



**Mansoura University**  
**Faculty of Computers and Information**  
**Department of Computer Science**  
**Second Semester: 2020-2021**



# **[MED-145] Genomics: Genome Indexing & Reads Mapping**

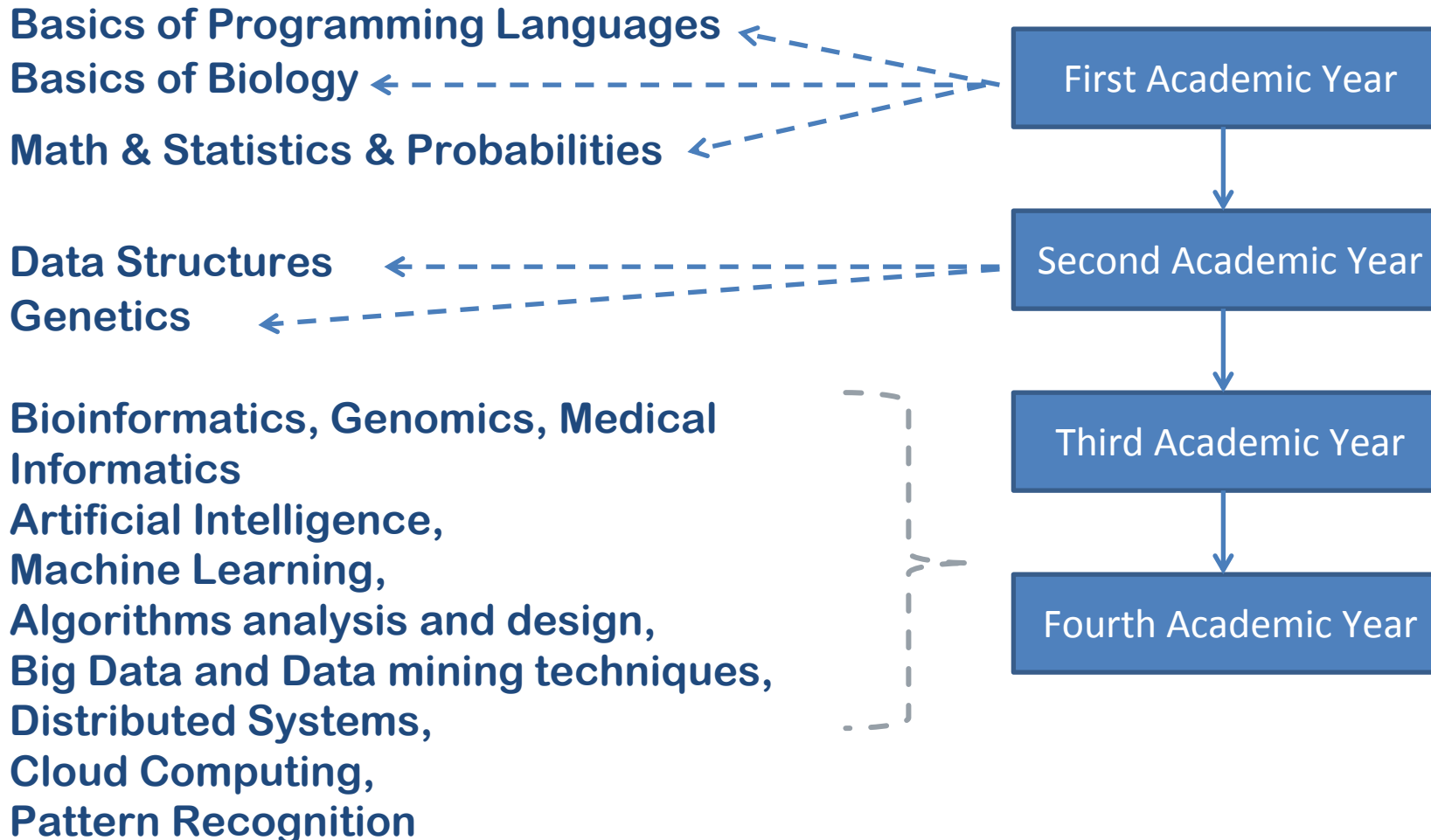
## **Grade: Third Year (Medical Informatics Program)**

**Sara El-Metwally, Ph.D.**  
**Faculty of Computers and Information,**  
**Mansoura University,**  
**Egypt.**

# COURSE OUTLINES

- **Course Meeting Time:** Monday 8:30-10:100, Wednesday 12:20-2.00.
- **Course Instructor:** Sara El-Metwally, PhD
- **Course TAs:** Eng. Nada El-Madah, Eng. Ola Magdy
- **Course Labs:** Unix-based commands, Shell Scripting, Python.
- **Course Grading:**
  - Midterm: 10%
  - Oral: 10%
  - Practical: 10%
  - Project / Paper: 10%
  - Final: 60%

# MANSOURA FCIS COURSE DEPENDENCY



# COURSE PROJECTS

## ❖ Choice No. 1

- Teams (up to 5 students).
- Pick a project from the projects list or propose an idea for your project.
- Projects outcomes:
  - A Project Proposal that describes the problem, how to solve it, the technologies/tools that you will use, team members and their tasks, etc.
  - A Github page that includes a Readme file that describes your project idea, algorithm, how to run and use the code and any useful links etc. and your project source code with any dependency.
  - Your proposal should be added to your Github page.
  - A Video demo that describes your project (English), the link should be added to your Github page and the video should be uploaded to our course channel on YouTube!
  - Group Photo with a faculty logo.
  - There is a competition among different genomics projects; the top best five projects will awarded a genomics course certificate plus some other prizes in the case of extraordinary projects.
  - Think big drive forward!

# COURSE PROJECTS

## Genomics Course List Projects 2021

Project Name	Example of already existing tools in the field
Genome/Transcriptome Browsers	UCSC Genome Browser
Genome/Transcriptome Assemblers	Velvet, Canu, LightAssembler, Trinity, SPAdes
Multiple Sequence Aligners (DNA, RNA, Amino Acids)	Clustal Omega, MAFFT, MUSCLE
Local alignment tool	BLAST
Short/long reads aligner to a reference genome	Bowtie, BWA
Phylogenetic Trees Drawing/Analysis Tool	iTOL, phylot, PhyML
Variant Calling Programs	GATK, SAMtools, FreeBayes, DeepVariant
Errors Correction Programs	Bless, Musket, Lighter, Fiona, NanoReviser, MARVEL
Efficient kmers counting tools	Jellyfish, KMC, DSK
Quality control of sequencing reads (short/long)	FastQC, FASTX-Toolkit, LongQC
Assembly Evaluation Software	QUAST
Fastq compression tool	MZPAQ, fqzcomp, Spring
<b>Your Idea?</b>	<b>Your Program</b>

# COURSE PROJECTS

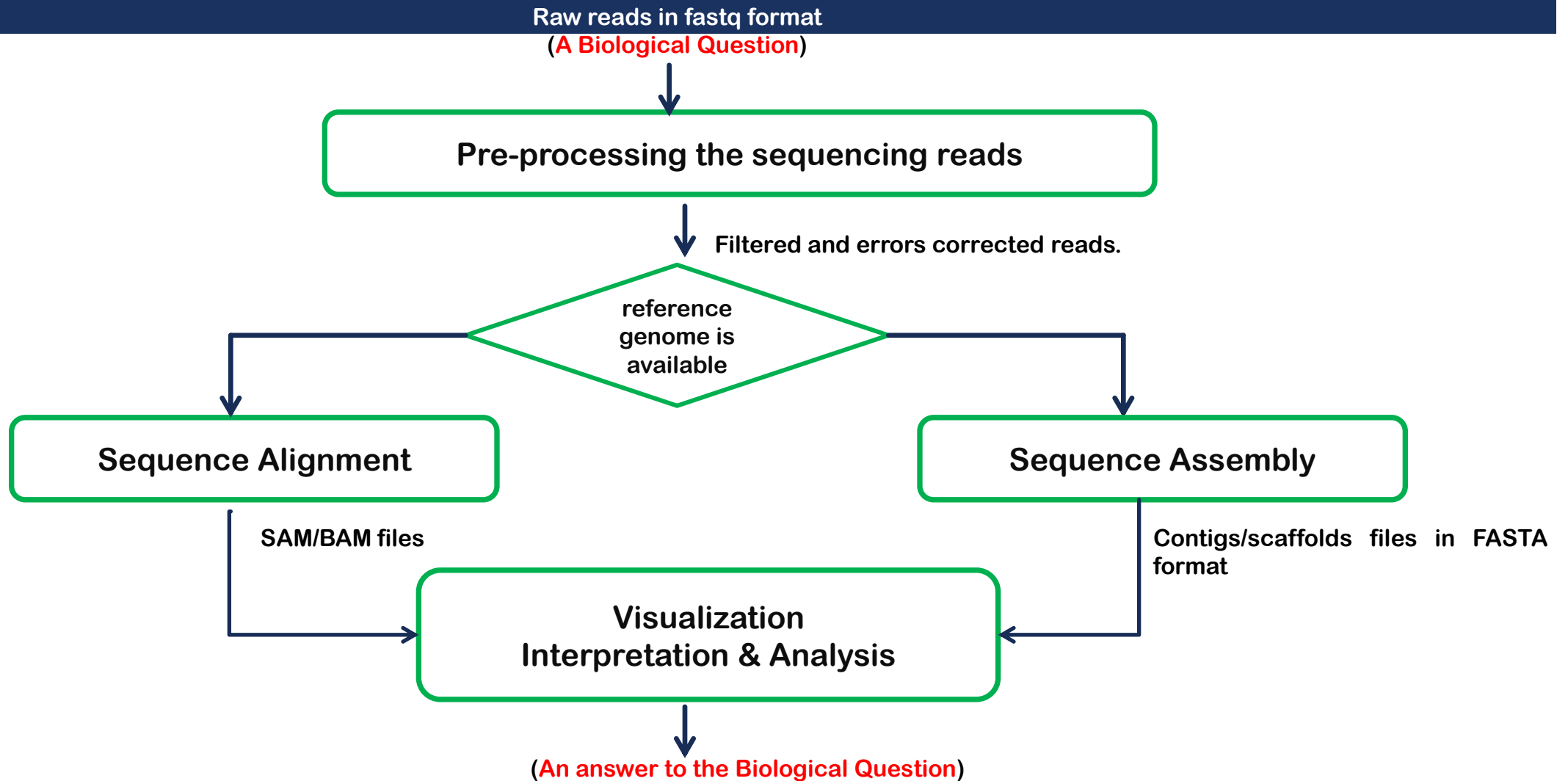
## ■ Choice No. 2

- ❑ One student.
- ❑ Pick a paper published in 2020/2021 in the genomics field in general or related to SARS-CoV-2 data analysis.
- ❑ Student outcomes:
  - ❑ A document that summarizes the paper, the analysis pipeline, the findings and your comment on the paper results (Max. 3 pages).
  - ❑ A Video that explains the paper idea using the presentation prepared by you.
  - ❑ Student Photo with a faculty logo.
  - ❑ A Github page that includes a Readme file that describes the paper idea, data analysis, etc. including the paper reference and your created video along with your photo.

# AGENDA

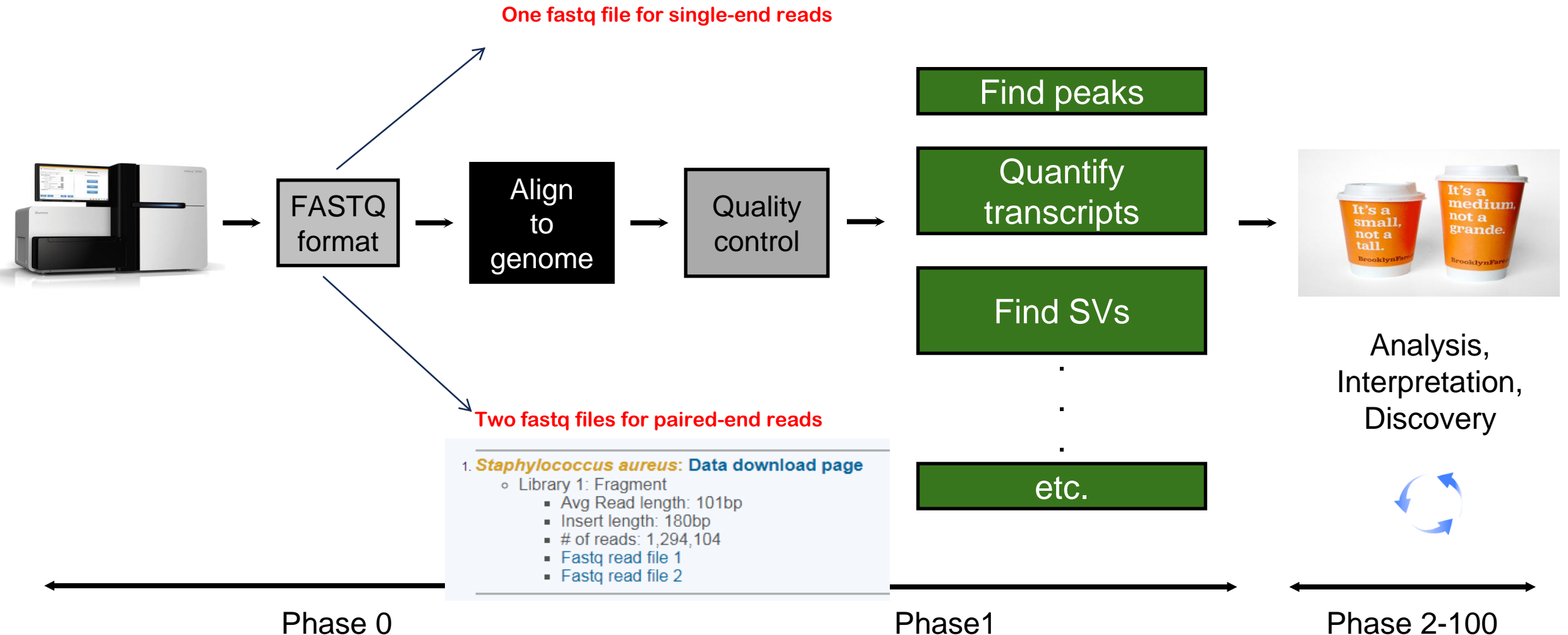
- A Typical Genomics Analysis Workflow.
- Alignment is central to most genomic research.
- Sequence mapping versus alignment.
- Reference based analysis mapping and challenges.
- Mapping Quality.
- Sequence alignment/Mapping software.
- Typical Mapping/Alignment Workflow.
- Hash-based mapping approach.
- Hash-based mapping approach drawbacks.

# A TYPICAL GENOMICS DATA ANALYSIS PIPELINE

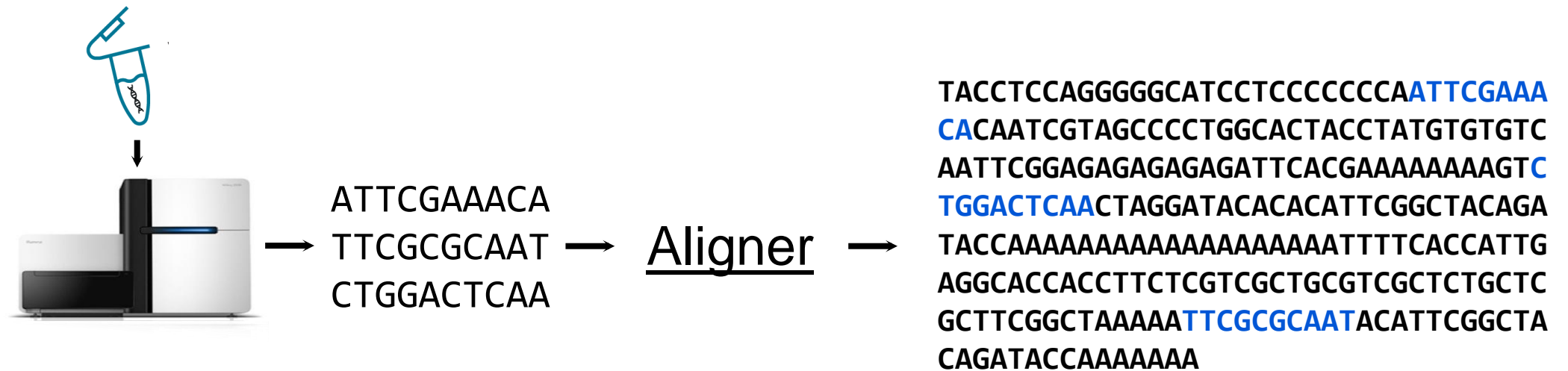




# ALIGNMENT IS CENTRAL TO MOST GENOMIC RESEARCH

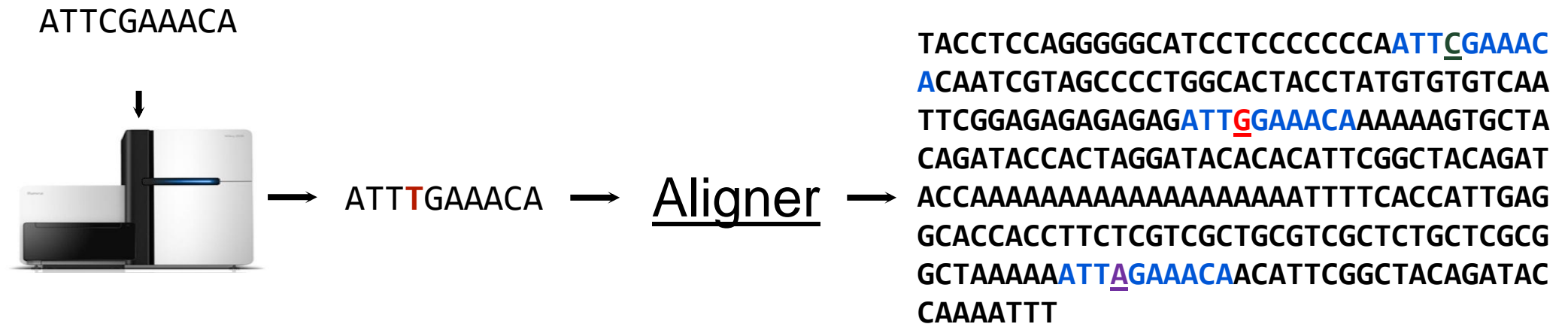


# BEST CASE SCENARIO: AN ERROR-FREE SEQUENCING TECHNOLOGY



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

# REALITY CHECK. ERRORS HAPPEN. FREQUENTLY.



“Fuzzy” matching is much more computationally expensive.

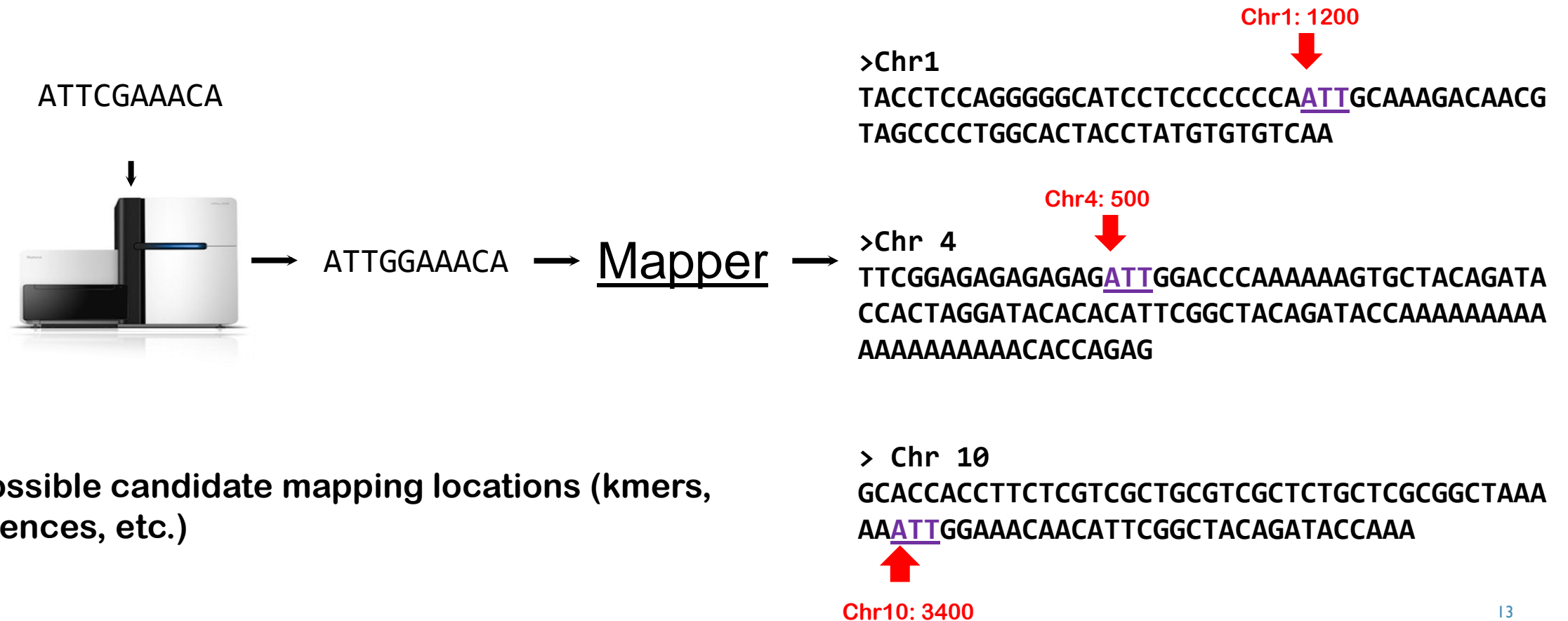
Think Google’s “Did you mean...”

# SEQUENCE MAPPING VERSUS ALIGNMENT

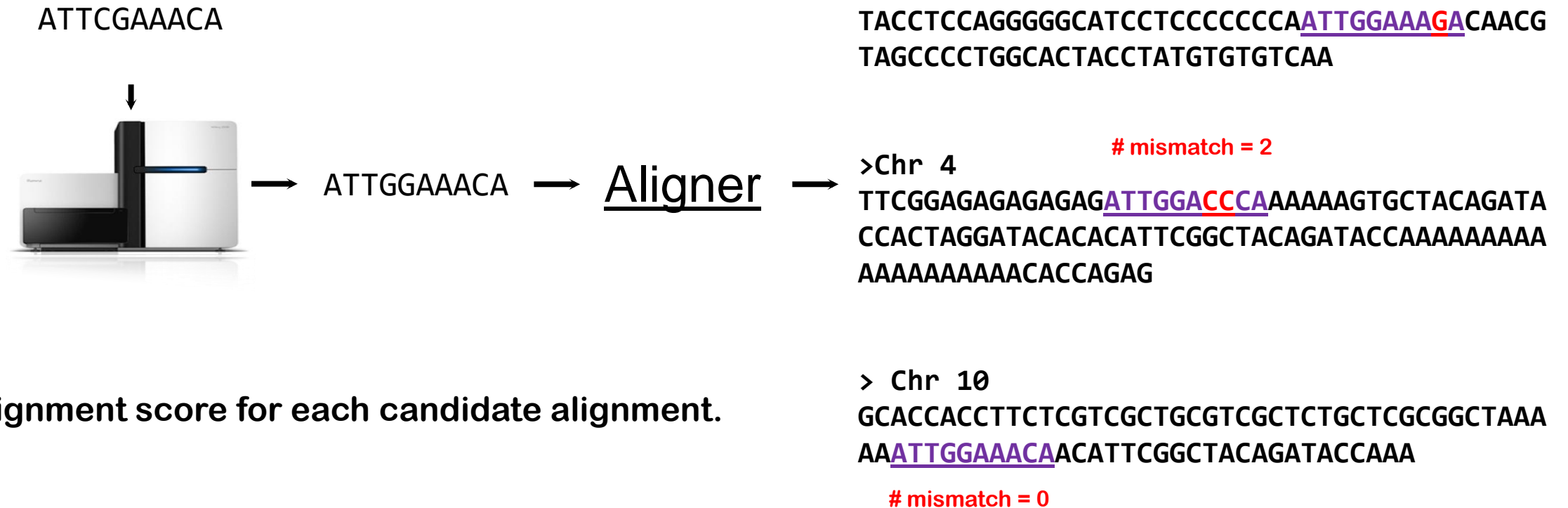
**Mapping:** (quickly) find the best possible loci to which a sequence could be aligned.

**Alignment:** for each locus to which a sequence can be mapped, determine the optimal base by base alignment of the query sequence to the reference sequence.

# SEQUENCE MAPPING VERSUS ALIGNMENT

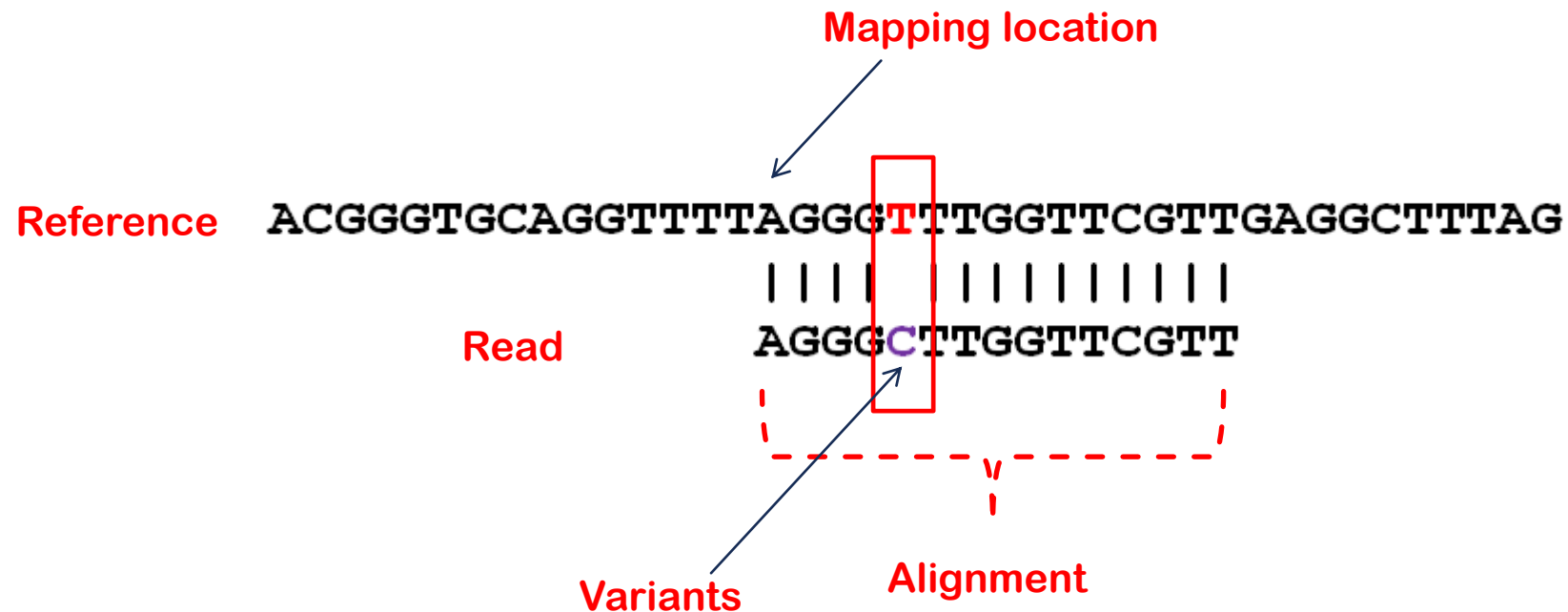


# SEQUENCE MAPPING VERSUS ALIGNMENT



- find alignment score for each candidate alignment.

# REFERENCE-BASED ANALYSIS



- ✓ Mapping for long reads, aligning for short reads, or used interchangeably.
- ✓ Discover genetic variations by comparing reads to a reference genome.
- ✓ To do this, the best mapping positions between reads and the reference should be identified (**Some Challenges will be here!**).

# REFERENCE MAPPING CHALLENGES

- ❖ Genomes are very large (3 billion bases in human) and have repetitive regions.
- ❖ Naïve algorithms would take too much time and memory to map reads to a reference

Reference



Text



Read

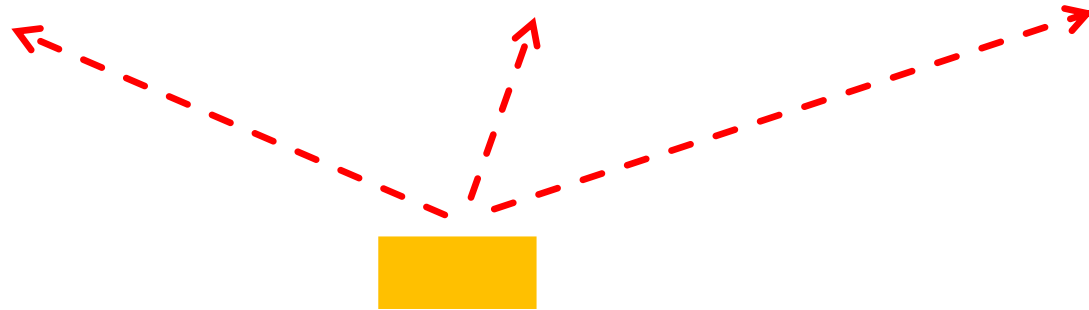
Pattern



# REFERENCE MAPPING CHALLENGES

- ❖ Genomes are very large (3 billion bases in human) and have repetitive regions.

Reference



Read

# REFERENCE MAPPING CHALLENGES

- ❖ Reads have sequencing errors (substitutions, insertions, and deletions).

**Reference**    ACGGGTGCAGGTTT TAGGG **T**TTGGTTCGTTGAGGCTTTAG  
                              | | | |    | | | | | | | |  
**Read**            AGGG **C**TTGGTTCGTT

**Reference**    ACGGGTGCAGGTTT TAGGG **G**TTTGG --- TTGAGGCT  
                              | | |    | | |        | |  
**Read**            AGG --- TTGG **TT**CGTT

- ✓ Mappers must be able to find inexact alignments by tolerating differences.
- ✓ Substitutions are more tolerance than Indels.

# REFERENCE MAPPING CHALLENGES

❖ Reads could be mapped to many locations across the genome, which one will be reported?

Chr10:1020

AGGGACCGGTTTCGTTTAGGGTTTGGTTCGTTGAGGCTTTAG

|||| |||||

mismatches : 1

AGGGATTGGTTCGTT

Start from the base/offset/position no. 2139 in Chr2

Chr2

Chr2:2139

AGGGACCGGTTTCGTTTAGGGTTTGGTTCGTTGAGGCTTTAG

|||| |

AGGGATTGGTTCGTT

mismatches : 2

# REFERENCE MAPPING CHALLENGES

- ❖ Reads could be mapped to many locations across the genome, which one will be reported?
- ❖ **MQ is an estimation of the probability that a mapping is incorrect (it encodes many factors such as the number of mismatches, type of mismatch, quality scores, etc.)**

Chr10:1020

AGGGACCGGTTTCGTTTAGGGTTTGGTTCGTTGAGGCTTTAG

|||| |||||

Mapping Quality : 10

AGGGATTGGTTCGTT

Chr2:2139

AGGGACCGGTTTCGTTTAGGGTTTGGTTCGTTGAGGCTTTAG

|||| |

AGGGATTGGTTCGTT

Mapping Quality : 1

# MAPPING QUALITY (MAPQ)

- What is the probability that the sequence should be mapped here and only here?
- MAPQ also uses the Phred (log) scale:

$$\text{MAPQ} = -10 \cdot \log_{10}(P_{\text{map\_loc\_wrong}})$$

$(P_{\text{map\_loc\_wrong}})$	$\log_{10}(P_{\text{map\_loc\_wrong}})$	MAPQ
1	0	0
0.1	-1	10
0.01	-2	20
0.001	-3	30
0.0001	-4	40

# REFERENCE MAPPING CHALLENGES

- ❖ Reads could be mapped to many locations across the genome, which one will be reported?
- ❖ Low quality mismatches are less important than high quality mismatches.

Chr10:1020

AGGGACCGGTTTCGTTTAGGGTTTGGTTCGTTGAGGCTTTAG

||||| |||||

AGGGATTGGTTCGTT

Q=15



Chr2:2139

AGGGACCGGTTTCGTTTAGGGTTTGGTTCGTTGAGGCTTTAG

||||| |||||

AGGGATCGGTTTCGTT

Q=40



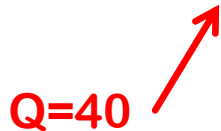
# REFERENCE MAPPING CHALLENGES

- ❖ Reads could be mapped to many locations across the genome, which one will be reported?
- ❖ When two mappings have the same exact alignment, the mapping is ambiguous. Two positions are equal? Which one is correct?

Chr10:1020

```
AGGGTTTGGTTCGTTT TAGGGTTTGGTTCGTTGAGGCTTTAG
||||| |||||
AGGGA TTGGTTCGTT
```

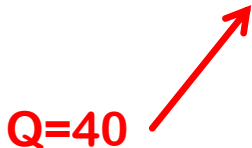
Q=40



Chr2:2139

```
AGGGA TCGGTTCGTTT TAGGGTTTGGTTCGTTGAGGCTTTAG
||||| |||||
AGGGA ACGGTTCGTT
```

Q=40



# REFERENCE MAPPING CHALLENGES


❖ Reads could be mapped to many locations across the genome, which one will be reported?

❖ Choose one mapping position at random, MQ=0

Chr10:1020

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


Q=40



Chr2:2139

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

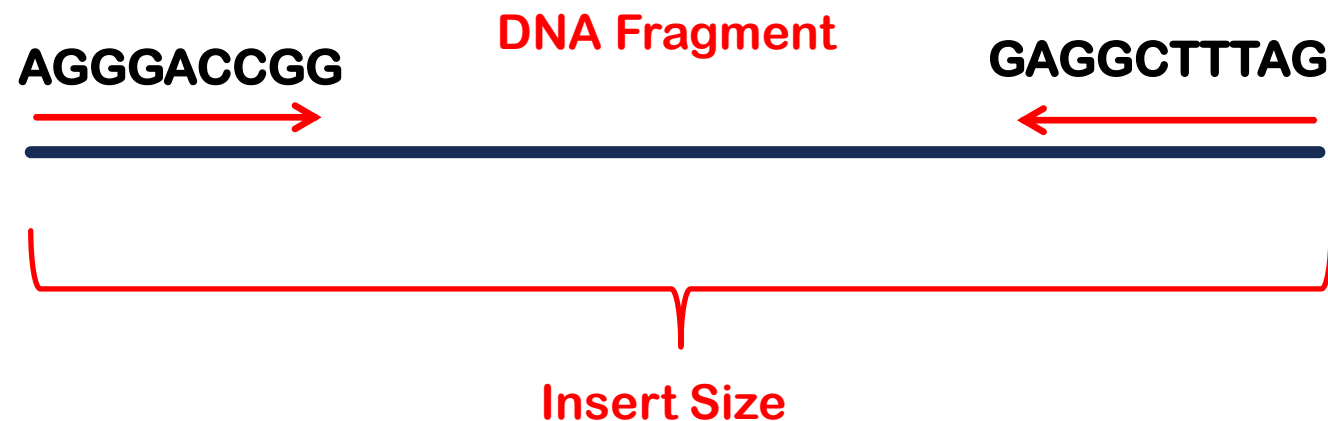
Q=40





# REFERENCE MAPPING CHALLENGES

- ❖ Reads could be mapped to many locations across the genome, which one will be reported?
- ❖ Paired-end reads can help in resolving ambiguous mappings.



# REFERENCE MAPPING CHALLENGES

❖ Reads could be mapped to many locations across the genome, which one will be reported?

Chr10:1020

```
AGGGTTGGTTCGTTT TAGGGTTGGTTCGTTGAGGCTTTAG
|||| | |||||
AGGGA TTGGTTCGTT
```

Chr2:2139

```
AGGGA TCGGTTCGTTT TAGGGTTGGTTCGTTGAGGCTTTAG
|||| | |||||
AGGGA ACGGTTCGTT
```

# REFERENCE MAPPING CHALLENGES

❖ Reads could be mapped to many locations across the genome, which one will be reported?

Chr10:1020

```
AGGGTTTGGTTCGTTT TAGGGTTTGGTTCGTTGAGGCTTTAG
|||| | ||||| ||||| |||||
AGGGAATTGGTTCGTT ——— GGTTCGTTGAGCTT
```

Chr2:2139

```
AGGGAACGGTTCGTTT TAGGGTTTGGTTCGTGAGGCTTTAG
|||| | ||||| ||||| || | | | |
AGGGAACGGTTCGTT ——— TCTTCGTTGATGCAG
```

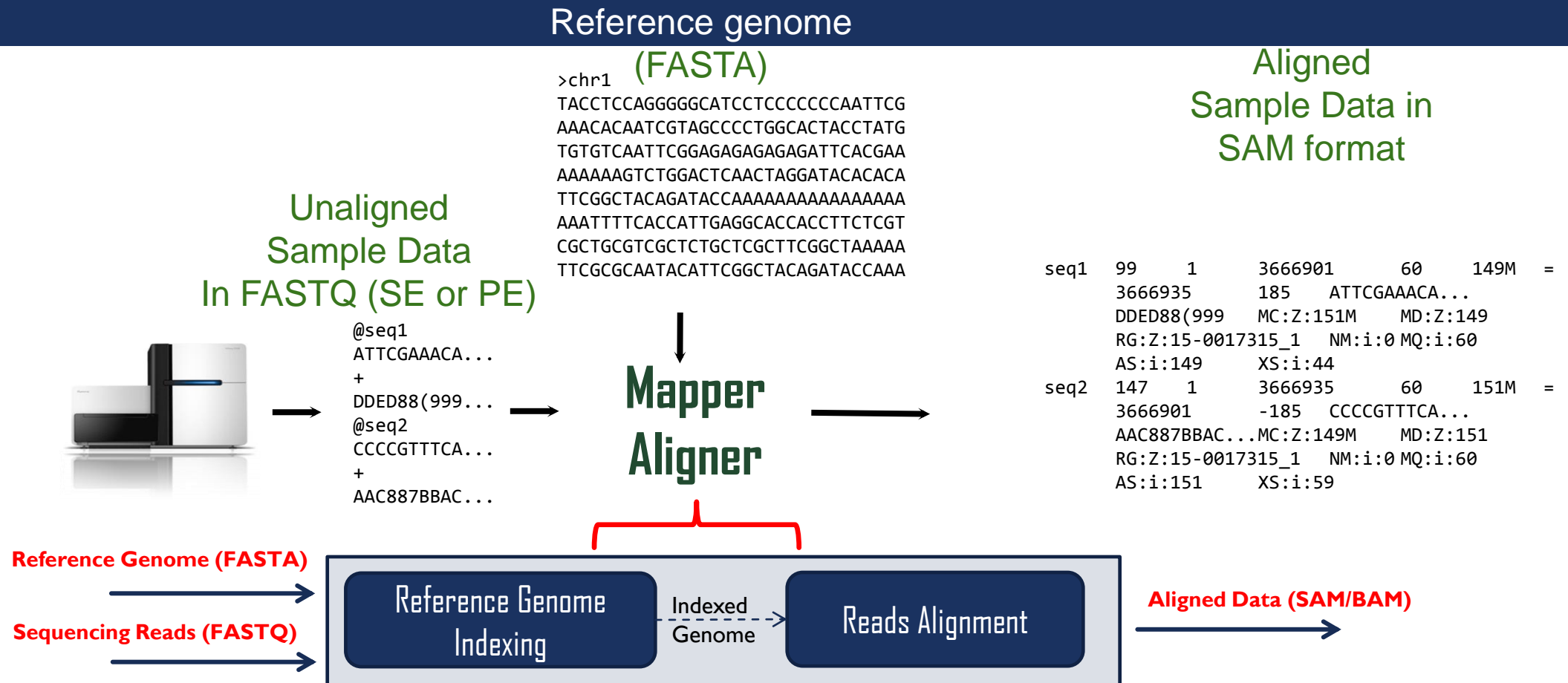
# SEQUENCE ALIGNMENT/MAPPING SOFTWARE

<u>Aligner</u>	<u>Approach</u>	<u>Applications</u>	<u>Availability</u>
BWA-mem	Burrows-Wheeler	DNA, SE, PE, SV	open-source
Bowtie2	Burrows-Wheeler	DNA, SE, PE, SV	open-source
Novoalign	hash-based	DNA, SE, PE	free for academic use
TopHat	Burrows-Wheeler	RNA-seq	open-source
STAR	hash-based (reads)	RNA-seq	open-source
GSNAP	hash-based (reads)	RNA-seq	open-source

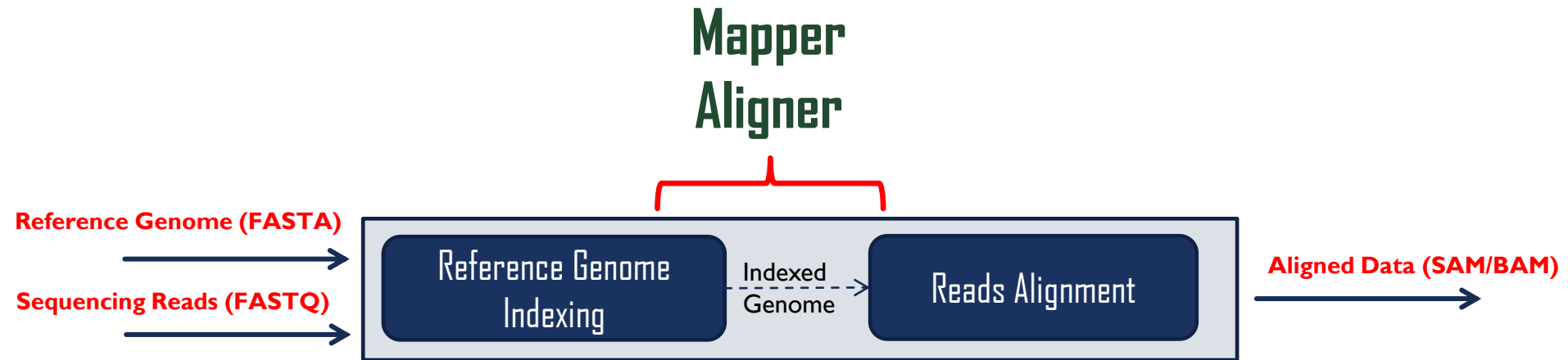
Slides curiosity from Aaron Quinlan: <https://github.com/quinlan-lab/applied-computational-genomics>

<https://academic.oup.com/bioinformatics/article/28/24/3169/245777>

# TYPICAL MAPPING/ALIGNMENT WORKFLOW



# TYPICAL MAPPING/ALIGNMENT WORKFLOW



## Genome Indexing and Mapping Approaches

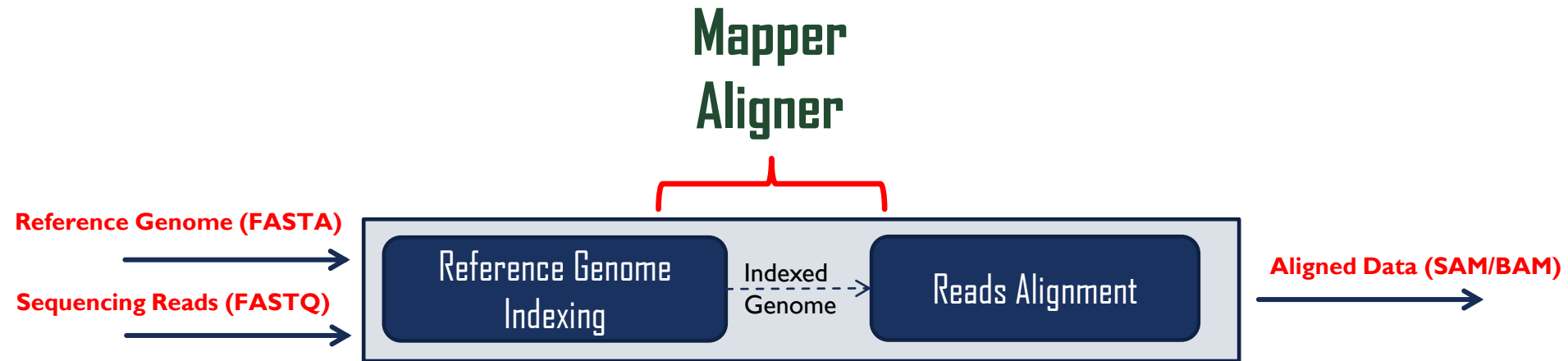


hash-based



Burrows-Wheeler

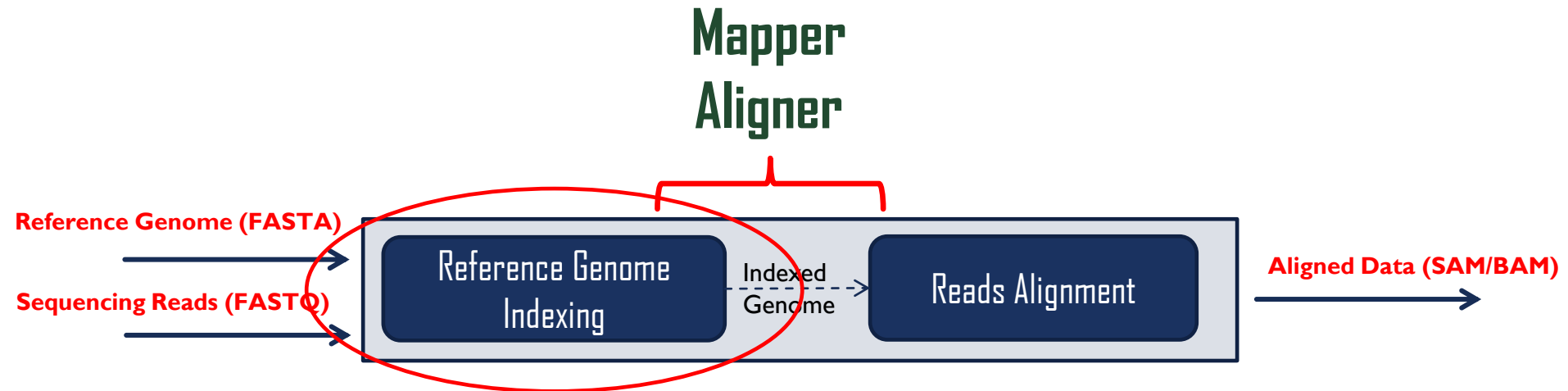
# TYPICAL MAPPING/ALIGNMENT WORKFLOW



## Genome Indexing and Mapping Approaches



# TYPICAL MAPPING/ALIGNMENT WORKFLOW



## Genome Indexing and Mapping Approaches





# HASH-BASED MAPPING:

## Step1: hash/index the genome

Toy  
genome  
(16 bp)

CATGGTCATTGGTTCC

# HASH-BASED MAPPING:

Step1: hash/index the genome

**CAT**GGTCATTGGTTCC

k = 3

Kmer/Hash

**CAT**

Genome Positions

**1**

Could be zero-based  
indexing

# 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

CATGGTCATTGGTTCC

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

# HASH-BASED MAPPING:

Step1: hash/index the genome

CATGGTCATTGGTTCC

k = 3	<u>Kmer/Hash</u>	<u>Genome Positions</u>
	CAT	1
	ATG	2
	TGG	3
	GGT	4

# HASH-BASED MAPPING:

Step1: hash/index the genome

CATG**GTC**ATTGGTTCC

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

# 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>	<b>6</b>

# HASH-BASED MAPPING:

Step1: hash/index the genome

CATGGT**CATT**GGTTCC

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



# HASH-BASED MAPPING:

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	A	T	G	G	T	C	A	T	T	G	G	T	T	C	C

Step1: hash/index the genome

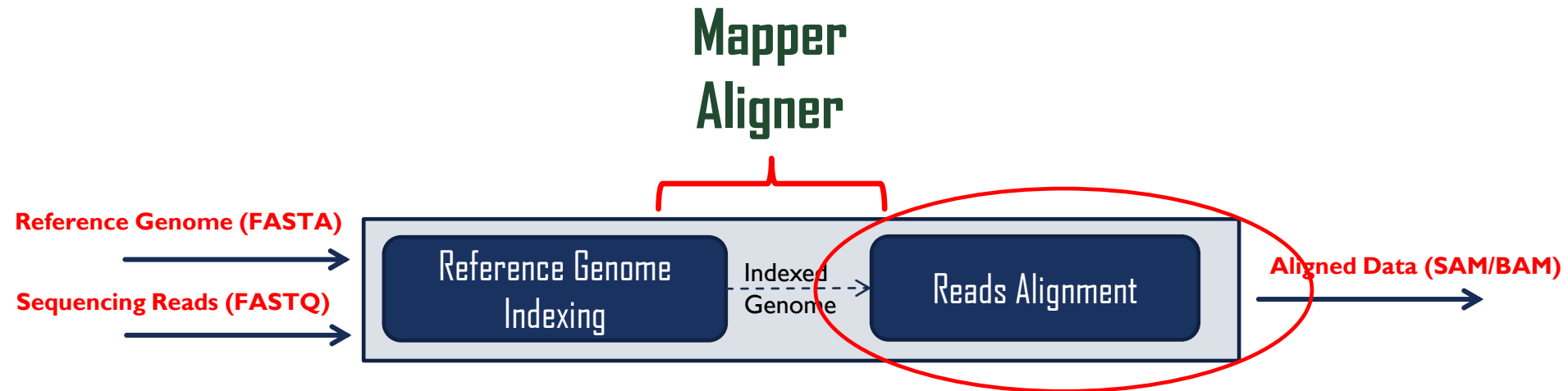
CATGGTCATTGGTTCC

*Complete hash/kmer index of our toy genome (forward strand only), k=3*

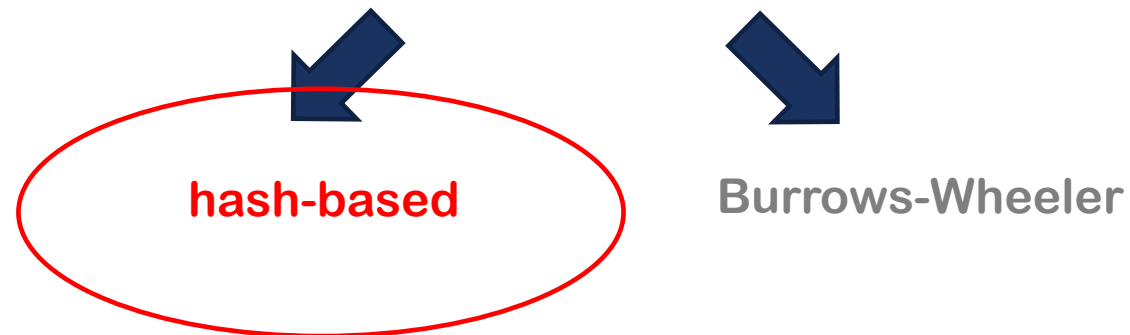
Kmer/hash	Genomic position
CAT	1,7
ATG	2
TGG	3,10
GGT	4,11
GTC	5
TCA	6
ATT	8
TTG	9
GTT	12
TTC	13
TCC	14

Genome Index

# TYPICAL MAPPING/ALIGNMENT WORKFLOW



## Genome Indexing and Mapping Approaches



# HASH-BASED MAPPING:

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

Genome Index



Read

TGGTCA

Kmer/hash	Genomic position
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

# HASH-BASED MAPPING:

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	A	T	G	G	T	C	A	T	T	G	G	T	T	C	C

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

Genome Index



Read **TGGTCA**

Hash match

TGG	3,10
-----	------

Kmer/hash	Genomic position
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

# HASH-BASED MAPPING:

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	A	T	G	G	T	C	A	T	T	G	G	T	T	C	C

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

Genome Index



Read **TGGTCA**

TGG	3,10
TCA	6

Hash match

Kmer/hash	Genomic position
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

# HASH-BASED MAPPING:

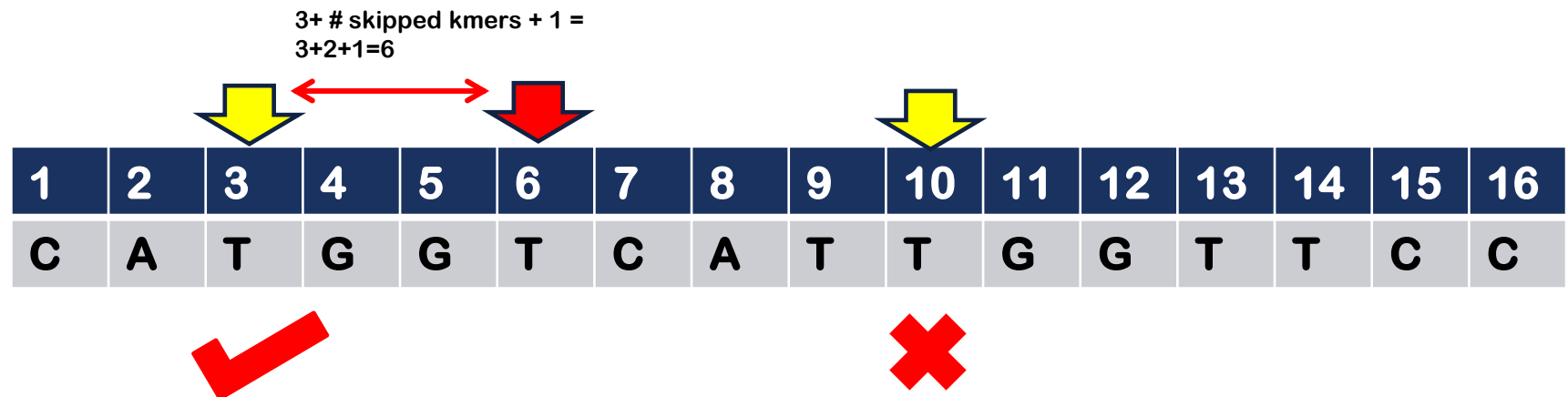
Step2: use the index to map (i.e., find alignment locations) reads.



→ Read TGGTCA

TGG	3,10
-----	------

TCA	6
-----	---



# HASH-BASED MAPPING:

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	A	T	G	G	T	C	A	T	T	G	G	T	T	C	C

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

Genome Index



Read

TGGTCT

**{3,10}** ?

T is an erroneous bases ?

Alignments will be tolerated to mismatches

Kmer/hash	Genomic position
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

# EDIT DISTANCE (LEVENSHTEIN DISTANCE)

How many edits (changes) must be made to a word or kmer to make it match (align) to another word or kmer?

What is the difference between Edit distance & Hamming distance?

CURLED → Edit distance = 1. Substitute C for H  
HURLED

SHORT → Edit distance = 1. Delete R  
SHO-T

TGTTACGG  
GGTTGACTA ?

TG-TT-ACGG  
-GGTTGACTA

Edit distance = 5

TGTT-ACGG  
GGTTGACTA

Edit distance = 4



# HASH-BASED MAPPING:

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	A	T	G	G	T	C	A	T	T	G	G	T	T	C	C

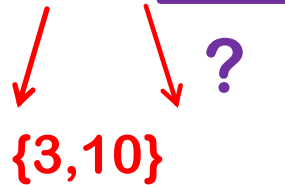
Step2: use the index to map (i.e., find alignment locations) reads

Genome Index



Read

**TGGTC****T**



Or :T is a True variation on the read and should be reported and studied.

Kmer/hash	Genomic position
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

## Note?

**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?

## Note?

**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?

Note?

**k=1?**

Note?

**k=3?** ( $4^3$  possibilities)  
AAA, AAC, AAG, ... , TTT

# Note?

**k=10?**

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

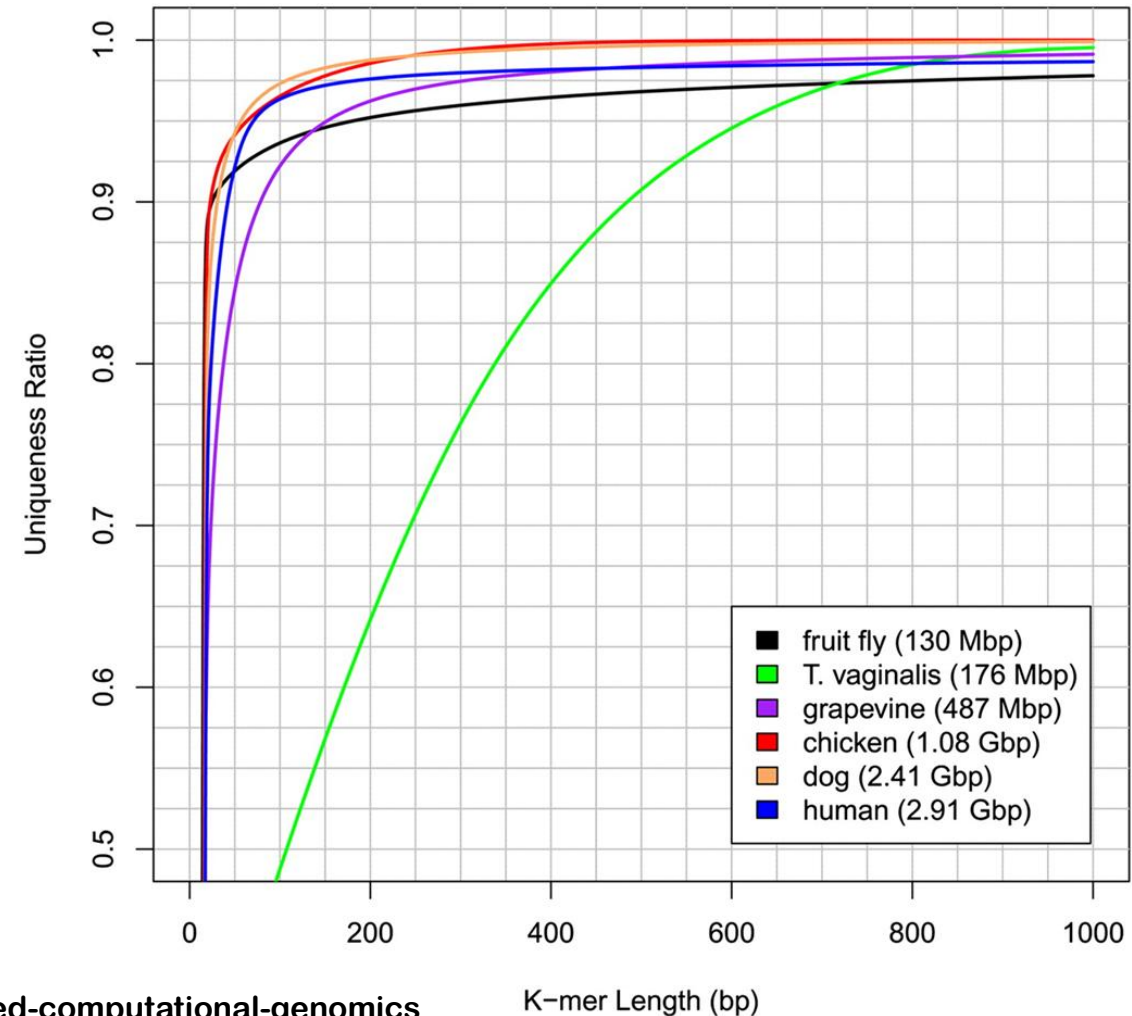
**Every one of these is present in the human genome at least once**

<http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-9-167>

# Note?

**It takes a very long k-mer  
to be unique in most genomes!**

<http://genome.cshlp.org/content/20/9/1165>



# Note?

- We have a sequencing run with over 100 million reads.
- After processing, the reads are between 20 and 25 nucleotide long.
- We would like to know **if these sequences are in the human genome**, and if so **where**.
- For a **20-mer** such as ACGTGTGACGTGATCTGAGC takes about **10 seconds**.
- Querying 100 million sequences would take more than 30 years.
- Without any indexing techniques, the whole genome will be scanned for every query.
- Hash-based indexing, **access time is fast**, does not depend on the text size.



## Note?

- Kmers length ranges from 20 to 25 chars, each character will be represented in two bits (A:00, C:01, G:10, T:11). 20 to 25 kmers will take 40 to 50 bits of storage.
- The human genome contains over 3.2 billion nucleotides, so we need at least 108 GB (in reality many 20 to 25-mers are repeated so this number would be lower).
- If the storage required for the locations and the overhead for the dictionary will be added, the total size will be over 200 GB.



**Thank you!**