

### Capacity of Shotgun Sequencing Channels

$G = 0$  at which  
message is transmitted  
with lowest

probability possible

$$c = \max_{P_X(x)} I(X; Y)$$

input

output

message

code

**THESE**

### Shotgun sequencing :-

Breaking the DNA into smaller fragments without any specification

Running the Sanger method/ Nasy method on the smaller fragments generated.

sequence the smaller segments individually

These fragments are known as contigs.

After getting the sequence from all those different fragments,

the length and data of the sequence generated from the fragments is provided to the

Computer, ~~to~~ The computer will

find the overlapping regions present in these fragments.

After finding the overlapping regions, complete stretch of DNA seq can be obtained.

AGGCT

CT

C C

	c	A
--	---	---

AC

1

1

1

1

→

During Reconstruction,

there is  $p(\text{error})$  as

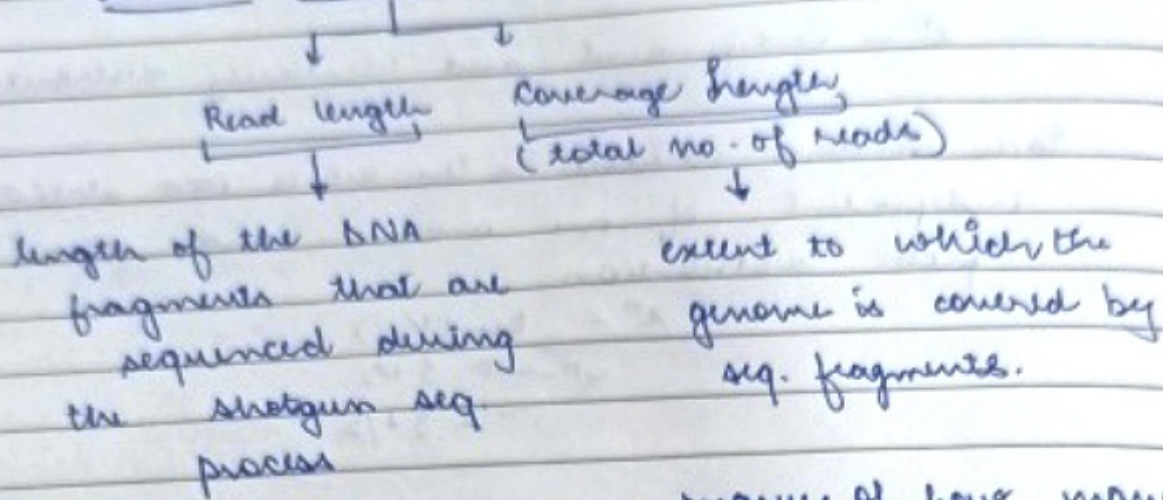
some sequences

obtained may get

erased.



## Reliable sequence reconstruction



Read length  $\uparrow$   
Contiguous coverage of the genome  $\uparrow$   
more challenging & costly to achieve in practice.

measure of how many times a specific position in a genome is covered by sequenced reads  
coverage length  $\uparrow$   
 $\downarrow$   
comprehensive seq of the genome

DNA mol - being seq - codeword from a predefined codebook

Utilizing coding techniques in DNA-based storage, number of reads required for reliable reconstruction of a binary seq. can be reduced from  $O(n)$  to  $O(n/\log n)$  making storage system more efficient & enabling higher capacity with fewer reads of shorter length



i.i.d. - Independent and Identically distributed

$x^n$  - independent and identically distributed

Each random variable for the set is ~~one~~ statistically independent of one another and follow same prob. distribution.

$$x^n \sim \text{Bern}(1/2)$$

$$x^n \rightarrow 0 \text{ } 1/2$$

$$\rightarrow 1 \text{ } 1/2$$

$$n \rightarrow \infty \quad L = L \log n$$

$$C = \frac{KL}{n}$$

↓

coverage depth

When  $L < 2 \rightarrow$  reconstruction is impossible

$L > 2 \rightarrow$  " " possible  $C_{LW} = \ln\left(\frac{n}{\epsilon}\right)$

$$C_{LW} = \ln\left(\frac{n}{\epsilon}\right)$$

all symbols in  $x^n$  are seq. with prob at least  $1-\epsilon$ .

### Feasibility region

$k \rightarrow$  factor that quantifies the fraction of the genome that is effectively covered by each read.

$O(n) / o(n/\log n) \rightarrow$  used to represent scaling behaviour or growth rate of certain quantities w.r.t size of input, which is denoted by  $n$



### Shuffling - Sampling Channel

The transmitted symbols are randomly rearranged or permuted during the comm. process.

no. of possible distinct ~~as~~ length sequences

$$\downarrow$$
$$2^{\bar{L} \log n} = n^{\bar{L}} = o(n / \log n) \Rightarrow o(k)$$

$$\frac{n^{\bar{L}} \log n}{n} \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\frac{n^{\bar{L}}}{k} \rightarrow 0 \text{ as } n \rightarrow \infty$$

~~$\frac{\log n}{n^{\bar{L}}} \rightarrow 0 \text{ as } n \rightarrow \infty$~~

$$n^{\bar{L}} = o(n / \log n)$$

$$\frac{n^{\bar{L}}}{k} \rightarrow 0 \text{ as } n \rightarrow \infty$$

when  $\bar{L} < 1$

probability that a given symbol in  $x^n$  is not seq. by any of the  $k$  reads

$$= \left(1 - \frac{k}{n}\right)^k$$

$$= \left(1 - \frac{L}{n}\right)^{cn/L} = \left(1 - \frac{L}{n}\right)^{\frac{c}{L} \frac{n}{n}}$$

$$= \lim_{n \rightarrow \infty} \left(1 - \frac{L}{n}\right)^{\frac{c}{L} \frac{n}{n}}$$

$$= \lim_{n \rightarrow \infty} \log e^{\frac{-c}{L} \frac{n}{n}}$$

$$= \boxed{e^{-c}}$$



$$C = \frac{K \bar{L} \log n}{n}$$

$$C' = \frac{K (\bar{L} - 1) \log n}{n}$$

$$C' = \frac{C}{\bar{L}} (\bar{L} - 1)$$

$$C' = C \left(1 - \frac{1}{\bar{L}}\right)$$

Even when  $\log n$  bits are removed, still the  $C$  of the channel would remain the same, which signifies the fact  $\log n$  is not providing any new information and is just helping in the ordering of the substrings.