# 1. Bioinformatic workflows.
# 2. Variant calling.
# 3. Cancer analysis.
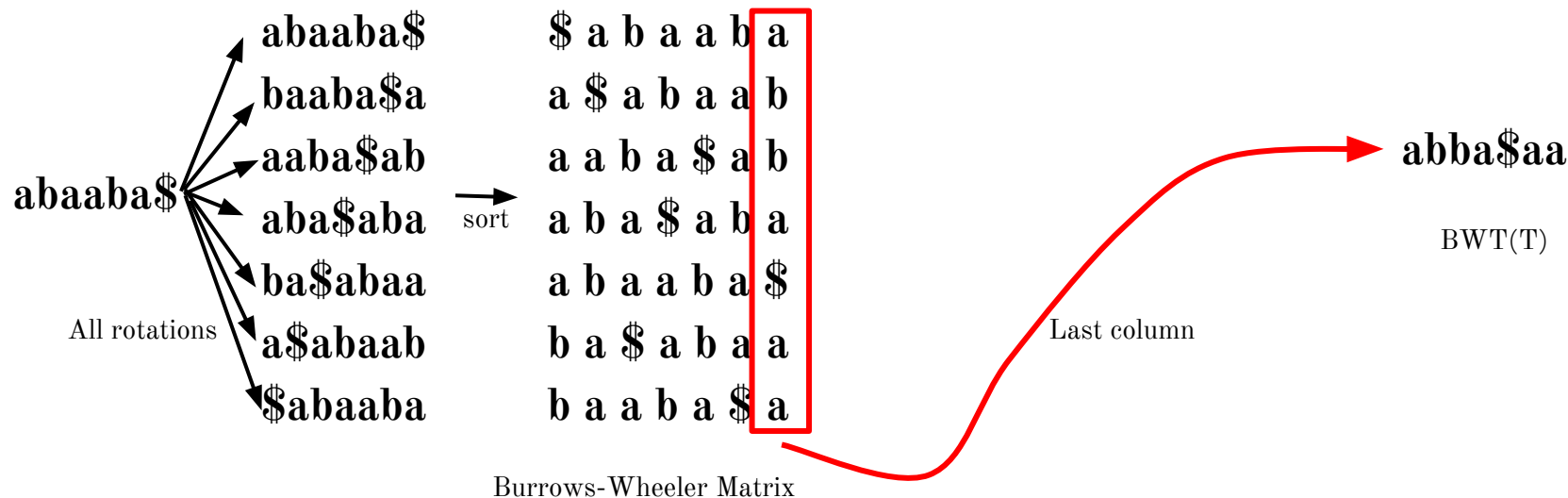
Lesson 05

# Recapitulation

Exact string matching algorithms

Multimap (sorted index)

on-line          off-line

Hash table

Suffix array

Boyer-Moore

Suffix trie

Suffix tree

Burrows-Wheeler
transform + FM index

# Burrows-Wheeler Transform

Reversible permutation of the characters of a string, used originally for compression

abaaba$

All rotations

| | | |
|---|---|---|
| abaaba$ | | $ a b a a b **a** |
| baaba$a | | a $ a b a a **b** |
| aaba$ab | | a a b a $ a **b** |
| aba$aba | sort | a b a $ a b **a** |
| ba$abaa | | a b a a b a **$** |
| a$abaab | | b a $ a b a **a** |
| $abaaba | | b a a b a $ **a** |

Burrows-Wheeler Matrix

Last column

abba$aa

BWT(T)

How is it useful for compression?          How is it reversible?          How is it an index?
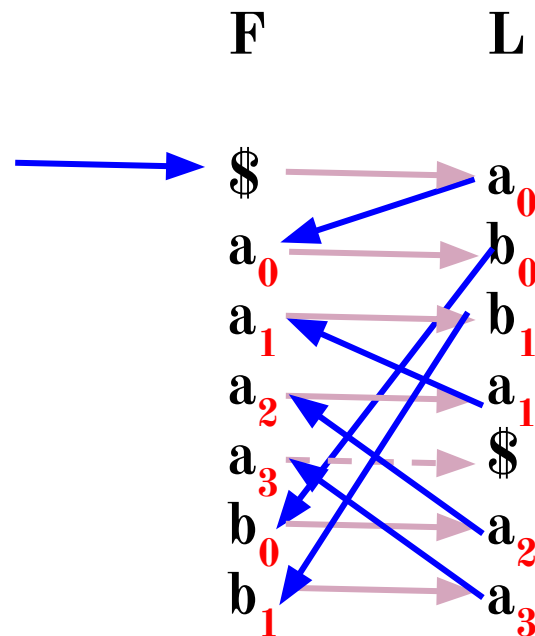
# Burrows-Wheeler Transform: reversing

Reverse BWT(T) starting at right-hand-side
of T and moving left

Start in first row. F must have \$.
L contains character just prior to \$: $a_0$
...

**F**          **L**

$\$$           $a_0$
$a_0$          $b_0$
$a_1$          $b_1$
$a_2$          $a_1$
$a_3$          $\$$
$b_0$          $a_2$
$b_1$          $a_3$

Reverse of chars we visited $= a_3\ b_1\ a_1\ a_2\ b_0\ a_0\ \$ =$ T

# FM Index: querying

Look for range of rows of BWM(T) with P as prefix
Do this for P's shortest suffix, then extend to successively longer suffixes
until range becomes empty or we've exhausted P

$$P = ab\textcolor{red}{a}$$

$$P = a\textcolor{red}{ba}$$

$$
\begin{array}{ccccccc}
\$ & a & b & a & a & b & a_0 \\
a_0 & \$ & a & b & a & a & b_0 \\
a_1 & a & b & a & \$ & a & b_1 \\
a_2 & b & a & \$ & a & b & a_1 \\
a_3 & b & a & a & b & a & \$ \\
b_0 & a & a & b & a & \$ & a_2 \\
b_1 & a & \$ & a & b & a & a_3 \\
\end{array}
$$

Look at those rows in L.
$b_0$, $b_1$ are b-s occurring
just to left.

Use LF Mapping. Let new
range delimit those b-s

$$
\begin{array}{ccccccc}
\$ & a & b & a & a & b & a_0 \\
a_0 & \$ & a & b & a & a & b_0 \\
a_1 & a & b & a & \$ & a & b_1 \\
a_2 & b & a & \$ & a & b & a_1 \\
a_3 & b & a & a & b & a & \$ \\
b_0 & a & a & b & a & \$ & a_2 \\
b_1 & a & \$ & a & b & a & a_3 \\
\end{array}
$$

# FM Index: querying

We have rows beginning with **ba**, now we seek rows beginning with **aba**

$$P = \text{a}\textcolor{red}{\text{ba}}$$

$$
\begin{array}{llllllll}
\$ & a & b & a & a & b & a_0 \\
a_0 & \$ & a & b & a & a & b_0 \\
a_1 & a & b & a & \$ & a & b_1 \\
a_2 & b & a & \$ & a & b & a_1 \\
a_3 & b & a & a & b & a & \$ \\
b_0 & a & a & b & a & \$ & a_2 \\
b_1 & a & \$ & a & b & a & a_3 \\
\end{array}
$$

Occurs just to the left

$$P = \textcolor{red}{\text{aba}}$$

**F**  **L**

$$
\begin{array}{llllllll}
\$ & a & b & a & a & b & a_0 \\
a_0 & \$ & a & b & a & a & b_0 \\
a_1 & a & b & a & \$ & a & b_1 \\
a_2 & b & a & \$ & a & b & a_1 \\
a_3 & b & a & a & b & a & \$ \\
b_0 & a & a & b & a & \$ & a_2 \\
b_1 & a & \$ & a & b & a & a_3 \\
\end{array}
$$

Use LF mapping

Now we have the rows with prefix **aba**

# FM Index

1. L = BWT(T)
2. First column (number of appearances of each character)
3. Suffix Array (or SA Sample)
4. Tally (rank, occurrences) matrix

# FM Index: Example



Search for:
GAGA

# Bioinformatic workflows and cloud computing

—

Lesson 05.1

# What is a workflow?

- Acyclic graph of tools connected to perform some analysis
- Workflow's nodes are:
  - Inputs (file or parameter)
  - Tools
  - Outputs
  - Workflow

```
bwa mem ref.fa read1.fq read2.fq >
aln.sam
sam2bam aln.sam > aln.bam
```

Fasta

BWA-MEM

SAM2BAM

SAM

FASTQ

# Why we need a workflow?

# Common Workflow Language

- Define inputs and outputs of a software, runtime and requirements
- Define how to connect software, creating a workflow
- Ensure reproducibility and portability
- Think of CWL as a detailed recipe!

# Common Workflow Language

- **Reproducible** analyses (standard)
- Scalable execution
- Metadata & file registry integration
- **Portability** - deployable on multiple platforms
- Revision management and versioning
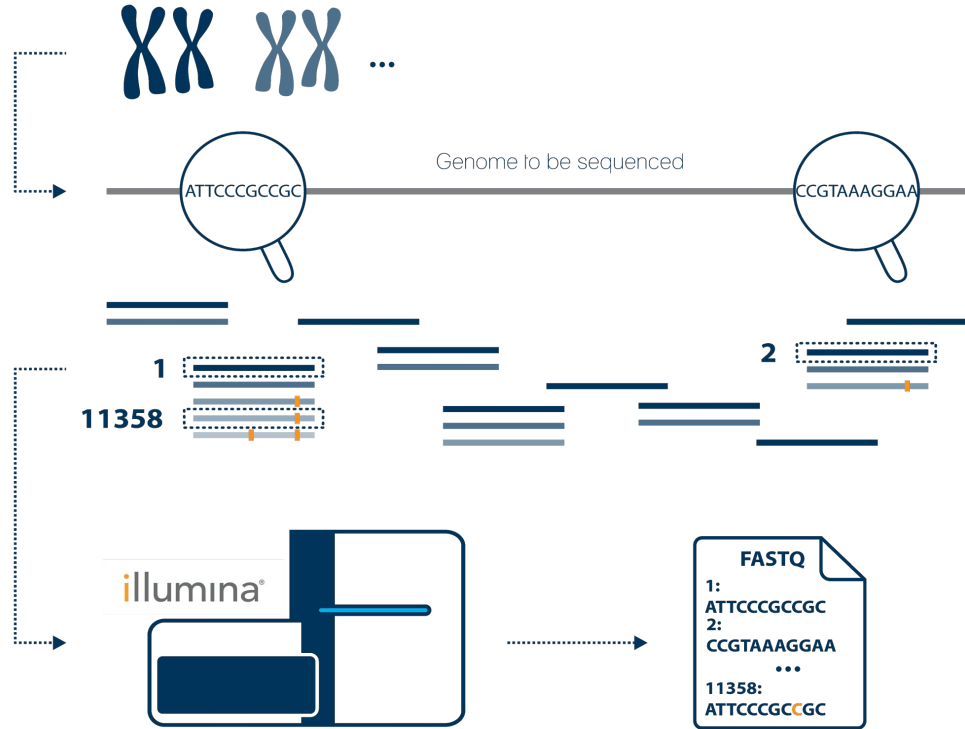  User management / permissions

# CWL @ Cloud
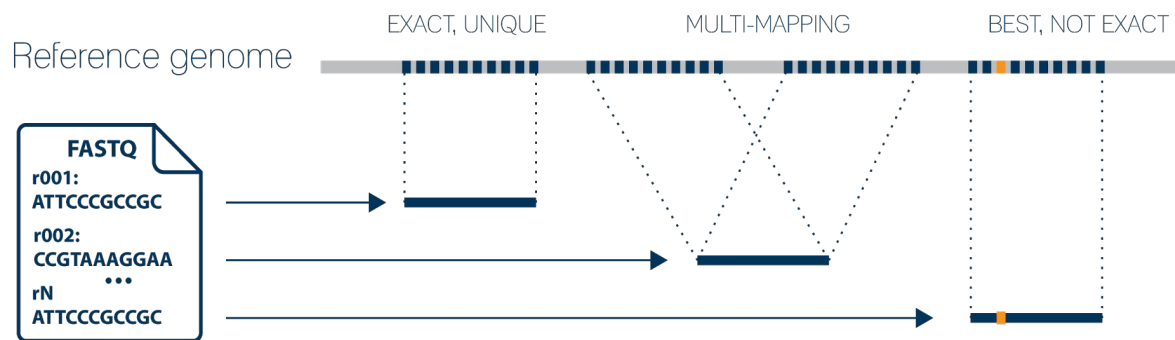
# Variant Calling

—

Lesson 05.2

# Reminder: DNA Sequencing

We got a FASTQ files with the "reads" – little pieces of the genome.

# Reminder: Alignment

# Introduction to Variant calling

- Variant calling is the process of finding differences between reference genome and observed sample

- We need aligned reads to the reference genome so we can find – "call" variants

- Different types of genomic variants

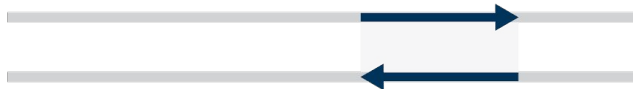# Genomic Variants

Single nucleotide variant

Copy number variant
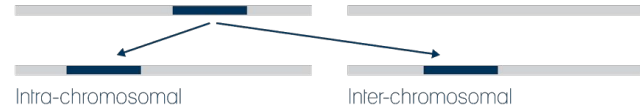
Deletion

Translocation

Intra-chromosomal          Inter-chromosomal

Insertion

Whole genome duplication
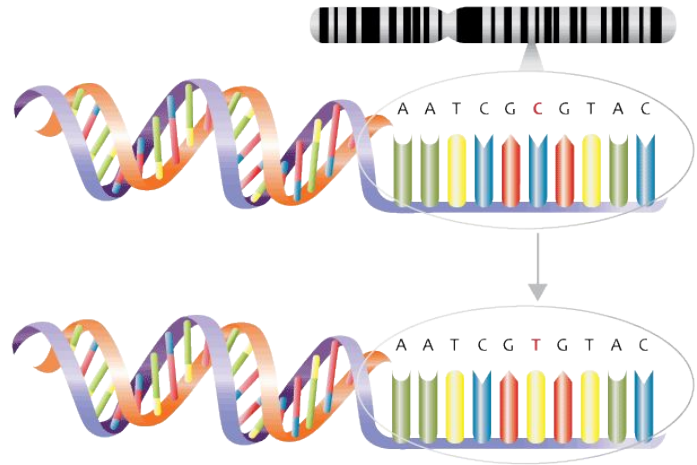
Inversion

Duplication

Tandem          Interspersed

# Genomic Variants

- SNV ( Single Nucleotide Variant)

  Simple ones - not a big change on the first look, but...
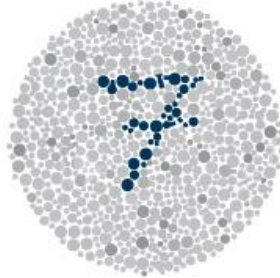
# Genomic Variants

Each of those characteristics causes one SNV

LONGER EYELASHES          DALTONISM          LESS SLEEPING          SUPER STRENGTH

# Genomic Variants

**Breast Cancer**

**BRCA2** gene (TS)

**SNV id : rs1799954**

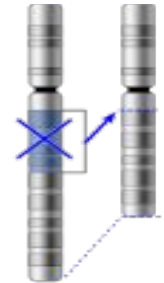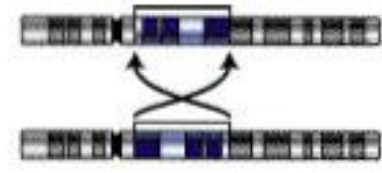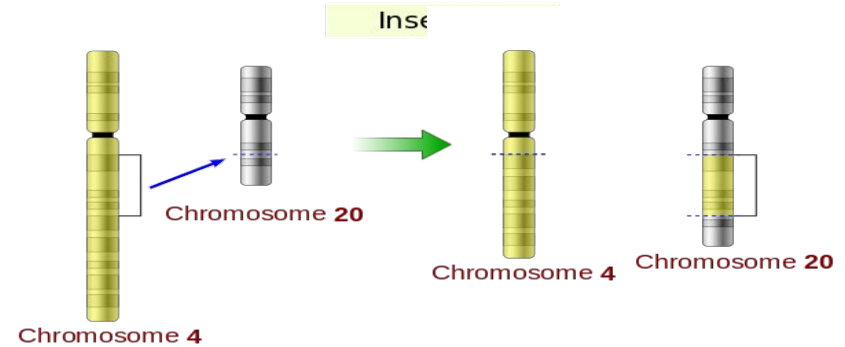Chromosome 13
Position 32,340,455

Cancer genotypes: **CC, CT and TT**

http://www.eupedia.com/genetics/cancer_related_snp.shtml
https://www.snpedia.com/index.php/Rs1799954

# Genomic Variants

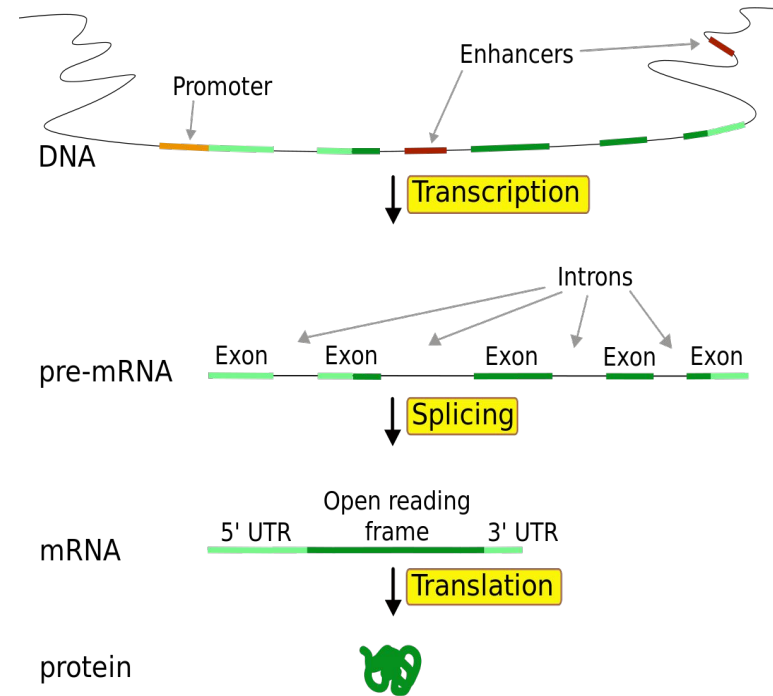Deletions, Insertions, Translocations, Inversions, and some others...

# Genomic Variants

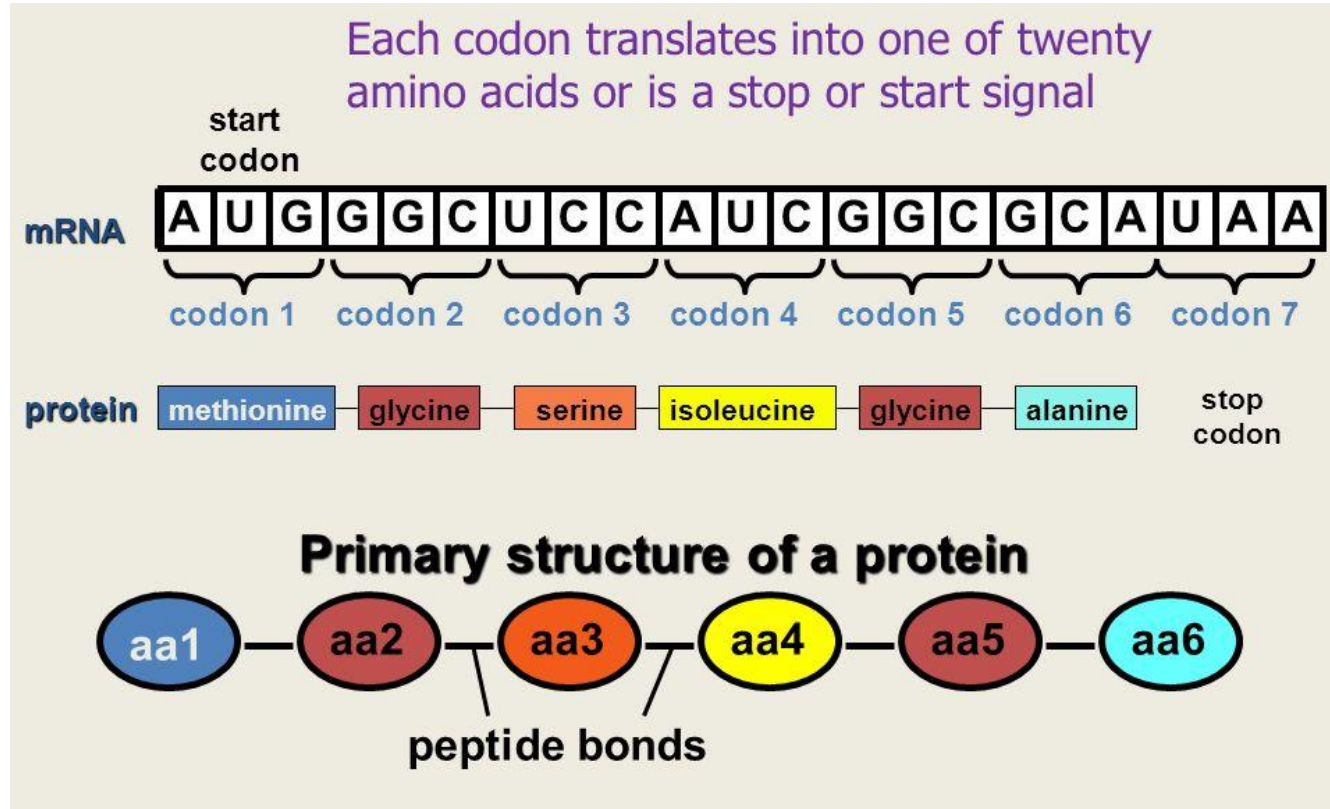Based on the variant location, we can predict if mutation will have impact.
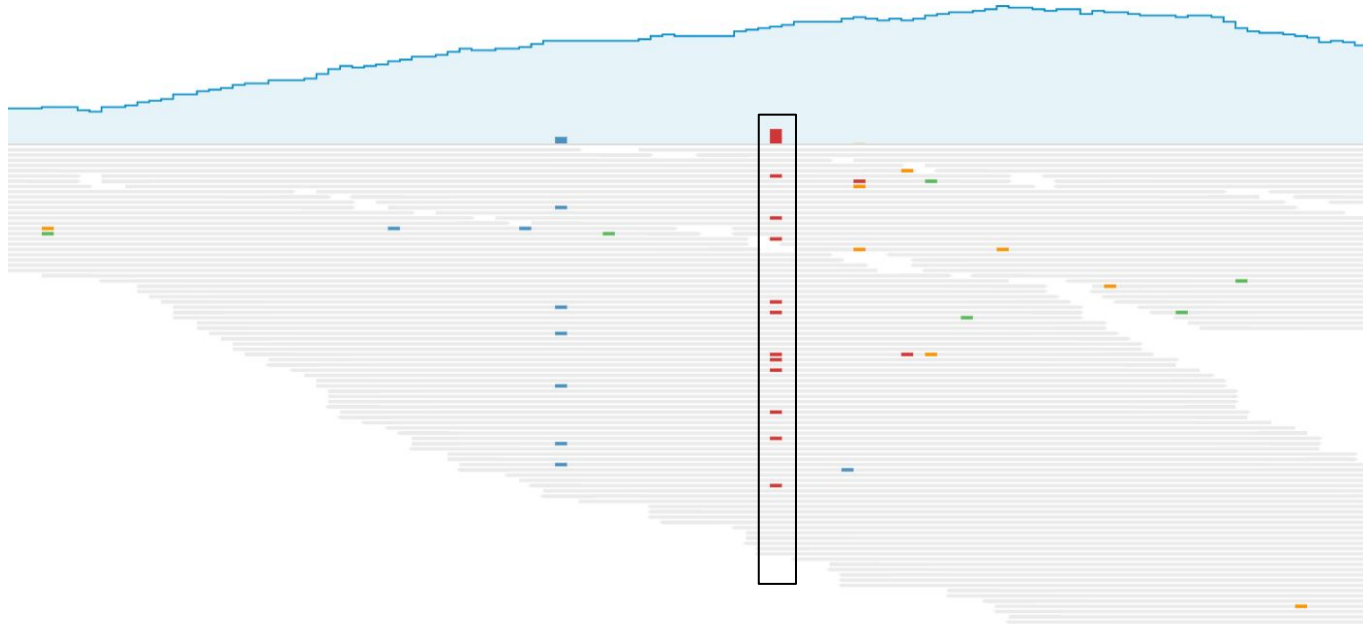


- Central dogma

# RNA to Protein

# Genomic Variants

- Variants can have different impact on human cells and organism
- Single Nucleotide Variants(**SNV**):
  - Harmless
    - **Silent** – Usually no effect
  - Harmful:
    - **Missense** – Amino acid change
    - **Nonsense**(Start/Stop Gain/Lost) – AUG / UAG, UAA, UGA
  - Depends on the location
    - **Noncoding regions** ( Promoter, Enhancer, lncRNA, miRNA...)
- Insertions/Deletions – **INDELS**
  - **In frame**
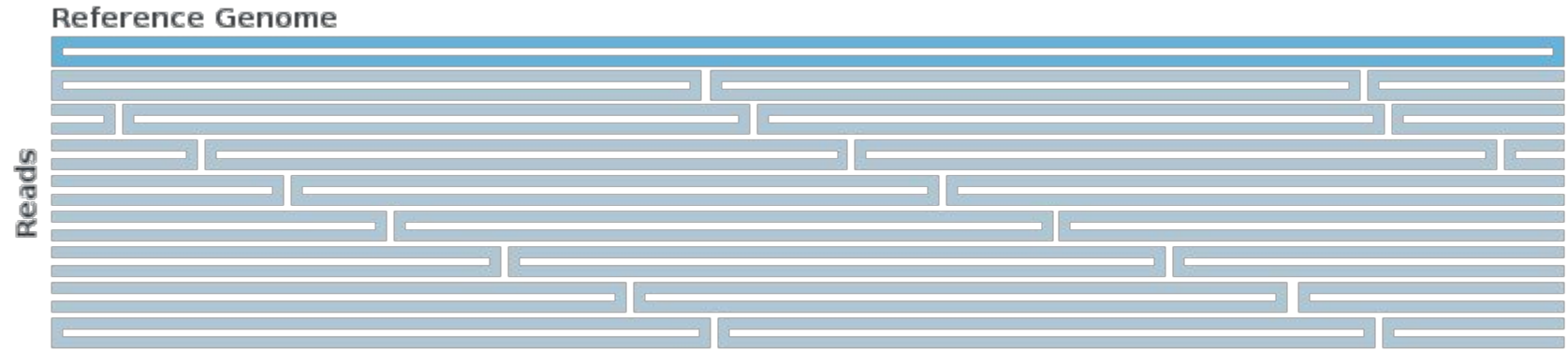  - **Out of frame (Frameshift)**

# What is the pileup?

# Ideal Variant Calling
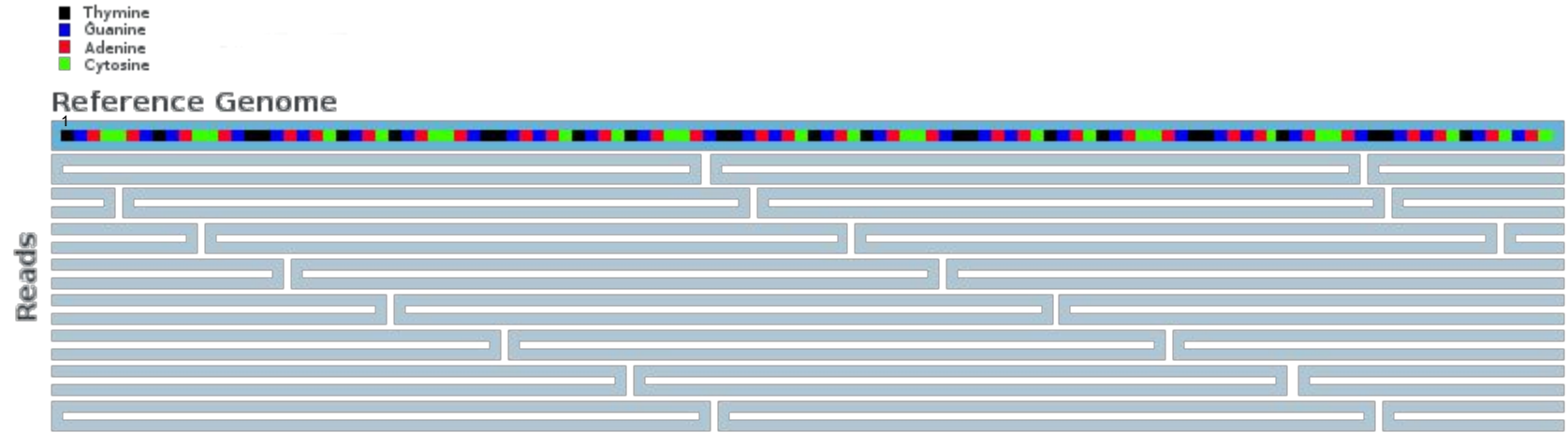
# Ideal Variant Calling



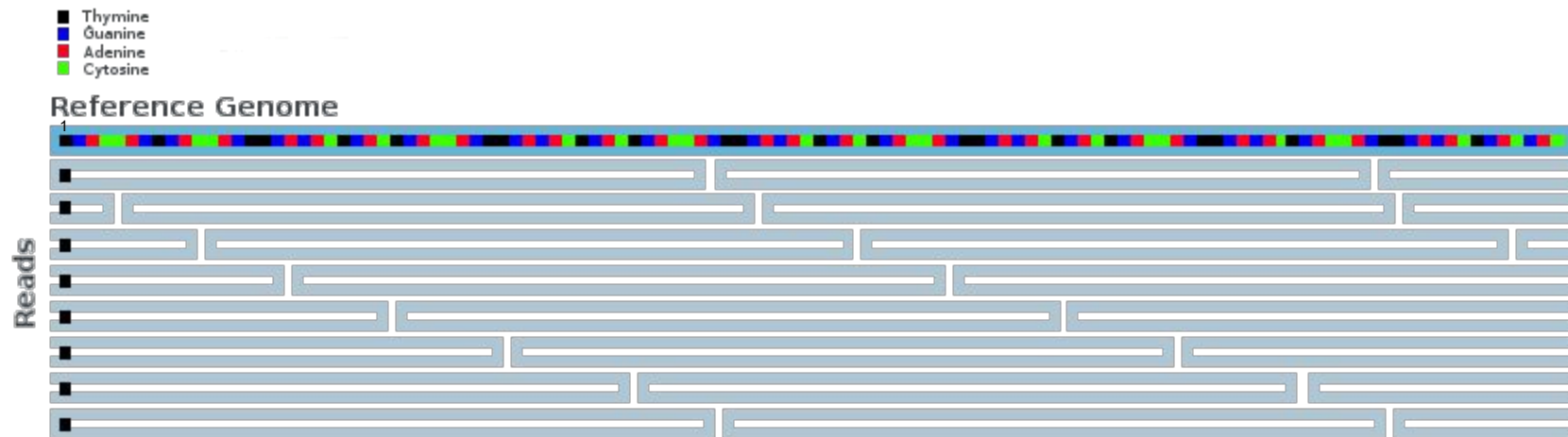Ideally we will have uniform distribution of reads.

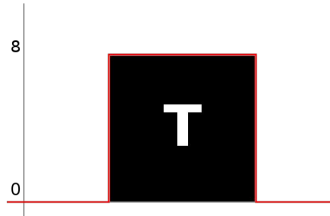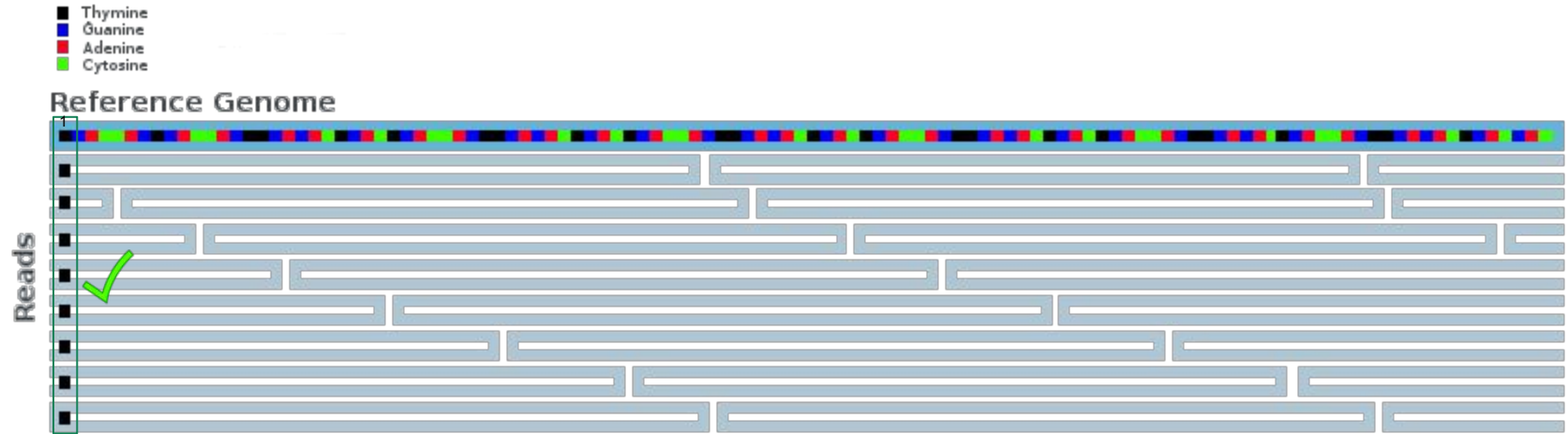# Ideal Variant Calling

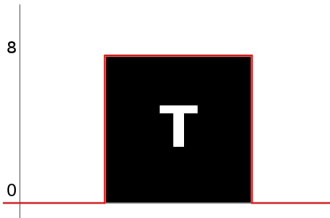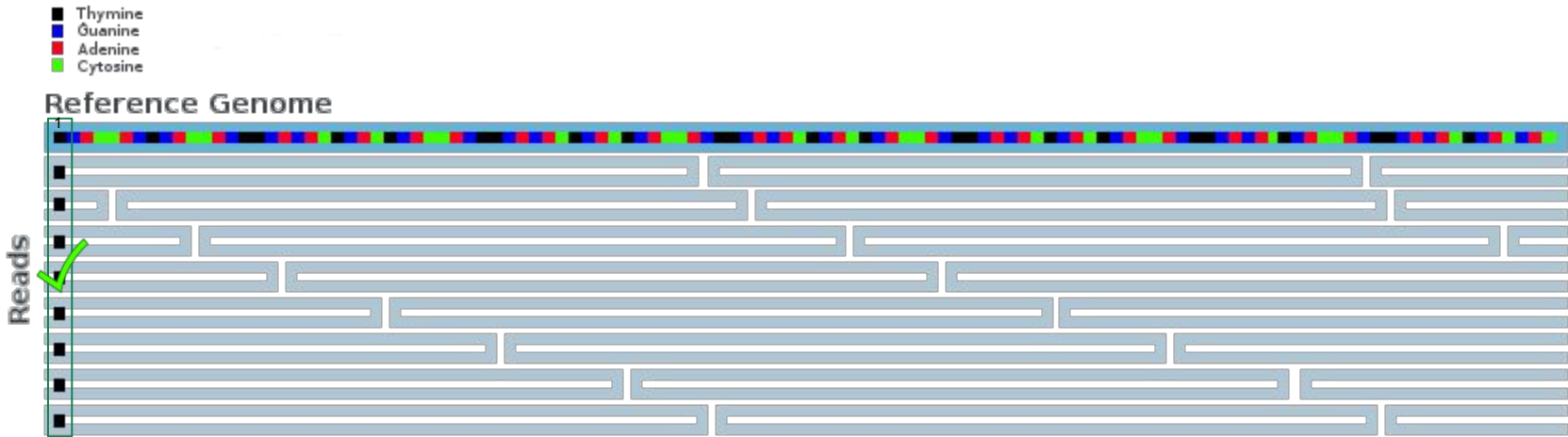# Ideal Variant Calling

# Ideal Variant Calling

# Ideal Variant Calling



- We have "T" in the all reads covering that position
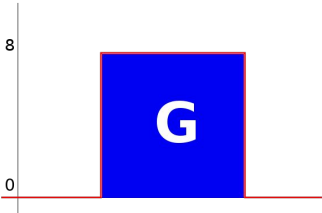
# Ideal Variant Calling



**How can we represent what we have observed?**

| CONTIG | POS | REF | ALT | GT |
| --- | --- | --- | --- | --- |
| X | 1 | T | - | 0/0 |

# Ideal Variant Calling



| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 1 | T | - | 0/0 |
| X | 2 | G | - | 0/0 |

# Ideal Variant Calling



| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 1 | T | - | 0/0 |
| X | 2 | G | - | 0/0 |
| X | 3 | A | - | 0/0 |

# Ideal Variant Calling
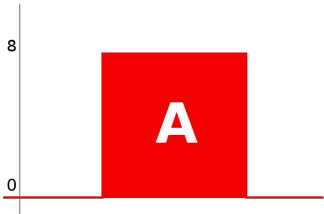


| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 2 | G | - | 0/0 |
| X | 3 | A | - | 0/0 |
| X | 4 | **C** | **T** | **1/1** |

# Ideal Variant Calling



| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 3 | A | - | 0/0 |
| X | 4 | **C** | **T** | **1/1** |
| X | 5 | **C** | **A** | **1/1** |

# Ideal Variant Calling



| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 4 | **C** | **T** | **1/1** |
| X | 5 | C | A | **1/1** |
| X | 6 | A | - | 0/0 |

# Ideal Variant Calling



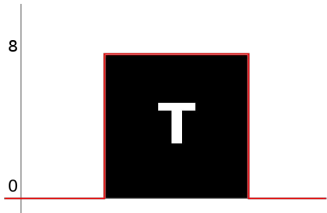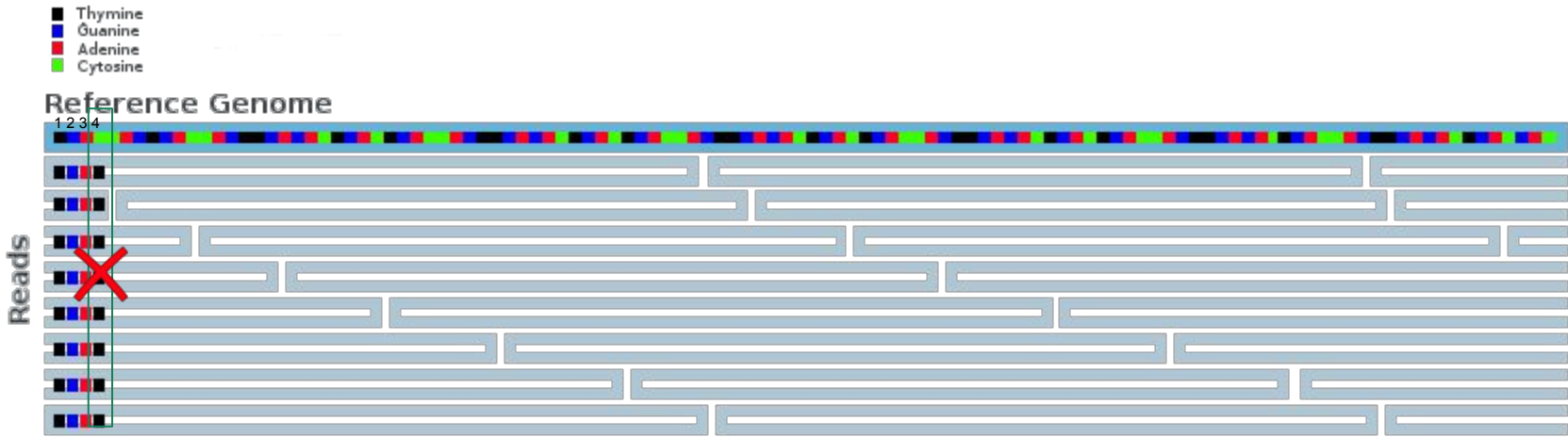| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 5 | **C** | **A** | **1/1** |
| X | 6 | A | - | 0/0 |
| X | 7 | G | - | 0/0 |

# Ideal Variant Calling



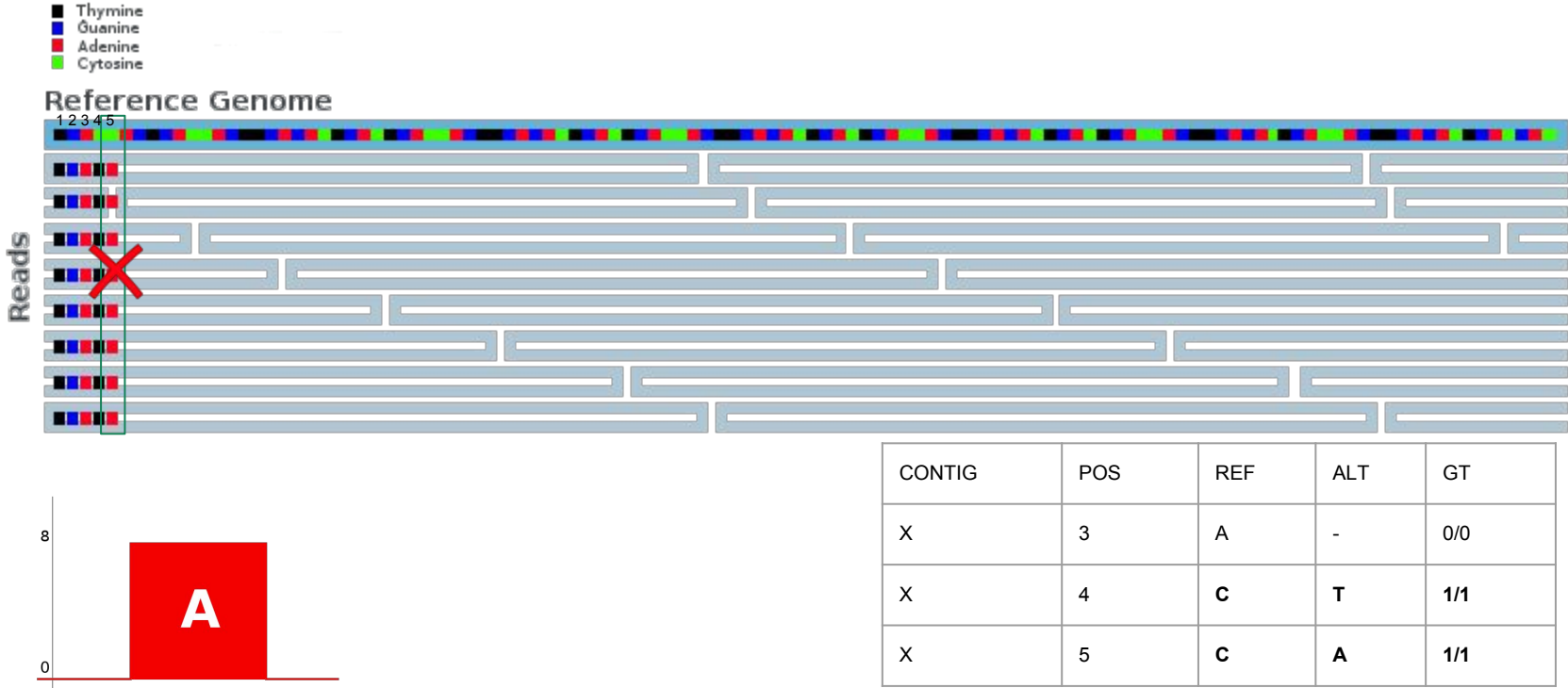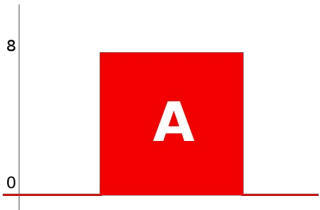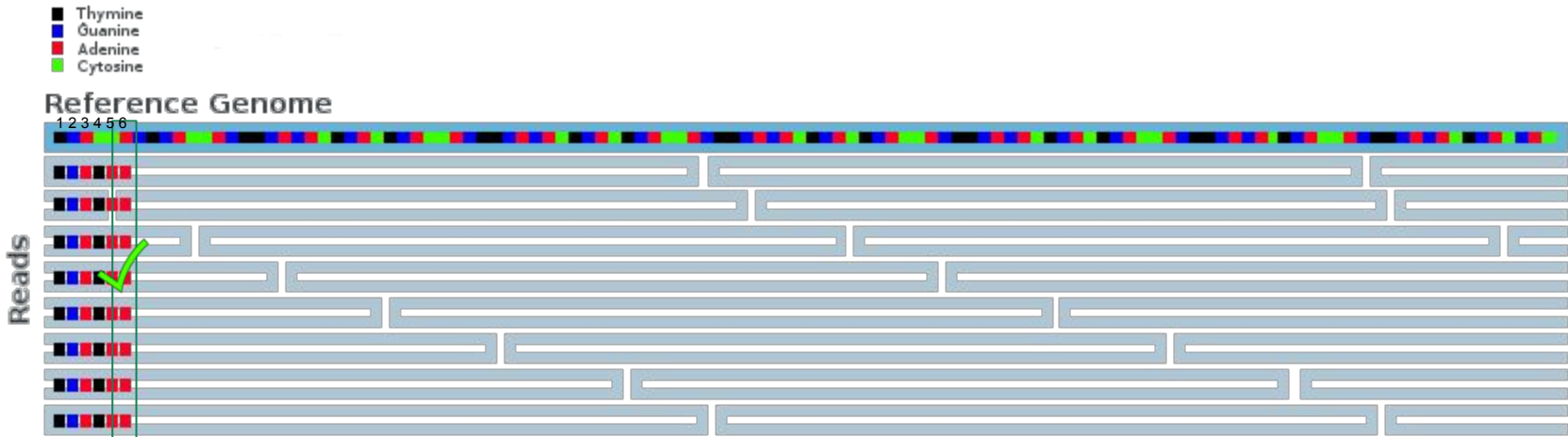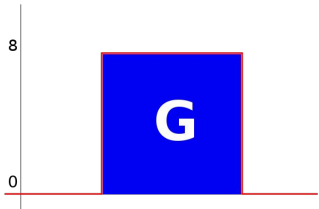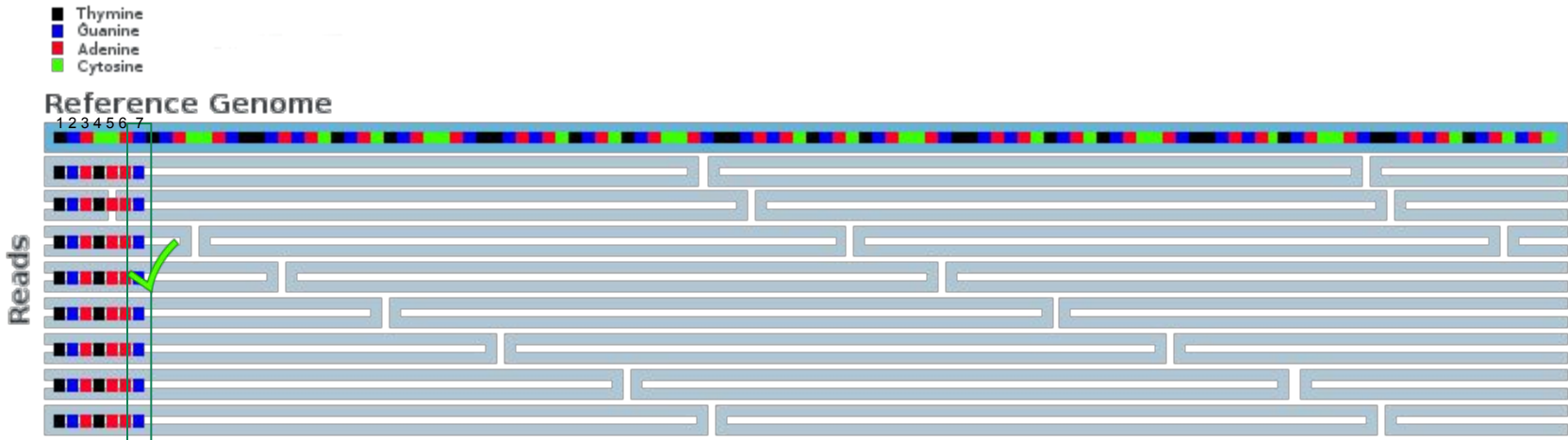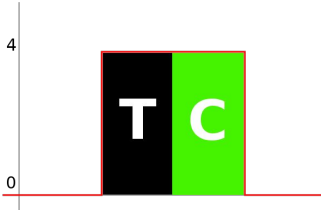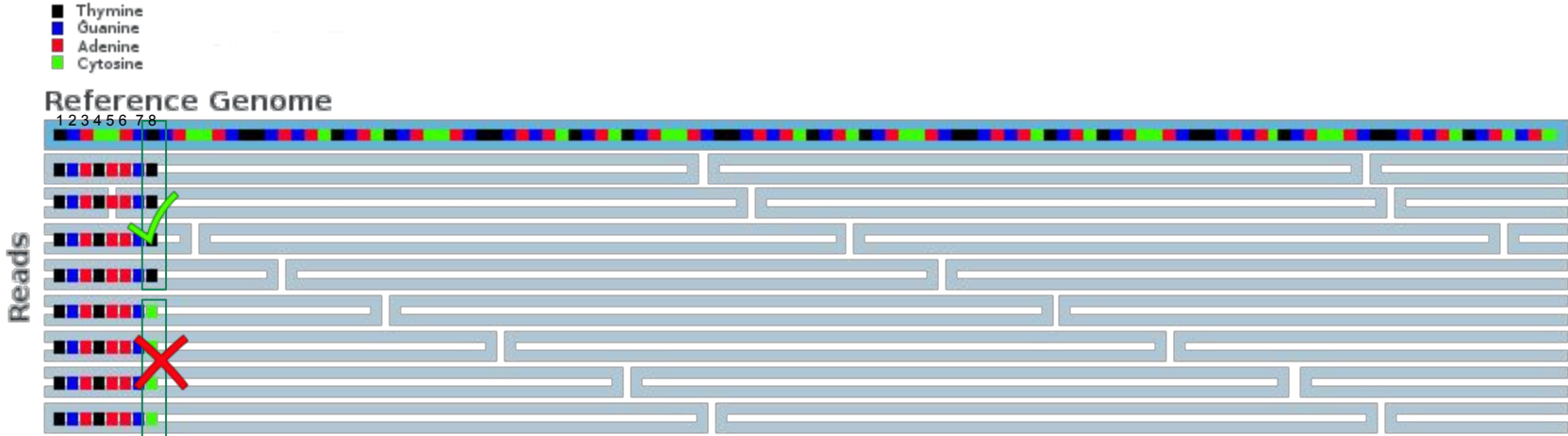| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 6 | A | - | 0/0 |
| X | 7 | G | - | 0/0 |
| X | 8 | **T** | **C** | **0/1** |

# Ideal Variant Calling



| CONTIG | POS | REF | ALT | GT |
|--------|-----|-----|-----|-----|
| X | 7 | G | - | 0/0 |
| X | 8 | **T** | **C** | **0/1** |
| X | 9 | **G** | **A,T** | **1/2** |

# Variant Calling

- Two possible cases:

1. All of the bases in pileup are the same nucleotide [A,T,C,G]

2. Different nucleotides exist in the pileup

- In the simplest case, assume diploidy
  - There can be only two alleles at a site
  - If there are more than two different letters in the pileup we will only consider the most common two
    (assume others are errors and discard them)

Normal Human Karyotype

Hypothetical - traffic light problem



?

Stevan

# Binomial distribution

- Models the number of successes in a sequence of yes/no experiments
- Parameters:
  - n - number of trials
  - p - probability of a success in a single trial
  - Probability that K out of n trials will be success

$$f(k; n, p) = \Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

# Binomial distribution

# Variant Calling

- So, when we have two letters in the pileup, what should we call?
  - Let's call the two "letters" **b** and **b'** *(b, b' ∈ [A, C, T, G])*
  - Let **n** be the total number of bases, and **k** number of b' bases
  - Three possible explanations for the pileup:
    - Genotype is bb; k bases are errors, n-k are correct
    - Genotype is b'b'; n-k bases are errors, k are correct
    - Genotype is bb'; all n bases are correct
  - Now we need to find the probabilities of these three cases
    - Will pick the most probable one!

# Variant Calling – advance

- We assumed a flat error rate
  - But we have Base qualities from the sequencer
  - Machine-specific error profiles
- We can look at mapping qualities
  - Mapping errors are a big source of errors
- We can look at haplotypes
  - Errors don't segregate nicely
- Population-based methods
  - Separate variant calling from genotyping

# Variant calling results – check out BAM file



| CHR | POS | REF | ALT | FORMAT | NA12877 |
|-----|-----|-----|-----|--------|---------|
| 1 | 14125 | T | C | GT, VAF | 0/1, 0.6 |

# Variant calling results

- The result of Variant Calling is a file in VCF format, which contains mutations

- A plain text file format for storing variant data

- A number of line starting with ## -the header

- Main header line:
  #CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE1

- This is followed by the actual variant data, one entry per line
  22 10001 . A C 40 PASS DP=14 GT 0/1

- More than one sample can be in one line

- For details: http://samtools.github.io/hts-specs/VCFv4.2.pdf

# Variant calling results

- Example of VCF format
- Each row represents one mutation

| CHR | POS | REF | ALT | FORMAT | NA12878 |
|-----|-----|-----|-----|--------|---------|
| 1 | 14300 | A | G | GT, VAF | 0/1, 0.4 |
| 2 | 15367 | A | C | GT, VAF | 1/1, 0.9 |
| 3 | 25612 | C | G,A | GT, VAF | 1/2, ? |
| 5 | 5632 | TA | T | GT, VAF | 0/1, 0.5 |
| 7 | 7824 | T | TA | GT, VAF | 1/1, 0.8 |

# Variant Calling Format File

# Computational Cancer Analysis

—

Lesson 05.3

# DNA replication

# What is cancer?

Mutation during DNA replication can fall to:
1. Intron (no change)
2. Important gene (cell dies, organism lives)
3. Gene that stops cell division (cell lives, organism...)

What causes cancer (increases probability of mutation)?
1. EM radiation
2. Chemical agents
3. Free radicals
4. Genetic factors
5. Infections (viruses)

A dividing lung cancer cell.
Credit: National Institutes of Health

# Genetic factors

A typical pedigree from a family with a mutation in the BRCA1 (tumor suppressor) gene

Fathers can be carriers and pass the mutation onto offspring

Not all people who inherit the mutation develop the disease, thus patterns of transmission are not always obvious



Breast Cancer

Ovarian Cancer

Unaffected Carrier

# Viruses and retroviruses

# What is metastasis?

Body's cells begin to divide without stopping and spread into surrounding tissues
Cancer cells - ignore signals that normally tell cells to stop dividing or that begin a process known as programmed cell death, or **apoptosis**, which the body uses to get rid of unneeded cells
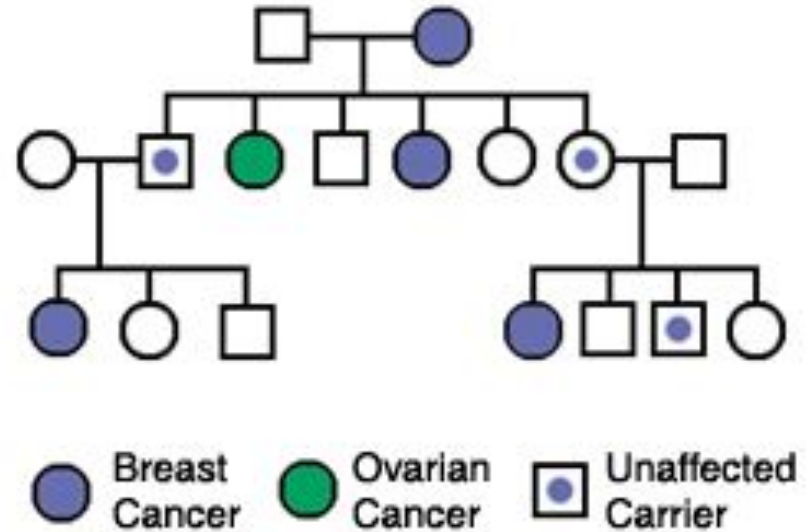


**Metastasis**

Brain metastasis

Cancer spreads to other parts of the body

Lung metastasis

Primary cancer

Metastatic tumor

Cancer cells in lymph system

Cancer cells in the blood

Primary cancer

© 2014 Terese Winslow LLC
U.S. Govt. has certain rights

# "Drivers" of Cancer

Cancer is a genetic disease that is caused by changes to genes which control the way our cells function, especially how they grow and divide:

1.  **Abnormal growth (proto-oncogenes)** Cellular growth mechanism is damaged and cell starts to multiply uncontrollably
2.  **Damaged control mechanism** (tumor suppressors) - Cells with certain alterations in tumor suppressor genes may divide in an uncontrolled manner (TP53 - Apoptosis)
3.  **Damaged DNA repair mechanism** (Accumulated errors in this group of genes can lead to uncontrollable proliferation)

# Cancer cells

**Our body develops thousands cancer cells every day! OMG OMG**



IDENTIFYING **THE ENEMY**

Damaged or infected cell alerts the immune system by holding out a fragment of protein for inspection

Each of your millions of T-cells has a slightly different receptor. Each detects different proteins

Fragment of protein

T-cell receptor

T-cell detects and destroys damaged or infected cell

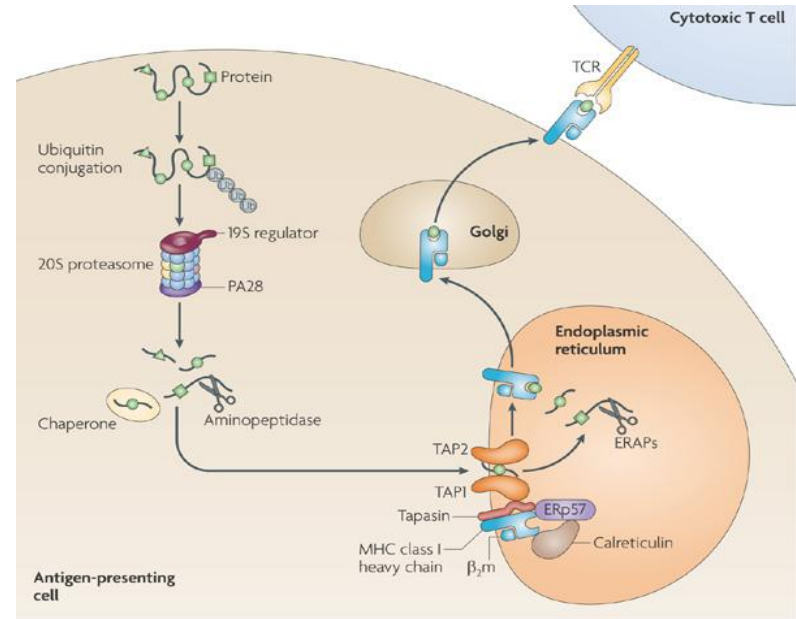T-cell with the right receptor recognises the enemy protein

# MHC Complex

MHC is a set of cell surface proteins essential for the acquired immune system to recognize foreign molecules (translated from HLA regions from the genome for humans)

MHC molecules bind to **protein fragments available in the cell**

MHC molecule with **antigen** (MHC complex) is "presented" outside of the cell to cytotoxic T cells and helper T cells
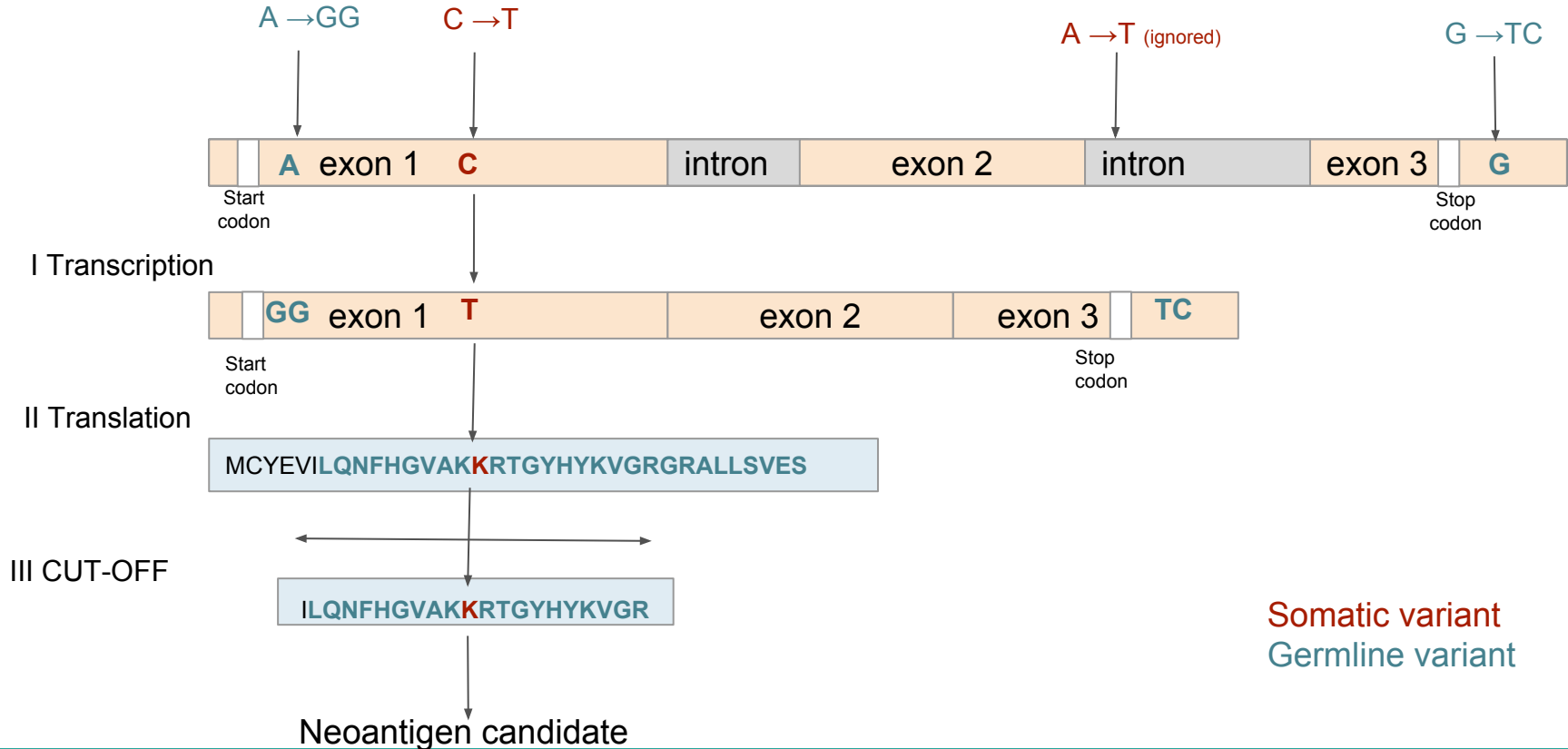


Nature Reviews | Immunology

# So, what can be done there?

1. Identify NEOANTIGENS - proteins presented only by cancer cells
2. "Program" T-cells to recognize neoantigens

Compare DNA from Tumor and Normal tissue
Mutations present in tumor - somatic mutations

# From DNA somatic mutation to neoantigen

# References

**How the Immune System Works - Lauren Sompayrac**