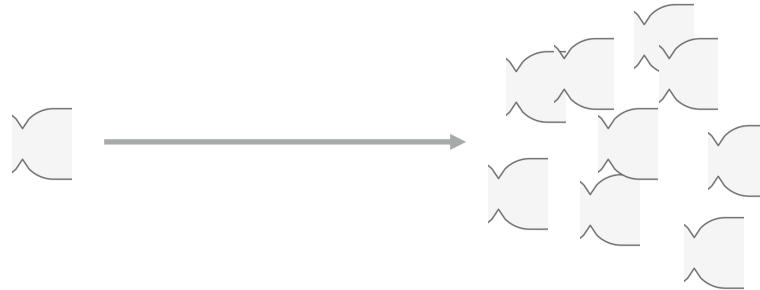


# Generate sequence reads

Fragment a genome → DNA library

PCR amplification

Sequence reads (ends of DNA fragment for mate pairs)



# Generate sequence reads

Fragment a genome → DNA library

PCR amplification

Sequence reads (ends of DNA fragment for mate pairs)



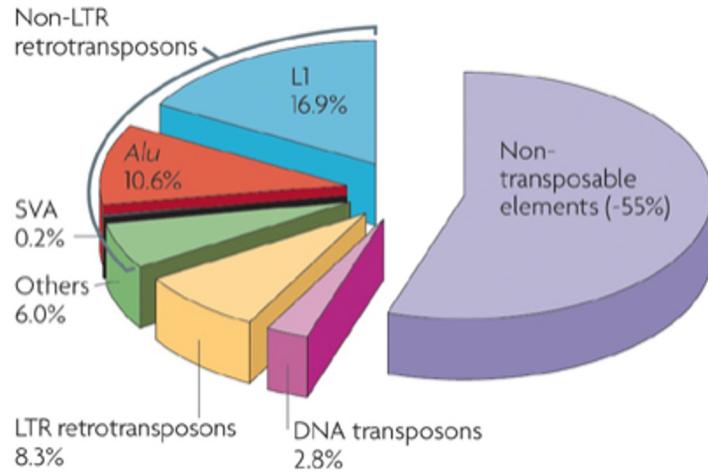
# We have FASTQ files. Now what?

- Need to find a home for every read in the file.
- Must get the alignment just right. Or else.....problems.
- Must choose the right tool for the experiment.

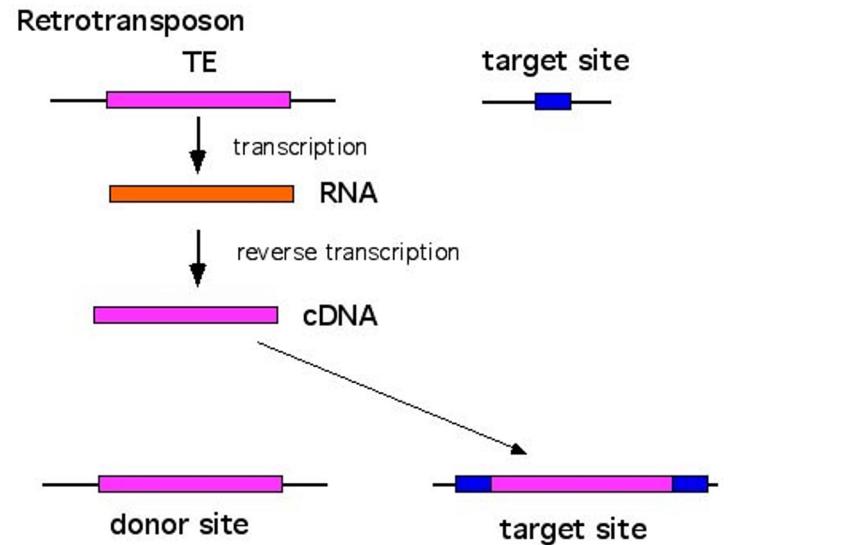


# Problem: Half of the human genome is comprised of repeats

a



McClintock's  
"jumping  
genes" in  
maize



Retrotransposons use a "copy/paste" mechanism

DNA transposons use a "cut/paste" mechanism

<http://www.nature.com/nrg/journal/v10/n10/pdf/nrg2640.pdf>

Slide courtesy of Aaron Quinlan

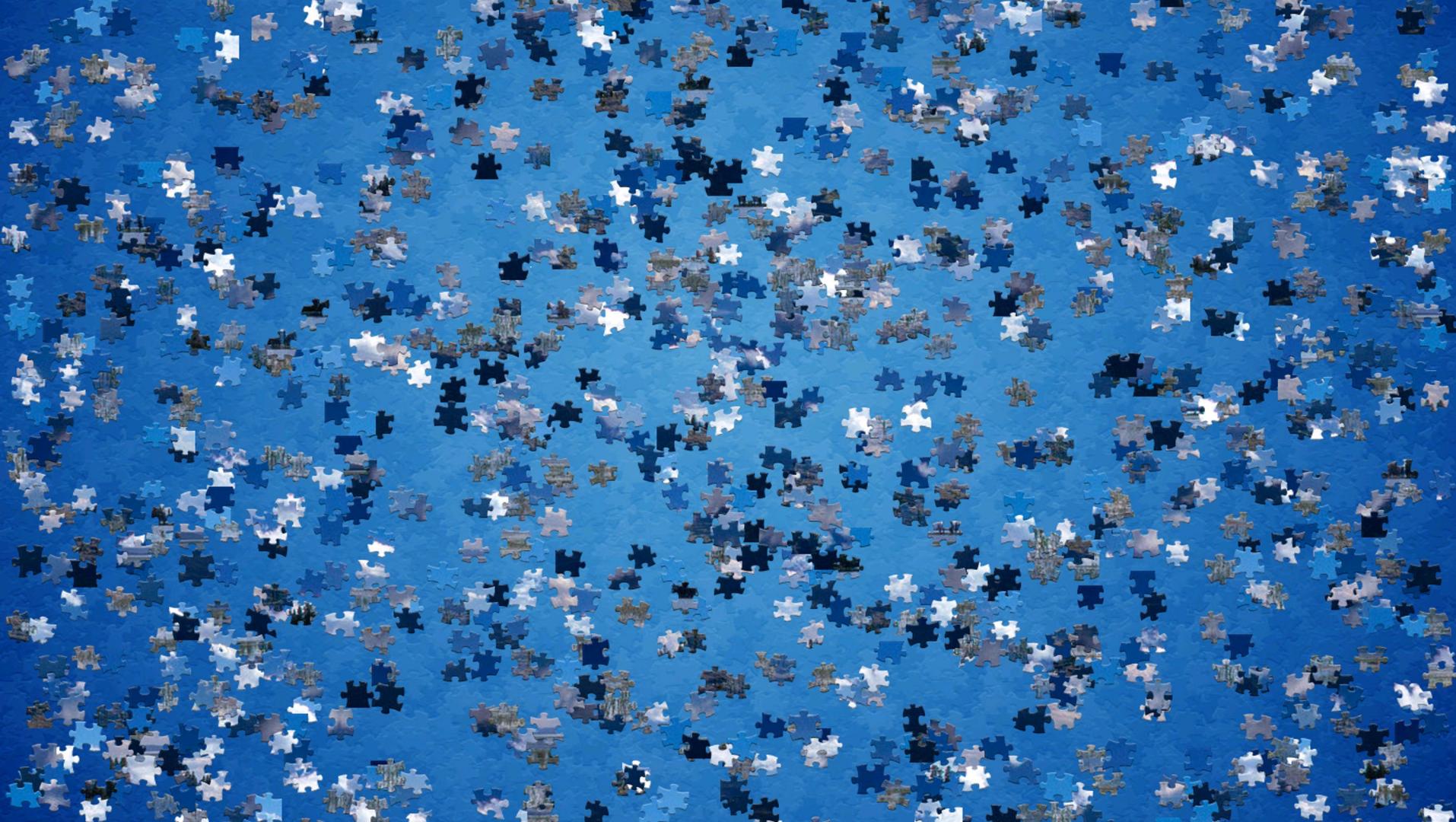


## Problem: Half of the human genome is comprised of repeats

( first bit of human chromosome 1 )







# Best case scenario: an error-free sequencing technology

ATTCGAAACA  
TTCGCGCAAT  
CTGGACTCAA



ATTCGAAACA  
TTCGCGCAAT  
CTGGACTCAA

→ Aligner →

TACCTCCAGGGGGCATCCTCCCCCA**ATTC**  
**GAAACA**CAATCGTAGCCCCTGGCACTACCTA  
TGTGTGTCAATTGGAGAGAGAGAGATTACAC  
GAAAAAAAAGT**CTGGACTCAA**CTAGGATACA  
CACATTGGCTACAGATACCAAAAAAAAAAA  
AAAAAAAATTTACCATTGAGGCACCACT  
TCTCGTCGCTCGTCGCTTGCTCGCTTCGG  
CTAAAAA**ATTCGCGCAAT**ACATTGGCTACAG  
ATACCAAAAAAA

Computers are rather good at finding *exact* matches.  
Think Google.



# Reality check. Errors happen. Frequently.

ATTCGAAACA



→ ATTTGAAACA

→ Aligner →

TACCTCCAGGGGGCATCCTCCCCCAATTC  
GAAACACAATCGTAGCCCTGGCACTACCTA  
TGTGTGTCAATTGGAGAGAGAGAGATTGGA  
AACAAAAAAGTGCTACAGATACCAACTAGGAT  
ACACACATTGGCTACAGATACCAACCAAAAAAAA  
AAAAAAAAAAATTTCACCATTGAGGCACCA  
CCTCTCGTCGCTGCGTCGCTGCTCGCGG  
CTAAAAAAATTGAAACAACATTGGCTACAG  
ATACCAAAATT

“Fuzzy” matching is much more computationally expensive.

Think Google’s “Did you mean...”



# There *are* optimal solutions.

Reference

cgggtatccaa

Read

cccttaggtcccc

What is the best alignment?



Reference

cgggtatccaa

Read

ccctaggtcccc

Reference

c~~ggg~~ta--t-ccaa

Read

c~~cc~~-taggt~~cc~~-a



Reference	cgggtatccaa
Read	ccctaggtcccc
Reference	c <span style="color:red">gggt</span> a--t-ccaa
Read	c <span style="color:red">cc-</span> taggt <span style="color:green">ccc-</span> a
Reference	c <span style="color:green">gggt</span> a---tccaa
Read	c <span style="color:red">c--</span> taggt <span style="color:green">cccc</span> a



Reference	cgggtatccaa
Read	ccctaggtcccc
Reference	c <span style="color:red">gggt</span> a--t-ccaa
Read	c <span style="color:red">cc-</span> taggt <span style="color:green">ccc-</span> a
Reference	c <span style="color:green">gggt</span> a---tccaa
Read	c <span style="color:red">c--c</span> taggtccca
Reference	c- <span style="color:green">gggt</span> a--tccaa
Read	c <span style="color:green">c--c</span> taggtccca



# Global: Needleman-Wunsch algorithm

# Local: Smith-Waterman algorithm

**A General Method Applicable to the Search for Similarities  
in the Amino Acid Sequence of Two Proteins**

SAUL B. NEEDLEMAN AND CHRISTIAN D. WUNSCH

*Department of Biochemistry, Northwestern University, and  
Nuclear Medicine Service, V. A. Research Hospital  
Chicago, Ill. 60611, U.S.A.*

(Received 21 July 1969)

*J. Mol. Biol.* (1981), **147**, 195–197

## Identification of Common Molecular Subsequences

The identification of maximally homologous subsequences among sets of long sequences is an important problem in molecular sequence analysis. The problem is straightforward only if one restricts consideration to contiguous subsequences (segments) containing no internal deletions or insertions. The more general problem has its solution in an extension of sequence metrics (Sellers 1974; Waterman *et al.*, 1976) developed to measure the minimum number of “events” required to convert one sequence into another.



# Local: Smith-Waterman algorithm

5' ACTACTAGATTACCTACGGATCAGGTACTTAGAGGCTTGCAACCA 3'

||||| ||||||| |||||||||||||

5' TACTCACGGATGAGGTACTTAGAGGC 3'

# Global: Needleman-Wunsch algorithm

5' ACTACTAGATTACCTACGGATCAGGTACTTAGAGGCTTGCAACCA 3'

|||||||||||      |||||||      |||||||||||||      |||||

5' ACTACTAGATT----ACGGATC--GTACTTAGAGGCTAGCAACCA 3'



# Local Alignment: dynamic programming

Scoring scheme:

Match: +3

Mismatch -3

Gap: -2

Adding an extra  
row and column  
allows us to  
align any base to  
any other  
positions in the  
other sequence

	T	G	T	T	A	C	G	G
G								
G								
T								
T								
G								
A								
C								
T								
A								



# Local Alignment: dynamic programming

## Scoring scheme:

Match: +3

Mismatch -3

Gap: -2

Initialize the matrix  
with a minimum  
score of 0.

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0							
G	0							
T	0							
T	0							
G	0							
A	0							
C	0							
T	0							
A	0							



# Local Alignment: dynamic programming

**Scoring scheme:**

Match: +3

Mismatch -3

Gap: -2

**Mismatch, but  
score can't go  
below 0.**

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0						
G	0							
T	0							
T	0							
G	0							
A	0							
C	0							
T	0							
A	0							



# Local Alignment: dynamic programming

## Scoring scheme:

Match: +3

Mismatch -3

Gap: -2

	T	G	T	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0	0
G	0	0	3						
G	0								
T	0								
T	0								
G	0								
A	0								
C	0								
T	0								
A	0								



# Local Alignment: dynamic programming

**Scoring scheme:**

Match: +3

Mismatch -3

Gap: -2

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0	3	1	0			
G	0	0	3	1	0			
T	0	3	1	6	4			
T	0	3	1	4	9	7		
G	0				7	6		
A	0							
C	0							
T	0							
A	0							



# Local Alignment: dynamic programming

## Scoring scheme:

Match: +3

Mismatch -3

Gap: -2

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3
G	0	0	3	1	0	0	0	3
T	0	3	1	6	4	2	0	1
T	0	3	1	4	9	7	5	3
G	0	1	6	4	7	6	4	8
A	0	0	4	3	5	10	8	6
C	0	0	2	1	3	8	13	11
T	0	3	1	5	4	6	11	10
A	0	1	0	3	2	7	9	8



# The traceback

Start at max score,  
traceback to next  
highest score, and  
so on. Stop at zero

C  
C

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3
G	0	0	3	1	0	0	0	6
T	0	3	1	6	4	2	0	1
T	0	3	1	4	9	7	5	3
G	0	1	6	4	7	6	4	8
A	0	0	4	3	5	10	8	6
C	0	0	2	1	3	8	13	11
T	0	3	1	5	4	6	11	10
A	0	1	0	3	2	7	9	8



# The traceback

**Start at max score,  
traceback to next  
highest score, and  
so on. Stop at zero**

A C  
A C

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3
G	0	0	3	1	0	0	0	3
T	0	3	1	6	4	2	0	1
T	0	3	1	4	9	7	5	3
G	0	1	6	4	7	6	4	8
A	0	0	4	3	5	10	8	6
C	0	0	2	1	3	8	13	11
T	0	3	1	5	4	6	11	10
A	0	1	0	3	2	7	9	8



# The traceback

Start at max score,  
traceback to next  
highest score, and  
so on. Stop at zero

Gap b/c that is the  
edit that led to this  
cell (9->7)

- A C  
G A C

	T	G	T	T	A	C	G	G
0	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3
G	0	0	3	1	0	0	0	3
T	0	3	1	6	4	2	0	1
T	0	3	1	4	9	7	5	3
G	0	1	6	4	7	6	4	8
A	0	0	4	3	5	10	8	6
C	0	0	2	1	3	8	13	11
T	0	3	1	5	4	6	11	10
A	0	1	0	3	2	7	9	8



# The traceback

Start at max score,  
traceback to next  
highest score, and  
so on. Stop at zero

G T T - A C  
G T T G A C

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3
G	0	0	3	1	0	0	0	3
T	0	3	1	6	4	2	0	1
T	0	3	1	4	9	7	5	3
G	0	1	6	4	7	6	4	8
A	0	0	4	3	5	10	8	6
C	0	0	2	1	3	8	13	11
T	0	3	1	5	4	6	11	10
A	0	1	0	3	2	7	9	8



This a "local" alignment.  
Subset of the full sequence.

**Start at max score,  
traceback to next  
highest score, and  
so on. Stop at zero**

G T T - A C  
G T T G A C

	T	G	T	T	A	C	G	G
T	0	0	0	0	0	0	0	0
G	0	0	3	1	0	0	0	3
G	0	0	3	1	0	0	0	3
T	0	3	1	6	4	2	0	1
T	0	3	1	4	9	7	5	3
G	0	1	6	4	7	6	4	8
A	0	0	4	3	5	10	8	6
C	0	0	2	1	3	8	13	11
T	0	3	1	5	4	6	11	10
A	0	1	0	3	2	7	9	8

# Aligning to a Reference Genome

- There are two major approaches:
  - Hashing the reference
  - Burrows Wheeler Transformation



# Hash-based mapping:

Step1: hash/index the genome

Toy  
genome      CATGGTCATTGGTTCC  
(16 bp)



# Hash-based mapping:

Step1: hash/index the genome

CATGGTCATTGGTTCC

k = 3

Kmer/Hash  
CAT

Genome Positions  
1



# Hash-based mapping:

Step1: hash/index the genome

CATGGTCATTGGTTCC

$k = 3$	<u>Kmer/Hash</u>	<u>Genome Positions</u>
	CAT	1
	ATG	2



# Hash-based mapping:

Step1: hash/index the genome

CAT**TGG**TCATTGGTTCC

k = 3

Kmer/Hash

CAT  
ATG  
**TGG**

Genome Positions

1  
2  
3



# Hash-based mapping:

Step1: hash/index the genome

CAT**GGT**CATTGGTTCC

k = 3

Kmer/Hash

CAT  
ATG  
TGG  
**GGT**

Genome Positions

1  
2  
3  
4



# Hash-based mapping:

Step1: hash/index the genome

CATGG**T**CATTGGTTCC

k = 3

Kmer/Hash

Genome Positions

CAT

1

ATG

2

TGG

3

GGT

4

**GTC**

5



# Hash-based mapping:

Step1: hash/index the genome

CATGG**TCA**TTGGTTCC

k = 3

<u>Kmer/Hash</u>	<u>Genome Positions</u>
CAT	1
ATG	2
TGG	3
GGT	4
GTC	5
<b>TCA</b>	6



# Hash-based mapping:

Step1: hash/index the genome

CATGGT**CAT**TGGTTCC

k = 3

<u>Kmer/Hash</u>	<u>Genome Positions</u>
CAT	1, 7
ATG	2
TGG	3
GGT	4
GTC	5
TCA	6



# Hash-based mapping:

Step1: hash/index the genome

CATGGTCATTGGTTCC

<u>Kmer/Hash</u>	<u>Genome Positions</u>
CAT	1, 7
ATG	2
TGG	3, 10
GGT	4, 11
GTC	5
TCA	6
ATT	8
TTG	9
GTT	12
TTC	13
TCC	14

*Complete hash/kmer index of our toy genome (forward strand only)*



# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads

Toy genome	CATGGTCATTGGTTCC	Kmer/Hash	Genome Positions
		CAT	1,7
		ATG	2
		TGG	3,10
		GGT	4,11
		GTC	5
		TCA	6
		ATT	8
		TTG	9
	→ Read	TGGTCA	
		GTT	12
		TTC	13
		TCC	14

*kmer index is used to quickly find candidate alignment locations in genome.*



# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads

Toy genome	CATGGTCATTGGTTCC	Kmer/Hash	Genome Positions
		CAT	1, 7
		ATG	2
		TGG	3, 10
		GGT	4, 11
		GTC	5
		TCA	6
		ATT	8
		TTG	9
		GTT	12
		TTC	13
		TCC	14

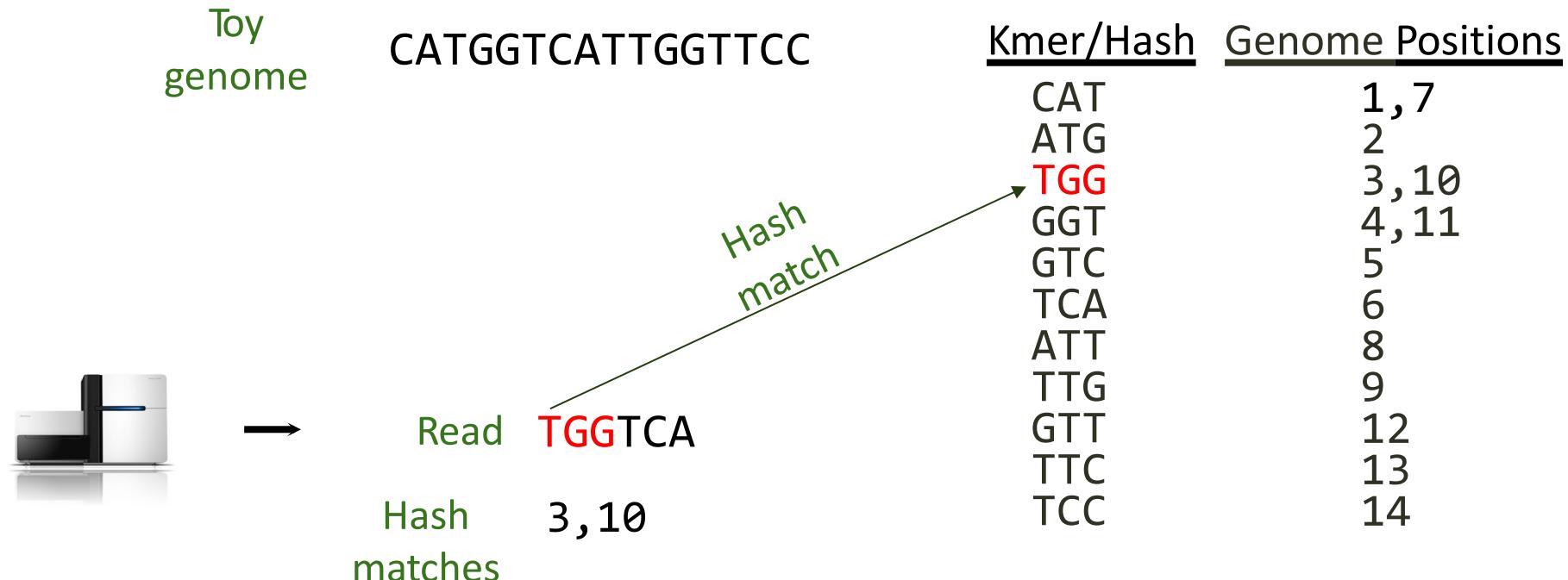


Read **TGG**TCA



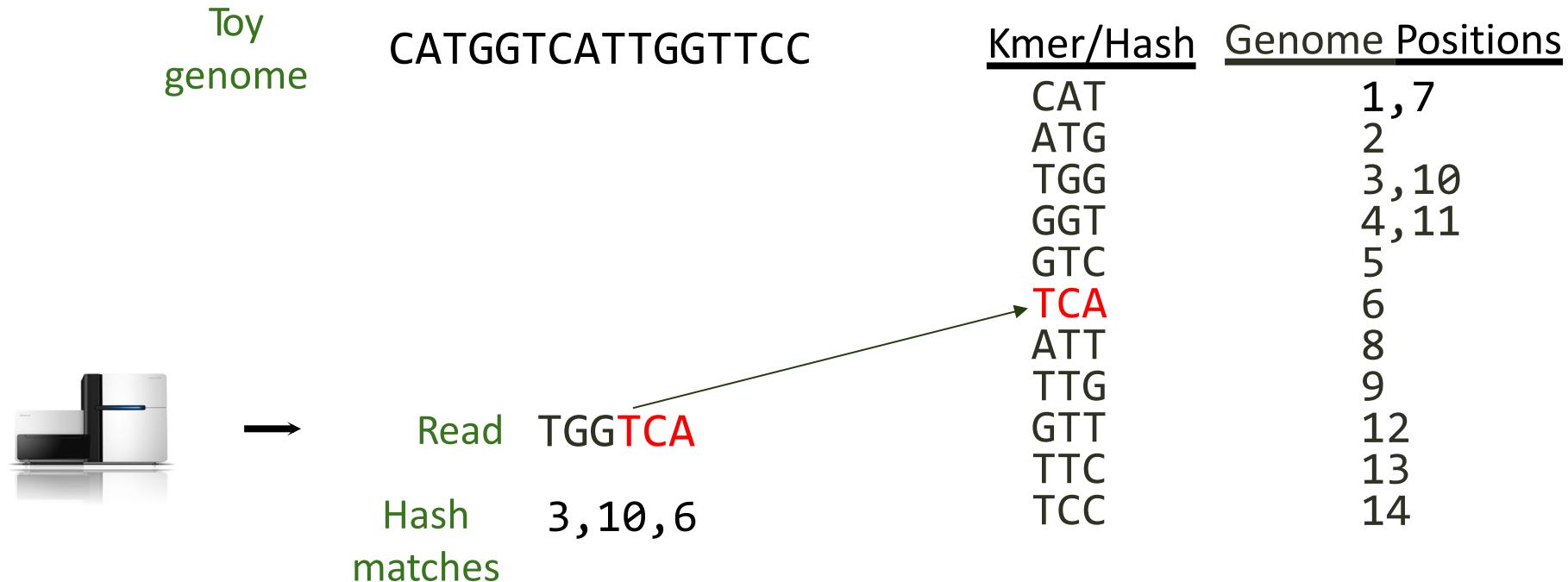
# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads



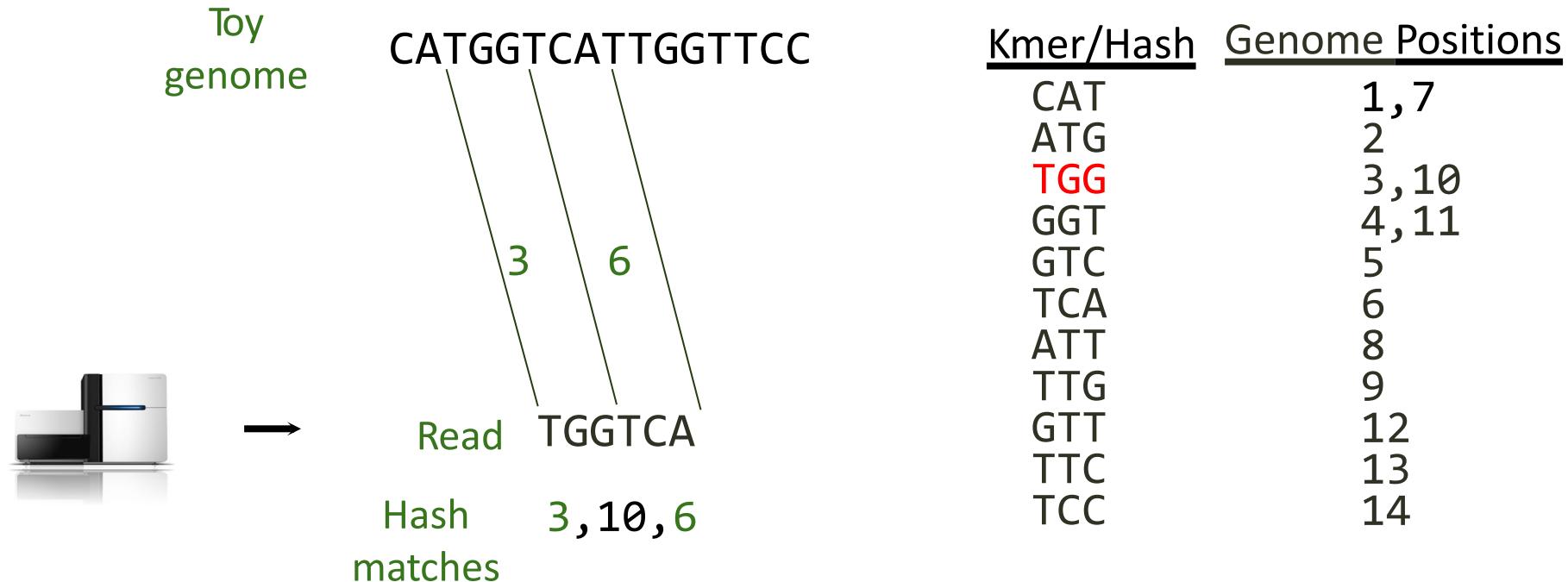
# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads



# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads



Okay, that was a bit easy because the read and the reference exactly matched. What about if there is a sequencing error or a genetic variant in the read?



# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads

Toy genome	CATGGTCATTGGTTCC	Kmer/Hash	Genome Positions
		CAT	1, 7
		ATG	2
		TGG	3, 10
		GGT	4, 11
		GTC	5
		TCA	6
		ATT	8
		TTG	9
		GTT	12
		TTC	13
		TCC	14



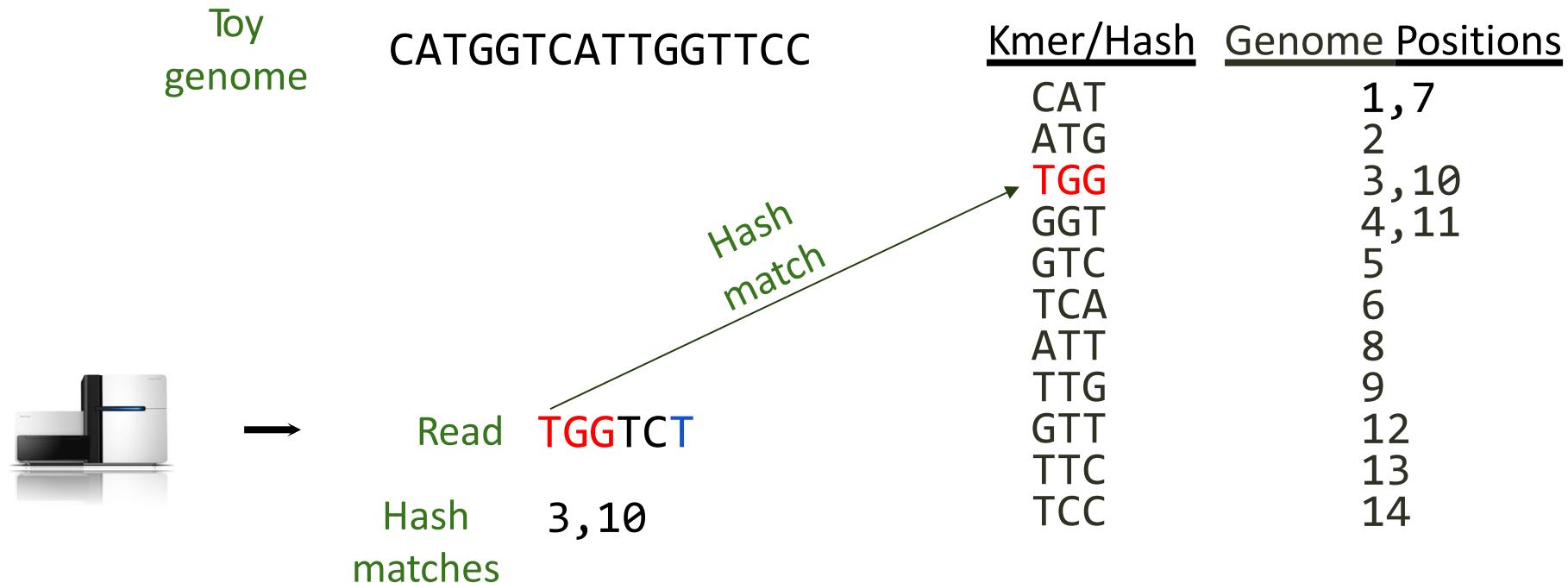
Read TGGTCT

*kmer index is used to quickly find candidate alignment locations in genome.*



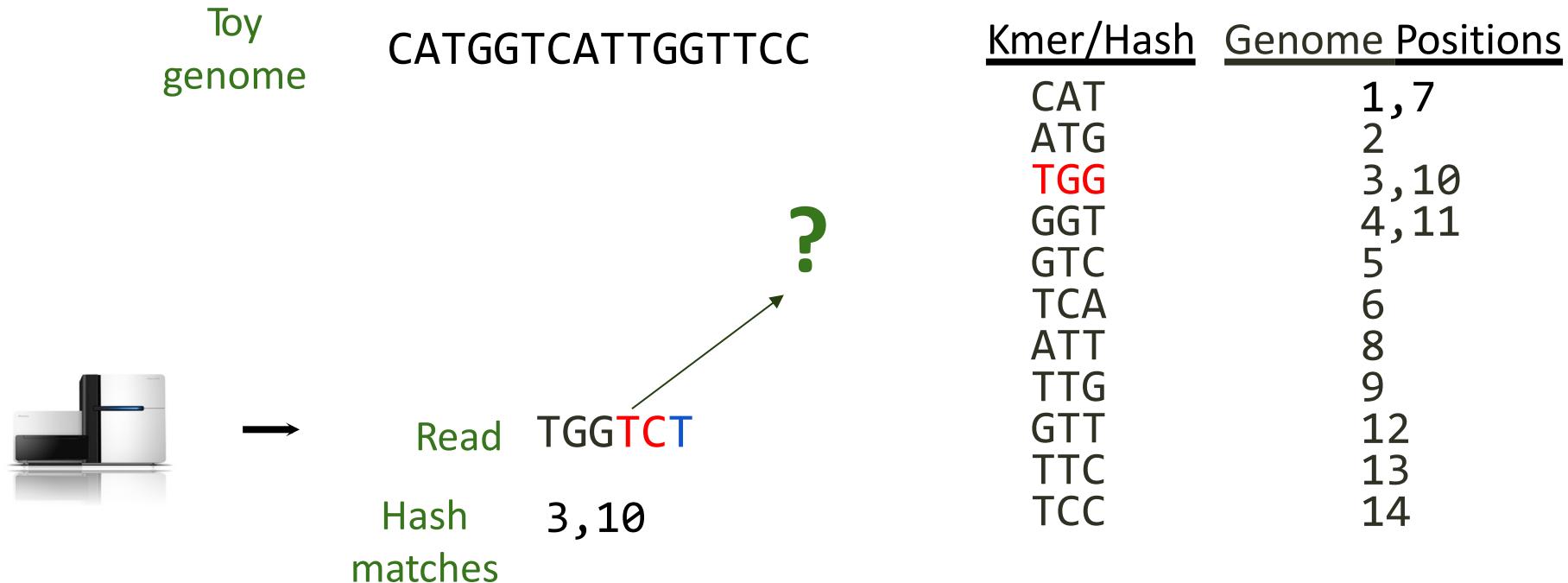
# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads



# Hash-based mapping:

Step2: use the index to map (i.e., find alignment locations) reads



**Thought experiment:** what is a good choice of hash size ( $k$  for  $k$ -mers) for building a hash table to facilitate sequence mapping to the human genome?



# k=1?



# k=3?

( $4^3$  possibilities)

AAA, AAC, AAG, ... ,

TTT



# k=10?

$4^{10}$  (1,048,576)

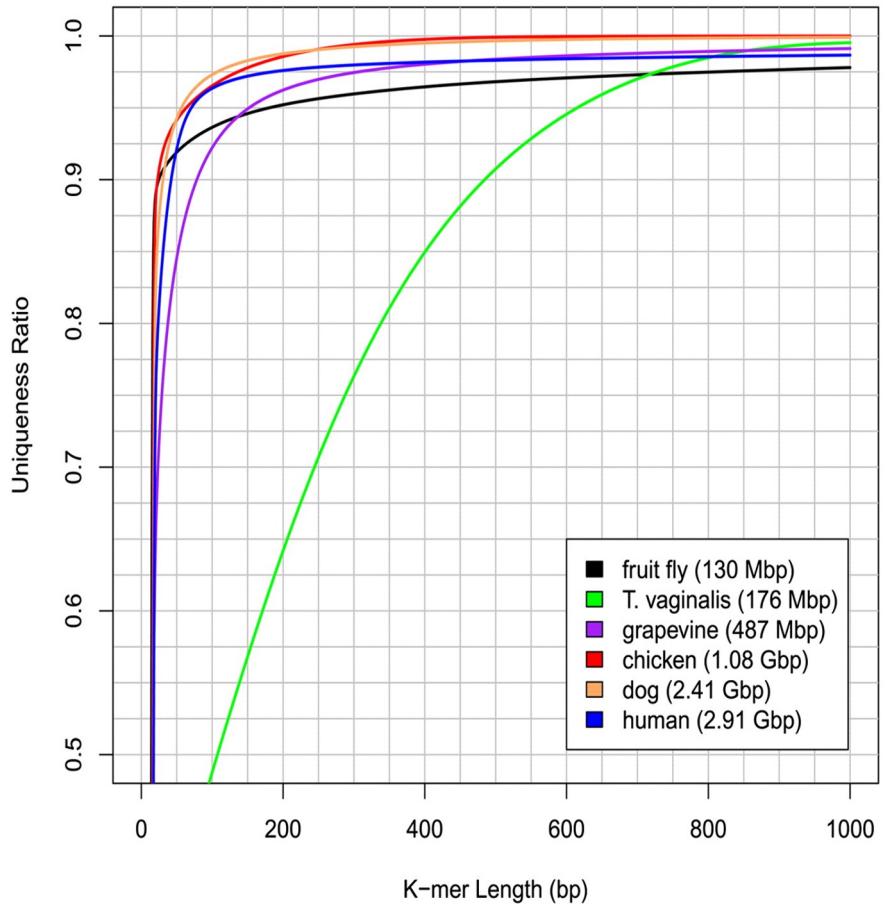
Every one of these is present in  
the human genome at least once

<http://bmcbioinformatics.biomedcentral.com/articles/10.1186/14>

71-2105-9-167

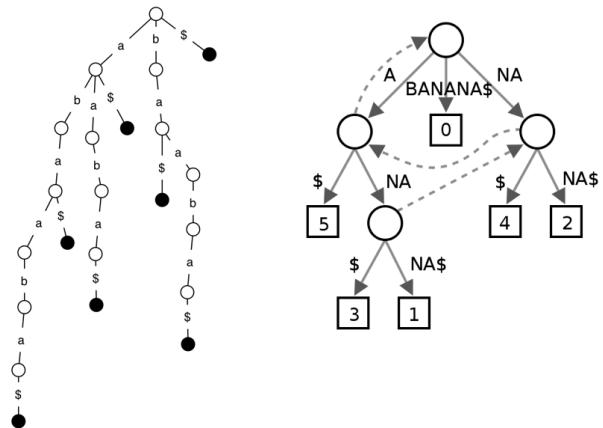


- What fraction of the k-mers occurs once and only once?
  - Proxy for how unique is the genome
- It takes a pretty long k-mer to be unique in most genomes
  - This example shows that you need K>50 for humans



# Burrows-Wheeler Transform

The latest fad. More involved computationally, but requires much less memory.



Bowtie1 / Bowtie2  
SOAP2  
BWA

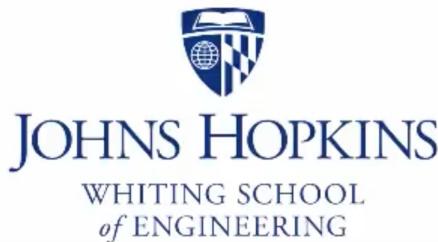
6	\$
5	A\$
3	ANAS\$
1	ANANA\$
0	BANANA\$
4	NA\$ BANA
2	NANA\$ BA

- Algorithm that underlies multiple sequence alignment strategies
  - It was originally designed for as a data compression strategy
- The BWT transforms an input string by permuting its characters in such a way that the inverse transform can recover the original string without using any information aside from the permuted string



# Burrows-Wheeler Transform and FM Index

Ben Langmead



You are free to use these slides. If you do, please sign the guestbook ([www.langmead-lab.org/teaching-materials](http://www.langmead-lab.org/teaching-materials)), or email me ([ben.langmead@gmail.com](mailto:ben.langmead@gmail.com)) and tell me briefly how you're using them. For original Keynote files, email me.

Ben Langmead

<https://genomebiology.biomedcentral.com/articles/10.1186/gb-2009-10-3-r25>

<https://www.youtube.com/watch?v=4n7NPk5lwbl>



# Pros and Cons of BWT

- Vast majority of sample sequence can be accurately placed
- Problems with:
  - Large scale differences - structural variation
  - Reference bias
  - Repetitive DNA



# Aligner      Strategy      Availability

---

<b>BFAST</b>	hash-based (genome)	open-source
<b>Bowtie2*</b>	Suffix array / BWT	open-source
<b>BWA*</b>	Suffix array / BWT	open-source
<b>MAQ</b>	hash-based (reads)	open-source
<b>Mosaik*</b>	hash-based (genome)	open-source
<b>Novoalign*</b>	hash-based (genome)	free for academic use
<b>RMAP</b>	hash-based (reads)	open-source
<b>SOAP2</b>	Suffix array / BWT	free for academic use
<b>Eland</b>	hash-based (reads)	commercial (Illumina)



# Unresolved problems

- Mapping tries to match the reference, so inherently introduces a bias towards the reference
- We have to modify parameters based on the read content (e.g. deletions)
- Mapping to repetitive DNA is still problematic
- What if there is no or an incomplete reference for the sequenced organism?

