# Last time:

* **Approximations & self consistency**

$$\epsilon x^2 + x - 1 = 0 \implies x \approx 1 \quad (\epsilon x^2 \approx \epsilon \ll 1)$$

<span style="color:blue">small?</span> <span style="color:red">dominant?</span>

$$\implies \epsilon \gg 1$$

* **Probability** $(x \sim p(x))$

$\implies$ **Generating functions**: $H_x(z) \equiv \langle e^{-zx} \rangle = \int e^{-zx} p(x) dx$

e.g. Poisson $\quad p(n) = \frac{\lambda^n}{n!} e^{-\lambda} \iff H_n(z) = e^{-\lambda(1-e^{-z})}$

$\implies$ **Central limit theorem**: $X_1, X_2, \ldots, X_n \sim p(x)$

as $n \to \infty \implies \frac{1}{n} \sum_{i=1}^{n} X_i \to \text{Gaussian}\left(\langle x \rangle, \frac{\text{Var}(x)}{n}\right)$
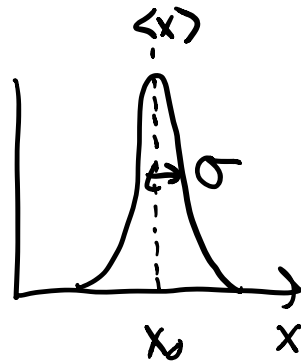
**Today:**
① Intuition about probability
② Biological background (#s & scales)
③ Simple model of evolution (if time permits)
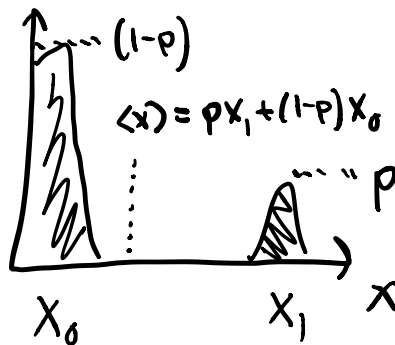
# "Average" vs "typical"

2 main classes of behavior:

**case 1:**



$\langle x \rangle$

$\sigma$

$x_0$  $x$

e.g. $\text{Binomial}(N, p)$

when $Np \gg 1$ $(N(1-p) \gg 1)$

$\Rightarrow x \approx x_0 \pm \sigma$

**case 2:**



$(1-p)$

$\langle x \rangle = p x_1 + (1-p) x_0$

$p$

$x_0$   $x_1$   $x$

e.g. $\text{Binomial}(N, p)$

when $Np \ll 1$

e.g. did a mutation occur?

$\Rightarrow$ becomes important if we do something w/ $x$:

e.g. $y = F(x) =$ future growth of $x$  # of mutations @ time 0.

$\Rightarrow$ in case 1: can use Taylor expansion:

$$y = F(x) = F\left(x_0 + (x - x_0)\right) \approx F(x_0) + F'(x_0)(x - x_0)$$

$$x_0 \Rightarrow \sigma$$

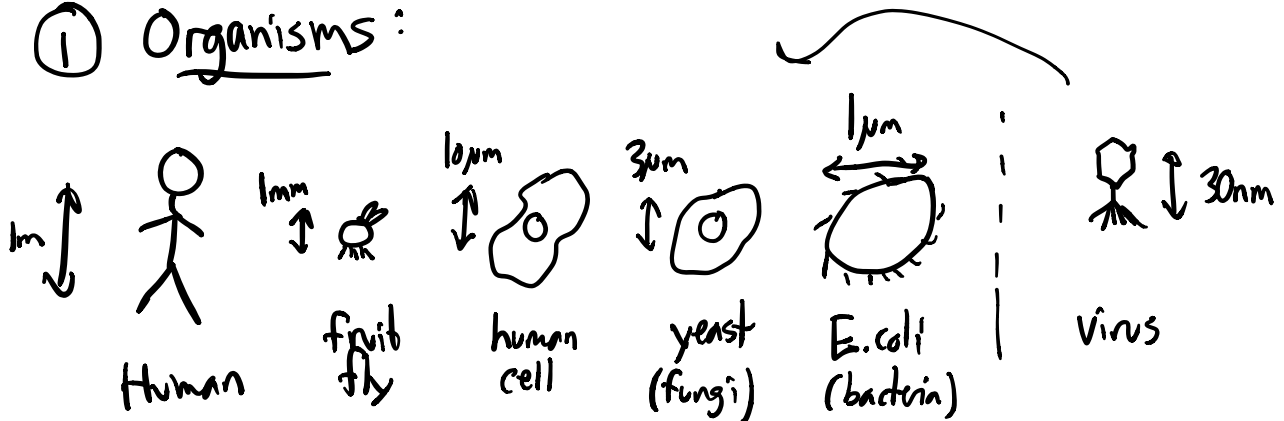$$\approx \boxed{F(x_0)} + F'(x_0)\, \sigma \quad \text{"spread"}$$

"deterministic"

case 2:

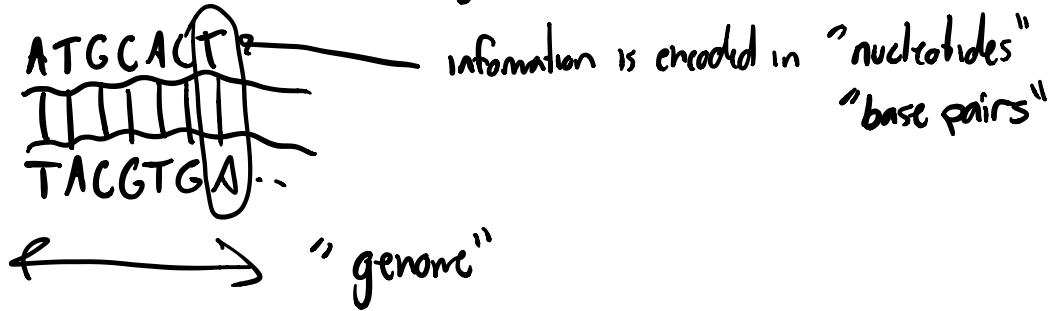$$y = \begin{cases} F(x_0) & \text{w/ prob } 1-p \\ \\ F(x_1) & \text{w/ prob } p \end{cases}$$

this can be "typical" case (most of time)

"rare event" happens repeatedly.

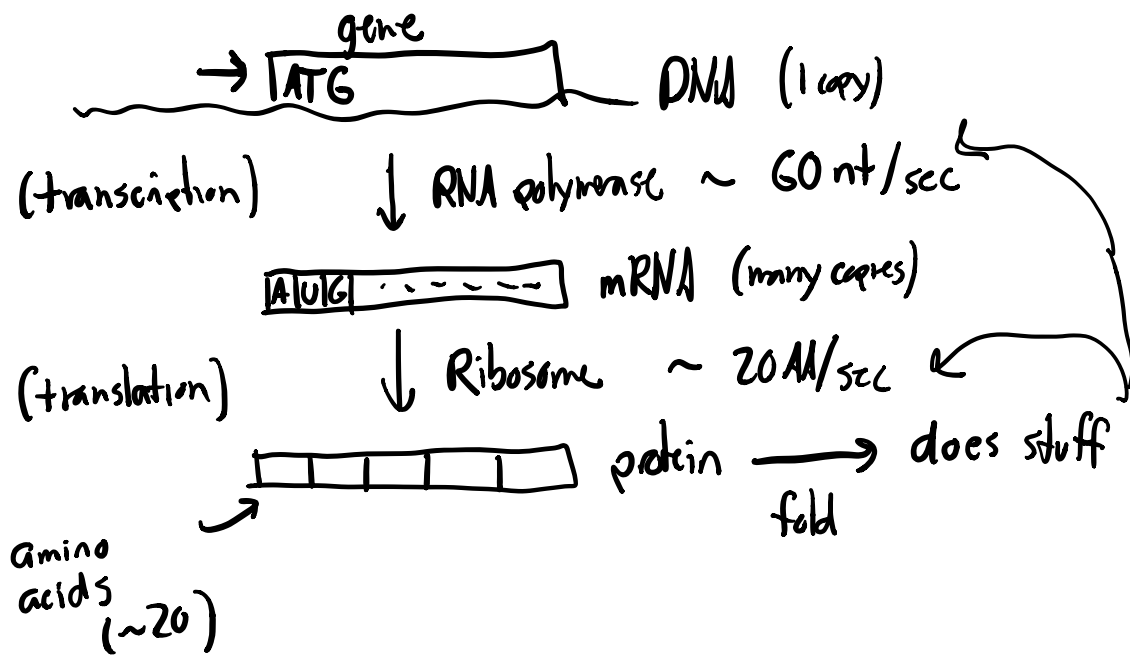# Biological Background (#s / scales)

① Organisms:



1m — Human

1mm — fruit fly

10μm — human cell

3μm — yeast (fungi)

1μm — E.coli (bacteria)

30nm — Virus

instructions encoded in a single molecule of DNA:

ATGCACT..

TACGTGA..

information is encoded in "nucleotides"
"base pairs"

$\longleftarrow \longrightarrow$ "genome"

Lengths of genomes vary widely

human: $\sim 10^9$ bp      yeast: $10^7$ bp      virus: $10^4 - 10^5$ bp

fruit fly: $\sim 10^8$ bp      bacteria: $10^6$ bp      ( 1Gbp , 1Mb    1kb
                                                              "          "     "
                                                         $10^9$ bp   $10^6$ bp  1000 bp.

information often encoded in genes (make proteins)

gene
$\rightarrow$ |ATG                    | DNA (1 copy)

(transcription)    $\downarrow$ RNA polymerase $\sim$ 60 nt/sec

|A|U|G| $\cdots$ | mRNA (many copies)

(translation)    $\downarrow$ Ribosome    $\sim$ 20 AA/sec

|   |   |   |   |   | protein $\longrightarrow$ does stuff
                              fold

amino
acids
(~20)

How does ribosome do it?

| A | T | T | = "codon"

⇩

1 amino acid
(isoleucine)

$4^3 = 64$ different codons ⟶ 20 amino acids
+ "start codon"
"genetic code" + "stop" codon

⟹ has degeneracy

⟹ typical protein ~ 300 AA (1000 bp of DNA)

⟹ # of genes varies widely across organisms:

humans: 20,000 genes    yeast: 6,000 genes
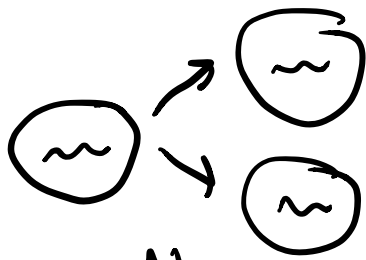
E-coli ~ 4,000 genes    viruses ~ 10 genes.

⟶ 1000x bigger genome ⟹ but 5x as many genes.

⟹ rest of genome is "noncoding" ⟶ regulation
("coding" = genes) ⟶ "junk"

cell makes a copy of itself!

① new cell wall, all other proteins
(including ribosomes!)

② needs to copy its DNA
( DNA polymerase )
(not usually limiting factor in growth)

1 doubling time
1 generation

| | |
|---|---|
| humans: ~20 yrs | E. coli ~ 20 mins - 1 hr |
| humans: ~1 day (HeLa) | Prochlorococus ~ 1 day (ocean bacteria) |

Virus: → HIV ~ 15 hrs
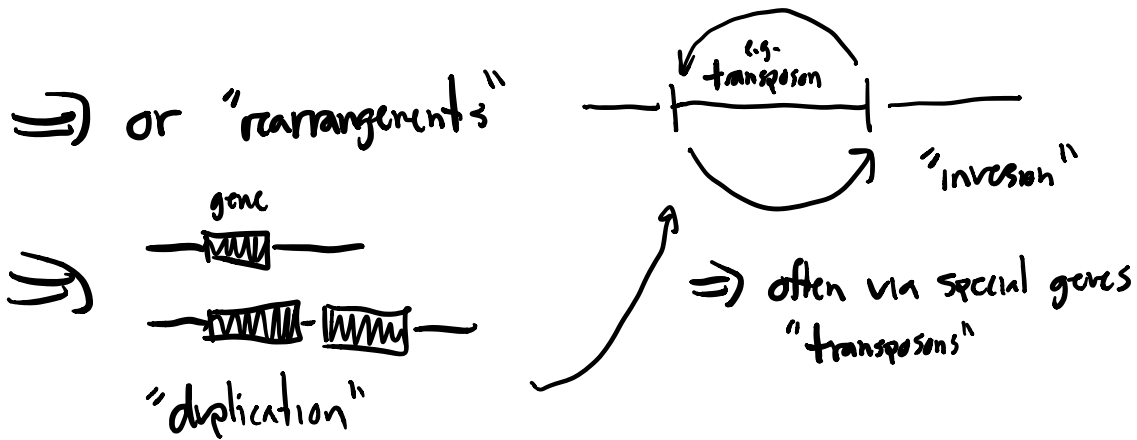SARS-CoV-2 ~ 10 hrs

⇒ when genome is copied, could introduce error ("mutations")

···· ATGCCA ···· parent
···· ATGTCA ···· offspring

⇒ simplest mutations are "point mutations" ($A \to T$, $T \to C$, ...)
"single nucleotide mut", "substitutions"

⇒ ⊕ also "insertion"  ATGTTTCA

                        ATGTTTTTTCA  "slippage of
                                      DNA pol"

⇒ or "deletion"  ATGTTTCA
                      ↓
                  ATGTCA

⇒ or "rearrangements"



"inversion"

⇒ often via special genes
  "transposons"

gene

"duplication"

⇒ cells have machinery for fixing errors!

⇒ mutation rates vary across organisms! ⟶ (N)

e.g. Humans: $N \sim 10^{-8}$/bp/gen.        E.coli ~ $N \sim 10^{-10}$/bp/gen

human      $N \sim 10^{-10}$/bp/division        viruses ~ $N \sim 10^{-5}$/bp/gen
cells                          "gen"
                                                (SARS-cov-2  $10^{-6}$/bp/gen)

**Humans:** genome is $L = 3\times10^9$ bp long $+$ $\mu \sim 10^{-8}$/bp/gen

$\Rightarrow$ so 30 mutations per genome/gen.

$\Rightarrow$ $\sim 10^{10}$ humans on earth $+$ $\mu \times 10^{-9}$/bp/gen/individual

$\Rightarrow$ $\sim$ 100 mutations produced @ every site in human genome every generation (in <u>some</u> individual)
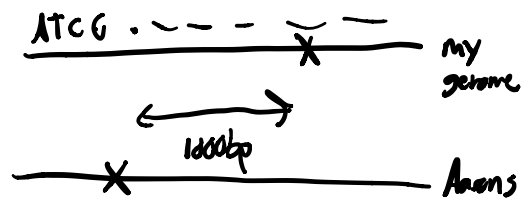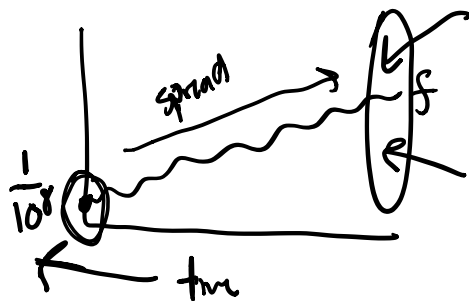
$\Rightarrow$ <u>but</u> if we pick 2 random people $\Rightarrow$ 

$\boxed{\text{differ @ 0.1\% of genome.}}$

$\boxed{\text{why?}}$ why not $10^{-4}$? $10^{-2}$?

<span style="color:red">wrong!</span>

$\boxed{\Rightarrow f = \dfrac{100}{10^{10}} \sim 10^{-8} \Rightarrow \text{differ @ } 10^{-8}?}$



$ATCG \cdot \texttt{----}\overset{\times}{\smile}\texttt{-}$ my genome

$1000\text{bp}$

$\underset{\times}{}$ Aaron's

$\boxed{Pr[\text{diff}] = 10^{-3} \neq f(1-f)}$

$10^{10}$ humans $\times 10^{-9} \times 10^{-8} = 10^{-6}$

$\Rightarrow$ Sequence space is big!

$\Rightarrow$ <u>but not all double mutants:</u>

E. coli:  genome is ~ $4 \times 10^6$ bp + $\Lambda \sim 10^{-10}$ /bp/gen.

$\Rightarrow 4 \times 10^{-4}$ mutations / genome / gen.

$\Rightarrow > 1000$ replications before single error!

in gut $10^9$ E. coli cells $\Rightarrow$ so almost bp mutated $\nearrow^{\text{w/in us.}}$ every day
$(10^{10})$

$\Rightarrow \times 10^{10}$ guts $\Rightarrow$ almost all double mutants present
in worldwide pop! $(10^{10} \times 10^{10} \times 10^{-10} \times 10^{-10} \sim 1)$

$\Rightarrow$ not triple: $\sim 10^{20} \times (10^{-10})^3 \ll 1.$        $\swarrow \begin{array}{l} 10^{82} \\ \text{atoms} \\ \text{in universe!} \end{array}$

$\Rightarrow$ more generally $L = 1000$ bp $\Rightarrow$ $4^L = 4^{1000} = 10^{602}$
                                                                    possible
                                                                    gene sequences

---

what do mutations do?      ( genotype $\Rightarrow$ phenotype map)

$\Rightarrow$ in general, we don't know (even for E. coli)

$\Rightarrow$ but in special cases, can guess based on genetic code.

e.g.   if mutation occurs in a gene:

⇒ changes a codon   ATC → ATT

① because of degeneracy, codons could code for same AA

⇒ doesn't change protein   "synonymous mutations"

② could change to something else   "nonsynonymous"

↳ other AA ("small change")   ("missense mutations")

↳ stop codon → °truncates the gene (big change)
("nonsense")                              ↓ loss-of-function