**Announcements:** Half way done!

**Recap:** Evolution at a single genetic locus



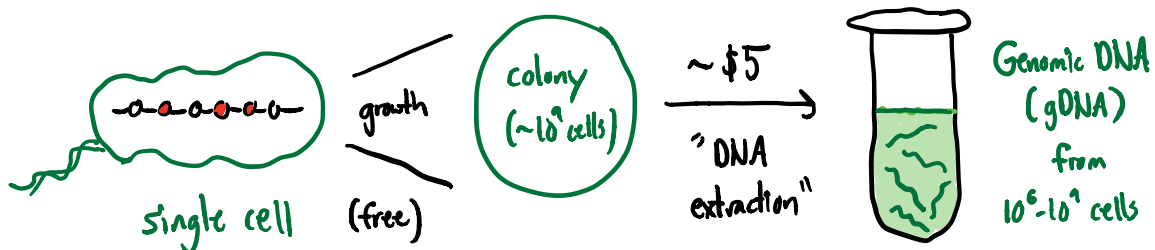$$\frac{df}{dt} = \underbrace{sf(1-f)}_{\text{selection}} + \underbrace{\sqrt{\frac{f(1-f)}{N}}\,\eta(t)}_{\text{genetic drift}}$$

$$+ \underbrace{\left[\mu(1-f) - \nu f\right]}_{\text{mutation}}$$

"Linear Branching Process"
$$\frac{df}{dt} = \underset{\text{mutation}}{\mu} + \underset{\text{selection}}{sf} + \underset{\text{drift}}{\sqrt{\frac{f}{N}}\,\eta(t)}$$

mutation freq, $f(t)$

$\tau_{est}$ $\quad$ $\tau_{fix}$

100% / 50% / $f^* \approx 1$ / dominated by selection / $\frac{1}{Ns}$ / dominated by genetic drift / $\frac{1}{N}$

time, $t$

**Where we're headed:** Evolution of longer genomes!

genome, $\vec{g}$ $\quad$ $\leftarrow L \gg 1 \rightarrow$

$$\frac{\partial f(\vec{g})}{\partial t} = \text{ ???}$$



**Last time:** How do we measure information in single DNA molecules?

**Step 1: Amplification** ("get a macroscopic amount of them")

single cell $\quad$ growth (free) $\quad$ colony ($\sim 10^9$ cells) $\quad$ $\sim \$5$ "DNA extraction" $\quad$ Genomic DNA (gDNA) from $10^6 - 10^9$ cells

<u>Problem</u>: $L = 10^6$ sized genomes too difficult do measure directly

⇒ most sequencing methods work w/ short sequences $(\lesssim 10^3 \, bp)$

Genome ....~~~~~ [ target sequence ] ~~~~~.... 

←——— ~1kb ———→

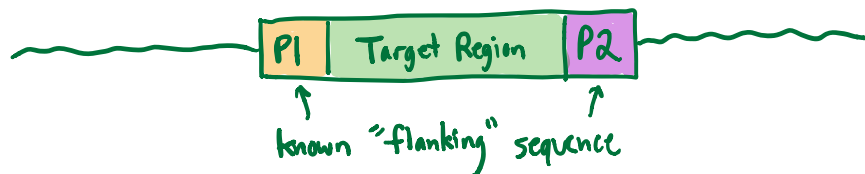How to get a <u>macroscopic</u> amount of just this region?

<u>Answer</u>: PCR ("polymerase chain reaction")

⇒ cell-free chemical reaction that's like but just for p<u>art</u> of genome.
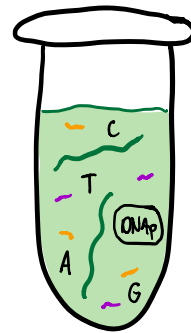
cell → colony

⇒ takes advantage of DNA replication machinery invented by bacteria

✳ <u>But</u> requires us to know some of the sequence near target

~~~~ [ P1 | Target Region | P2 ] ~~~~
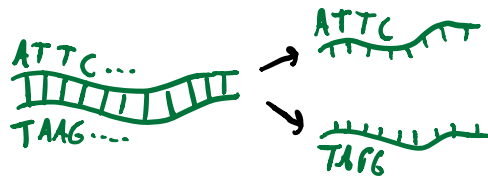
known "flanking" sequence

(1) can have company synthesize "primers" (short sequences of ~20bp)
that correspond to P1 & P2.   (~ $0.30/bp for ~ $10^3$ reactions)

(2) mix w/ <u>DNA polymerase</u>,

   <u>dNTP</u> (free A's, C's, T's, G's)

   + your <u>genomic DNA</u>.

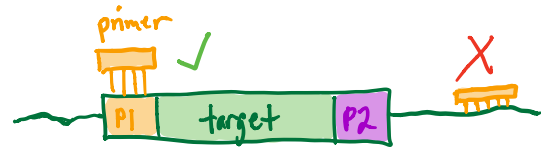(3) Heat sample so that
    DNA strands "<u>melt</u>" (~90°C)

    ATTC...        ATTC
    TAAG....   →   TAAG

(4) ~60°C
    Now cool sample so that
    primers "anneal" to input DNA

    primer

    P1  target  P2

    ⟹ melting & annealing is physics problem: $\dfrac{p(\text{bound})}{p(\text{unbound})} \sim e^{-\frac{\Delta E}{kT}}$

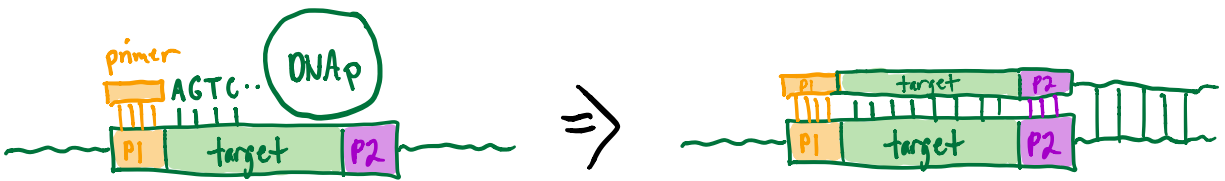    w/ $\Delta E \approx$ # matched bp
    (very roughly...)

⟹ want primers to bind to known region, but not anywhere else!



✻ we saw in PSET 1 that most of E.coli genome is uniquely identifiable w/ $\ell \approx 20bp$ sequence ⟹ ~20bp primers sounds ok ✻
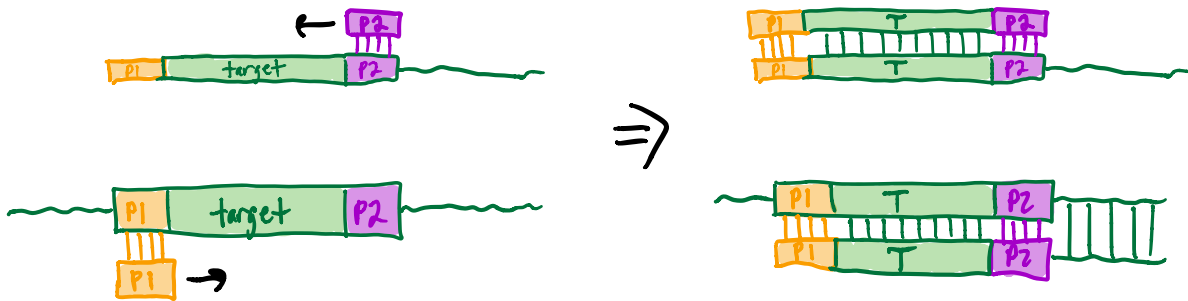
( more generally ⟹ information theory argument, see p.5 of notes...)

④ After primers are bound, DNA polymerase will start incorporating dNTPs onto *primer* to create complementary strand...



("extension phase")
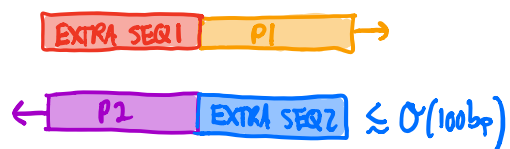
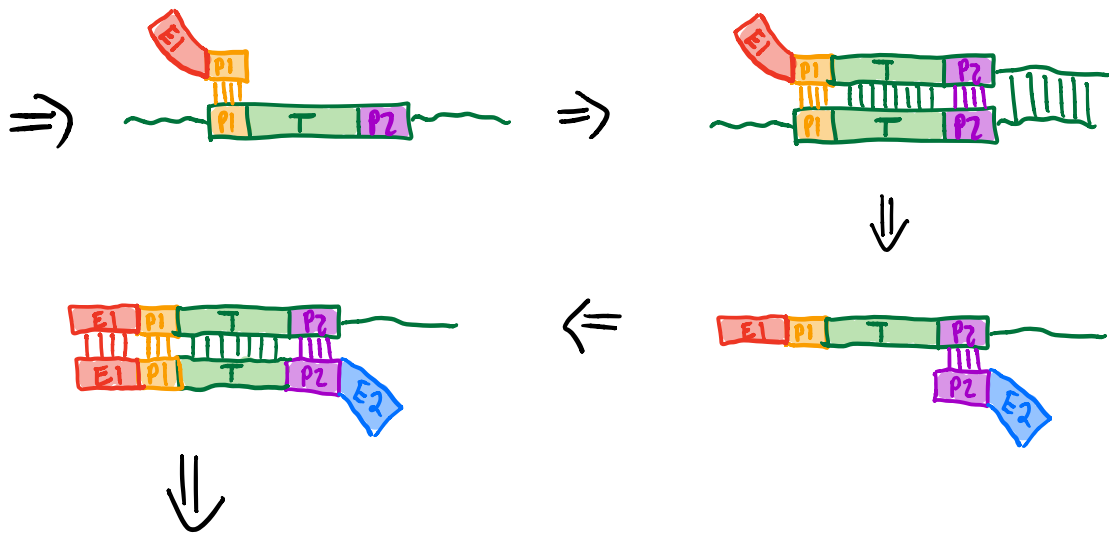(5) Melt, anneal, & extend again:



(6) Repeat for K~20-30 cycles...
⇒ exponential amplification of [P1 | T | P2] sequence!
("amplicon")

(7) "clean up" to remove leftover primers, etc. ⇒



macroscopic # of amplicons!

⇒ Note: can also use PCR to add extra bit of DNA sequence to your target region...

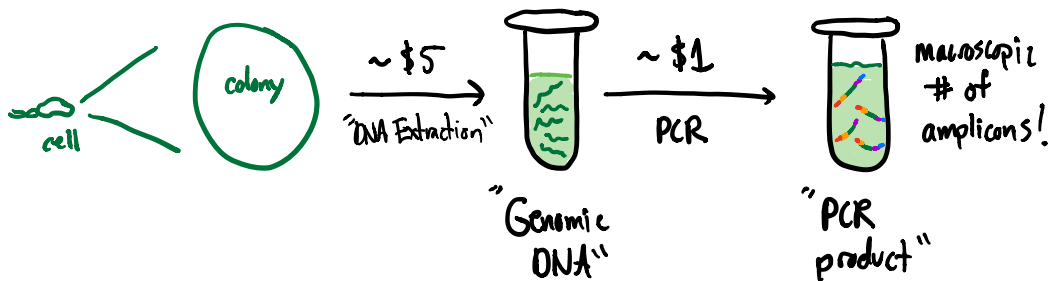⇒ order modified primers:    [EXTRA SEQ1 | P1] →
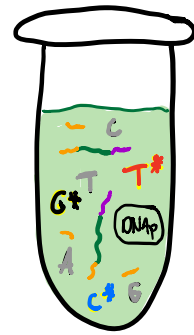← [P2 | EXTRA SEQ2] ≲ 𝒪(100bp)

Amplified sequence: 

# Recap:
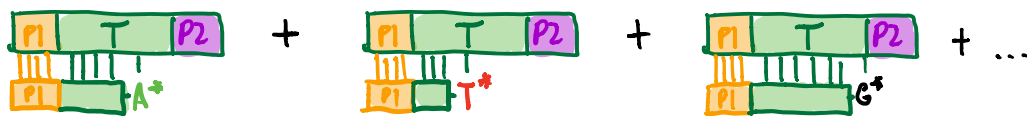


⇒ How do we read out information from PCR products?

Traditional approach: Sanger sequencing
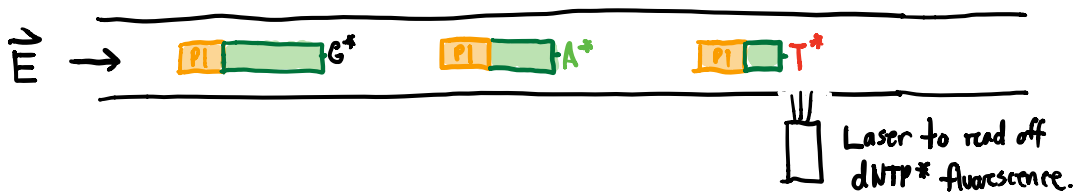
(like 1 extra round of PCR...)

Idea: mix PCR product w/ P1, DNAp, dNTPs,
+ special fluorescent dNTPs that block DNAp
(A*, G*, C*, T*)

⇒ After 1 round of extension, random mixture of

| P1 | T | P2 | + | P1 | T | P2 | + | P1 | T | P2 | + ...
P1 —A*          P1 —T*          P1 —G*

⇒ Flow in electric field (shorter fragments move faster)

E⃗ →   | P1 | G*      | P1 | A*      | P1 | T*

Laser to read off
dNTP* fluorescence.

⇓ "chromato-gram"

Fluorescence ↑

C T T C C C C G G G T T G C A C C A C A G A C T T G G G G A C A C G T T T G C G C T C G C T C T T C T G G C C
    270         280           290           300              310

Can read
off sequence!
(~1000 bp)

Time (~ bp) →

⇒ costs ~$5 (send away overnight:  )

⇒ so have seen how to go from:



cell — colony $\xrightarrow[\text{DNA Extraction}]{\sim \$5}$ gDNA $\xrightarrow[\substack{\text{PCR + Sanger} \\ \text{Sequencing}}]{\sim \$5}$ ATCGATTAGCC

sequence of ~1kb region

⇒ in Problem 1 of PSET 2, Layla Murray used sanger seq
   to sequence URA3 genes in ~ 300 yeast colonies
      ⇒ $1500 total

⇒ But expensive to sequence whole genomes!
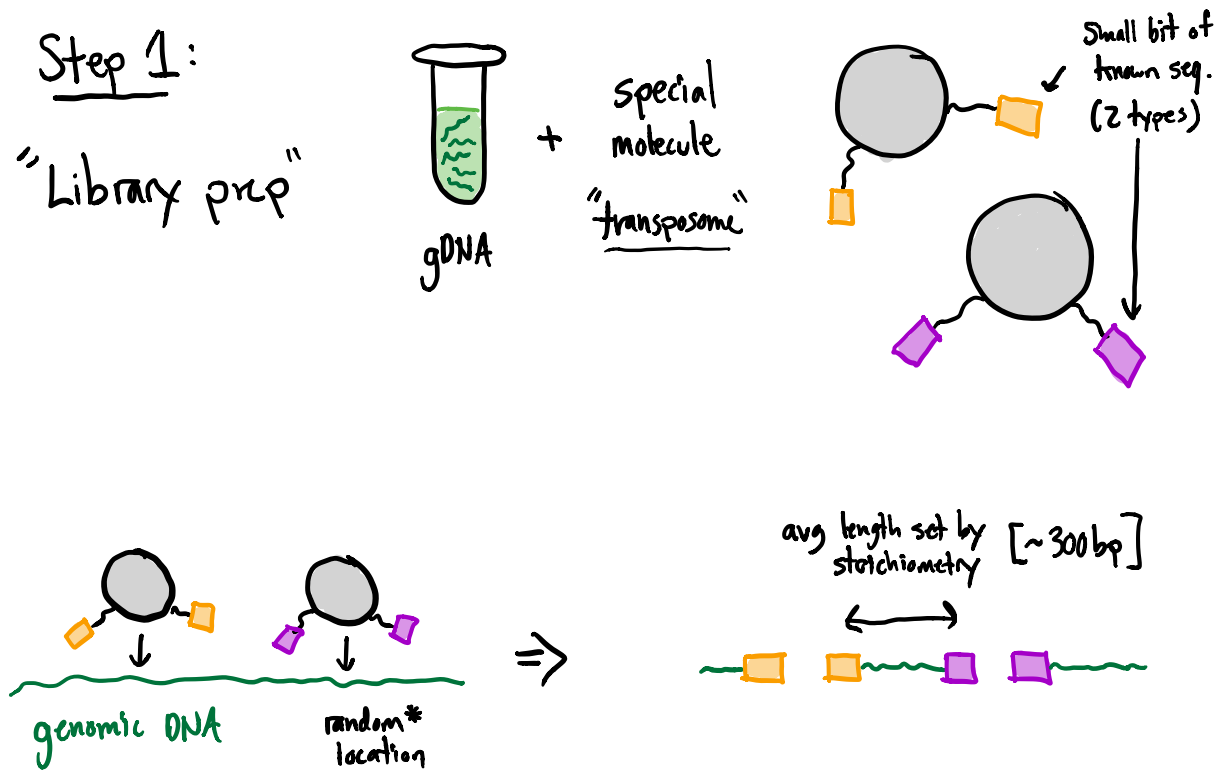   e.g. E. coli: 1 clone = $10^6 \sim 10^3$ sanger → $5k
                                          seq runs
   e.g. Humans: 1 person = $10^9$ bp ⇒ $10^6$ sanger ~ $5M
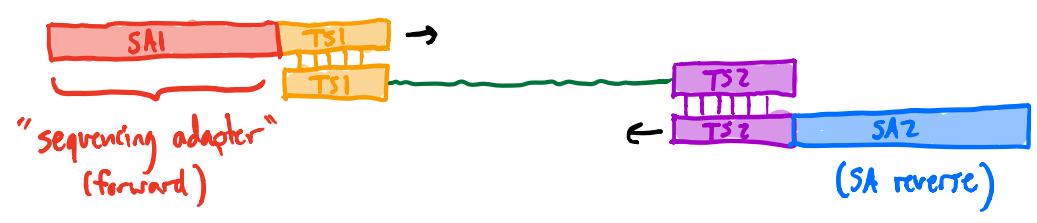                                               seq runs

⇒ Now things are much cheaper w/ <u>Next-gen</u> <u>sequencing</u>
("Illumina sequencing")

⇒ same idea, but higher throughput!

Step 1:

"Library prep"

gDNA  +  special molecule "transposome"

Small bit of known seq. (2 types)

genomic DNA      random* location

⇒

avg length set by stoichiometry  [~300 bp]

Step 2: Now in position to do PCR & add extra known sequence:

SA1   TS1  →

TS1

TS2

TS2   SA2

"sequencing adapter" (forward)

(SA reverse)

⇒ This gives us:

cell ⟶ colony ──~$5──→ gDNA ──~$10──→ "sequencing Library"
              "DNA Extraction"        "library prep"

macroscopic # of diverse amplicons!

Step 3: "Sequencing"
(Illumina Machine)

Flow library onto chip
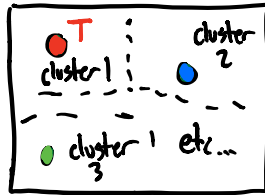w/ probes that
bind sequencing adapter


flow cell

← "flow cell"

Step 4: do more PCR on bound fragments to turn each
molecule into macroscopic cluster of identical* molecules

top down view:

cluster 1    cluster 2
cluster 3    etc...

**Step 5:** flow primer + fluorescent dNTPs that incorporate <u>once</u> (then stop)
(SA1)

⇒ if take picture of flow cell:



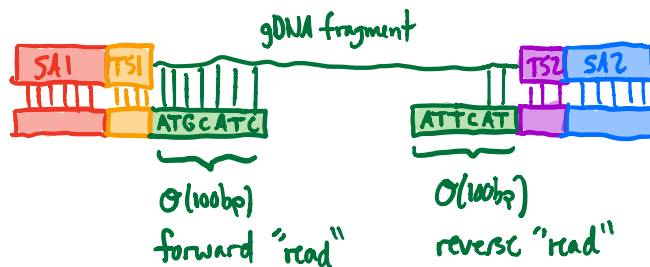⇒ different colors
= nucleic acid
@ first position

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

**Step 6:** remove fluorescent part (& block), & repeat w/ new round of dNTPs      ⇒ 2nd position.

**Step 7:** can repeat for Θ(100) cycles [until pictures ≈ desynchronized]

**Step 8:** can also repeat for reverse direction (SA2)

Net result:



gDNA fragment

SA1  TS1  ATGCATC          ATTCAT  TS2  SA2

Θ(100bp) forward "read"

Θ(100bp) reverse "read"

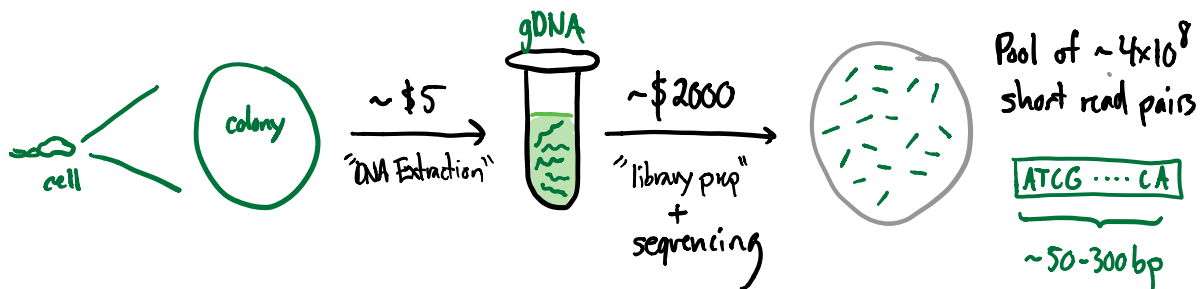⇒ get Θ(100 bp) read out from each end of <u>single</u> DNA fragment

⇒) w/ modern Illumina machines, process is very high throughput!

⇒ ~$4 \times 10^8$ read pairs in ~few days for ~$2000
(catch: can't do smaller batches)

Recap: Now have method for going from:



gDNA

cell — colony — $\xrightarrow[\text{"DNA Extraction"}]{\sim \$5}$ — $\xrightarrow[\substack{\text{"library prep"} \\ + \\ \text{sequencing}}]{\sim \$2000}$ — Pool of ~$4 \times 10^8$ short read pairs
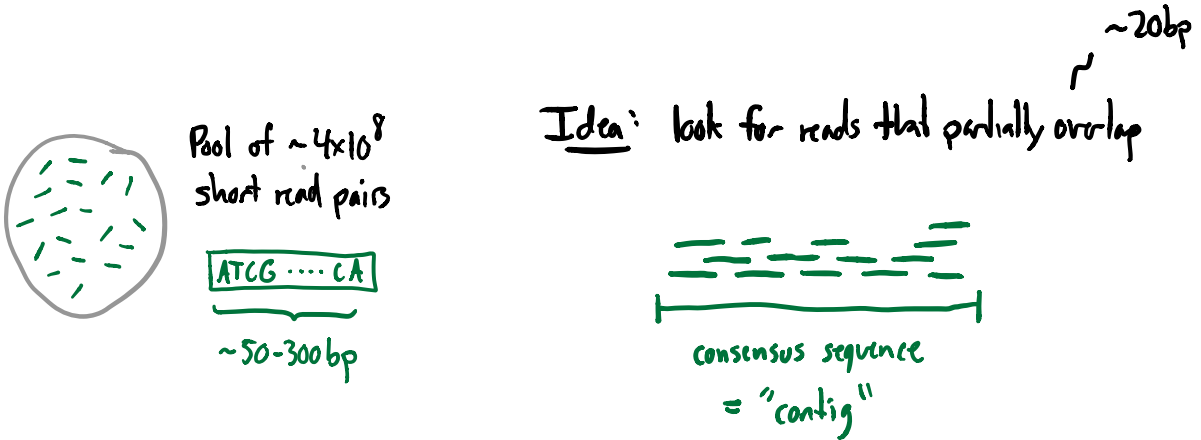
ATCG .... CA
~50-300 bp

Next: what can we do with this kind of data?

⇒ need to put puzzle back together... 2 main methods

# ① "De novo Genome Assembly"

(common programs: spades, Velvet, ...)

Pool of ~$4 \times 10^8$
short read pairs

```
ATCG .... CA
```

~50-300 bp

**Idea:** look for reads that partially overlap

~20 bp

consensus sequence
= "contig"

⇒ simple in principle, but lots of corner cases...

⇒ e.g. what if 2 regions of genome are identical
for ≥ 100 bp (or length of read?)

gene copy 1      gene copy 2

"fork in assembly"

⇒

gene sequence
(copy 1 ± 2)

$\Rightarrow$ for these & other reasons, assembly typically
yields collection of discrete contigs,

each $\sim 10^3 - 10^5$ bp long :  ━━━  ━━  ━━━  . . . .
                                    contig 1  contig   contig 3
                                              2

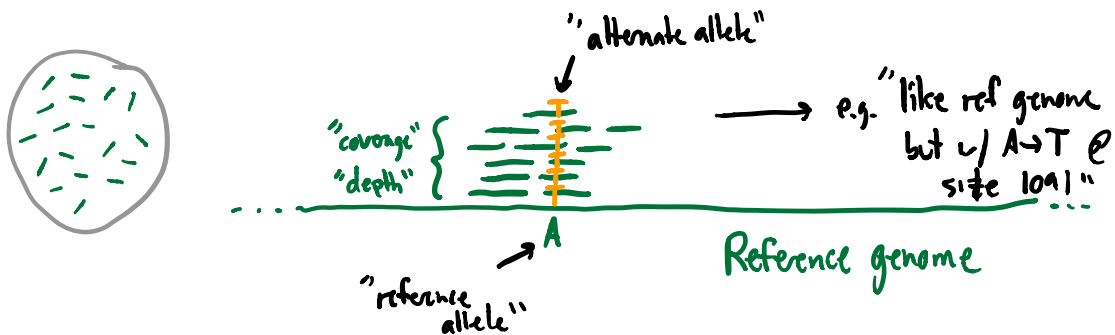$\Rightarrow$ much harder (+ manual effort)

to "finish" into complete genome

+ <u>HUGE</u> memory requirements ($\sim 32Gb - 1Tb$ depending on L)
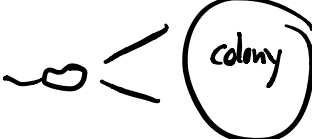(since need to compare all* pairs of reads to each other...)

② Alignment of reads to Reference Genome

⇒ if already have assembled genome from related strain...
can align reads to best matching place & look for changes
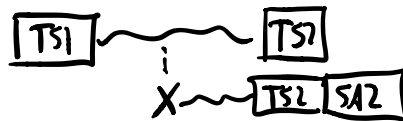(common programs: Bowtie2, BWA-MEM, + mpileup)



"alternate allele"

"coverage"
"depth"

→ e.g. "like ref genome
but w/ A→T @
site 1091"

A

Reference genome

"reference allele"

Major wrinkle: sequencing errors (A→T by chance, "shot noise")

where could these errors come from?

① initial grow up: ◦< (colony)   ⇒ but happen w/ low rate
$(\mu \sim 10^{-10}/\text{bp/gen})$

② library prep & PCR

TS1 ~~~ TS2
X ~~~ TS2 SA2

$\Rightarrow$ higher rates $\sim 10^{-6}$/bp/round $\rightarrow 10^{-4}$/bp total.

③ Errors on sequencing machine $\Rightarrow$ estimated $\sim 10^{-3}$/bp $\Big)$ larger

(cluster generation PCR errors)
wrong fluorescent base....

&ast; **but** can vary from site to site

(up to $10^{-2}$ sometimes)

$\Rightarrow$ dominant source of noise: $P_{err} \sim 10^{-3}$ (up to $\sim 10^{-2}$)

in special cases.

$\Rightarrow$ low rate, but $L \gg 1$!

$10^3$ errors per E.coli

$\Rightarrow$ # errors per genome $= L \times P_{err} \sim 10^6$ errors for Human

$\Rightarrow$ fortunately, can correct errors w/ consensus sequences!



$\rightarrow$ correct to "A"

ref genome

vs

$\rightarrow$ correct to "T"

$\Rightarrow$ higher coverage is helpful... how much?

$$\Pr\left(\text{consensus error}\right) = \sum_{A=\bar{D}/2}^{\infty} \frac{(P_{err} \times \bar{D})^A}{A!} e^{-P_{err}\bar{D}} \simeq \frac{(P_{err}\bar{D})^{\bar{D}/2} e^{-P_{err}\bar{D}}}{(\bar{D}/2)!}$$

\# errors in genome $= L \times \Pr(\text{consensus error})$

$$\simeq \exp\left[\log L - \frac{\bar{D}}{2}\left[\log\left(\frac{1}{2P_{err}}\right) - 1\right]\right] \lesssim \mathcal{O}(1)$$

$\Rightarrow$ solve for $\bar{D}$ $\Rightarrow$ $\bar{D} \simeq \dfrac{2\log L}{\log\left(\frac{1}{2P_{err}}\right) - 1} \approx \begin{cases} 5 \text{ if } L \sim 10^{+6}, P_{err} \sim 10^{-3} \\ 8 \text{ if } L \sim 10^9 \end{cases}$

$\Rightarrow$ so need coverage of $\gtrsim 10x$ coverage to detect single muts.

$\Rightarrow$ How many reads is this?

E. coli: $10 \times 10^6$ bp $\Rightarrow 10^7$ bp $\Rightarrow 10^5$ reads
per site

$\Rightarrow$ but a single run of Illumina seq produces $4 \times 10^8$ reads

$\Rightarrow$ overkill! (wastes money)