# Sequencing & Genomics II

Last time: Next-gen / Illumina sequencing of bacterial isolates / "clones"
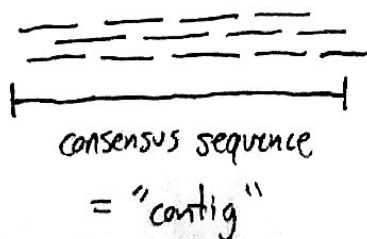
single cell → colony $\xrightarrow[\text{"DNA extraction"}]{\sim\$5}$ genomic DNA $\xrightarrow[\substack{\text{"library} \\ \text{prep"}}]{\sim\$10}$ sequencing library $\xrightarrow{\sim\$2000}$ → 1 read $\sim O(100bp)$

(~free)

$\sim 4 \times 10^8$ ~~~~ read pairs ("short reads")

genomic DNA molecule

| SA1 | TS1 | ———— | TS2 | SA2 |

→ R1    ← R2

ATCGATT···GC

← $\sim 50-300bp$ →

What can we do w/ this kind of data?
  ⟹ need to put puzzle back together...  2 main methods

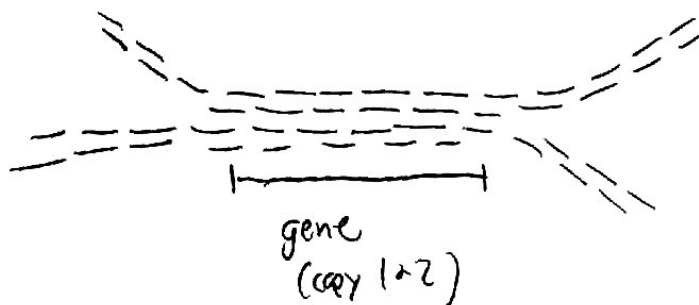① ⬤ De Novo Genome Assembly  (common programs: Spades/Velvet ...)

  ⟹ look for overlapping reads ($\gtrsim 20bp$)

consensus sequence
= "contig"

Simple in principle, but lots of corner cases...

   e.g. what if 2 regions of genome are identical for
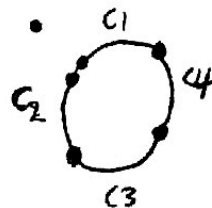     more than 100 bp (or 1 read)?

                     gene copy 1                gene copy 2

                                             &larr; fork in assembly

                gene
              (copy 1 & 2)

$\Longrightarrow$ for these & other reasons, assembly often results in
collections of different contigs                    
                            contig 1    contig 2    contig 3   etc...

each $\sim 1000 - 10^5$ bp long.

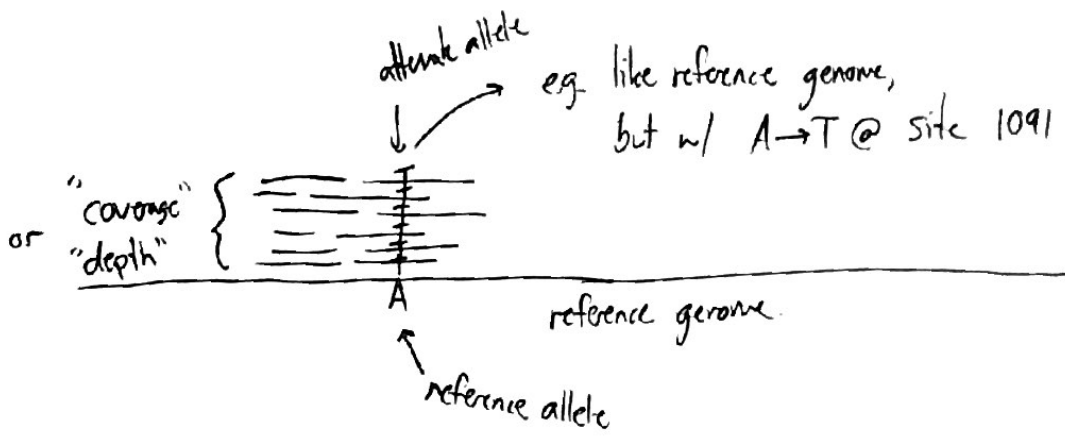       $\sqsupset\!\!\Longrightarrow$ much harder (+ manual effort) to "finish" into complete

        genome:
                         C1
                                    C4
                     C2
                            C3

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

# ② Alignment of Reads to Reference Genome

⟹ if already have assembled genome from closely related strain, can align reads to best matching place in genome & look for changes: ( ~~common~~ common programs: Bowtie2, BWA-MEM

+ mpileup (samtools) )



alternate allele
↓
eg. like reference genome, but w/ A→T @ site 1091

"coverage" or "depth" {

reference genome.

↑ reference allele

⟹ this is $O(\#reads)$ & much lower memory footprint. (laptop)

⟹ saw that ~20bp sufficient to locate most sites in E.coli, so ~100bp reads are mostly ok.
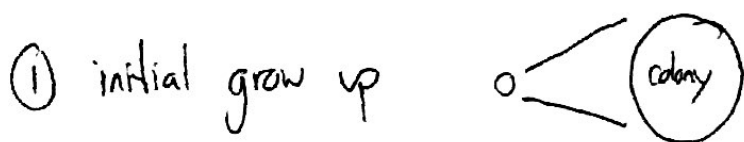
( like difference between putting puzzle together from scratch, & putting together a puzzle when the completed one is right next to you. )

⟹ still some corner cases ( but small fraction of genome ) & works best when ref genome is "close" θ to sample ( ~θ at most 1-2 diffs per read

## one winkle: sequencing errors

$\Rightarrow$ ability to sequence single molecules has its drawbacks:

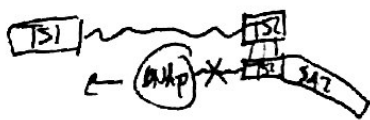random molecular errors ("shot noise") can lead to e.g. $A \to T$ by chance.

## where could these errors come from?

① initial grow up    o $<$ (colony)    e.g. mutation during first division, 2nd, etc.
$\Rightarrow$ Luria-Delbrück process

but happen @ low rate, $N \sim 10^{-10}$/bp/gen

② Library prep PCR    error during early rounds of PCR ("PCR errors")
also Luria-Delbrück-like process.

[TS1]   [TS2]
← (AMP) × [TS] [SA2]

higher rates, $\sim 10^{-6}$/bp/round $\sim 10^{-4}$ total. ← presumably this too

③ Error on sequencing machine    $\Rightarrow$ estimated to be $\sim 10^{-3}$/bp, but
(e.g. cluster generation PCR error
→ wrong fluorescent base, etc.)    can vary a lot from site-to-site.

$\Rightarrow$ dominant source of error

$\Rightarrow$ still low rate per site, but lots of sites in genome, so expect
@ least $L \times P_{err} \sim 10^{3}$ in Mb-sized genome.

(big problem for detecting single mutations...)

Fortunately, can correct many of these errors by taking consensus across independent reads covering same site:

```
—A—
—A—
—A—   →  "A"        vs
—T—
=T=
————
   A
```

```
—T—
=T=
—T—   →  "T"
—A—
————
   A
```

⟹ higher coverage is helpful.

⟹ how much coverage necessary before we expect $\leq 1$ consensus errors in whole genome?

error in consensus require $\geq \frac{1}{2}$ of all reads to have an error.

⟹ if ~~coverage~~ coverage is ~~Poisson~~ Poisson process w/ mean $\bar{D}$

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

$$\# \text{errors} = L \times \Pr\left(\text{consensus error}\right) = \exp\left[\log L + \frac{\bar{D}}{2}\log(P_{err}\bar{D}) - P_{err}\bar{D} - \frac{\bar{D}}{2}\log\left(\frac{\bar{D}}{2}\right) + \frac{\bar{D}}{2}\right]$$

$$\approx \exp\left[\log L - \frac{\bar{D}}{2}\left[\log\left(\frac{1}{2P_{err}}\right) - 1\right]\right]$$
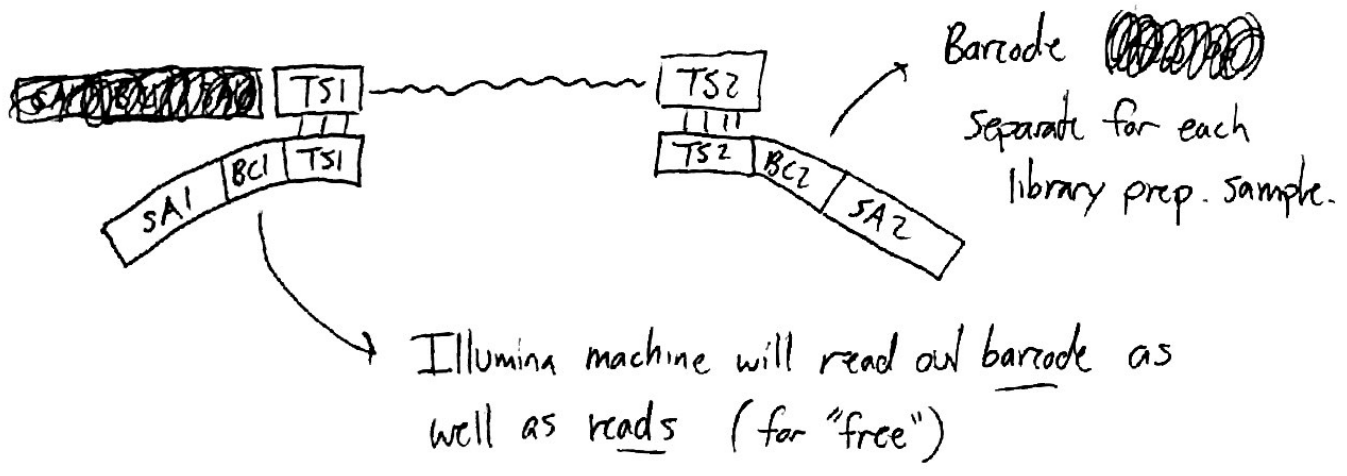
⟹ crosses 1 when $\bar{D} \approx \frac{2\log L}{\log\left(\frac{1}{2P_{err}}\right) - 1} \approx \begin{cases} 5 \text{ if } L \approx 10^6, P_{err} \approx 10^{-3} \to 10 \text{ if } P_{err} \approx 10^{-2} \\ 8 \text{ if } L \approx 10^9 \to \bullet \end{cases}$

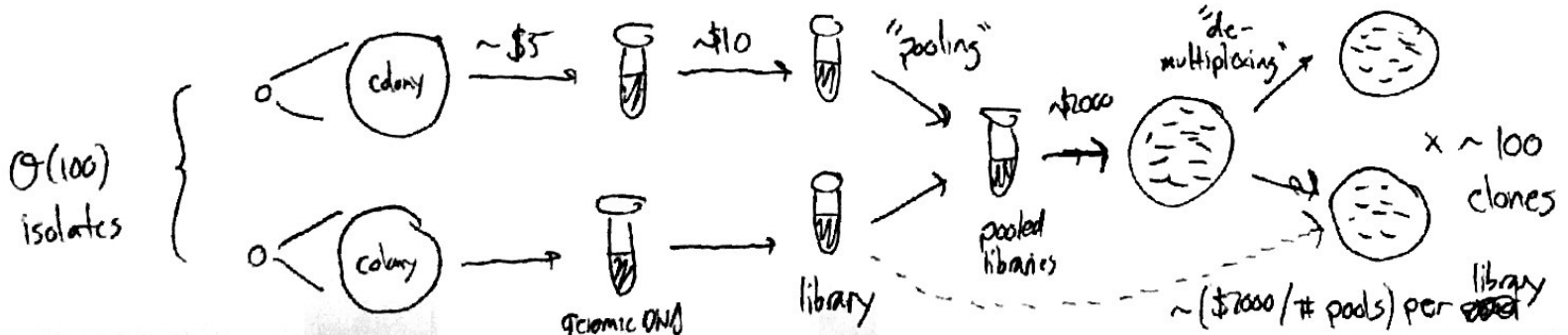So need coverage $\geq 10$ to ~~eliminate~~ eliminate errors (though still weird locs)

$\Rightarrow$ how many reads is this?    ~~~~ $10 \times 10^6$ bp sequenced. $= 10^7$ bp.

$\approx \rightarrow 10^5$ reads per genome.

$\Rightarrow$ but single run generates $\sim 10^8$ reads! ~~~~ overkill. (waste money)
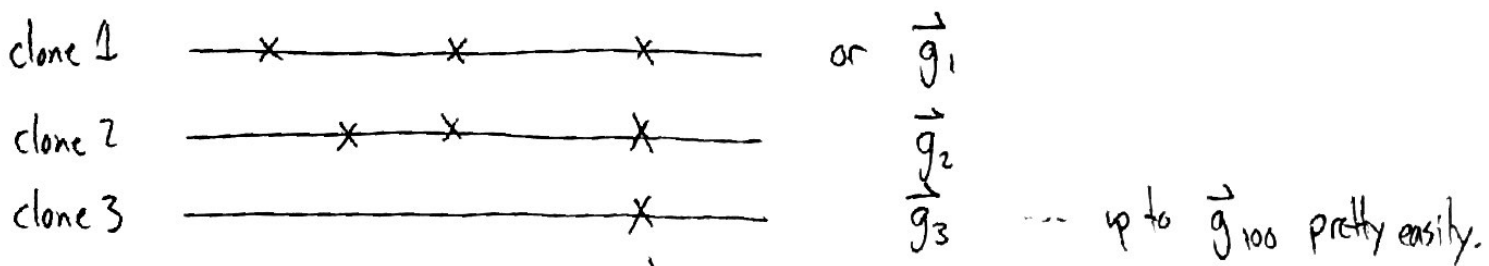
Solution: "multiplexing" w/ DNA barcodes. can design sequencing adapter w/ special sequence of letters (barcode) used to distinguish different samples.



Barcode ~~~~ separate for each library prep. sample.

Illumina machine will read out barcode as well as reads (for "free")

$\Rightarrow$ can then tell which library each ~~~~ read pair came from.



$\theta(100)$ isolates

$\Rightarrow$ after aligning reads, detecting "true" mutations, get sequences

clone 1 ——×————×————×———— or $\vec{g}_1$

clone 2 ——×—×————×———— $\vec{g}_2$

clone 3 ————————×———— $\vec{g}_3$ ... up to $\vec{g}_{100}$ pretty easily.

e.g. single nucleotide variants
(but also indels, deletions, etc.)

How are these related to
the distribution of genotypes
in our population?

$\Rightarrow$ ~~xxxxxxxxxxxx~~ if clones are sampled randomly, then

$$Pr\left[\{n_{\vec{g}}\} \mid n, \{f(\vec{g})\}\right] = \underset{\text{dist'n}}{\text{multinomial}} \propto \prod_{\vec{g}} f(\vec{g})^{n_{\vec{g}}} \frac{1}{n_{\vec{g}}!}$$

where $f(\vec{g})$ = fraction of population w/ genotype $\vec{g}$ (random from evolution)

$n_{\vec{g}}$ = # of ~~clones~~ in sample w/ ~~genotype $\vec{g}$~~ (random from sampling)
                clones                     genotype $\vec{g}$

$\Rightarrow$ genotype space is huge, so often reduce to summary <u>statistics</u>.

e.g. # of mutations separating 2 genomes

Since depends on length of genome, often normalize by L:

$$\frac{\text{\# mutations between 2 randomly sampled clones}}{L} = \begin{cases} \text{"heterozygosity" } (\pi) \text{ if from same pop'n} \\ \\ \text{"divergence" } (d) \text{ if from diff. "species"} \\ \qquad (\text{\# isolated pop'ns}) \end{cases}$$

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

divergence between humans & chimps is $\sim 10^{-2}$

heterozygosity among E.coli from different humans is $\sim 10^{-2}$

to relate $\pi$ to genotype freqs, $f(\vec{g})$, note that

$$\pi = \frac{1}{L} \sum_{\ell=1}^{L} \left[ \left[ g_{1\ell}(1-g_{2\ell}) \right] + \left[ (1-g_{1\ell})g_{2\ell} \right] \right] \longrightarrow \text{two ways of sampling genomes that differ at site } \ell.$$

$\searrow \in \{0,1\}$

So on average,

$$\langle \pi \rangle = \frac{1}{L} \sum_{\ell=1}^{L} \left[ \langle g_{1\ell}(1-g_{2\ell}) \rangle + \langle (1-g_{1\ell})g_{2\ell} \rangle \right]$$

$$= \frac{1}{L} \sum_{\ell=1}^{L} 2f_\ell(1-f_\ell) \longrightarrow \text{fraction of population w/ mutation @ site } \ell.$$

the $f_e$'s are themselves random (from evolution)

so technically, have only calculated $\langle \pi | \{f_e\} \rangle$

averaging over these, we have

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

If genome is collection of <u>neutral sites</u>, then $p(f_e) = p(f) = \frac{2N_e\nu}{f}$, and

$$\langle \pi \rangle = \langle 2f(1-f) \rangle = \int 2f(1-f) \cdot \frac{2N_e\nu}{f} df = 2N_e\nu$$

$\Rightarrow$ thus, connection between $\langle \pi \rangle$ and $N_e$

(sometimes people even call $\langle \pi \rangle$ $N_e$, as if $N_e$ were an empirical property of the data. this is really bad, and we should stop doing it )
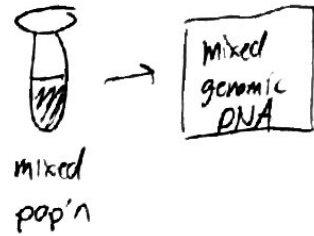
Note that Variance of $\pi$ is much more complicated, since it depends on correlations between $g_{1e}$ and $g_{1e'}$ @ different sites.

However, related summaries that are <u>linear</u> in sites, but involve bigger samples can still be calculated.
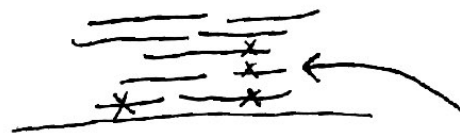
(will see examples later on )

So far, have focused on sequencing clones, but
lots of other things you could put in your library prep.

⟹ common one is <u>population</u> of bacteria

("pooled sequencing" or "metagenomic sequencing")

mixed pop'n → mixed genomic DNA

⟹ in this case, assembly can be very hard (since dealing w/ mixture of different genomes)

payoff comes from
⟹ reference mapping :

for typical coverages
& input pop'ns, each
read samples from a
different cell in pop'n.

e.g. # of reads w/ mutation @ site $\ell$

is $Pr(A \mid D, f_\ell) \approx Binomial(D, f_\ell) +$ sequencing error

e.g. if want to calculate $\langle \pi \mid f_\ell \rangle$,

$$\langle \pi \mid f_\ell \rangle = \left\langle \frac{A(D-A)}{\binom{D}{2}} \right\rangle + \text{sequencing error} = 2 f_\ell (1 - f_\ell) + \text{seq error}$$
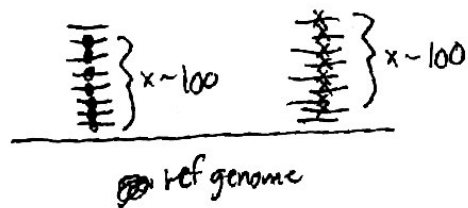
$\Rightarrow$ Since can get $\sim 100x$ coverage for $\sim 100$ E.coli genomes in 1 run of Illumina seq, can effectively sample $\sim 100$ clones $\sim 100x$ more cheaply by sequencing pools. ( _much_ cheaper way to ~~do~~ measure $\langle \Pi \rangle$ or to track frequencies of individual mutations.)

Downsides: also _much_ harder to distinguish low freq mutations from sequencing errors. ~~$\Pi \lesssim$~~ $\langle \Pi \rangle \gtrsim P_{err} \sim 10^{-3}$

If $P_{err} \sim 10^{-3}$, not even possible theoretically to measure freqs below this (even w/ infinite coverage) unless you do fancier things.

$\Rightarrow$ also lose information about which mutations are in same cells ("linkage information") unless you catch them on same sequencing read.

$\Rightarrow$ Sometimes can make progress w/ "pigeonhole principle", e.g. :

$\underset{\text{ref genome}}{\underline{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}}$  $\Big\} x \sim 100$  $\Big\} x \sim 100$

$\Rightarrow$ probably many cells w/ $\leftarrow\!\!*$ , since genotypes must add to 1:

$$1 = f(-) + f(\leftarrow) + f(-*) + f(\cdot *) = f(-) \cdot -f(\cdot *) + f_\bullet + f_x$$

$\underset{\text{must be} > 0!}{\nearrow}$  $\underset{\gg 1}{\underbrace{\phantom{xxx}}}$

$\Rightarrow$ show example data from Lenski's LTEE

Don't have to sequence mixed population of same species

$\Rightarrow$ nothing keeping you from extracting DNA from mixed community of bacteria in native community (e.g. fecal samples in gut microbiome)

("shotgun metagenomic sequencing")

$\Rightarrow$ in this case, since don't have to grow the bacteria, can work even when bacteria hard to grow in lab / frozen / dead / etc.

"culture independent sequencing" $\longrightarrow$ if genomes in sample sufficiently close (or sufficiently far), can use for de novo assembly to discover new bacteria/genes

often contrasted w/ amplicon sequencing (i.e. add Illumina sequencing adapters to PCR product)

$\quad$ $\llcorner\rightarrow$ commonly used target in microbiology is 16S Ribosomal RNA gene. all bacteria have it, & there are a few regions that are highly conserved across tree of bacterial life $\rightarrow$ good targets for primers.

$\Rightarrow$ amplicon metagenomic sequencing $\approx$ distribution of species* abundances in sample.