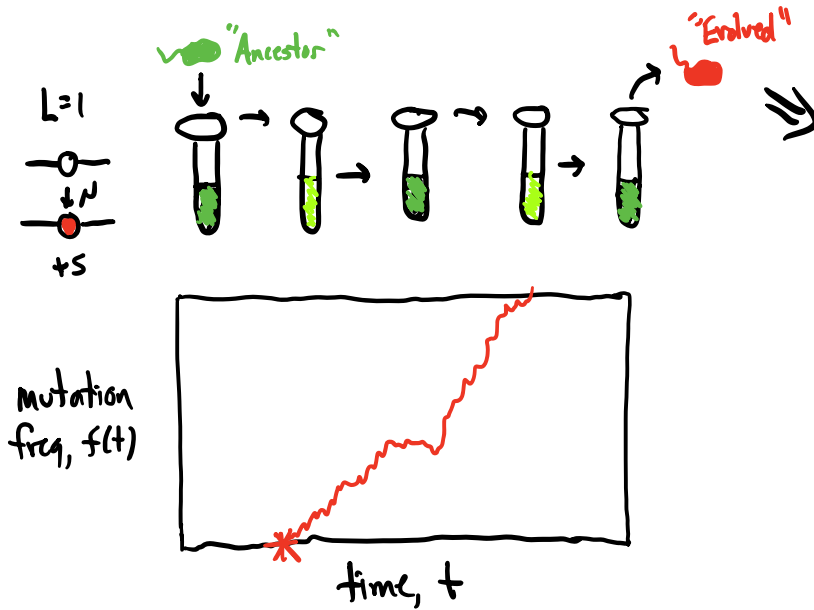


## **Chapter 8**

# **DNA sequencing & genomics**

# DNA Sequencing + Genomics

So far....



Qualitative phenotype  
(e.g. red vs green)  
that could be detected  
by colony counting

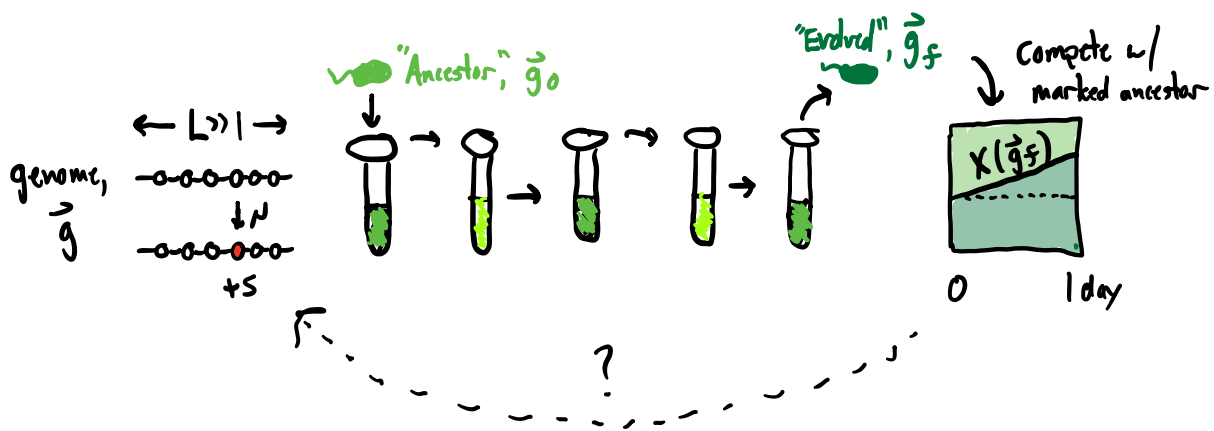


$\Rightarrow$  In practice, genomes contain many sites

$\Rightarrow$  don't know what phenotypes mutations  
@ these sites produce or how to  
measure them w/ colony counting assay...

$$\left( \begin{array}{l} L \sim 10^4 - 10^5 \text{ for viruses} \\ L \sim 10^6 - 10^7 \text{ bacteria} \\ L \sim 10^9 \text{ for humans} \end{array} \right)$$

Historically, experimental evolution relied on competitive fitness

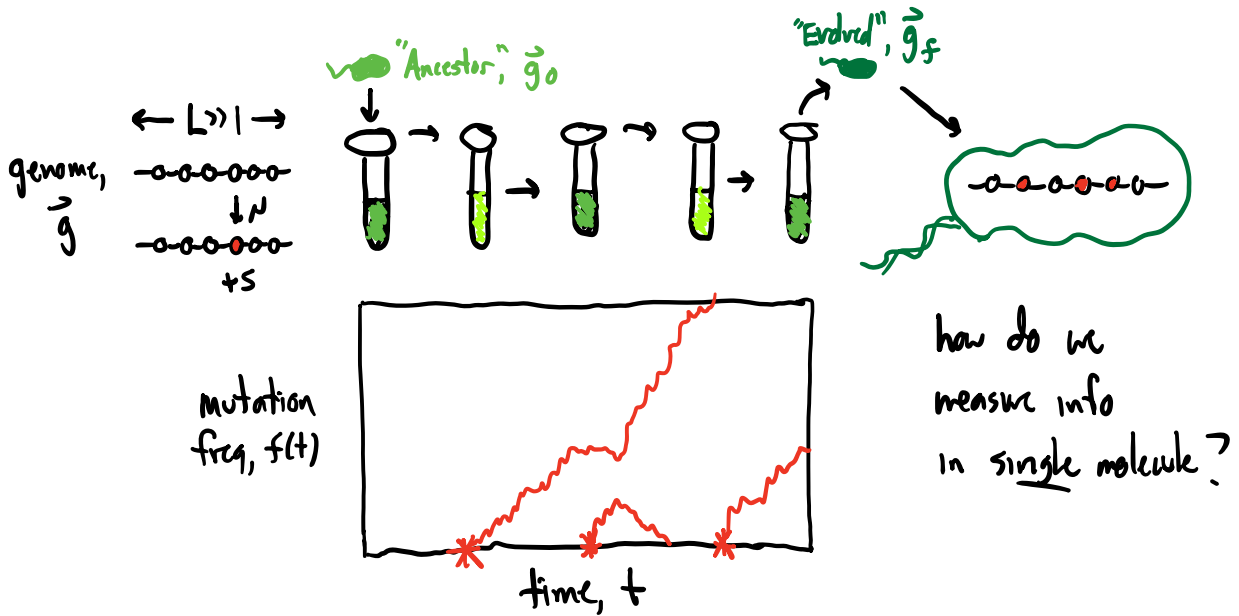


$\Rightarrow$  statistics of  $X(\vec{g}_f)$  w/in + between populations  
tell us something about evolutionary dynamics of  $\vec{g}$

$\Rightarrow$  downside: indirect! many different dynamics of  $\vec{g}$   
consistent w/ same dynamics of  $X(\vec{g})$ ...

+ mapping from  $\vec{g} \rightarrow X(\vec{g})$  poorly understood...


Now: DNA sequencing allows us to measure genomes directly\*



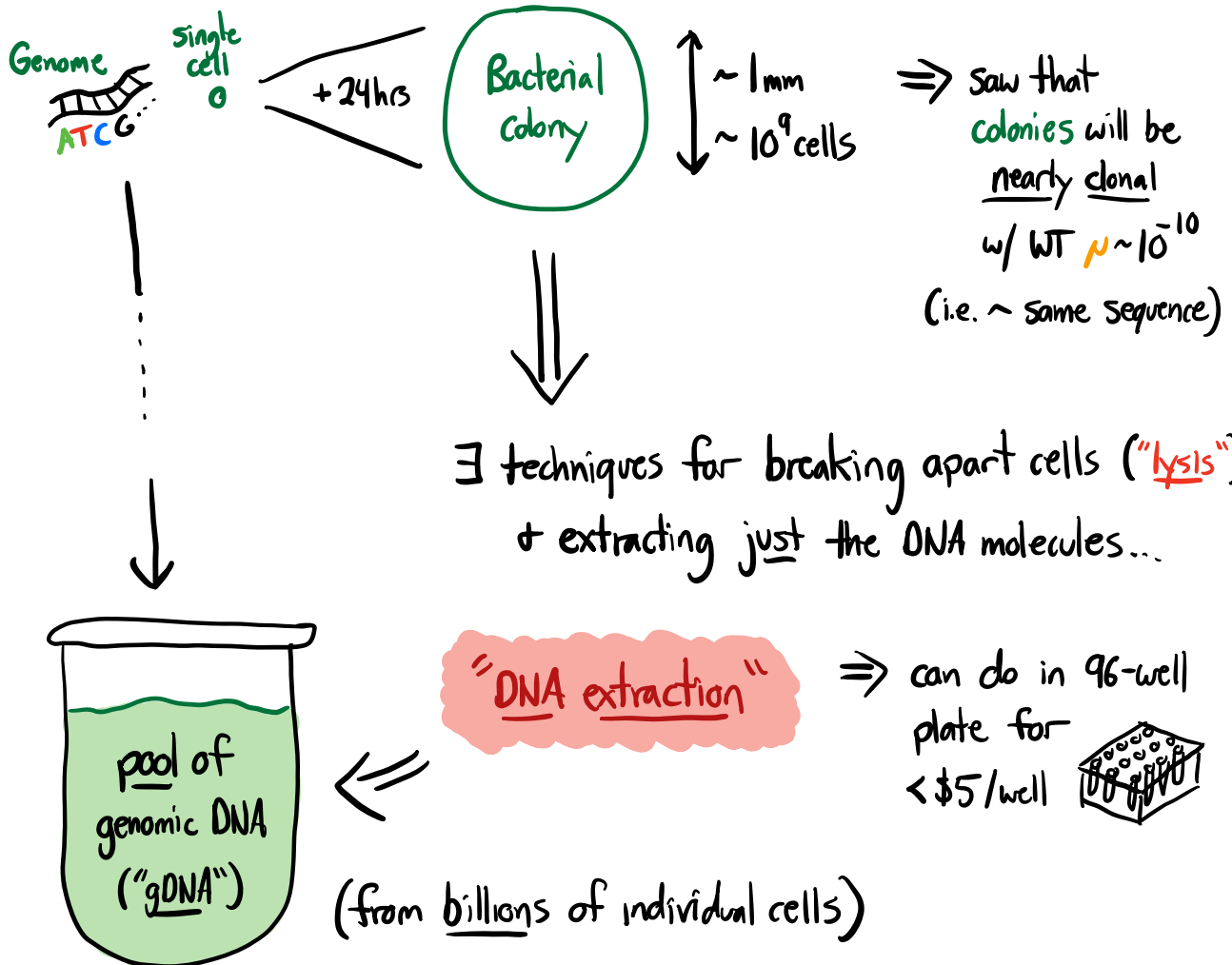
Recall: genome =  $\begin{matrix} \text{ATCGTA} \dots \\ \text{TA} \dots \end{matrix}$  complementary strands

$L \sim 10^{4-5}$  viruses  
 $L \sim 10^6$  bacteria  
 $L \sim 10^9$  humans

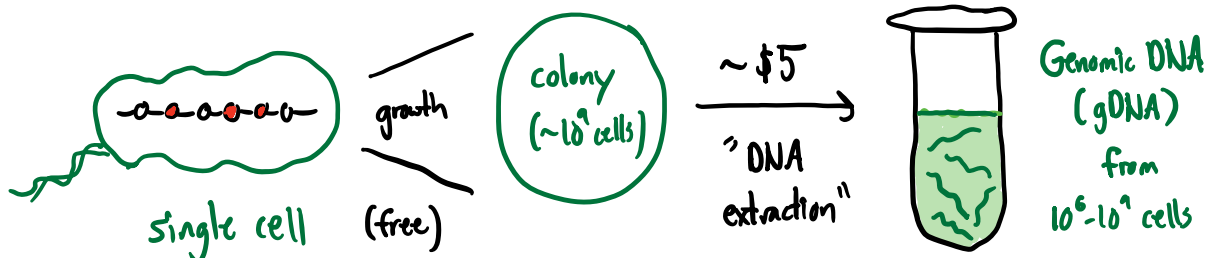
$\Rightarrow$  Step 1 for reading genomes: amplification!

$n=1$  molecule  $\ll$   need macroscopic quantities of our DNA molecule (w/ same DNA sequence) to work with...

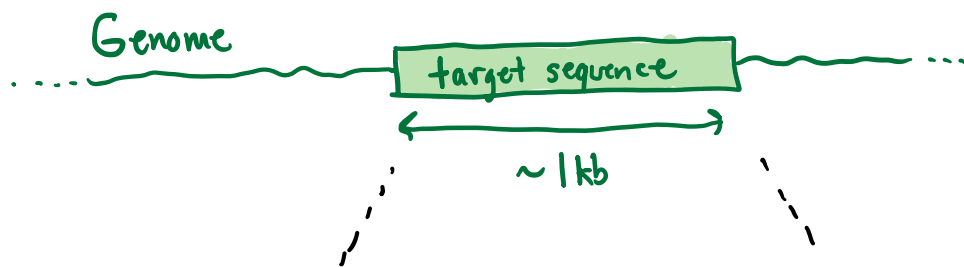
⇒ easy for lab bacteria! use built-in ability to grow exponentially:



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



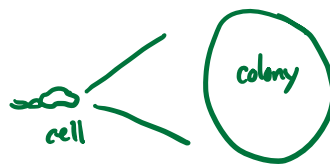
Problem:  $L=10^6$  sized genomes too difficult to measure directly  
 $\Rightarrow$  most sequencing methods work w/ short sequences ( $\approx 10^3$  bp)



How to get a macroscopic amount of just this region?

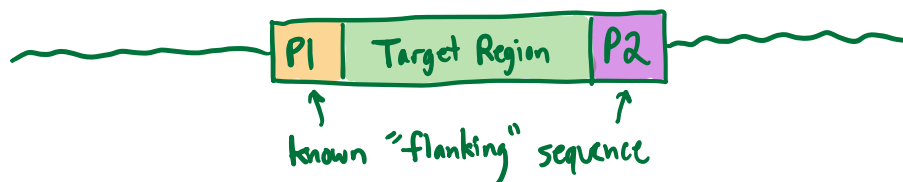
Answer: PCR ("polymerase chain reaction")

$\Rightarrow$  cell-free chemical reaction that's like  
but just for part of genome.



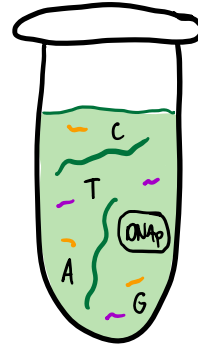
$\Rightarrow$  takes advantage of DNA replication machinery invented by bacteria

\* But requires us to know some of the sequence near target

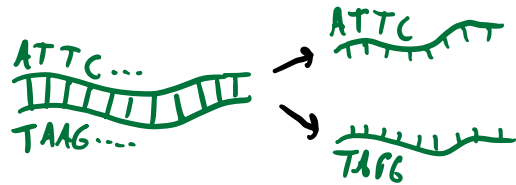


① can have company synthesize "primers" (short sequences of  $\sim 20$ bp) that correspond to P1 + P2. ( $\sim \$0.30/\text{bp}$  for  $\sim 10^3$  reactions)

② mix w/ DNA polymerase,  
dNTP (free A's, C's, T's, G's)  
+ your genomic DNA.



③ Heat sample so that DNA strands "melt" ( $\sim 90^\circ\text{C}$ )



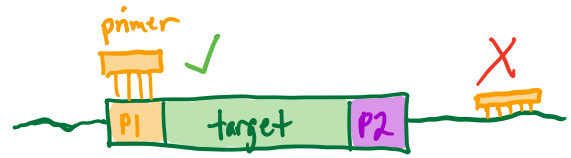
④ Now cool sample ( $\sim 60^\circ\text{C}$ ) so that primers "anneal" to input DNA



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

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

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



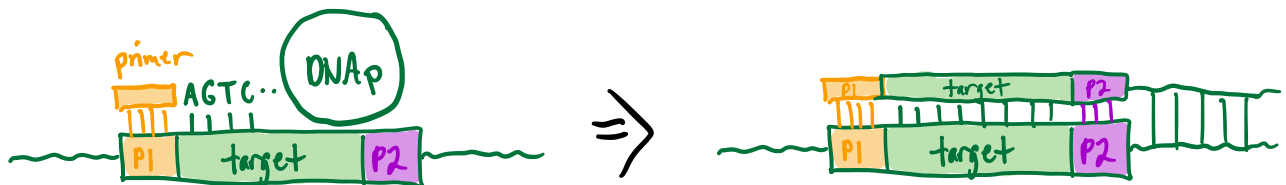
\* we saw in PSET 1 that most of E.coli genome is uniquely identifiable w/  $l \approx 20$ bp sequence ⇒  $\sim 20$ bp primers sounds ok \*

(more generally ⇒ "information theory argument" :

$$\# \text{ unique locations } (L) \lesssim \# \text{ unique sub-sequences } (4^l)$$

$$\Rightarrow l^* \geq \log_4(L) \left. \begin{array}{l} \sim 10 \text{ bp (E.coli)} \\ \sim 15 \text{ bp (Humans)} \end{array} \right)$$

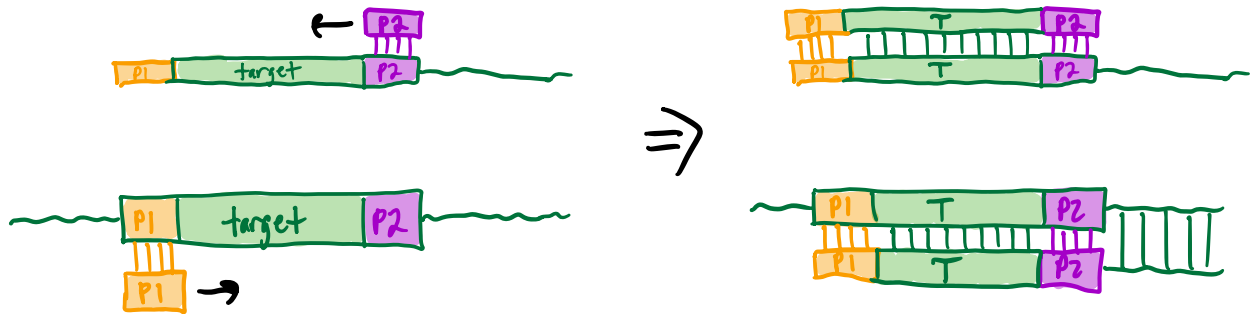
⑤ After primers are bound, DNA polymerase will start incorporating dNTPs onto primer to create complementary strand...




("extension phase")



⑥ Melt, anneal, + extend again:





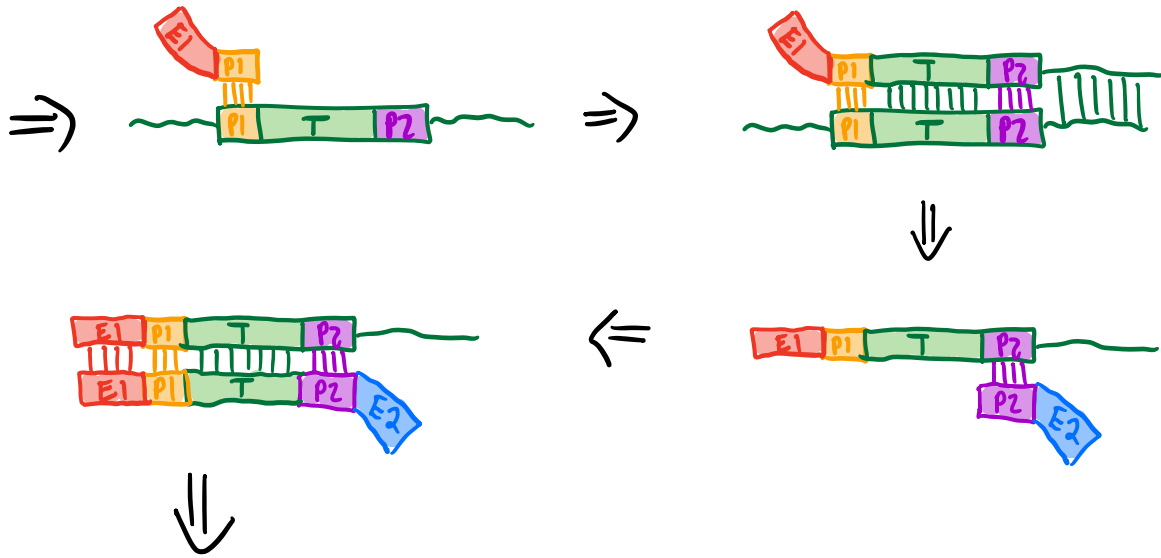
⑦ Repeat for  $k \sim 20-30$  cycles...

$\Rightarrow$  exponential amplification of  sequence!  
("amplicon")

⑧ "clean up" to remove leftover primers, etc.  $\Rightarrow$   macroscopic # of amplicons!

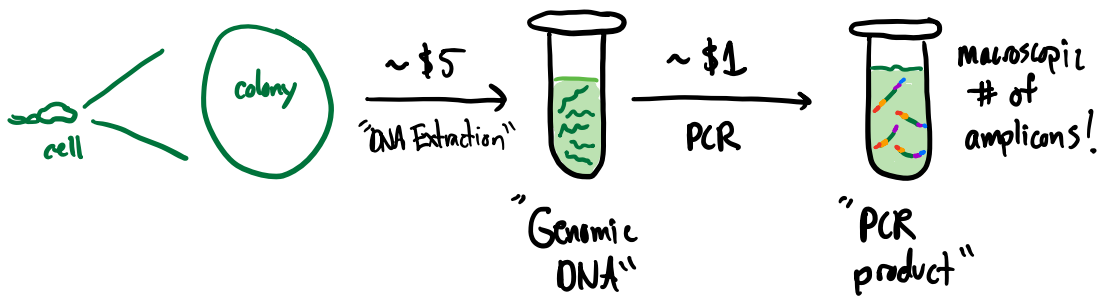
$\Rightarrow$  Note: can also use PCR to add extra bit of DNA sequence to your target region...

$\Rightarrow$  order modified primers:   $\rightarrow$   
 $\leftarrow$    $\approx O(100bp)$



Amplified sequence:

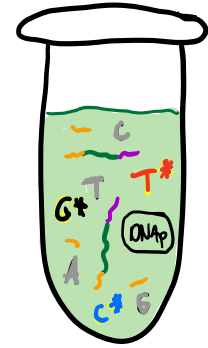
Recap:



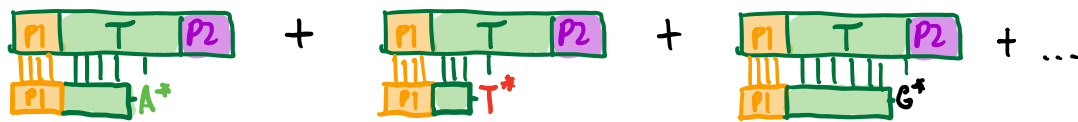
⇒ How do we read out information from PCR products?

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

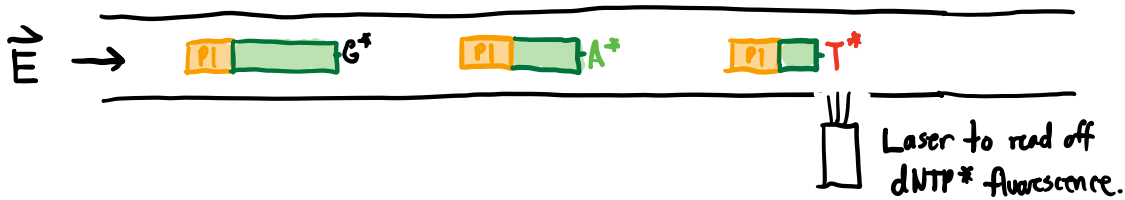
Idea: mix PCR product w/ **PI**, **DNAp**, **dNTPs**,  
 + special fluorescent dNTPs that block DNAp  
 ( $A^*$ ,  $G^*$ ,  $C^*$ ,  $T^*$ )



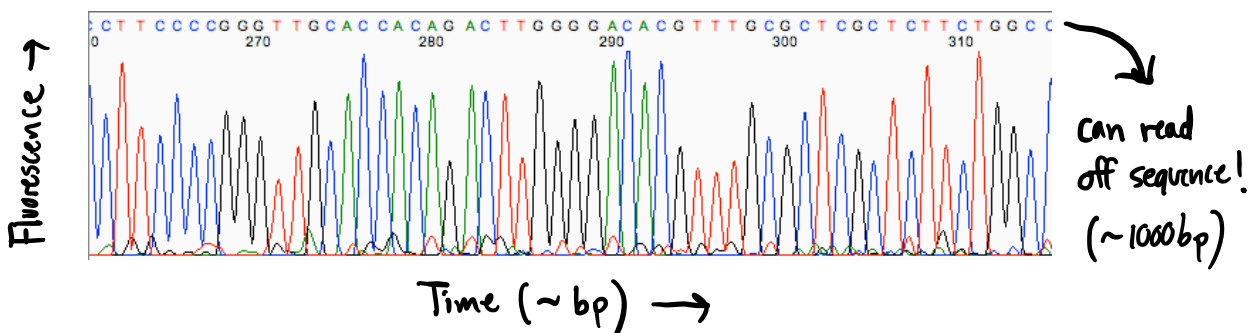
⇒ After 1 round of extension, random mixture of



⇒ Flow in electric field (shorter fragments move faster)



⇓ "chromato-gram"

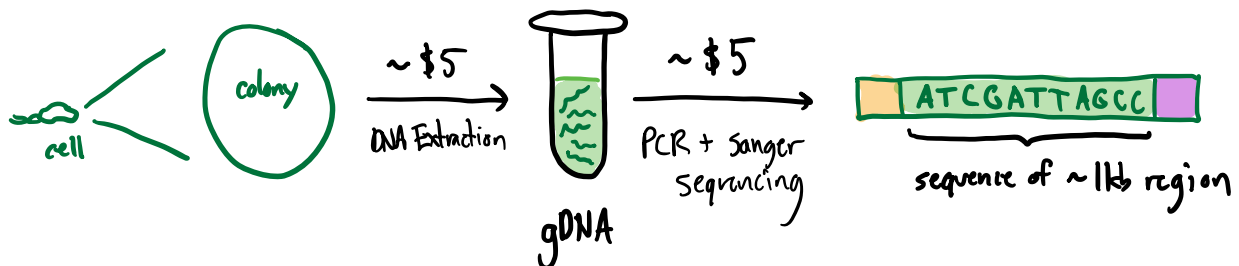


⇒ costs ~\$5 (send away overnight:



)

⇒ so have seen how to go from:



⇒ in Problem 1 of PSET 2, Lang & Murray used sanger seq to sequence URA3 genes in ~300 yeast colonies

⇒ \$1500 total

⇒ But expensive to sequence whole genomes!

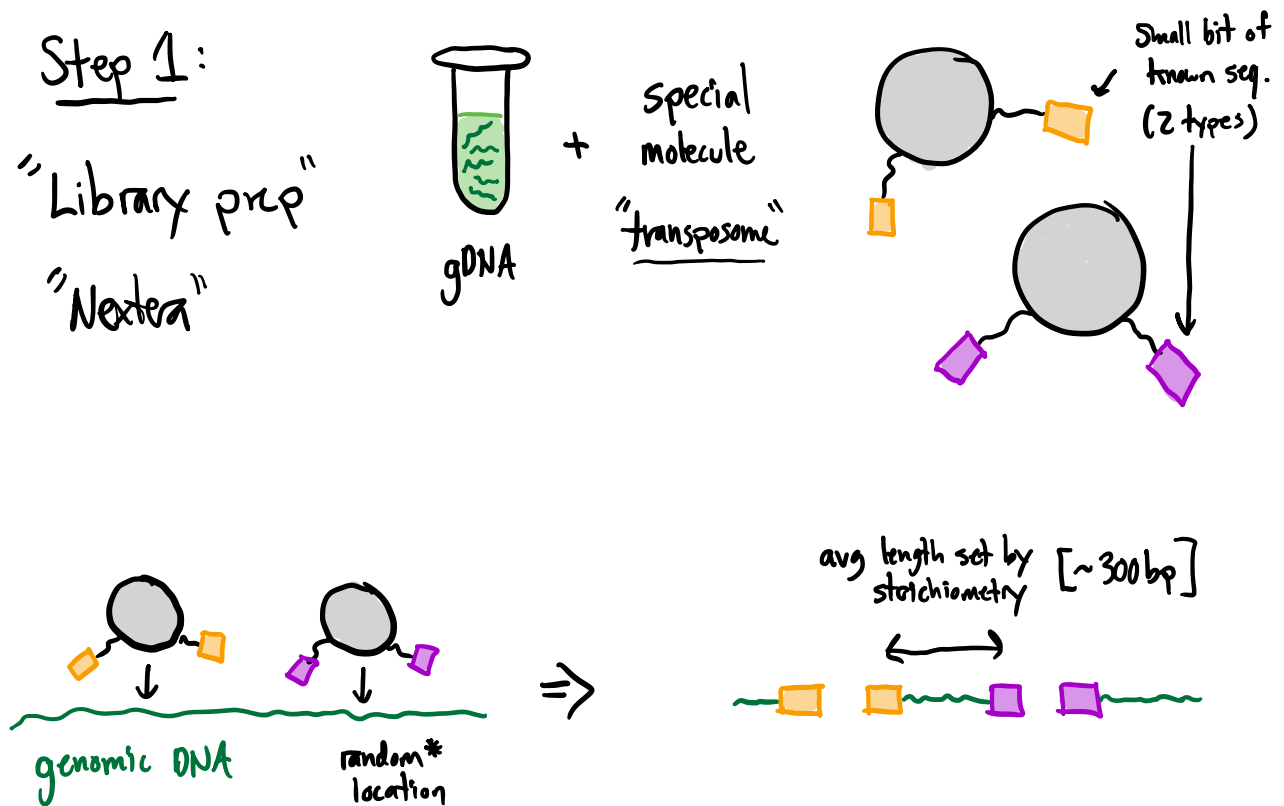
e.g. E. coli: 1 clone =  $10^6$  bp ⇒  $10^3$  Sanger seq runs → \$5k

e.g. Humans: 1 person =  $10^9$  bp ⇒  $10^6$  Sanger seq runs → \$5M

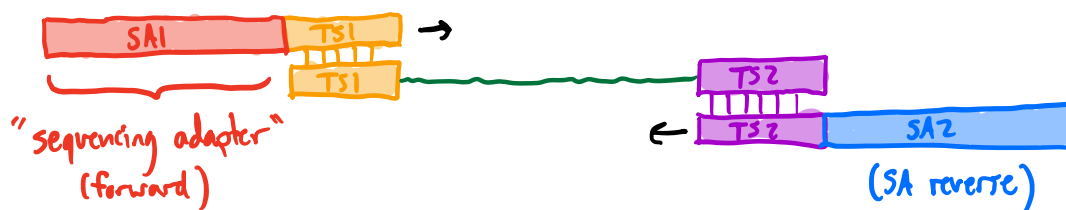
⇒ Now things are much cheaper w/ Next-gen sequencing  
("Illumina sequencing")

⇒ same idea, but higher throughput!

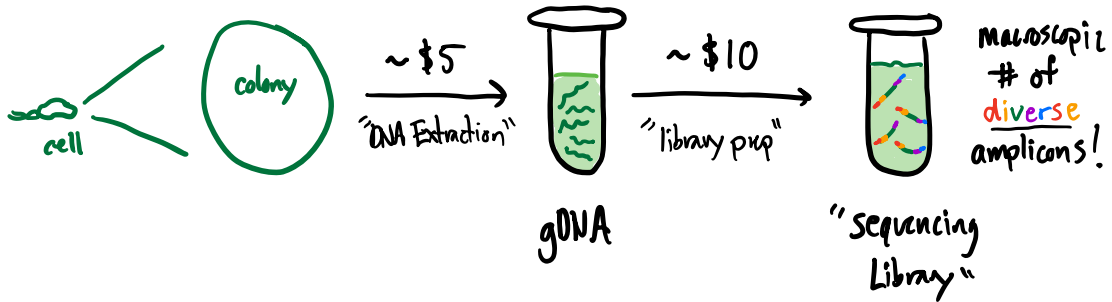
Step 1:  
"Library prep"  
"Nextera"



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

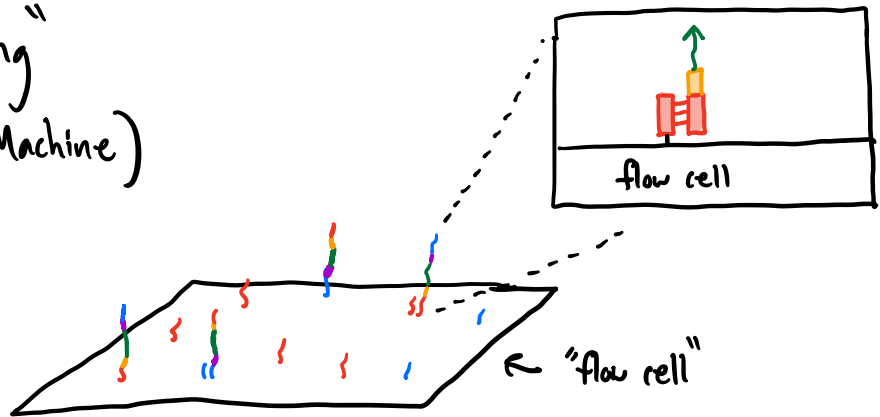


⇒ This gives us:



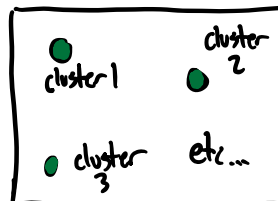
Step 3: "Sequencing"  
(Illumina Machine)

Flow library onto chip  
w/ probes that  
bind sequencing adapter



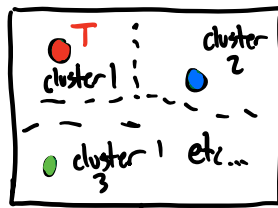
Step 4: do more PCR on bound fragments to turn each molecule into macroscopic cluster of identical\* molecules

top down  
view:



Step 5: flow primer + fluorescent dNTPs that incorporate once (then stop)  
(SA1)

⇒ if take picture of flow cell:



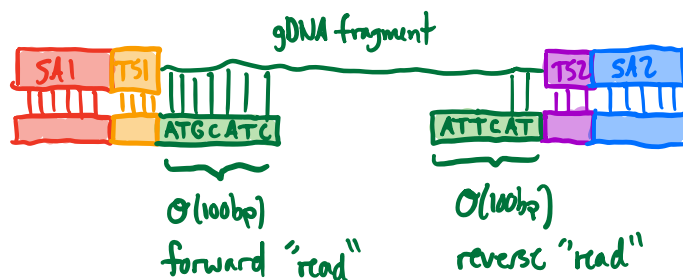
⇒ different colors = nucleic acid @ first position

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

Step 7: can repeat for  $\mathcal{O}(100)$  cycles [until pictures  $\approx$  desynchronized]

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

Net result:

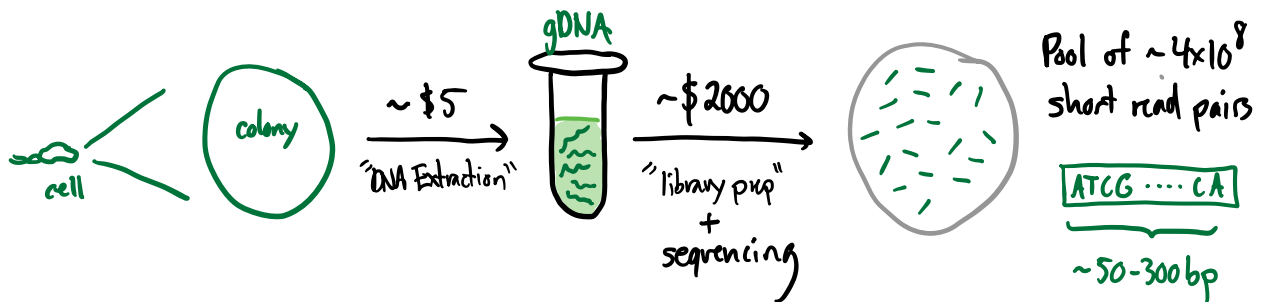


⇒ get  $\mathcal{O}(100\text{ bp})$  read out from each end of single DNA fragment

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

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

Recap: Now have method for going from:



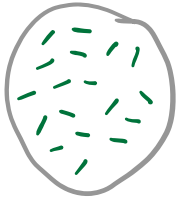
Next: what can we do with this kind of data?

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



# ① "De novo Genome Assembly"

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



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

ATCG ... CA

$\sim 50-300$ bp

Idea: look for reads that partially overlap

(e.g.  $\approx 20$ bp)



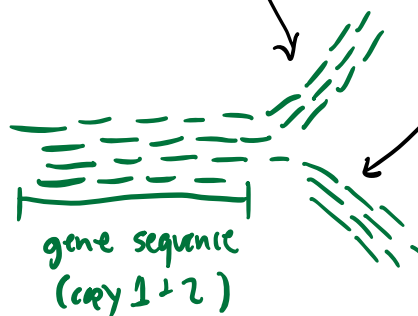
consensus sequence  
= "contig"

$\Rightarrow$  simple in principle, but lots of corner cases ...

$\Rightarrow$  e.g. what if 2 regions of genome are identical  
for  $\geq 100$ bp (or length of read?)



$\Rightarrow$

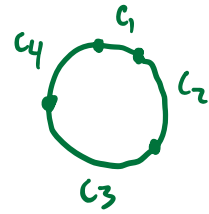


"fork in assembly"

⇒ for these & other reasons, assembly typically yields collection of discrete contigs,

each  $\sim 10^3 - 10^5$  bp long: 

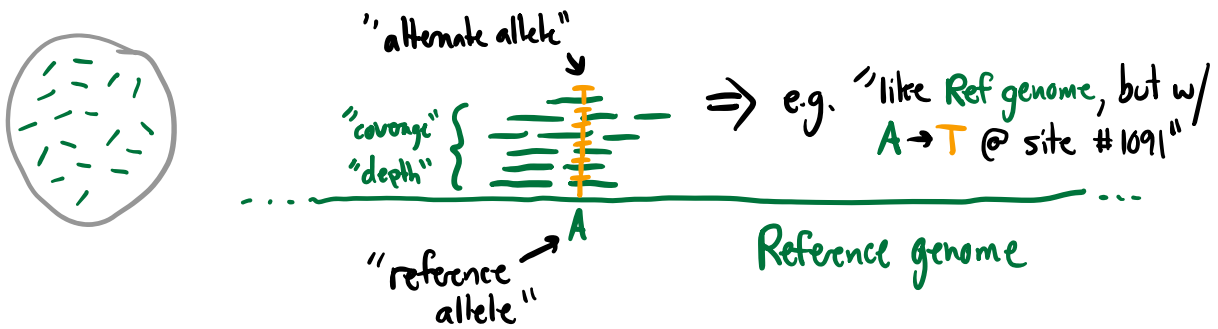
⇒ much harder (+ manual effort)  
to "finish" into complete genome



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

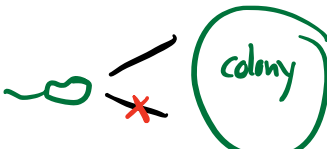
## ② Alignment of reads to Reference Genome

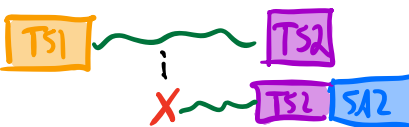
⇒ if already have assembled genome from related strain...  
can align reads to best-matching place & look for changes  
(common programs: BLAST, Bowtie2, BWA-MEM, + mpileup)



\* Major wrinkle: sequencing errors! (A → T by chance; "shot noise")

⇒ where could these errors come from?

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

② library prep + PCR: 

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

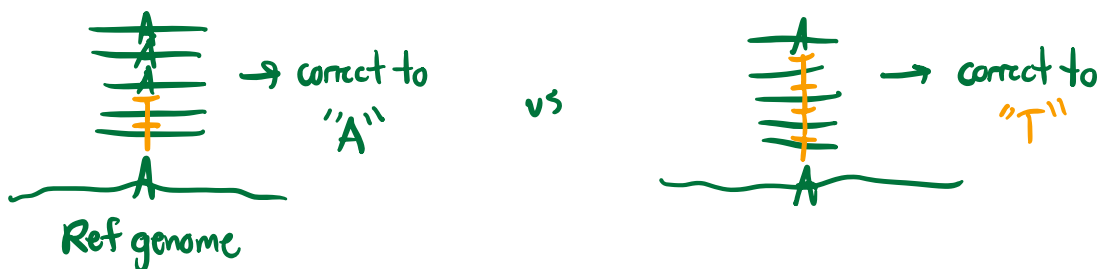
③ Errors on sequencing machine ⇒ estimated  $\sim 10^{-3}$  / bp  
(e.g. PCR errors during cluster formation, wrong fluorescent base, etc...)  
*\* but* varies from site-to-site  
(as high as  $\sim 10^{-2}$  @ some sites...)

⇒ often dominant source of noise ⇒  $P_{err} \sim 10^{-3}$  (up to  $10^{-2}$  in special cases...)

⇒ Upshot: low rate per site, but # sites  $L \gg 1$ !

⇒ # errors per genome =  $L \times P_{err} \sim 10^3$  errors / E. coli genome  
 $\sim 10^6$  errors / human genome

⇒ fortunately, can correct errors w/ consensus sequences!



⇒ higher coverage is helpful... how much do we need?

Fermi problem :

$$P_{\Gamma}(\text{consensus error}) = \sum_{A=\bar{D}/2}^{\infty} \frac{(\text{Perr} \times \bar{D})^A}{A!} e^{-\text{Perr} \cdot \bar{D}} \approx \frac{(\text{Perr} \times \bar{D})^{\bar{D}/2}}{(\bar{D}/2)!} e^{-\text{Perr} \cdot \bar{D}}$$

" errors in >50% of reads

(Poisson distribution)

↳ dominated by 1<sup>st</sup> term

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

$$\approx \exp\left[\log L - \bar{D}/2 \log\left(\frac{1}{2e \cdot \text{Perr}}\right)\right] \Rightarrow \text{want} \lesssim \mathcal{O}(1)$$

$$\Rightarrow \text{solve for } \bar{D} \Rightarrow \bar{D} \sim \frac{2 \log L}{\log\left(\frac{1}{2e \cdot \text{Perr}}\right)} \approx \begin{cases} 6 & \text{if } L \sim 10^6, \text{Perr} \sim 10^{-3} \\ 8 & \text{if } L \sim 10^9 \end{cases}$$

Upshot: need  $\gtrsim 10\times$  coverage to detect single mutations...  
(e.g. Problem 2 of HW 3)

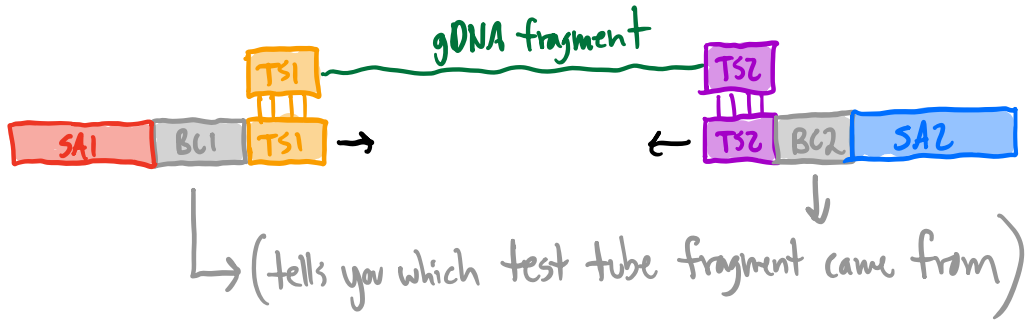
$\Rightarrow$  How many reads is this?

$$\text{E. coli: } 10 \frac{\text{reads}}{\text{site}} \times 10^6 \text{ bp} \Rightarrow 10^7 \text{ bp} \Rightarrow 10^5 \text{ reads}$$

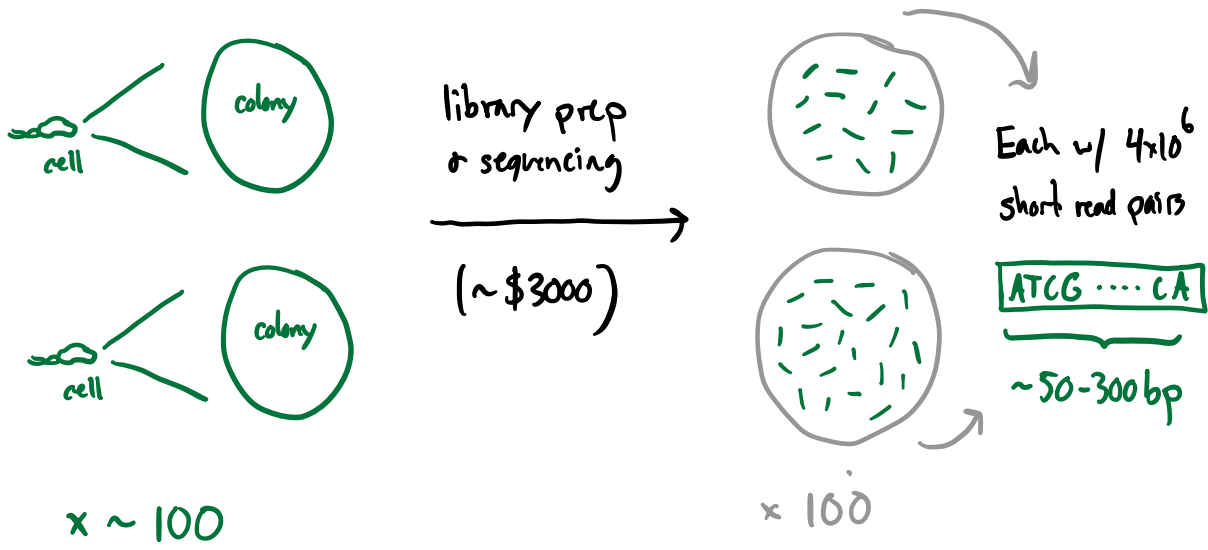
$\Rightarrow$  but a single run of Illumina seq produces  $4 \times 10^8$  reads...

$\Rightarrow$  overkill! (wastes money...)

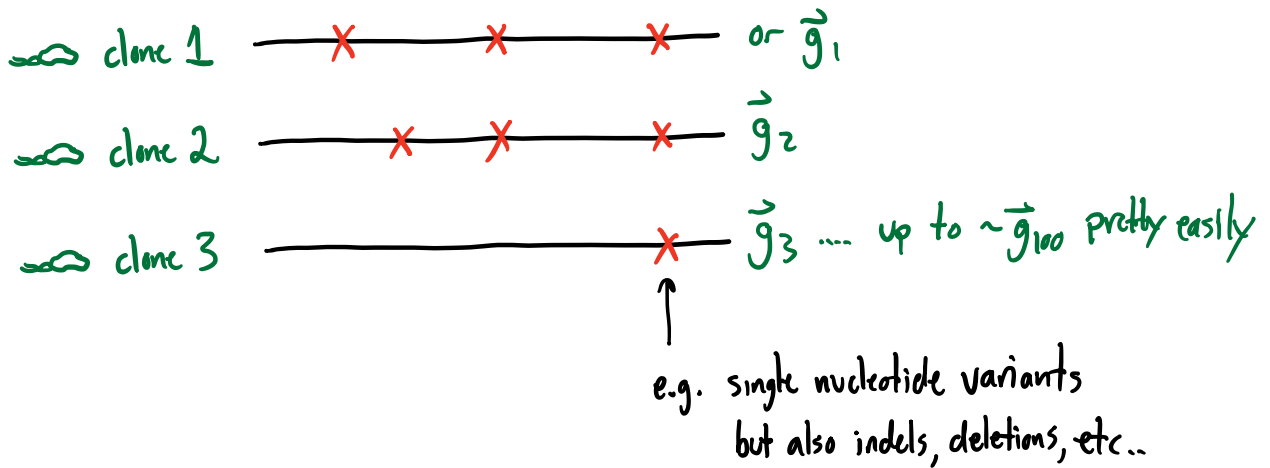
Solution: "multiplex": Add sample specific "barcode" sequence during library prep step



Upshot: can sequence ~100 E. coli libraries on one flow cell & get 300-fold coverage of E. coli genome



⇒ After aligning reads & detecting "true" mutations, get sequences of genomes:



Question:

How are sampld genomes related to dist'n of genomes in pop'n?

⇒ let  $n_{\vec{g}} \equiv$  # sampled clones w/ genome  $\vec{g}$  (random from sampling)

$f(\vec{g}) \equiv$  frequency of genome  $\vec{g}$  in population (random from evolution)

⇒ then sampling process  $\sim$  multinomial distribution

$$Pr[\{n_{\vec{g}}\} | n, \{f(\vec{g})\}] \propto \prod_{\vec{g}} \frac{f(\vec{g})^{n_{\vec{g}}}}{n_{\vec{g}}!}$$

↑  
 total # clones sampled

⇒ Genotype space is huge! (e.g.  $\vec{g} \in \{0,1\}^L$ )

⇒ often coarse-grain using summary statistics:

Eg 1:  $n_e \equiv$  # individuals w/ mut'n @ site  $e \equiv \sum_{\vec{g}} g_e n_{\vec{g}}$

↳  $\begin{cases} 1 & \text{if mut @ site } e \\ 0 & \text{else.} \end{cases}$

⇒ can show  $\Pr[n_e | n, \{f(\vec{g})\}] = \text{Binomial}(n, f_e)$

where  $f_e \equiv \sum_{\vec{g}} g_e f(\vec{g}) =$  frequency of mut'n @ site  $e$

Eg 2: total # of mutations separating 2 genomes

⇒ since depends on genome length, often normalized by  $L$ :

$\frac{\text{\# mutations between 2 random clones}}{L} = \begin{cases} \text{"heterozygosity"} (\pi) & \text{if from same pop'n} \\ \text{"divergence"} (d) & \text{if from diff "species" (or isolated sub-pops)} \end{cases}$

e.g. heterozygosity ( $\pi$ ) in humans is  $\sim 10^{-3}$

divergence ( $d$ ) between humans + chimps is  $\sim 10^{-2}$

heterozygosity ( $\pi$ ) between E-coli in different humans is  $\sim 10^{-2}$



Can we relate  $\pi$  to the genotype distribution,  $f(\vec{g})$ ?

$\Rightarrow$  Note that for specific pair of genomes,  $\vec{g}_1 + \vec{g}_2$ :

$$\pi(\vec{g}_1, \vec{g}_2) = \frac{1}{L} \sum_{e=1}^L \left[ \underbrace{g_{1e} \cdot (1 - g_{2e})}_{\vec{g}_1 \text{ has mut'n @ } e \text{ \& } \vec{g}_2 \text{ doesn't}} + (1 - g_{1e}) \cdot g_{2e} \right]$$

$\Rightarrow$  Average over sampling process ( $\vec{g}_1, \vec{g}_2 | \{f(\vec{g})\}$ ):

$$\langle \pi | \{f(\vec{g})\} \rangle_{\substack{\text{sampling} \\ \vec{g}_1, \vec{g}_2}} = \frac{1}{L} \sum_{e=1}^L \left[ \underbrace{\langle g_{1e} (1 - g_{2e}) \rangle}_{f_e (1 - f_e)} + \underbrace{\langle (1 - g_{1e}) g_{2e} \rangle}_{(1 - f_e) f_e} \right]$$

$$= \frac{1}{L} \sum_{e=1}^L 2 f_e (1 - f_e) \quad \text{Remember: } f_e \equiv \sum_{\vec{g}} g_e f(\vec{g})$$

$\Rightarrow$  Average over evolution ( $f_e$ ):

$$\langle \pi \rangle = \frac{1}{L} \sum_{e=1}^L \langle 2 f_e (1 - f_e) \rangle = \frac{1}{L} \sum_{e=1}^L \int 2 f_e (1 - f_e) p(f_e) df_e$$

What to choose for  $p(f_e)$ ?

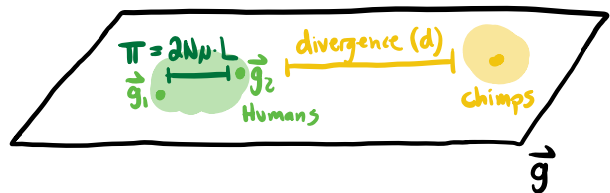
$\Rightarrow$  e.g. if genome is collection of neutral sites ( $s=0$ )

$$\Rightarrow p(f_e) \approx \frac{2N_e\mu}{f_e} \quad (\text{"quasi-stationary" dist'n from last lecture})$$

$$\Rightarrow \langle \pi \rangle = \int_0^1 2s(1-s) \cdot \frac{2N_e\mu}{s} \cdot ds = 2N_e\mu \cdot \int_0^1 2(1-s) ds$$

$$\Rightarrow \text{Neutral heterozygosity: } \langle \pi \rangle = 2N_e\mu \quad ***$$

population can only spread out so far 'in genotype space!



$\Rightarrow$  Note: doesn't agree w/ human data if  $N=10^9$ ...

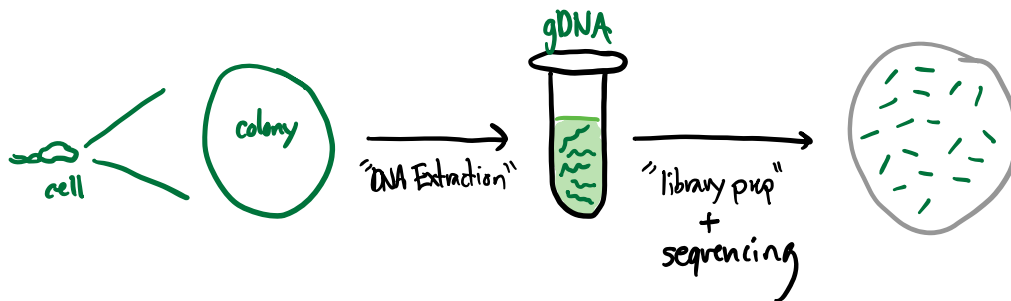
$$\Rightarrow \text{often used to "fit" } N_e \equiv \frac{\langle \pi \rangle}{2\mu} \xrightarrow{\text{humans}} \frac{10^{-3}}{2 \times 10^{-8}}$$

(\* beware! \*)

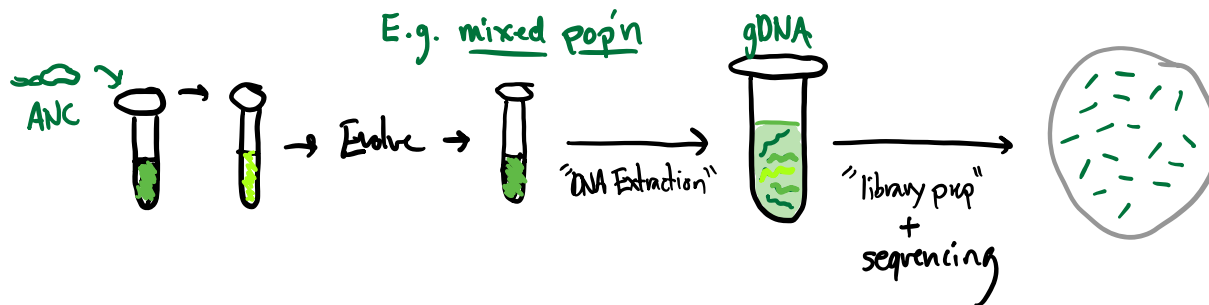
$\Rightarrow$  Variance (or dist'n) of  $\pi$  is much more complicated!

$\Rightarrow$  correlations between  $g_e + g_{e'}$   $\Rightarrow$  will see more later!

So far, have focused on clones



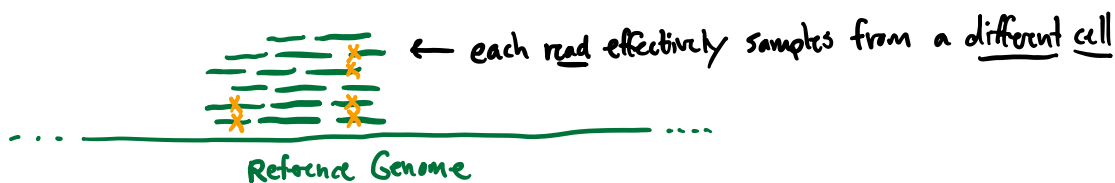
⇒ lots of other things we could put in our library prep...



⇒ known as "pooled sequencing" / "metagenomic sequencing"

⇒ in this case, assembly much harder!

⇒ Payoff comes from reference mapping:



e.g. if  $A_\ell \equiv$  # reads w/ mutations @ site  $\ell$

$$\Rightarrow \Pr[A_\ell | D_\ell, \{f(\hat{g})\}] = \text{Binomial}(D_\ell, f_\ell) + \text{sequencing/PCR errors}$$

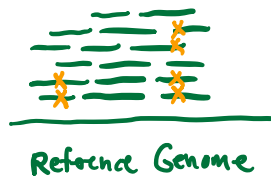
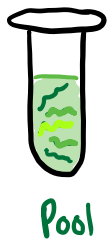
↑  
total coverage  
@ site  $\ell$

$\Rightarrow$  since we can sequence  $\sim 100$  E. coli genomes  
@  $>100x$  coverage in 1 run of Illumina sequencing

$\Rightarrow$  can effectively sample  $\sim 100$  clones  
 $\sim 100x$  more cheaply by sequencing pools!

(much cheaper way to track freqs of individual mut'ns)

### Downsides:

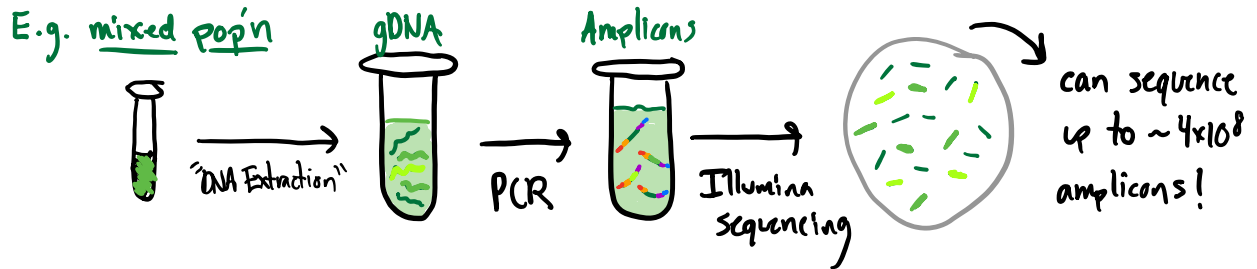


① sequencing errors!

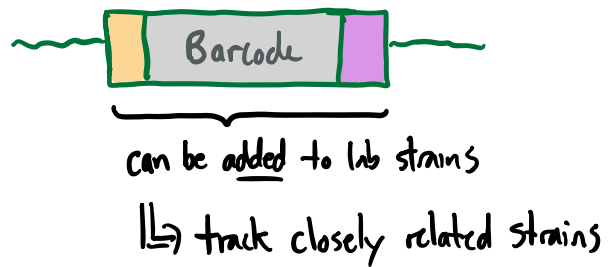
② which mutations are in same cells? ("linkage information")



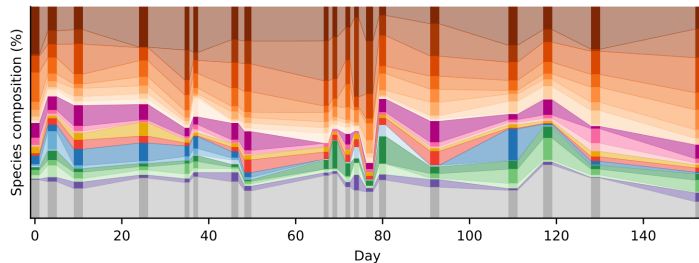
Can also sequence pools of amplicons :



Two common targets :



E.g. 1 person's microbiome over time

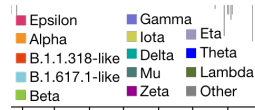
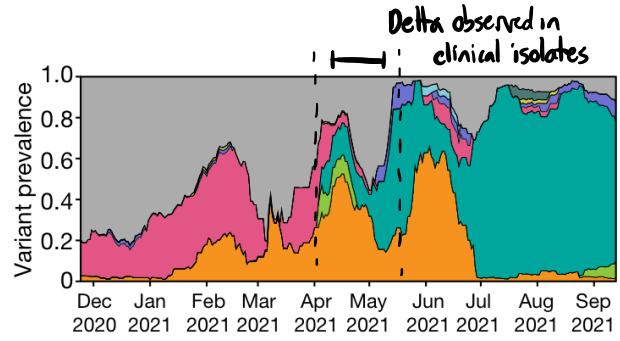
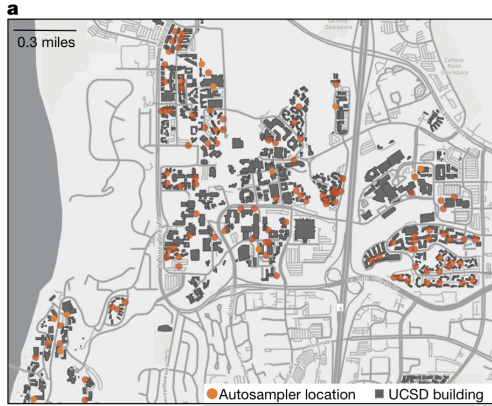


- *Bacteroides vulgatus*
- *Bacteroides coprocola*
- *Bacteroides uniformis*
- *Bacteroides cellulosilyticus*
- *Bacteroides eggerthii*
- *Bacteroides faecis*
- *Bacteroides massiliensis*
- *Bacteroides caccae*
- *Alistipes sp*
- *Alistipes onderdonkii*
- *Alistipes finegoldii*
- *Parabacteroides distasonis*
- *Paraprevotella clara*
- *Butyrivibrio crossotus*
- *Coprococcus sp*
- *Coprococcus comes*
- *Eubacterium rectale*
- *Eubacterium siraeum*
- *Eubacterium eligens*
- *Phascolarctobacterium sp*
- Other

E.g. Homework Problem

# Wastewater sequencing reveals early cryptic SARS-CoV-2 variant transmission

Published online: 7 July 2022



MGX sequencing + read alignment