



Molecular Biology (3)

The human genome

Mamoun Ahram, PhD

Resources



- This lecture
- Cooper, Ch. 6, pp. 157-160, 195-205, 209-212

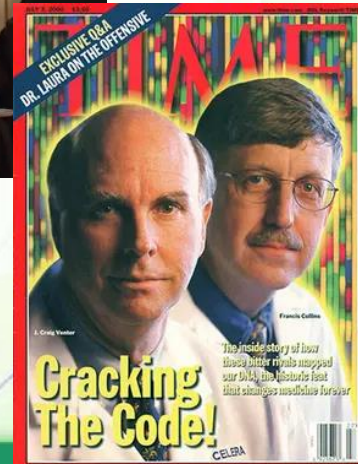
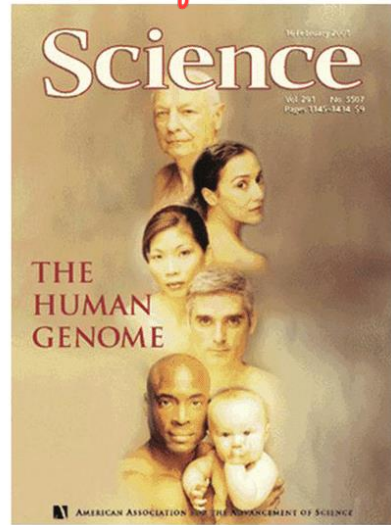
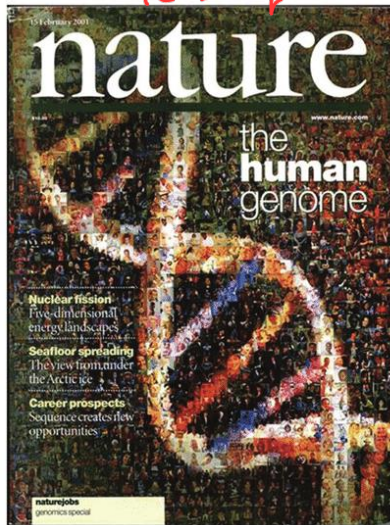
The human genome project



↳ total collection of DNA

- A \$3 billion, 13-year, multi-national project launched in 1990 led by the US government to (know the) **sequence the human genome** and to map and identify the genes (a draft was published in 2001 and 92% was completed in 2004).

↳ the sequence of Nitrogen bases



Major outcomes



- Determination of the number of human genes
- Development of major technologies
- Completed sequences of other genomes
- Open discussion of legal and ethical issues

absolutely not important
↳ (now we can know the number of gens in about 24 hours and it cost about 900 \$)



| SPECIES | BASE PAIRS (estimated) | GENES (estimated) | CHROMOSOMES |
|--|---------------------------|----------------------|-------------|
| Human (<i>Homo sapiens</i>) | 3.2 billion | X ~ 25,000 | 46 |
| Mouse (<i>Mus musculus</i>) | 2.6 billion | X ~ 25,000 | 40 |
| Fruit Fly (<i>Drosophila melanogaster</i>) | 137 million | 13,000 | 8 |
| Roundworm (<i>Caenorhabditis elegans</i>) | 97 million | 19,000 | 12 |
| Yeast (<i>Saccharomyces cerevisia</i>) | 12.1 million | 6,000 | 32 |
| Bacteria (<i>Escherichia coli</i>) | 4.6 million | 3,200 | 1 |
| Bacteria (<i>H. influenzae</i>) | 1.8 million | 1,700 | 1 |

the doctor did not said memorise it (I will not even if he said)
and it is old numbers too that 😂

Nucleotides per genomes



He said not for memorizing

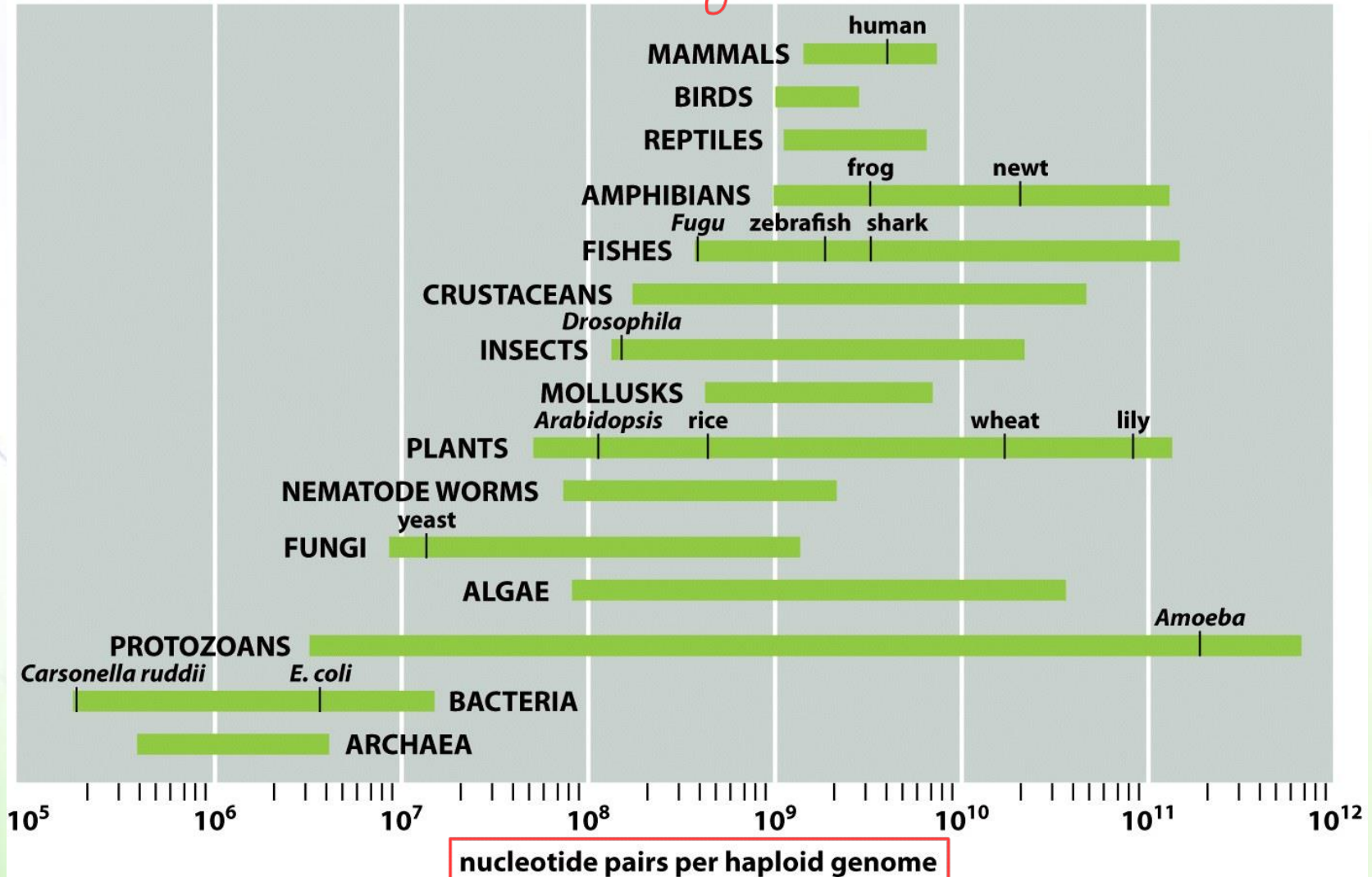


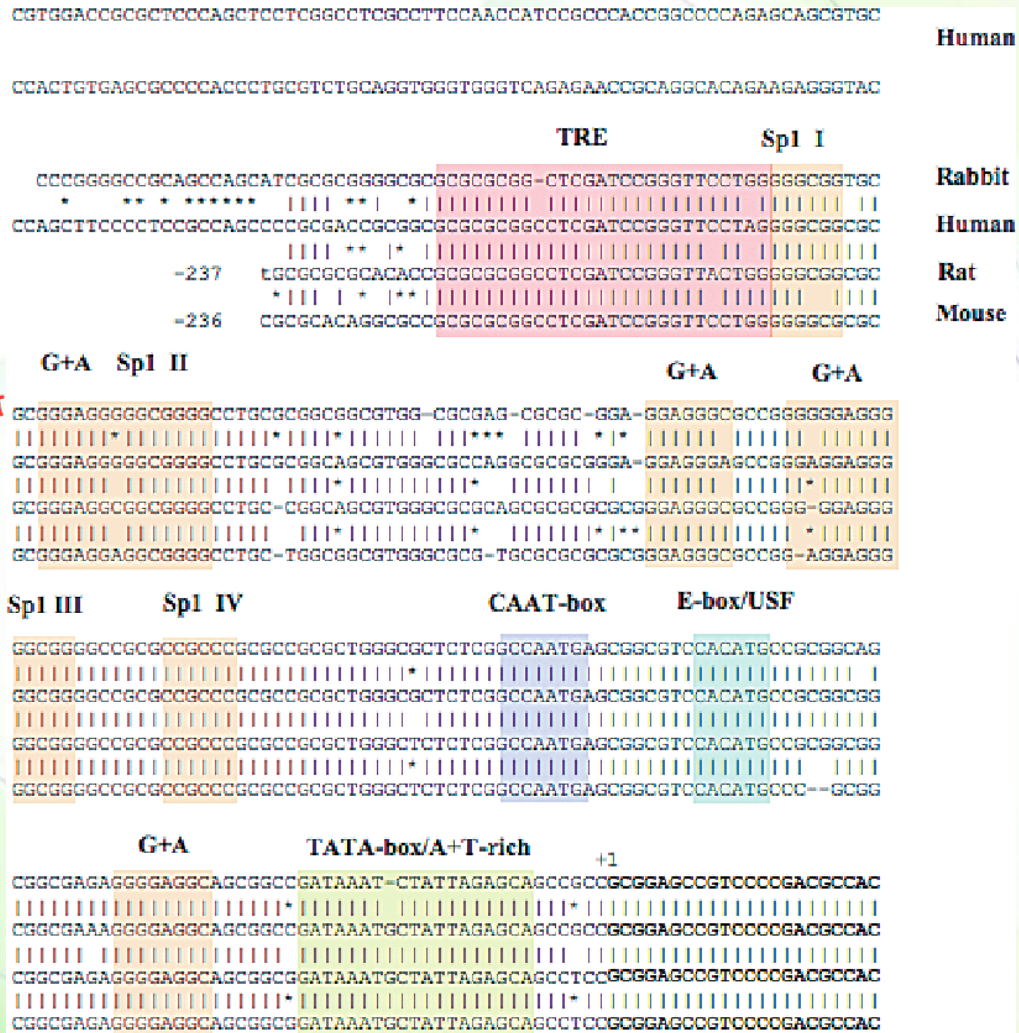
Figure 1-41 Essential Cell Biology 3/e (© Garland Science 2010)



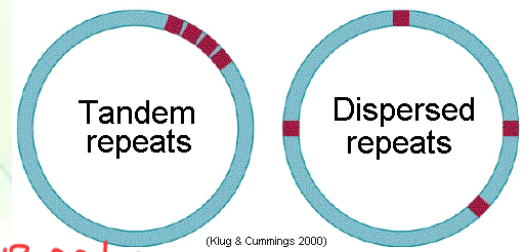
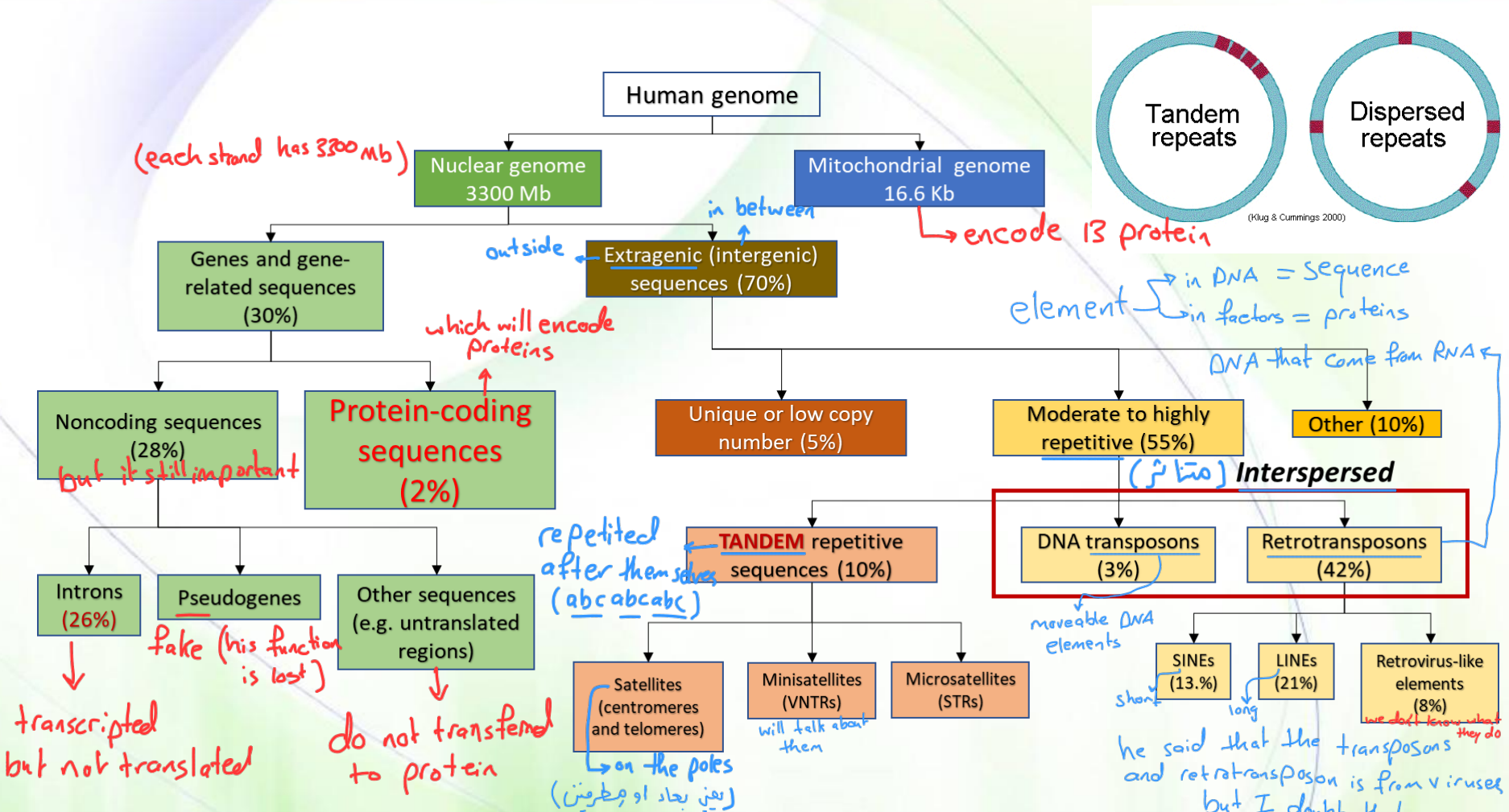
DNA homology

(i.e. sequence similarity)

There are similar to each other in some regions which is most important region



Components of the human genome



~5% of the genome contains sequences of **noncoding DNA** that are **highly conserved** indicating that they are **critical to survival**.

Note: All numbers are approximate

an again he mention the evolutionary method
 So I advise you to watch 8 الدكتور اباد لقيني رحلة , ليقين

The ENCODE project (2003-on)



all the numbers are approximate (don't memorize significant numbers)
but memorize it as approximately number so it may be asked in the exam

- ENCODE: Encyclopedia of DNA Elements (ENCODE)
- 80% of the entire human genome is relevant (either transcribed, binds to regulatory proteins, or is associated with some other biochemical activity).

Summary of ENCODE Results

| | |
|----------------------|--------|
| Protein-coding genes | 20,687 |
| Short noncoding RNAs | 8801 |
| Long noncoding RNAs | |
| Pseudogenes | 11,224 |

| | | |
|---|-------|----------------------------|
| Percentage of genome <u>transcribed into RNA</u> | 74.7% | but 2% transfer to protein |
| Percentage of genome-binding <u>transcription factors</u> will talk about it. | 8.1% | |

→ to transfer Data to information

On March 31, 2022...

remember in 2004 they have finished 92% of human genome



A gene: a region of DNA that is transcribed.

A transcript: a RNA molecule that is produced by transcription

Gene annotation

| | |
|---------------------------------|---------|
| Number of genes | 63,494 |
| Protein coding | 19,969 |
| Number of exclusive genes | 3,604 |
| Protein coding | 140 |
| Number of transcripts | 233,615 |
| Protein coding | 86,245 |
| Number of exclusive transcripts | 6,693 |
| Protein coding | 2,780 |

RESEARCH ARTICLE

HUMAN GENOMICS

The complete sequence of a human genome

to understand the repetitive sequence of DNA

Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.055 billion-base pair sequence of a human genome, T2T-CHM13, that includes gapless assemblies for all chromosomes except Y, corrects errors in the prior references, and introduces nearly 200 million base pairs of sequence containing 1956 gene predictions, 99 of which are predicted to be protein coding. The completed regions include all centromeric satellite arrays, recent segmental duplications, and the short arms of all five acrocentric chromosomes, unlocking these complex regions of the genome to variational and functional studies.

after processes they could be took



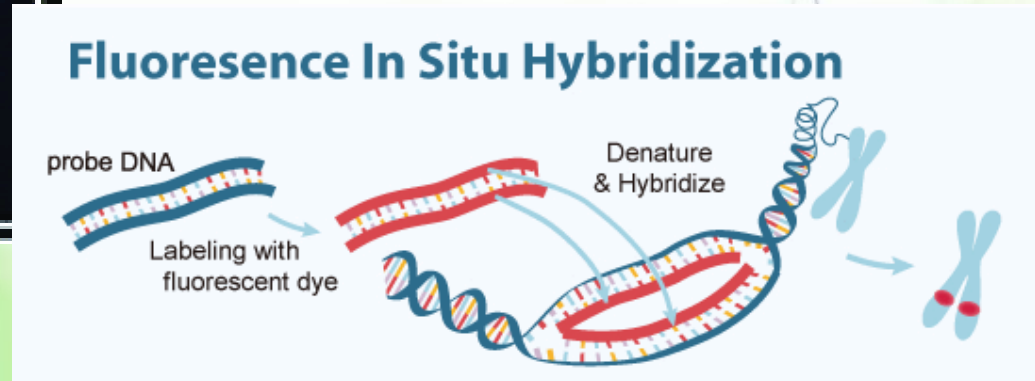
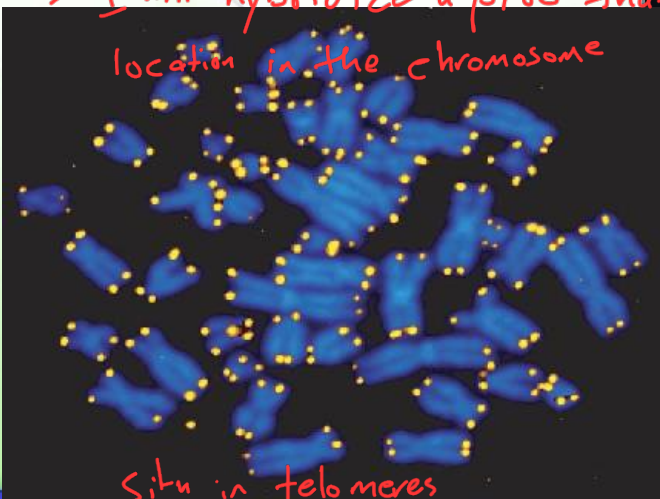
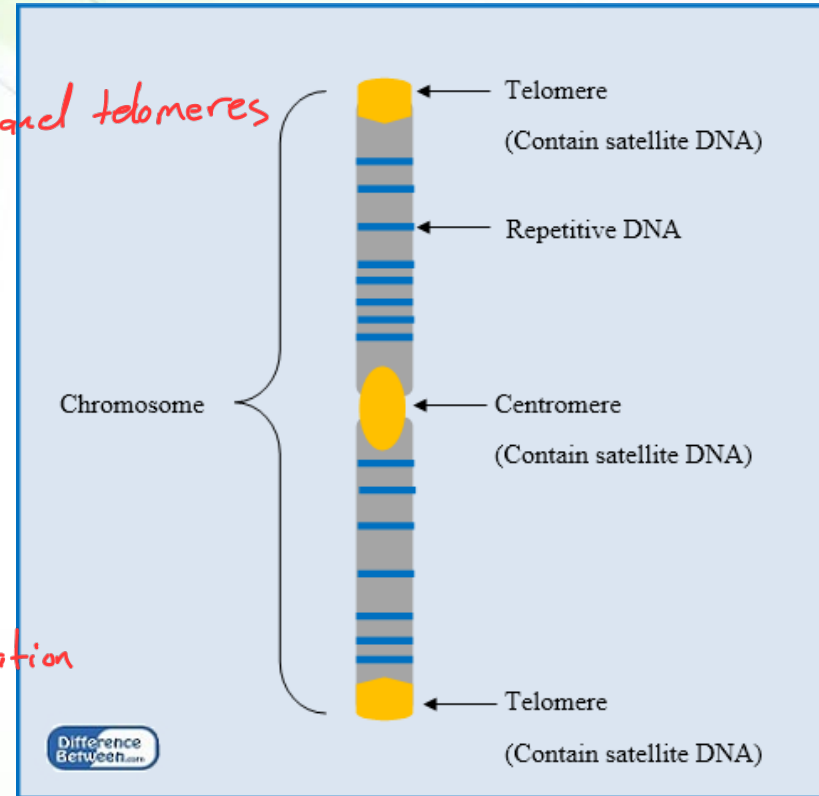
Tandem repeats

Satellite (macro-satellite) DNA



- Regions of 5-300 bp repeated 10^6 - 10^7 times *and found in centromeres and telomeres*
- Centromeres and telomeres
- Centromeric **A/T-rich repeats** (171 bp) called **α -satellite** **unique to each chromosome** (you make chromosome-specific probes) by **fluorescence in situ hybridization (FISH)**. *in specific location*

I will hybridize a prob that move to specific



Telomeric repeats

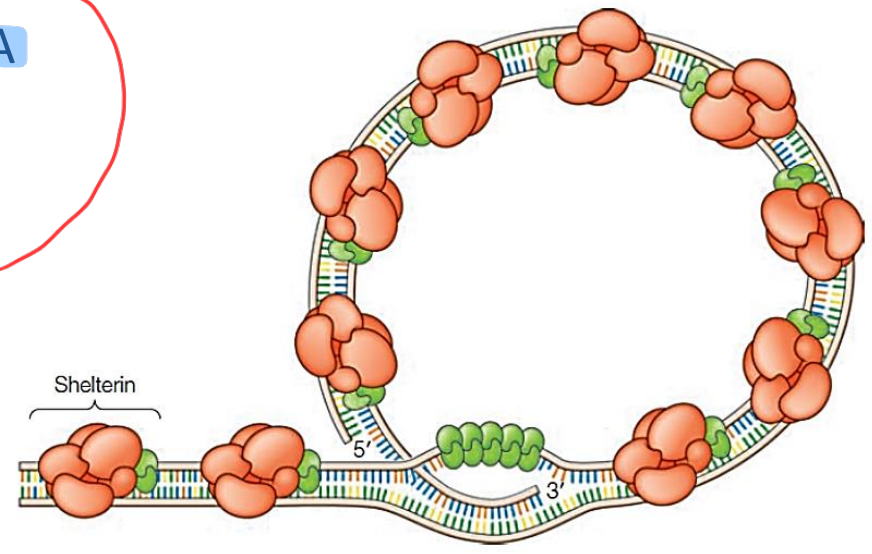


- (TTAGGG) is repeated hundreds ^{مردد} to thousands of times at the termini of human chromosomes with a 3' ^{result of elongation of strand} overhang of single-stranded DNA.

- The repeated sequences form loops that bind a protein complex called **shelterin**, which **protects the chromosome termini from degradation**.

- **Telomeric repeat-containing RNA (TERRA):** a **long non-coding RNA** transcribed from telomeres and **functions in:**

- maintaining the integrity of chromosome termini,
- regulating telomerase activity,
- maintaining the heterochromatic state of telomeres,
- protecting DNA from deterioration or fusion with neighboring chromosomes

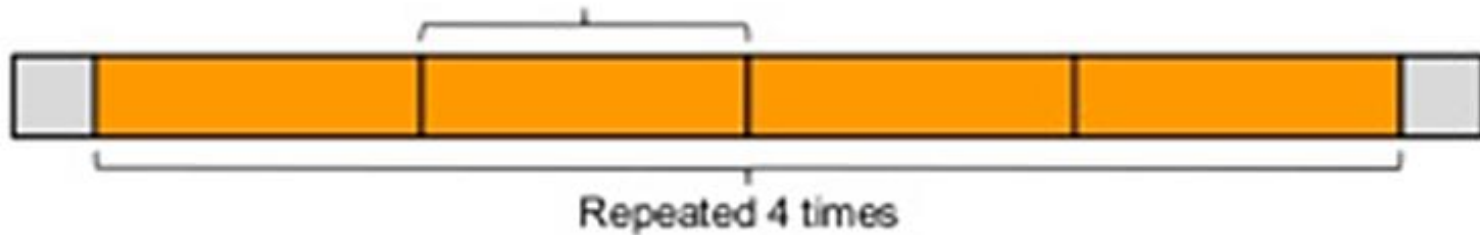


Mini- and Micro-satellite DNA



↙ another name ↘

Minisatellite: Variable Number Tandem Repeats (VNTR)



Microsatellite: Short Tandem Repeats (STR) – Simple Sequence Repeats (SSR)



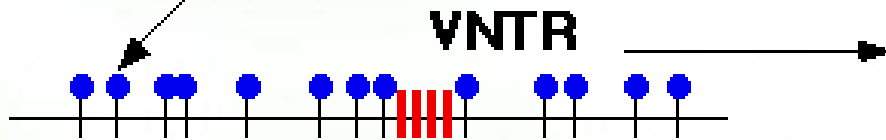
the only different is the size

Mini-satellite DNA

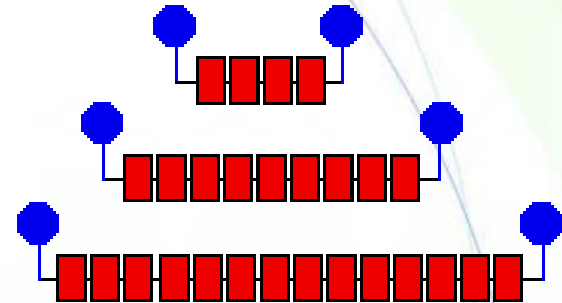


- Mini satellite sequences or VNTRs (variable number of tandem repeats) of 20 to 100 bp repeated 20-50 times

Restriction Endonuclease
Cutting Sites



Variable Number of
Tandem Repeats



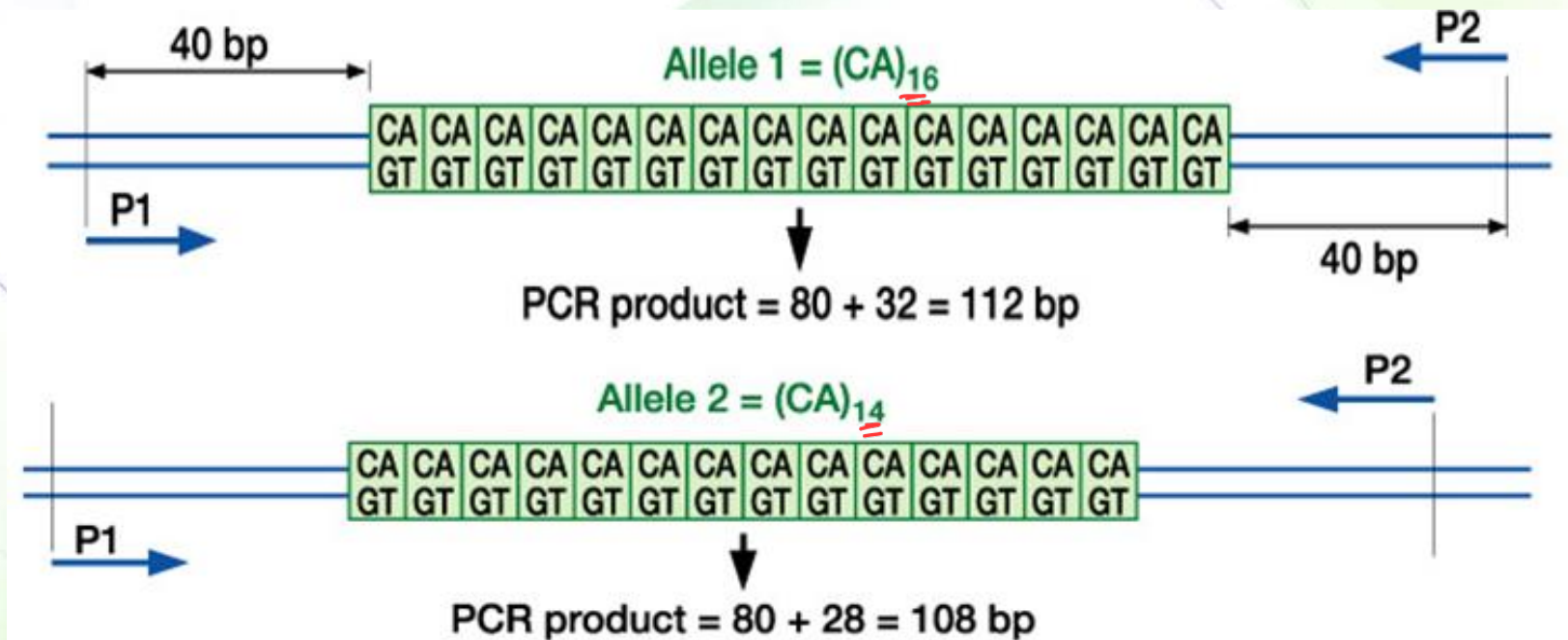
* So I have the same Repeat location with you but different sequence

Micro-satellite DNA



Same the mini-satellite \Rightarrow Same location but different sequence

- STRs (short tandem repeats) of 2 to 10 bp repeated 10-100 times



maybe I have it 16 times on one chromosome and 14 on another because it is diploid chromosomes

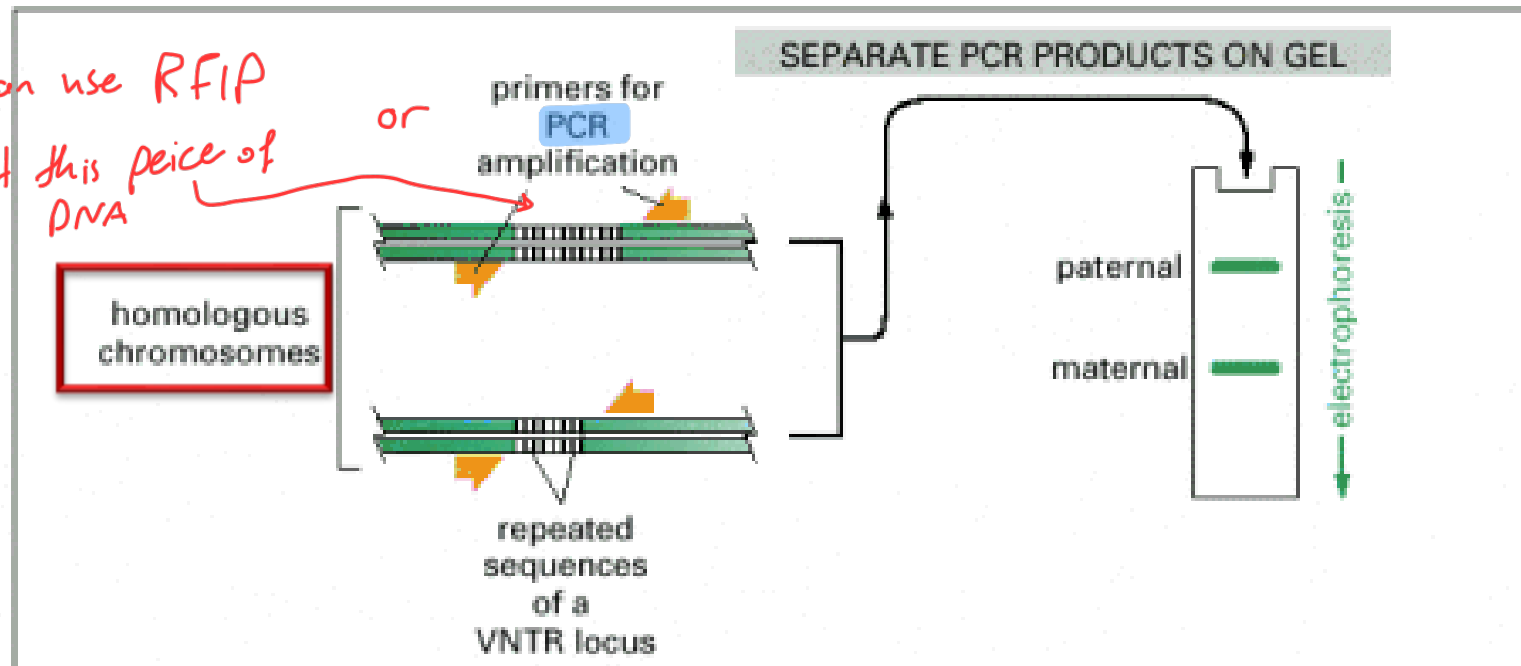
Polymorphisms of VNTR and STR



differences in DNA

- STRs and VNTRs are highly **variable among individuals** (polymorphic).
- They are useful in DNA profiling for forensic testing.

*we can use RFLP
to cut this piece of
DNA*



STRs and VNTRs as DNA Markers



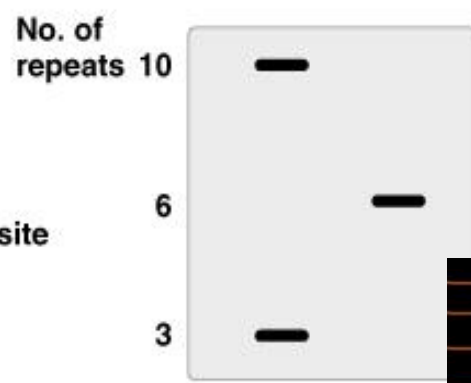
may be
RFLP then southern
blotting

or

PCR then
Gel electrophoresis

Cut with restriction enzyme
and analyze by gel electrophoresis,
Southern blotting, and probing with
a monolocus probe

↓ = Restriction site

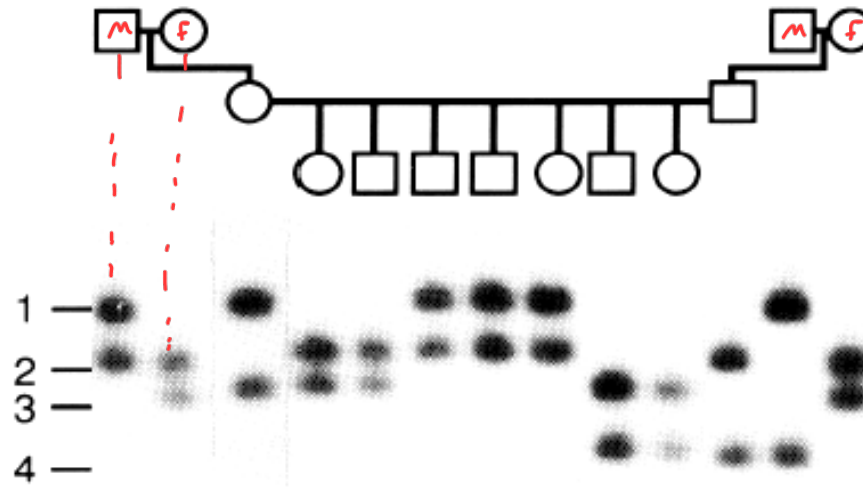


The likelihood of 2 unrelated individuals
having same allelic pattern is **extremely
improbable.**

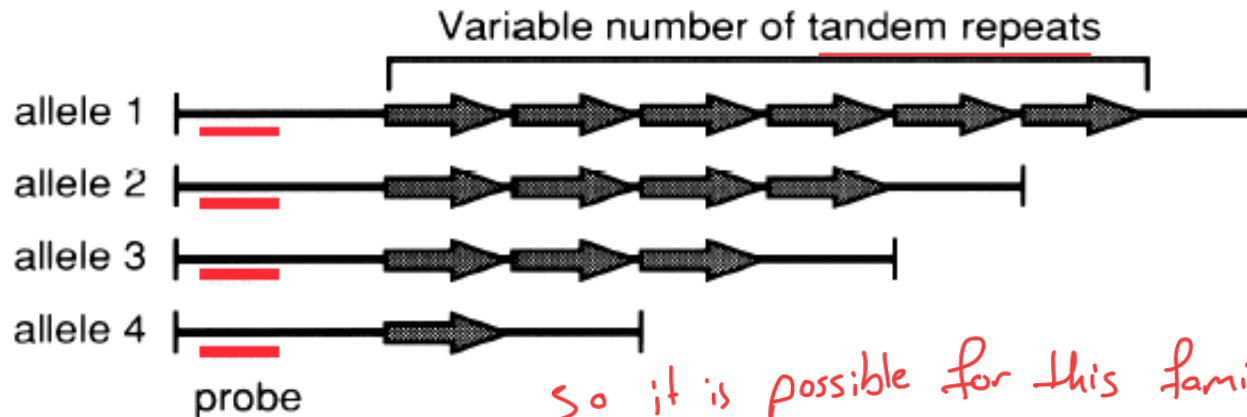


Real example

STR is more used than VNTR

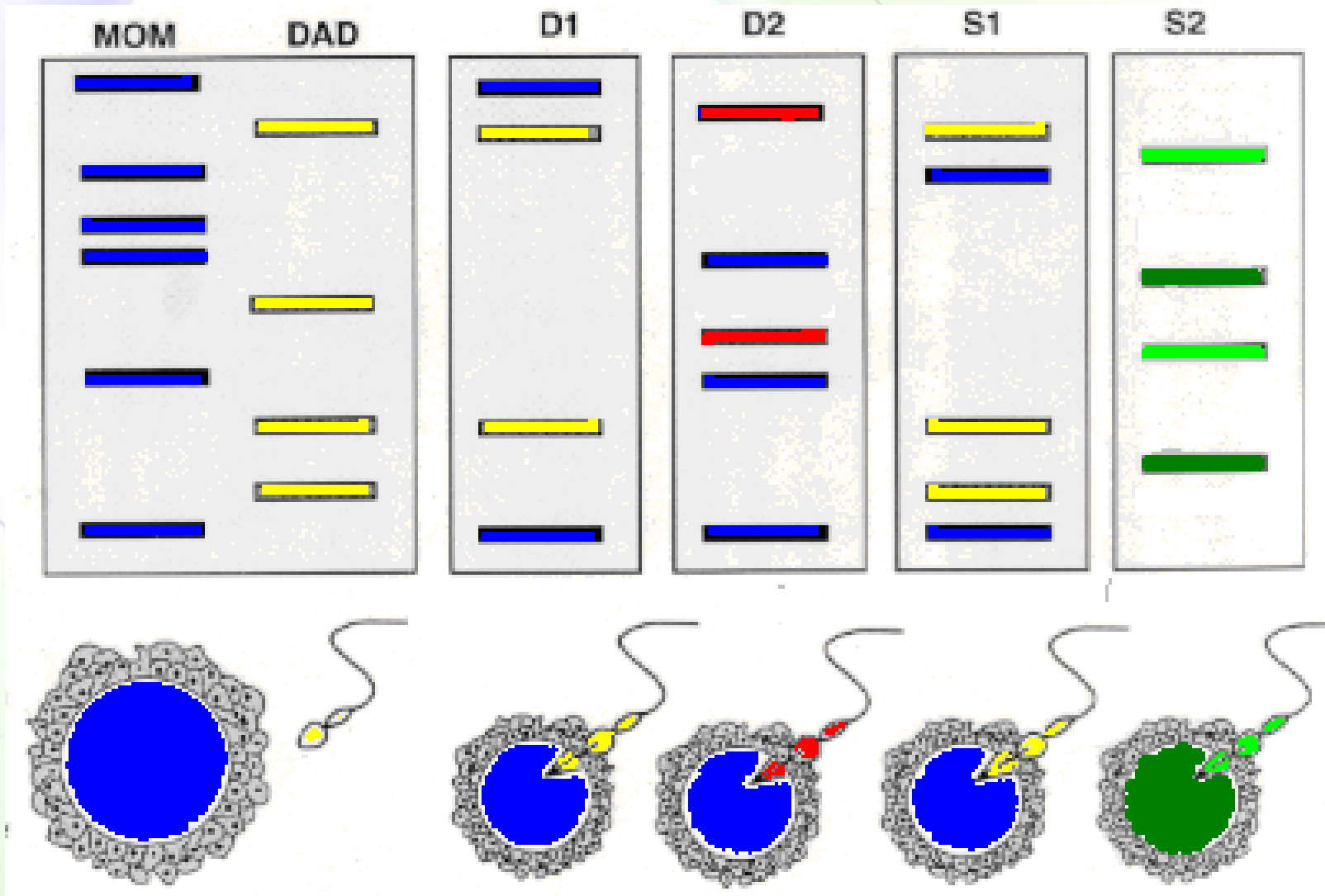


single-locus probe but multiple alleles



So it is possible for this family to repeat the tandem 1, 3, 4, 6 times

Paternity testing



Single nucleotide polymorphism (SNPs)



Another source of polymorphism

- Another source of genetic variation
- Single-nucleotide substitutions of one base for another
- Two or more versions of a sequence must each be present in at least one percent of the general population
- SNPs occur throughout the human genome - about **one in every 300 nucleotide** base pairs.
 - **~10 million SNPs** within the 3-billion-nucleotide human genome
 - **Only 500,000 SNPs are thought to be relevant**

that connected with disease or something important

Examples



| | Homozygous SNP | | Heterozygous SNP | | |
|--------------------------|-------------------------------|----------|-------------------------------|----------|-----------|
| Paternal allele | AACTGGACTT | G | AAGCATCTACGTT | A | TCCATGAAG |
| Maternal allele | AACTGGACTT | G | AAGCATCTACGTT | C | TCCATGAAG |
| Frequency in population: | G 51% T 49% (minor allele) | | A 90% C 10% (minor allele) | | |

the percent of of SNPs must be more 1% of population otherwise we call it mutation

2 strands for 1 chromosome

Heterozygous SNPs

Individual 1

Chr 2 copy1 ...CGATATTCC**T**ATCGAATGTC...
 ...GCTATAAGG**A**TAGCTTACAG...
 Chr 2 copy2 ...CGATATTCC**C**ATCGAATGTC...
 ...GCTATAAGG**G**TAGCTTACAG...
Both strands for Both chromosome

Individual 4

Chr 2 copy1 ...CGATATTCC**T**ATCGAATGTC...
 ...GCTATAAGG**A**TAGCTTACAG...
 Chr 2 copy2 ...CGATATTCC**C**ATCGAATGTC...
 ...GCTATAAGG**G**TAGCTTACAG...

Homozygous SNPs

Individual 2

Chr 2 copy1 ...CGATATTCC**C**ATCGAATGTC...
 ...GCTATAAGG**G**TAGCTTACAG...
 Chr 2 copy2 ...CGATATTCC**C**ATCGAATGTC...
 ...GCTATAAGG**G**TAGCTTACAG...

Individual 5

Chr 2 copy1 ...CGATATTCC**C**ATCGAATGTC...
 ...GCTATAAGG**G**TAGCTTACAG...
 Chr 2 copy2 ...CGATATTCC**T**ATCGAATGTC...
 ...GCTATAAGG**A**TAGCTTACAG...

Individual 3

Chr 2 copy1 ...CGATATTCC**T**ATCGAATGTC...
 ...GCTATAAGG**A**TAGCTTACAG...
 Chr 2 copy2 ...CGATATTCC**T**ATCGAATGTC...
 ...GCTATAAGG**A**TAGCTTACAG...

Individual 6

Chr 2 copy1 ...CGATATTCC**C**ATCGAATGTC...
 ...GCTATAAGG**G**TAGCTTACAG...
 Chr 2 copy2 ...CGATATTCC**T**ATCGAATGTC...
 ...GCTATAAGG**A**TAGCTTACAG...

Categories of SNPs

do not forget we talk about 500k SNPs

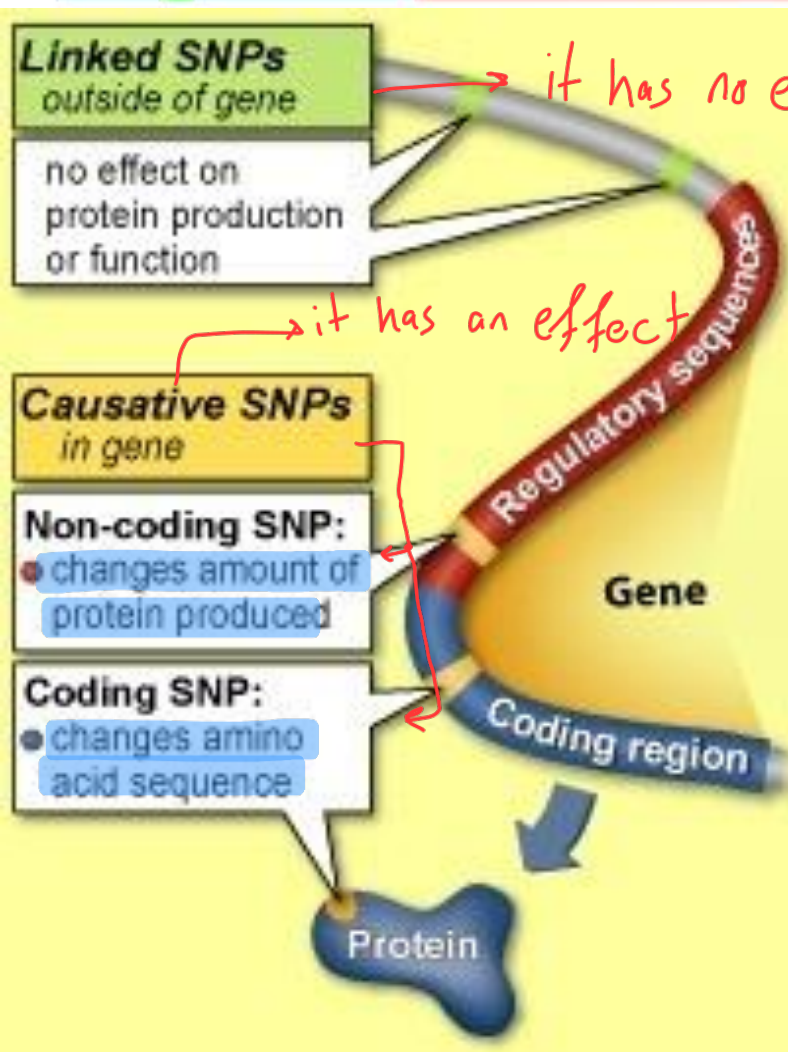


Linked SNPs

Causative SNPs

Regulatory Sequence

Coding Region

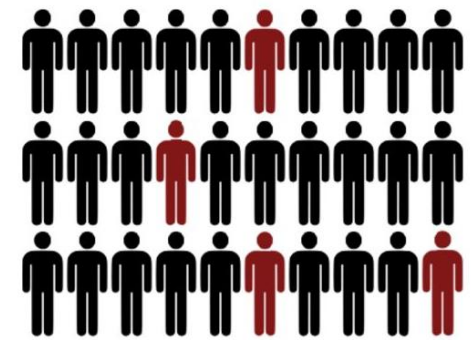
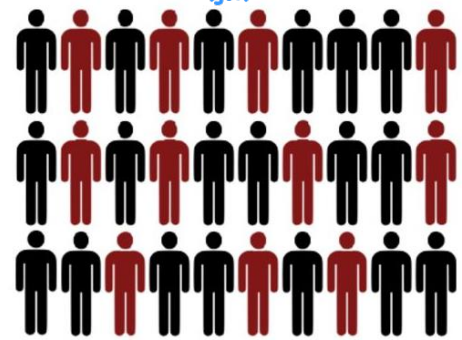


it has no effect But linked with phenotype

but if I removed this gene the phenotype will Not removed

Some of them has C but little

So some of them has T but little



healthy

TTGGCCAGCTGGACGAGGGGCGATGAC

TTGGCCAGCTGGATGAGGGGCGATGAC

that does NOT mean T is the reason of the disease



Interspersed repeats

Repeated Regions but with spaces between them
(مجاورة)

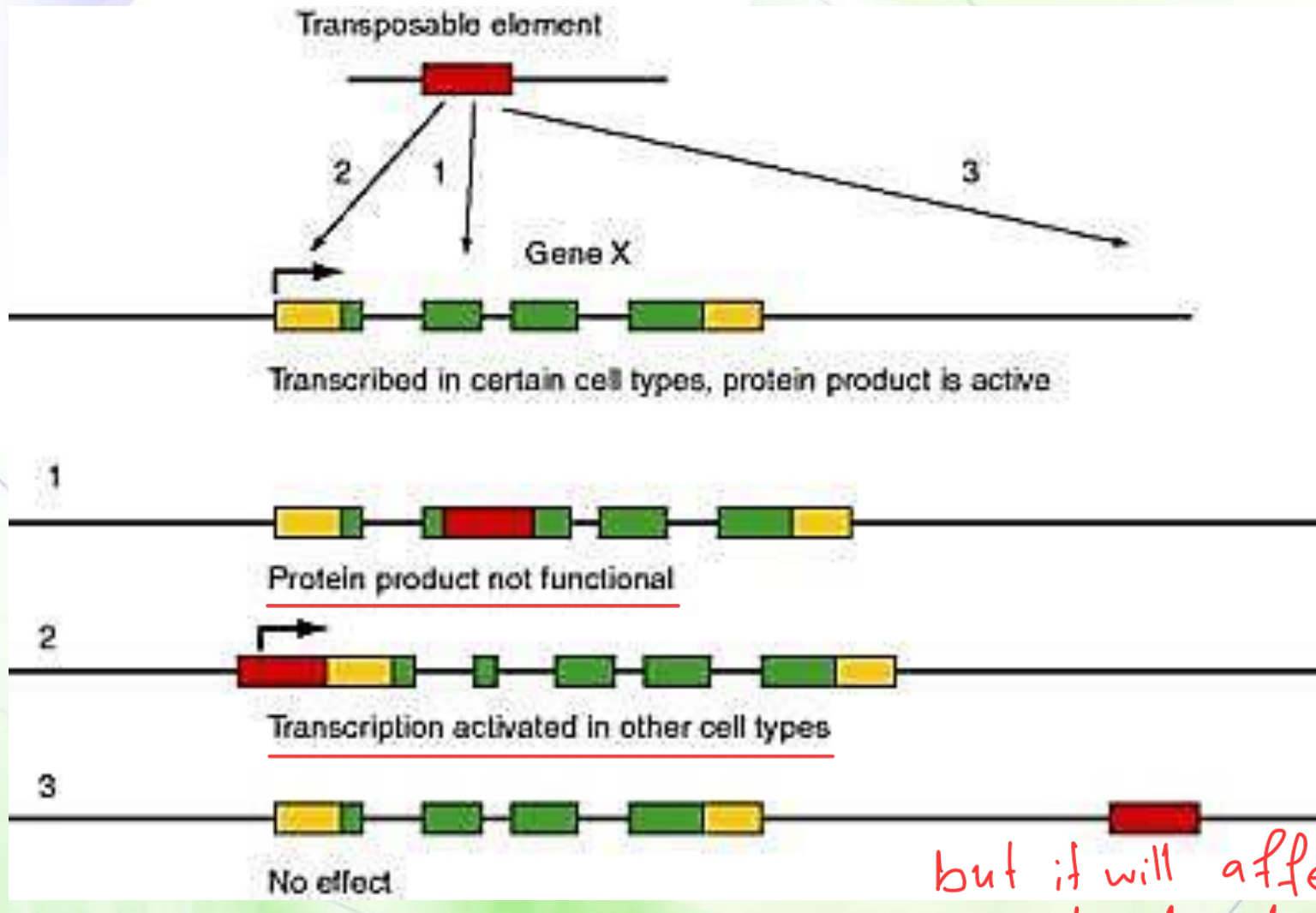
Transposons (jumping genes)



Just another name

↳ can move (but in human it mostly lost ability to move)

- They are segments of DNA that can move from their original position in the genome to a new location.
- Two classes:
 - DNA transposons (3% of human genome) *from DNA of viruses*
 - RNA transposons or retrotransposons (42% of human genome).
 - Long interspersed elements (LINEs, 21%)
 - Short interspersed elements (SINEs, 13%)
 - An example is Alu (300 bp) *↳ an restriction endo clease called Alu can cut this element*
 - Retrovirus-like elements (8%)
- Over 99% of the transposons in the human genome lost their ability to move, but we still have some active transposable elements that can sometimes cause disease.
 - Hemophilia A and B, severe combined immunodeficiency, porphyria, predisposition to cancer, and Duchenne muscular dystrophy.



but it will affect other non coding elements