# Approaches to Gene Discovery

Bruce R. Korf, MD, PhD

- The Human Genome

- Genetic Variation

- Gene Identification

# Gene Regulation

# Splicing
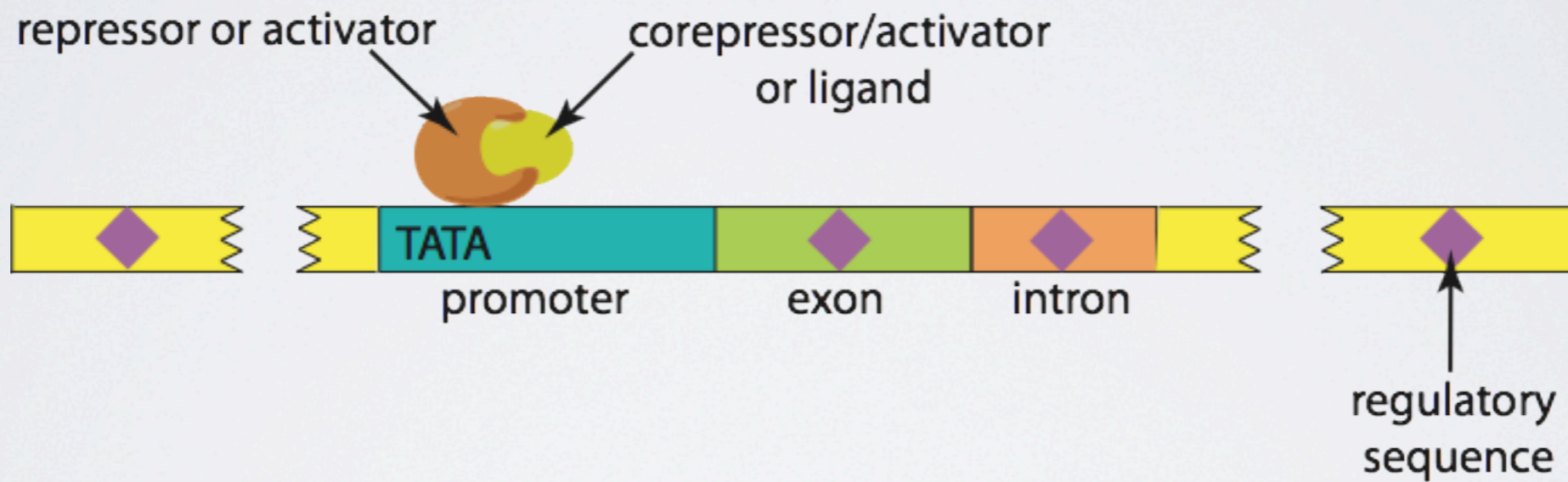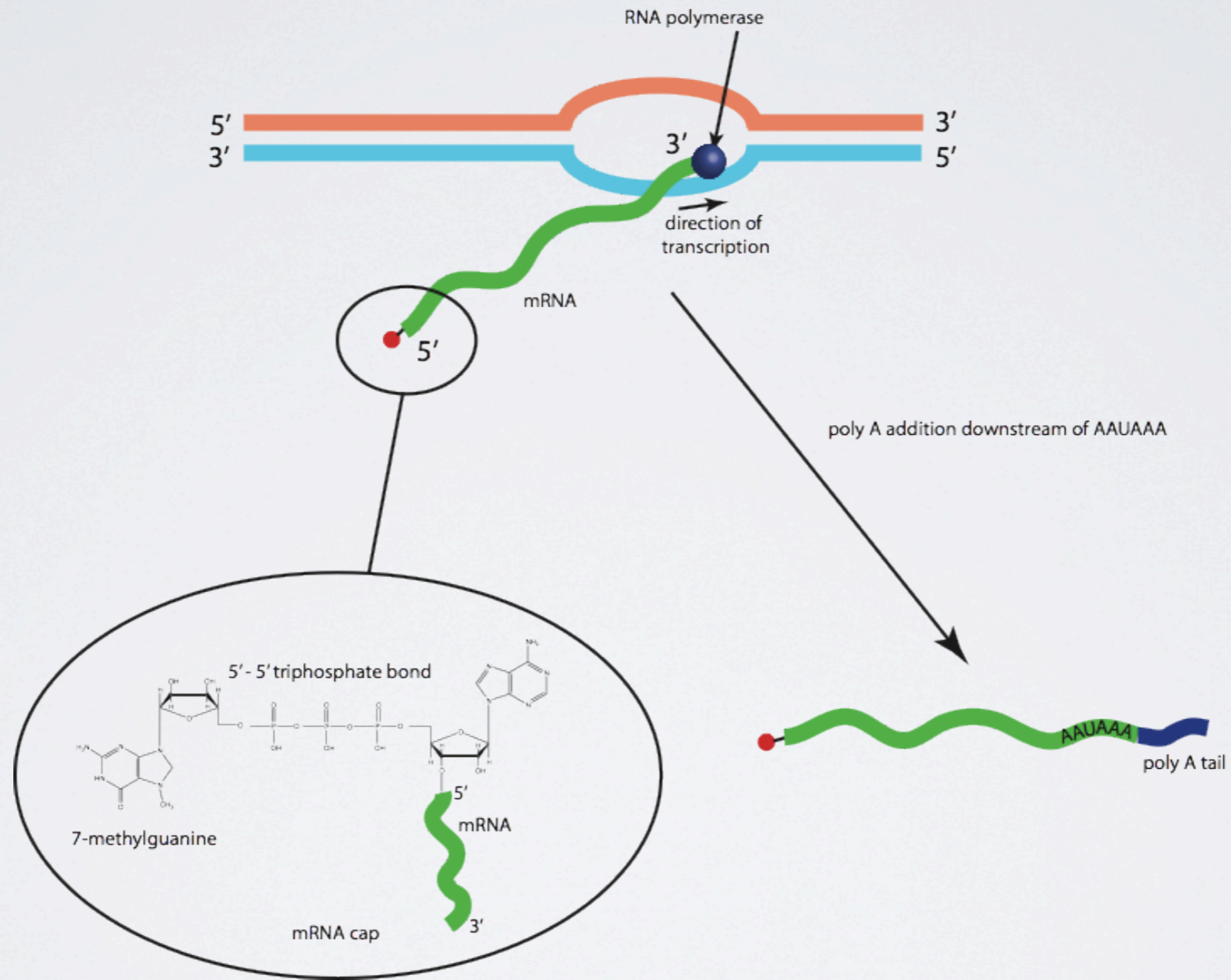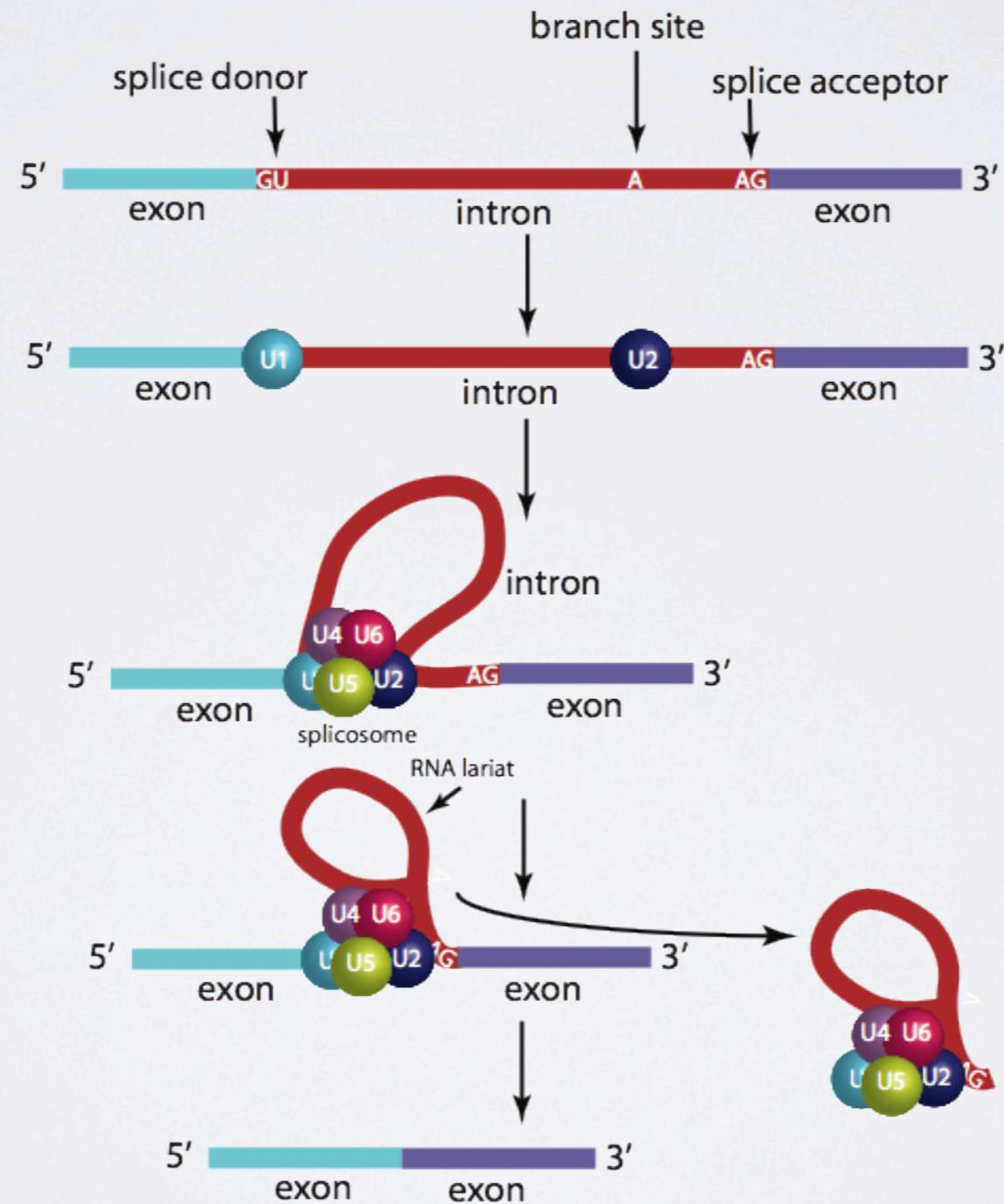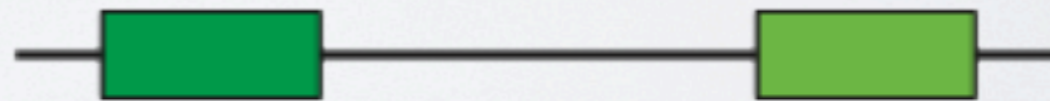
# Repeated Sequences

simple sequence repeat

...GCGACACACACACACACAGT...

variable number tandem repeat
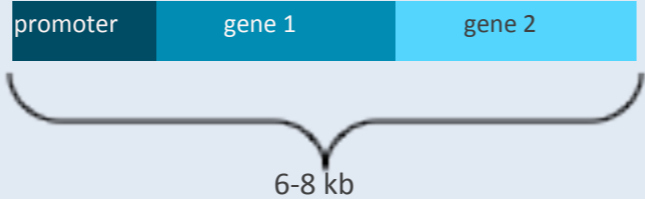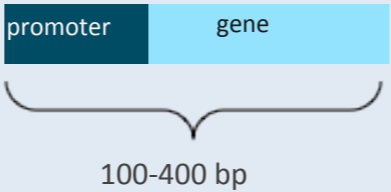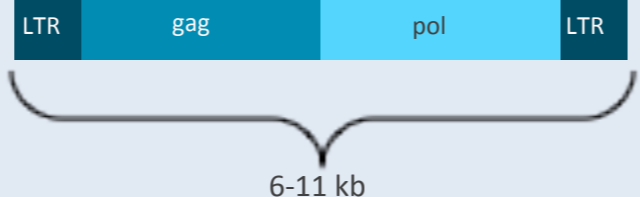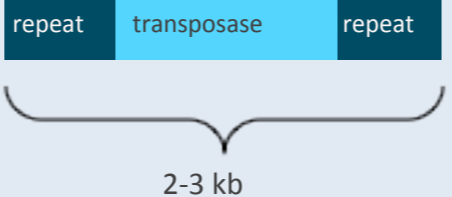
14-100 base pair repeat unit

highly repeated sequences at
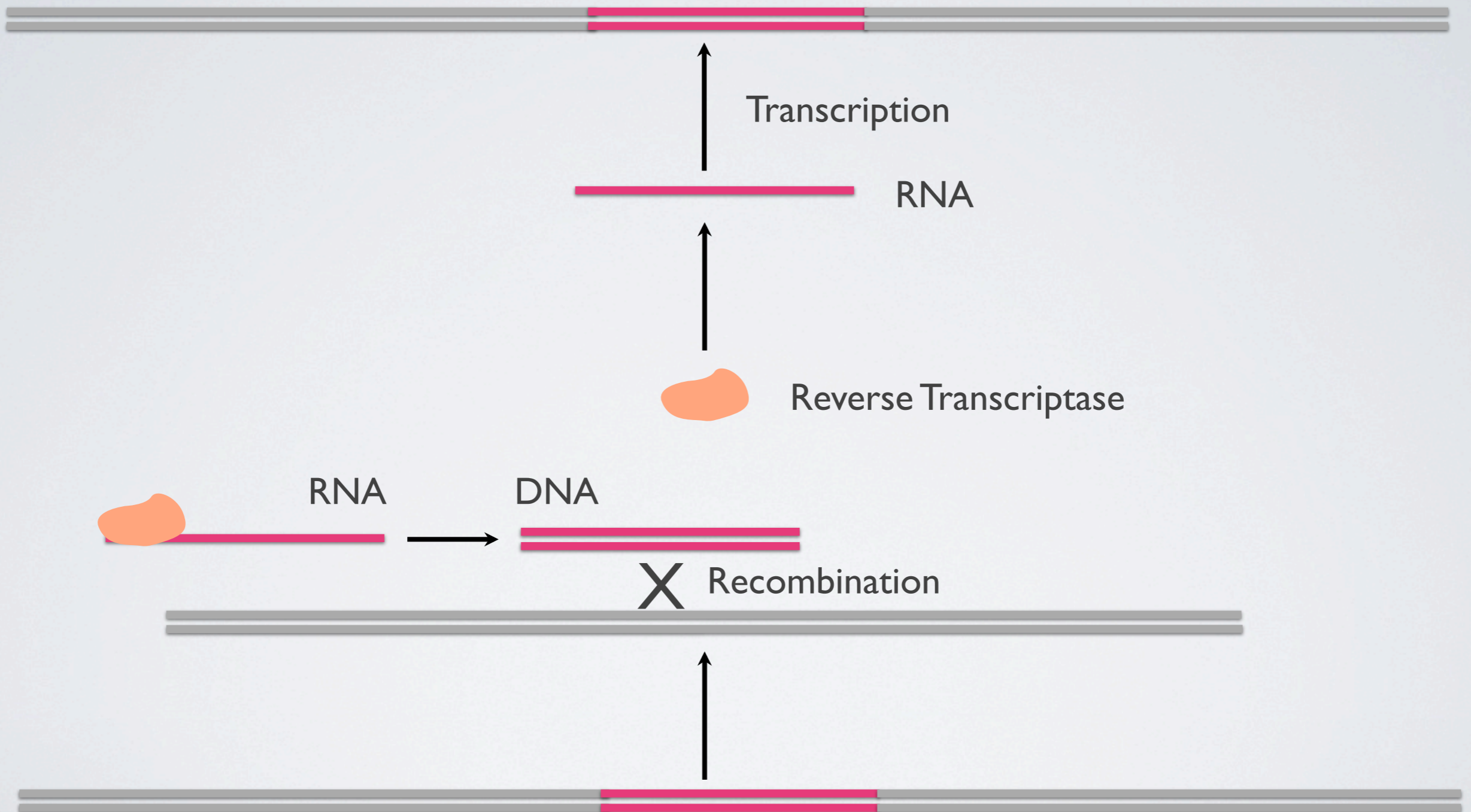centromeric and subtelomeric regions

segmental duplications

# Transposable Genetic Elements

| Type | Structure | Copy Number | Percent |
|------|-----------|-------------|---------|
| LINE | promoter / gene 1 / gene 2 — 6-8 kb | 850,000 | 21 |
| SINE | promoter / gene — 100-400 bp | 1,500,000 | 13 |
| Retroviral-like | LTR / gag / pol / LTR — 6-11 kb | 450,000 | 8 |
| Transposon | repeat / transposase / repeat — 2-3 kb | 300,000 | 3 |

# LINE "Life Cycle"

Transcription

RNA

Reverse Transcriptase

RNA          DNA

X Recombination

# ENCODE project

- annotated 20,687 protein-encoding genes

- average 6.3 alternatively spliced isoforms per gene

- 8,801 small RNAs; 9,640 long non-coding transcripts

- >80% genome transcribed in some cell type
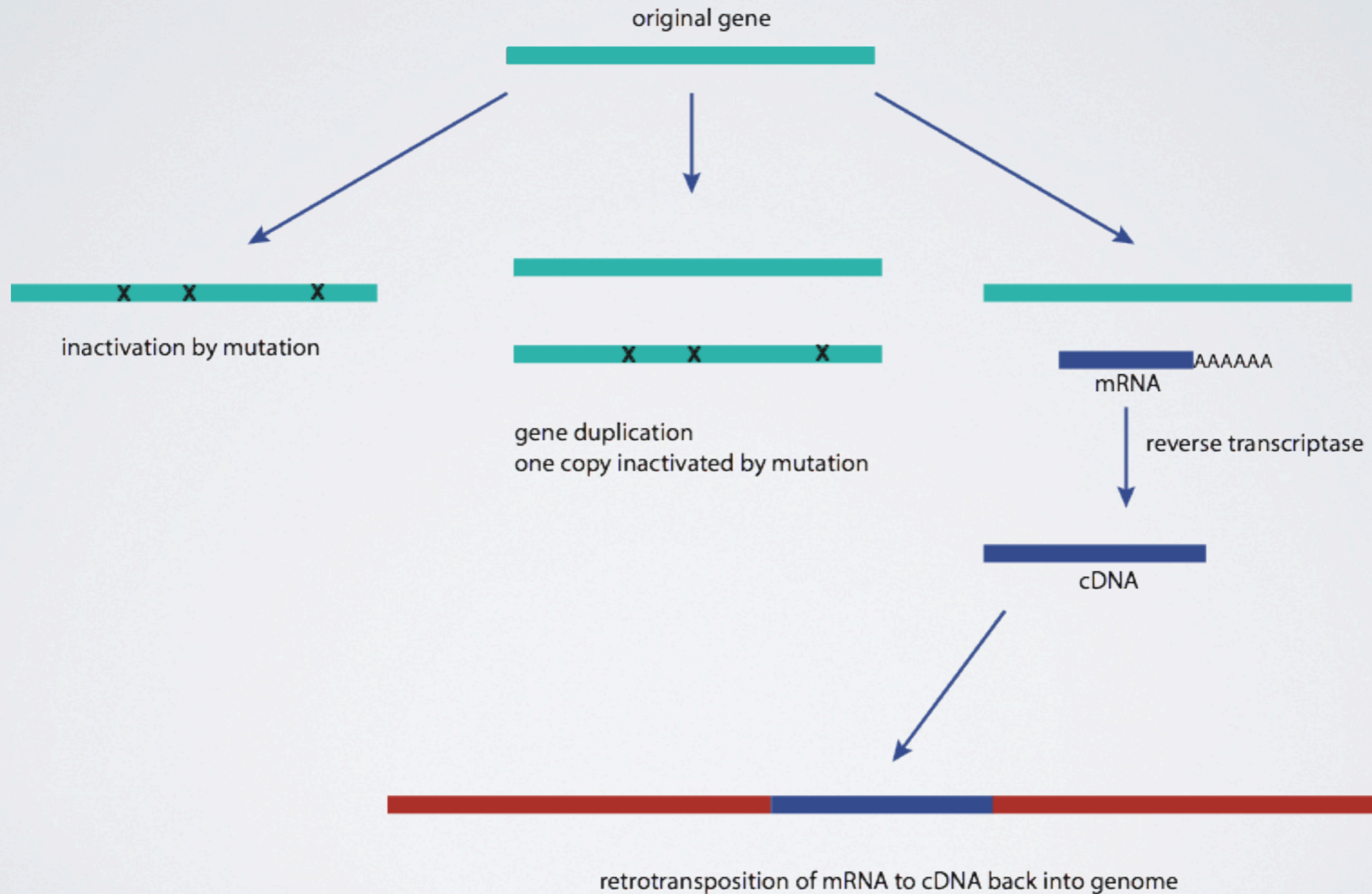
- >400,000 enhancers and 70,000 promoters

# Non-Coding RNAs

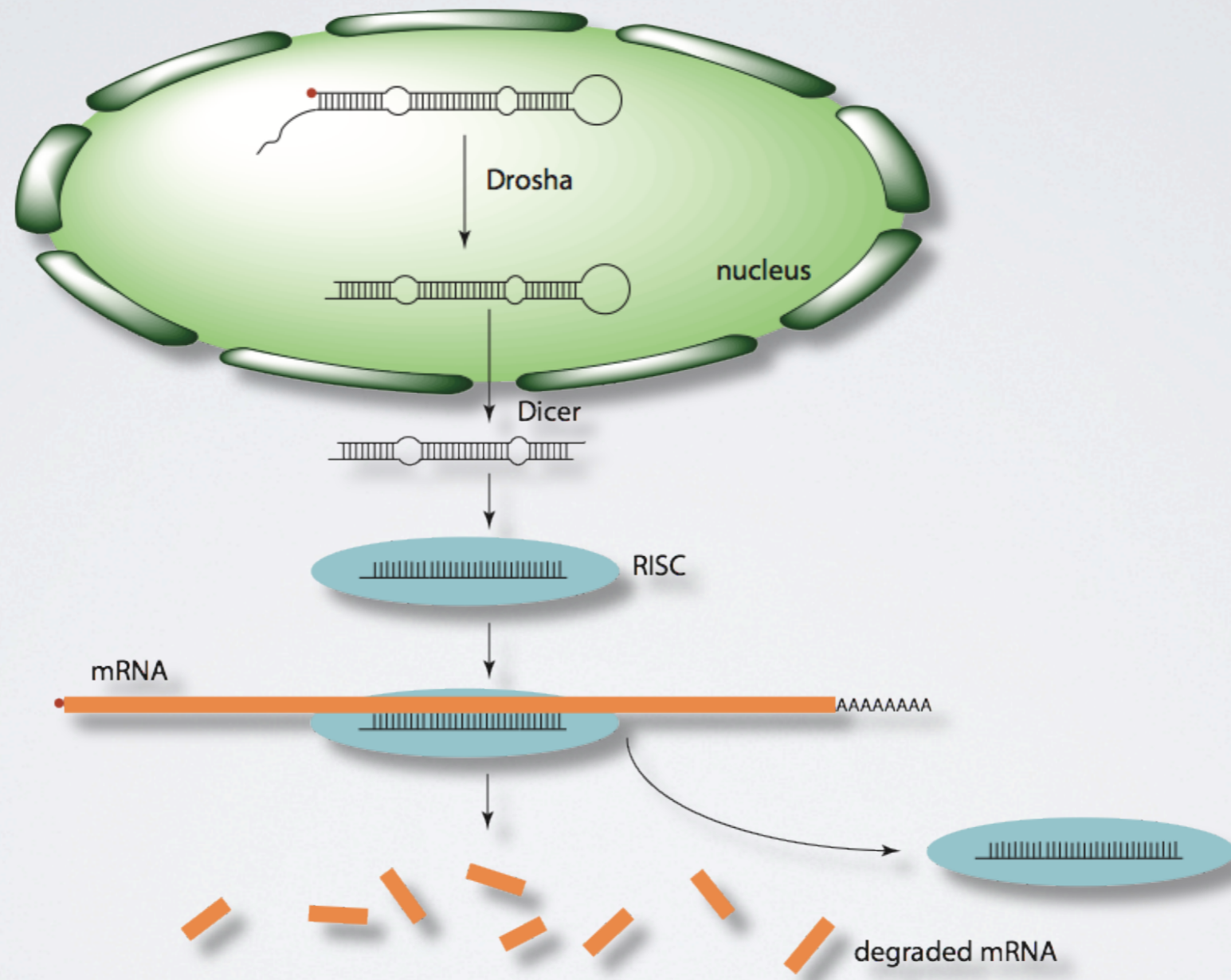| | | |
|---|---|---|
| tRNA | transfer RNA | protein synthesis |
| rRNA | ribosomal RNA | protein synthesis |
| snRNA | small nuclear RNA | splicing |
| snoRNA | small nucleolar RNA | RNA modification |
| miRNA | micro RNA | gene regulation |
| siRNA | small interfering RNA | viral defense |
| lncRNA | long non-coding RNA | gene regulation/unknown |

# Long Non-Coding RNAs

- antisense

- intergenic

- sense overlapping
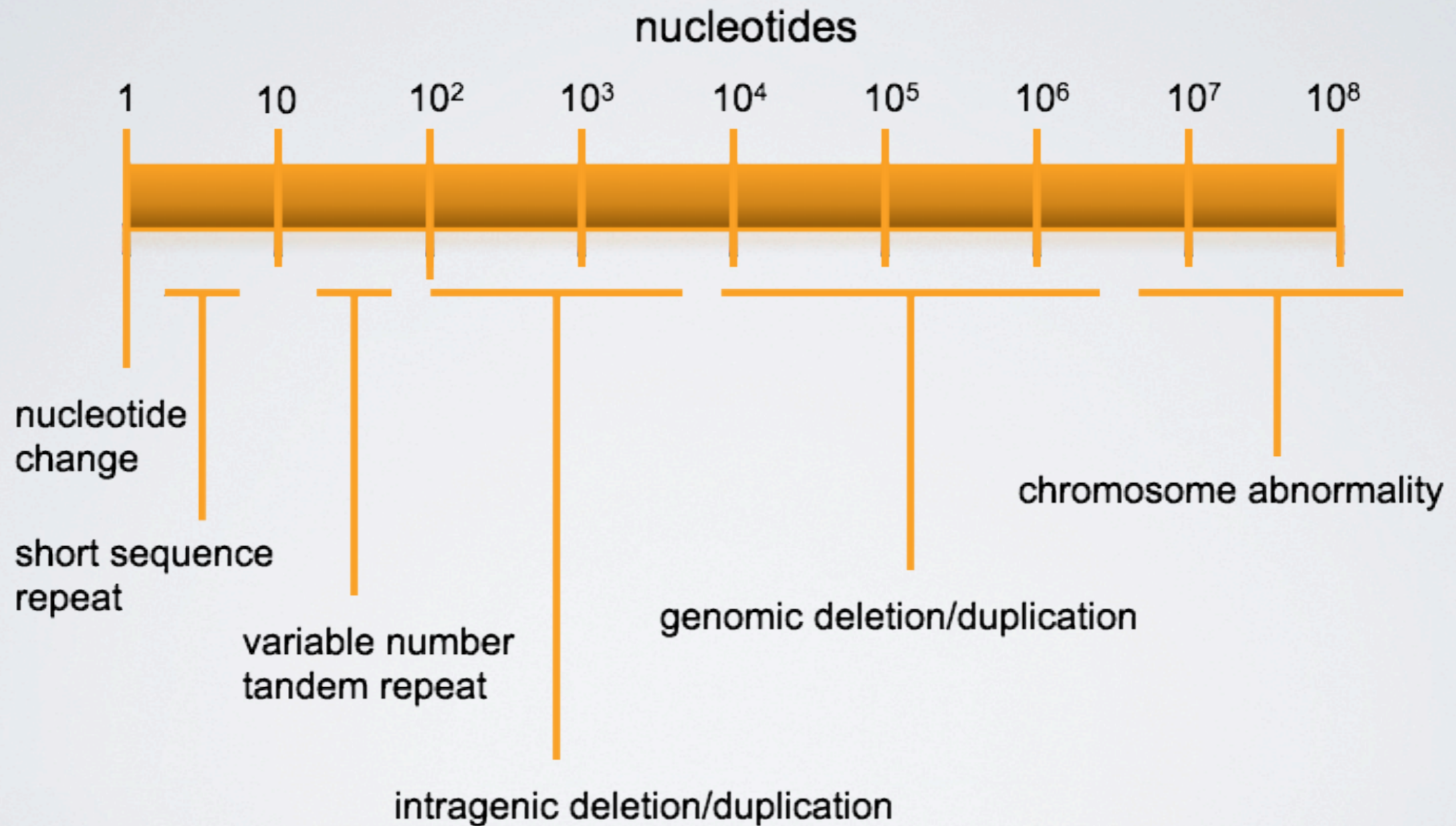
- sense intronic

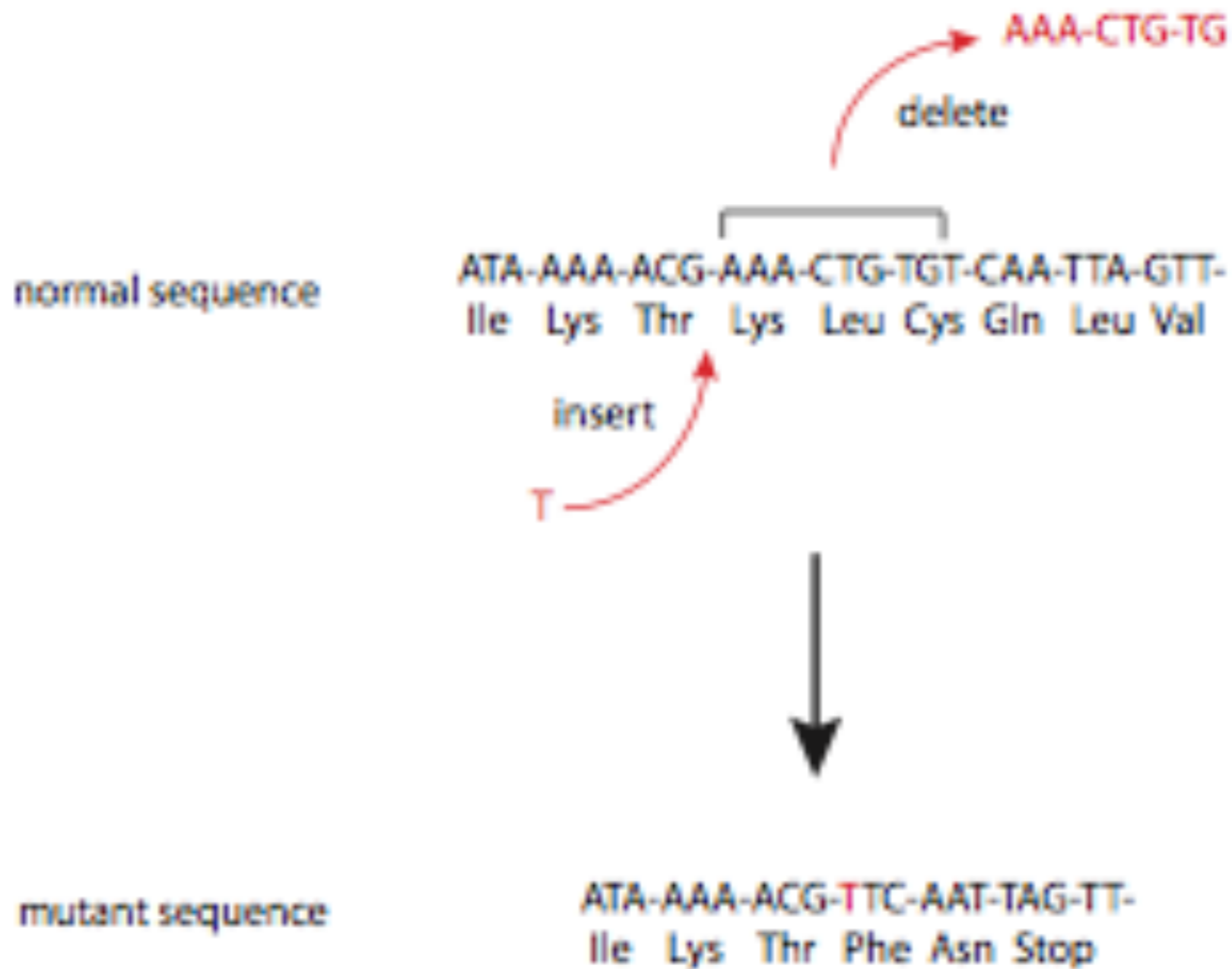- processed transcript
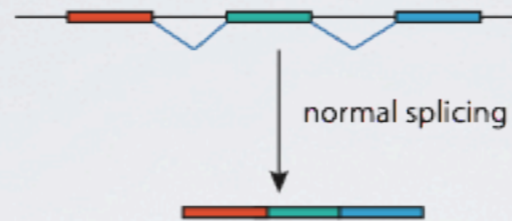
# Pseudogenes



original gene

inactivation by mutation

gene duplication
one copy inactivated by mutation

mRNA        AAAAAA

reverse transcriptase

cDNA

retrotransposition of mRNA to cDNA back into genome

# MicroRNA

# Genetic Variation

# Point Mutations

TCC CAA ATC GTC CCT CGA GTT
  ser  gln  ile  val  pro  arg  val
                                                wild type sequence

TCC CAG ATC GTC CCT CGA GTT
  ser  gln  ile  val  pro  arg  val
                                                silent mutation

TCC CAA ATC CTC CCT CGA GTT
  ser  gln  ile  leu  pro  arg  val
                                                conservative mutation

TCC CAA ATC GTC GCT CGA GTT
  ser  gln  ile  val  ala  arg  val
                                                non-conservative mutation

TCC CAA ATC GTC CCT TGA GTT
  ser  gln  ile  val  pro  stop
                                                stop mutation

TCC CAG AAT CGT CCC TCG AGT T
  ser  gln  asn  arg  pro  ser  ser
                                                frameshift mutation
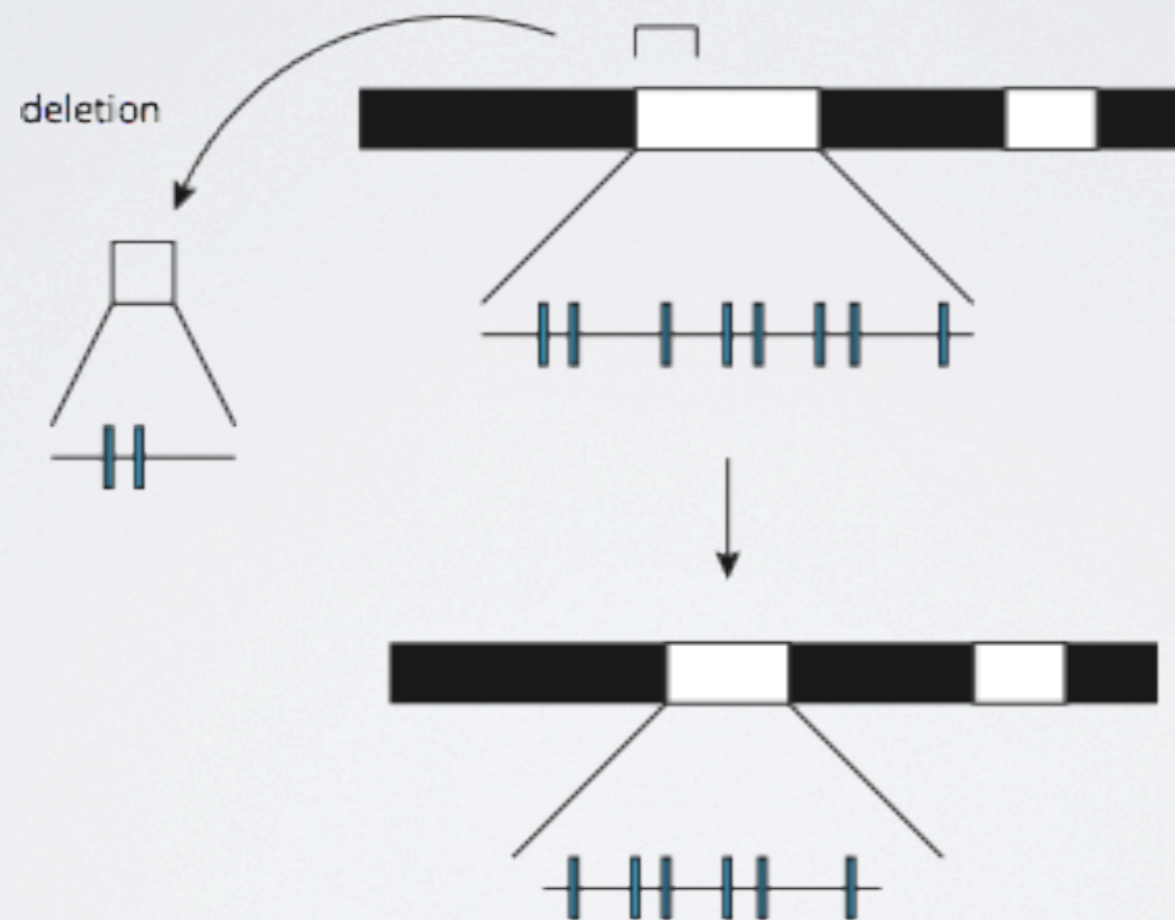
# Indel

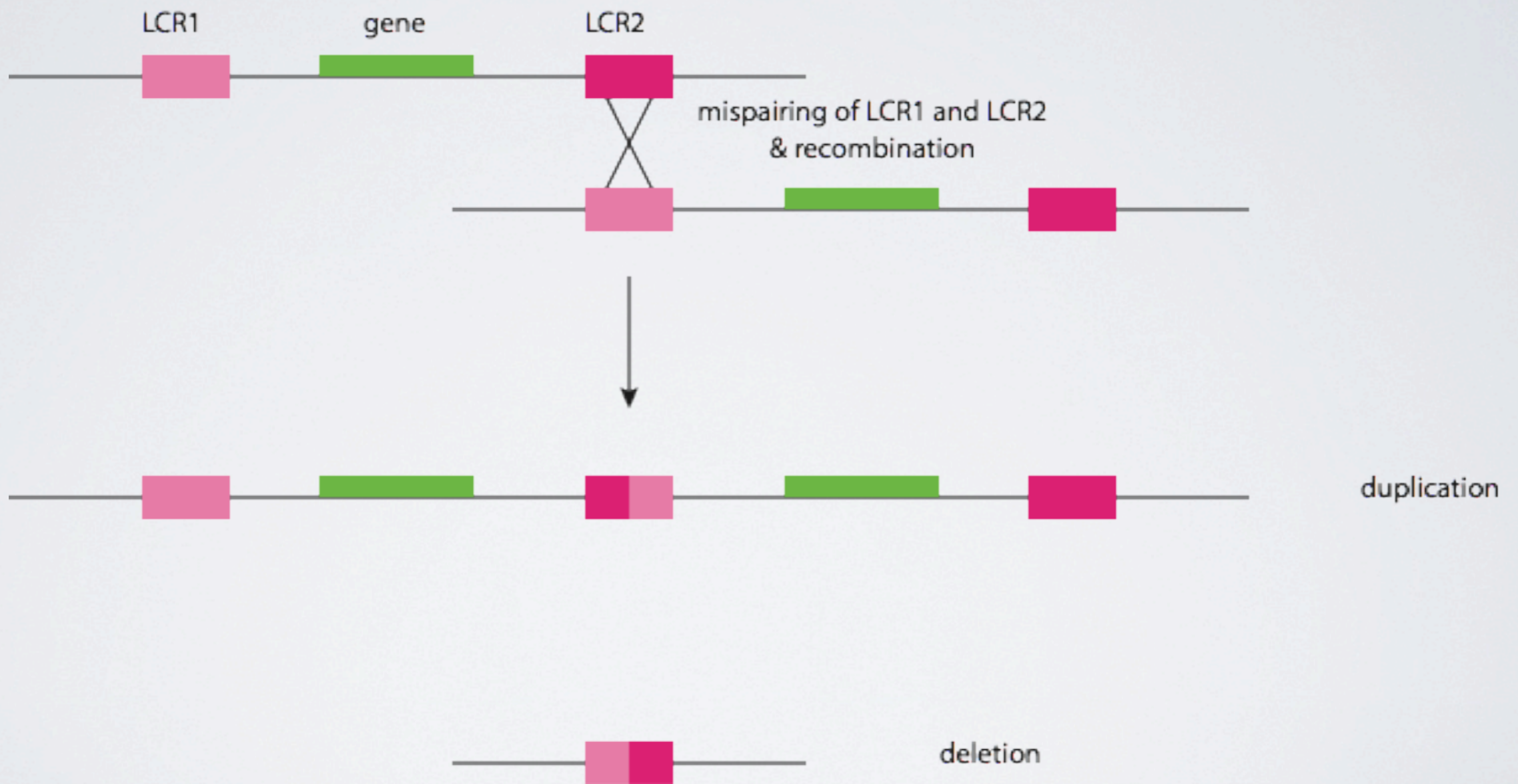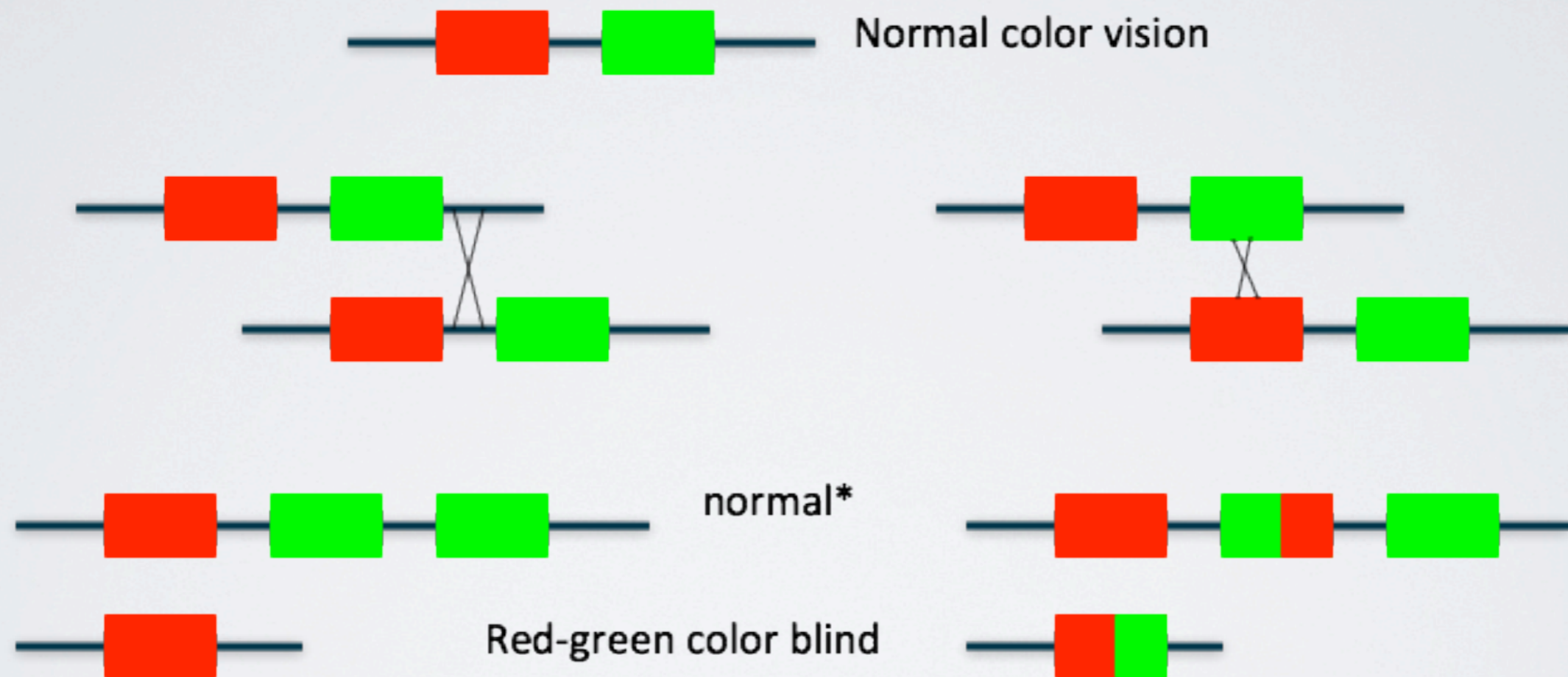# Splicing Mutations

# Triplet Repeat Expansions

# Multiexon Deletion

# Chromosome Microdeletion



deletion

# LCR Mispairing

# Red-Green Color Blindness



Normal color vision

normal*
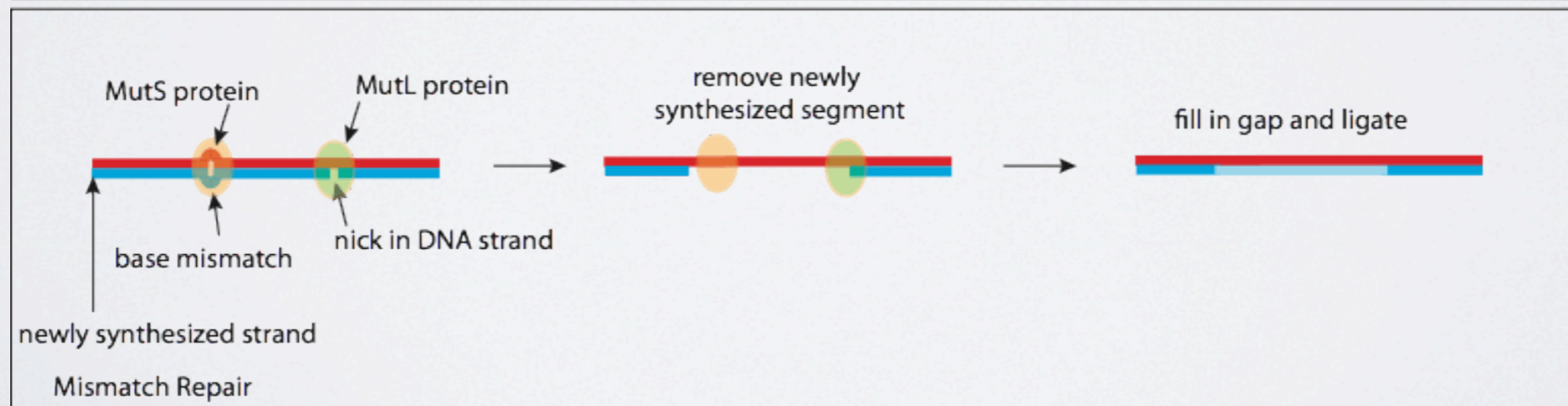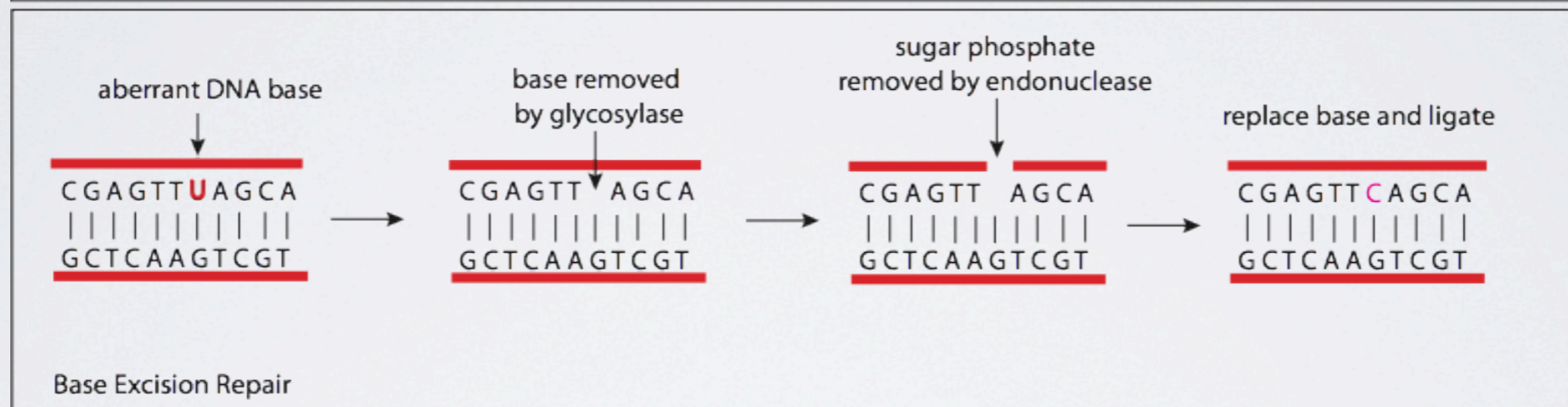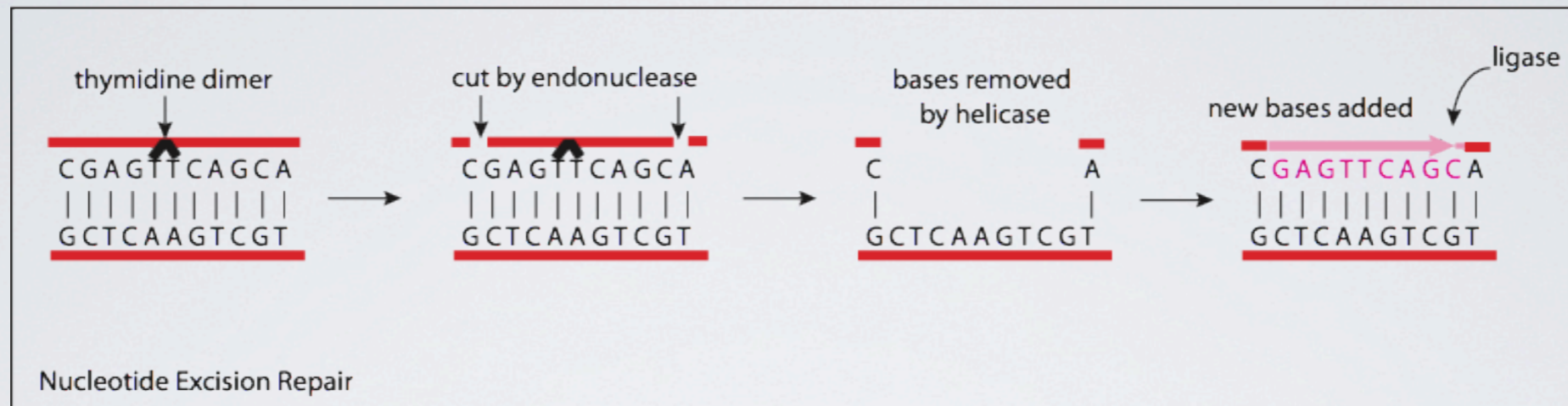
Red-green color blind

\* Color vision may be abnormal if green gene not expressed

# DNA Repair

# Frequency of Mutation

## A map of human genome variation from population–scale sequencing

The 1000 Genomes Project Consortium*

The 1000 Genomes Project aims to provide a deep characterization of human genome sequence variation as a foundation for investigating the relationship between genotype and phenotype. Here we present results of the pilot phase of the project, designed to develop and compare different strategies for genome-wide sequencing with high-throughput platforms. We undertook three projects: low-coverage whole-genome sequencing of 179 individuals from four populations; high-coverage sequencing of two mother–father–child trios; and exon-targeted sequencing of 697 individuals from seven populations. We describe the location, allele frequency and local haplotype structure of approximately 15 million single nucleotide polymorphisms, 1 million short insertions and deletions, and 20,000 structural variants, most of which were previously undescribed. We show that, because we have catalogued the vast majority of common variation, over 95% of the currently accessible variants found in any individual are present in this data set. On average, each person is found to carry approximately 250 to 300 loss-of-function variants in annotated genes and 50 to 100 variants previously implicated in inherited disorders. We demonstrate how these results can be used to inform association and functional studies. From the two trios, we directly estimate the rate of *de novo* germline base substitution mutations to be approximately $10^{-8}$ per base pair per generation. We explore the data with regard to signatures of natural selection, and identify a marked reduction of genetic variation in the neighbourhood of genes, due to selection at linked sites. These methods and public data will support the next phase of human genetic research.

If there are $10^8$ sperm per ejaculate, in principle every base
could be mutated in at least one sperm cell and each germ cell
has around 10 mutations
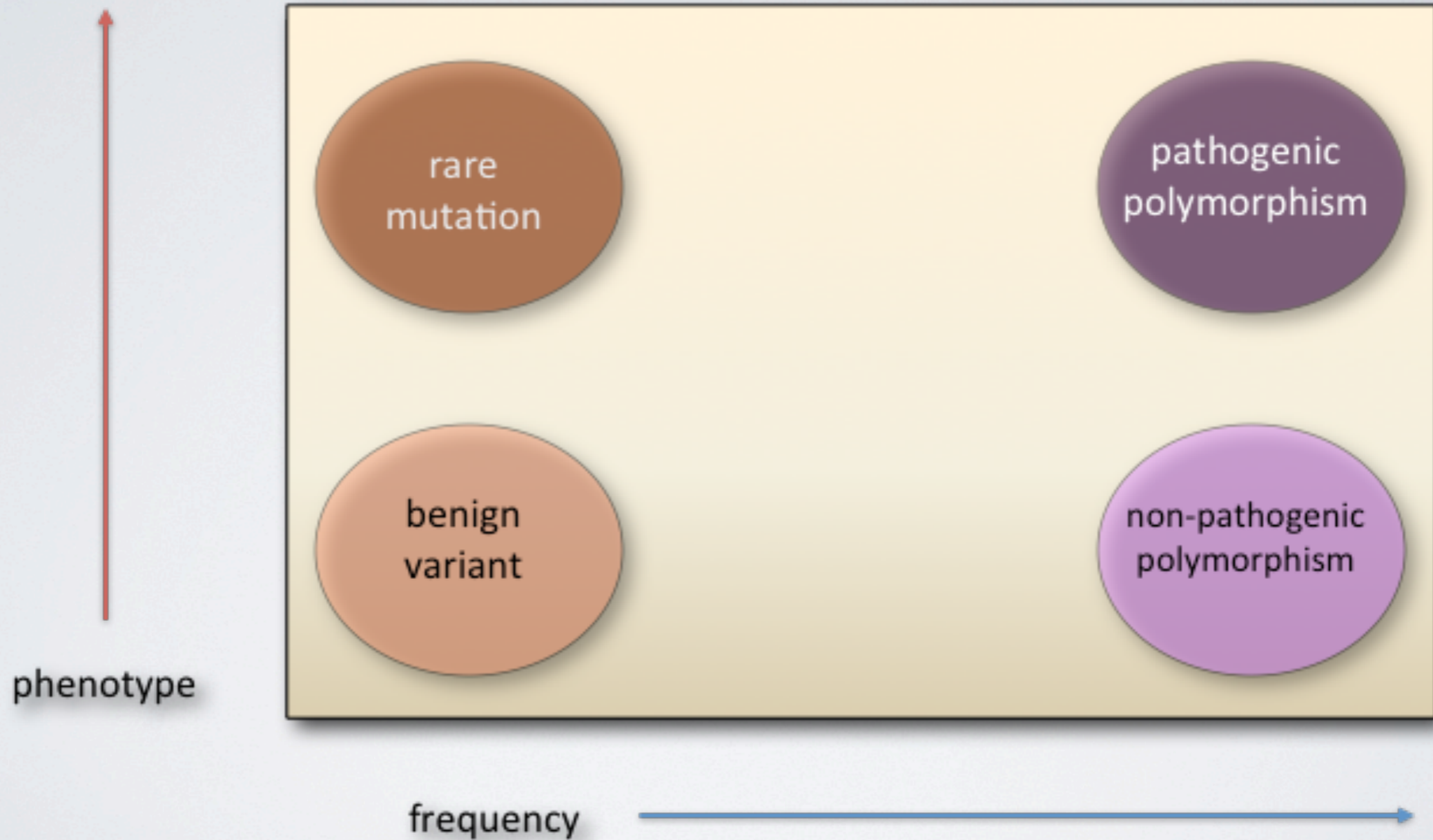
# Human Mendelian Phenotypes

**OMIM Entry Statistics:**
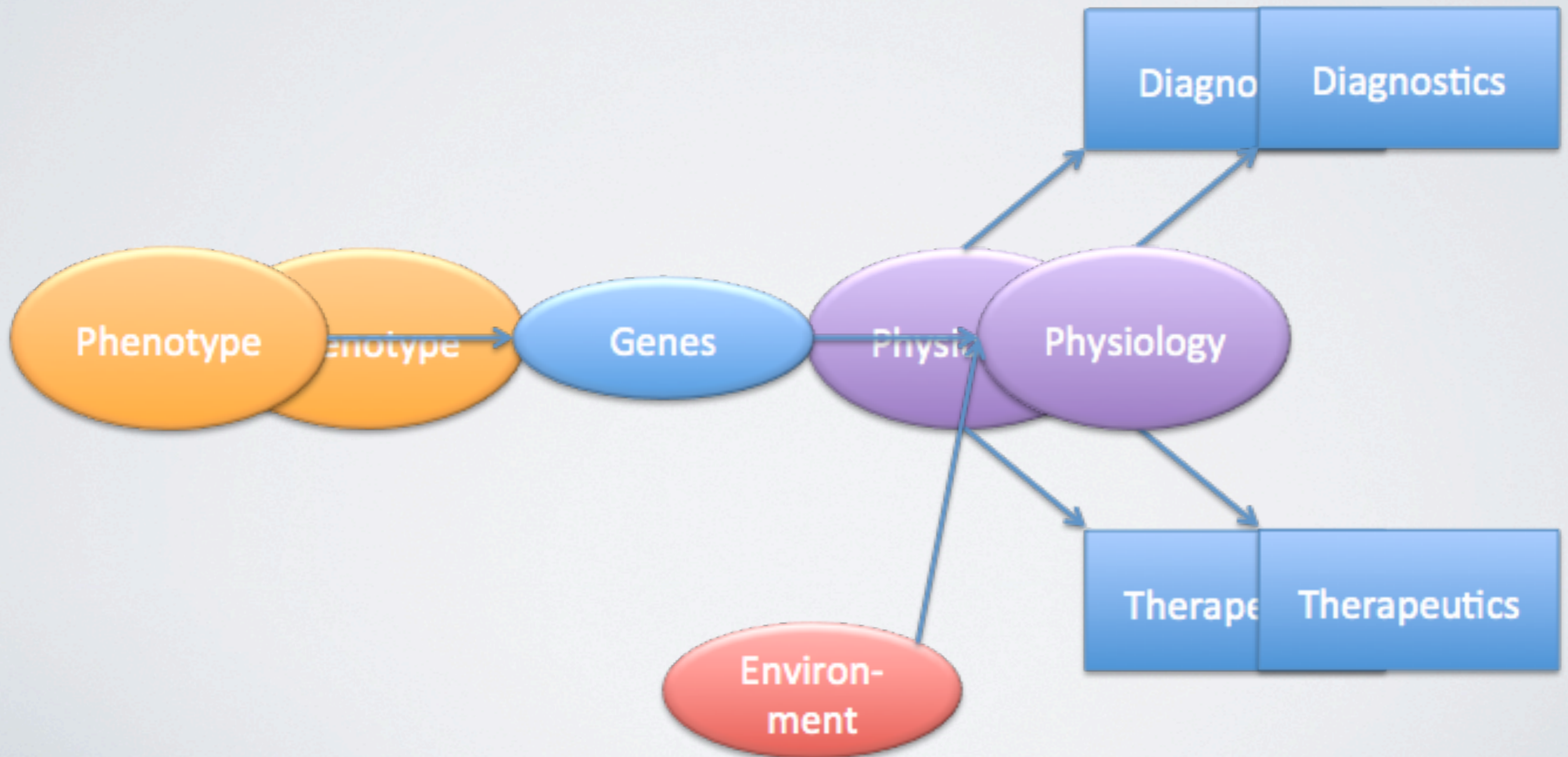
**Number of Entries in OMIM** (1 January 2012) :

| Prefix | Autosomal | X Linked | Y Linked | Mitochondrial | Totals |
|---|---|---|---|---|---|
| \* Gene description | 13,041 | 640 | 48 | 35 | 13,764 |
| + Gene and phenotype, combined | 161 | 6 | 0 | 2 | 169 |
| # Phenotype description, molecular basis known | 3,064 | 258 | 4 | 28 | 3,354 |
| % Phenotype description or locus, molecular basis unknown | 1,654 | 136 | 5 | 0 | 1,795 |
| Other, mainly phenotypes with suspected mendelian basis | 1,799 | 129 | 2 | 0 | 1,930 |
| Totals | 19,719 | 1,169 | 59 | 65 | 21,012 |



Cumulative Pace of Gene Discovery 1981-2003[1]

http://www.genome.gov/Pages/News/PaceofDiseaseGeneDiscovery.pdf

# Types of Variants

# Approach to Genetic Disorders

# Genetic Linkage

# Polymorphism

Polymorphism: occurrence of at least two alleles at a locus having a frequency of at least 1%

| Type | Description |
|------|-------------|
| VNTR | 14-100 bp repeat unit with variable number of repeats |
| STR | di, tri, tetranucleotide repeats |
| SNP | Single base change |
| CNV | Copy number variation |

11

# Linkage



Independent Assortment

Complete Linkage

10% Recombination

Likelihood Ratio

$$\text{odds ratio} = \frac{(1-\theta)^n (\theta)^r}{(1/2)^{n+r}}$$

n = number non-recombinants
r = number recombinants

# LOD Analysis



Family 1

Family 2

Family 3

Family 4

Family 5

lod score

θ

| Family | Sibs | Recombinants | Nonrecombinants | 0 | 0.1 | 0.2 | 0.3 | 0.4 |
|--------|------|--------------|-----------------|-----|------|------|------|------|
| 1 | 12 | 2 | 10 | - ∞ | 1.15 | 1.25 | 1.02 | 0.60 |
| 2 | 9 | 2 | 7 | - ∞ | 0.39 | 0.96 | 0.58 | 0.36 |
| 3 | 8 | 2 | 6 | - ∞ | 0.13 | 0.43 | 0.43 | 0.28 |
| 4 | 10 | 2 | 8 | - ∞ | 0.64 | 0.84 | 0.73 | 0.44 |
| 5 | 7 | 1 | 6 | - ∞ | 0.83 | 0.83 | 0.65 | 0.38 |
| Total | 46 | 7 | 39 | - ∞ | 3.14 | 4.31 | 3.41 | 2.06 |

# Haplotype Analysis

# Positional Cloning
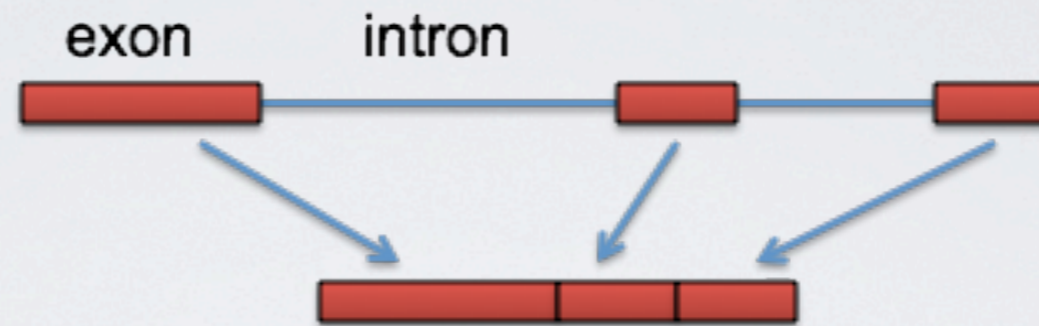


Duchenne
muscular dystrophy

# Genome Browser

Cost per Genome
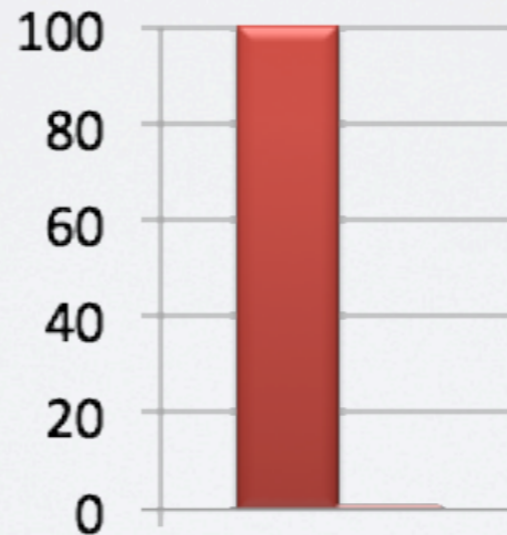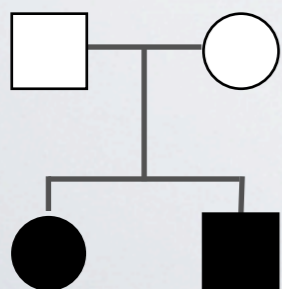
# Massively Parallel Sequencing
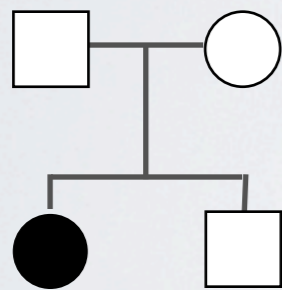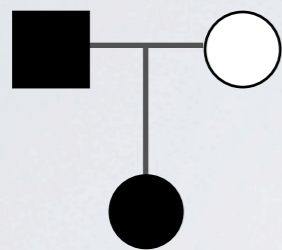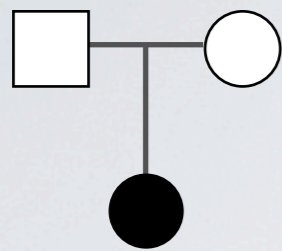
# Exome vs. Genome Sequencing



exon    intron

Genome

Exome

# Gene discovery

variants

↓

not in database of benign variants
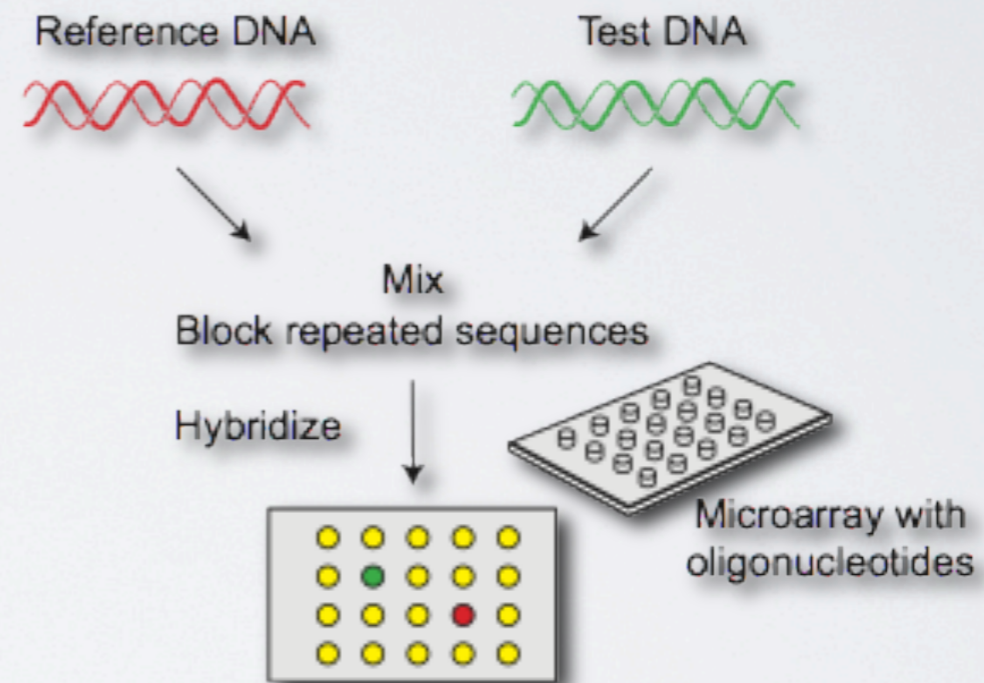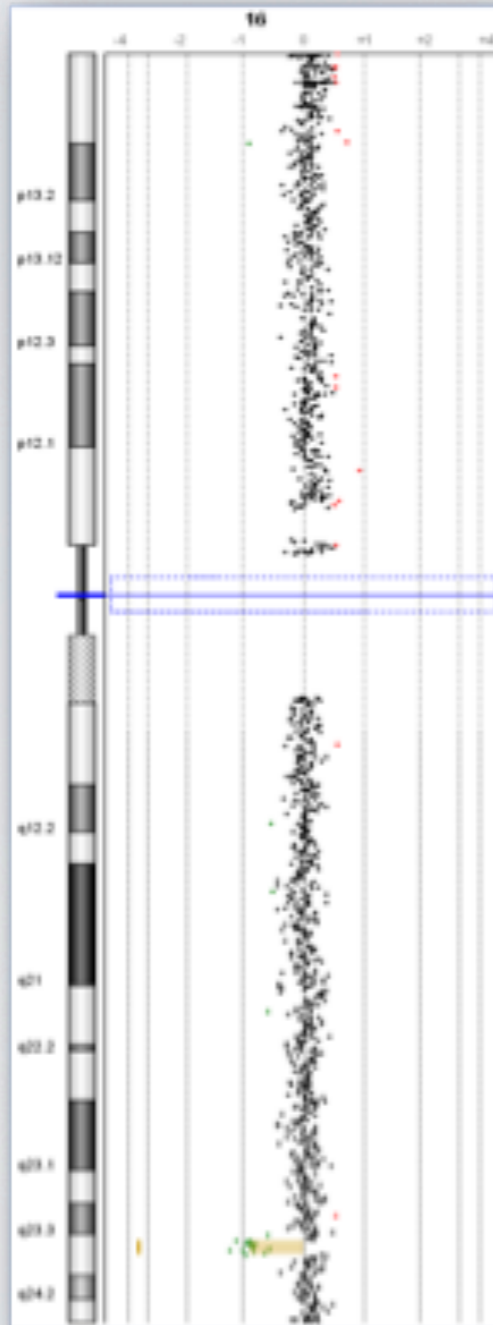
↓

predicted damaging

↓

affects one or both alleles

↓

shared by affected relatives

# Cytogenomics

# Diagnostic Odyssey