

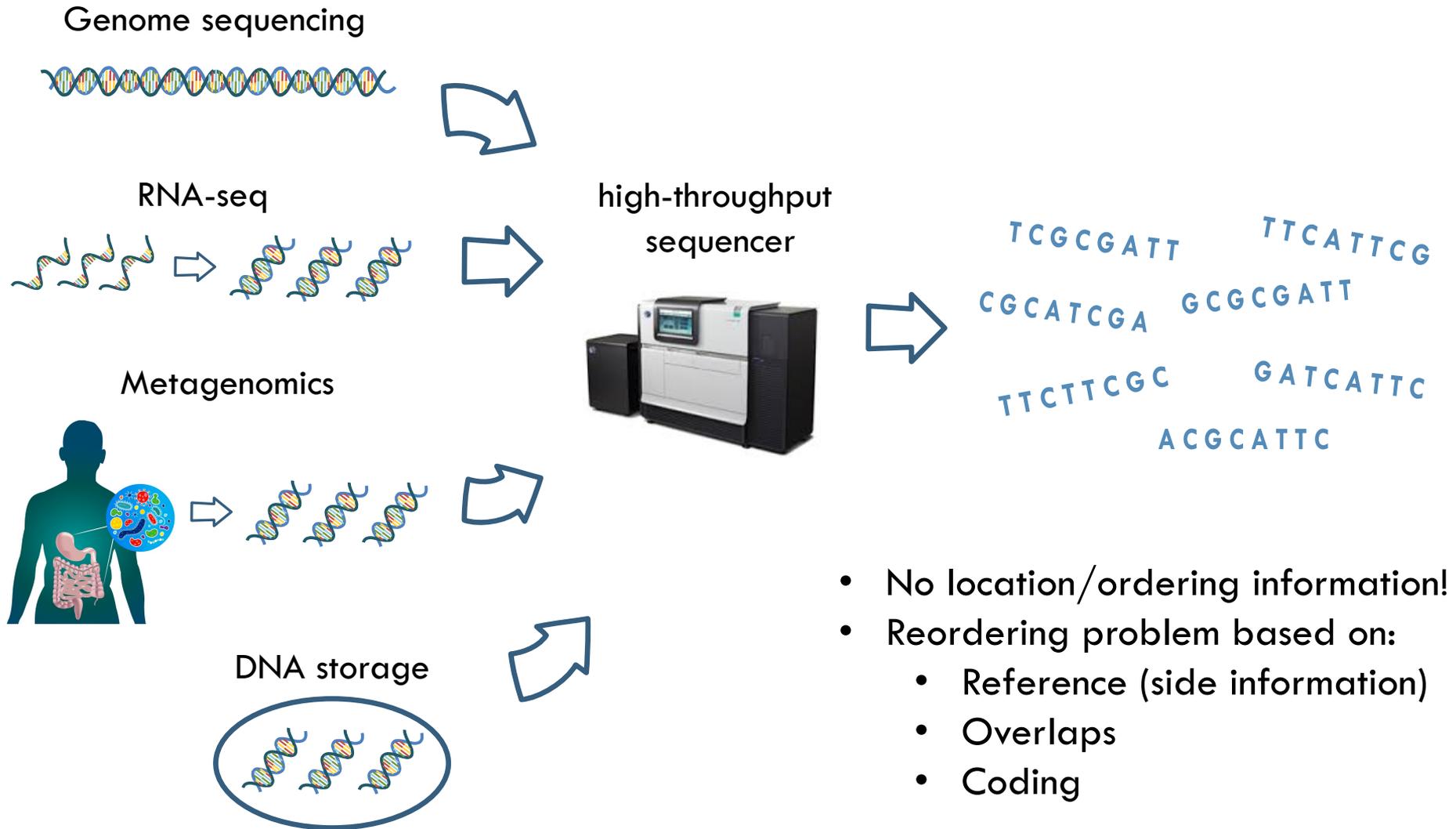
An Information Theory for Out-of-Order Information: DNA Data Storage and Beyond



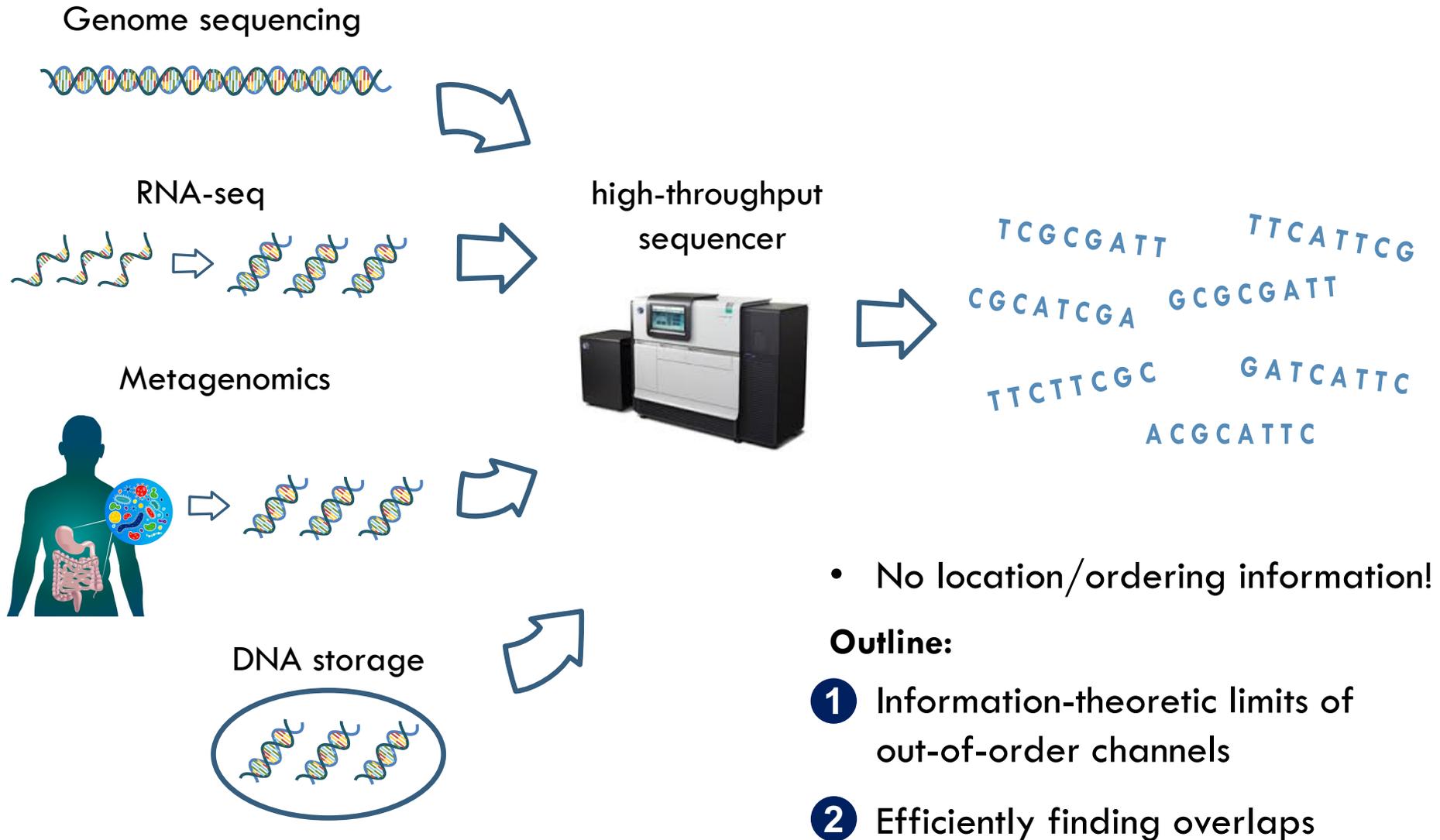
Ilan Shomorony
University of Illinois at Urbana-Champaign

Stanford ISL Colloquium – 03/16/23

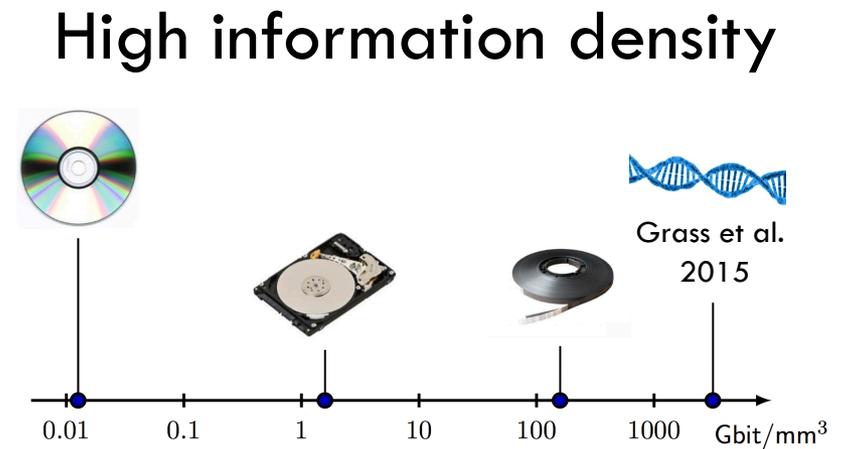
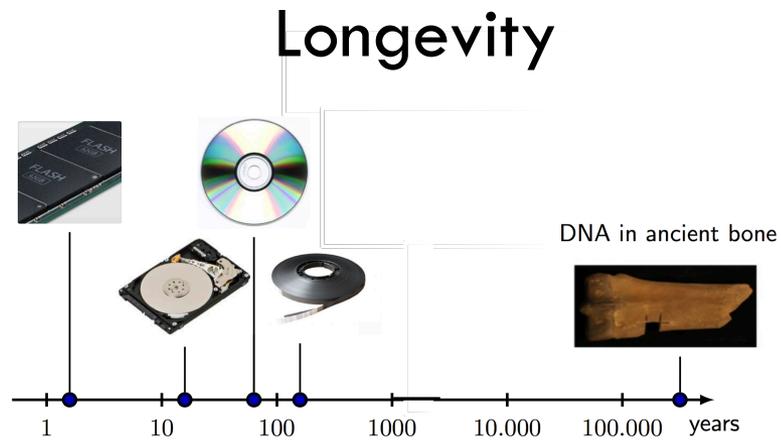
The wonders of high-throughput sequencing



The wonders of high-throughput sequencing

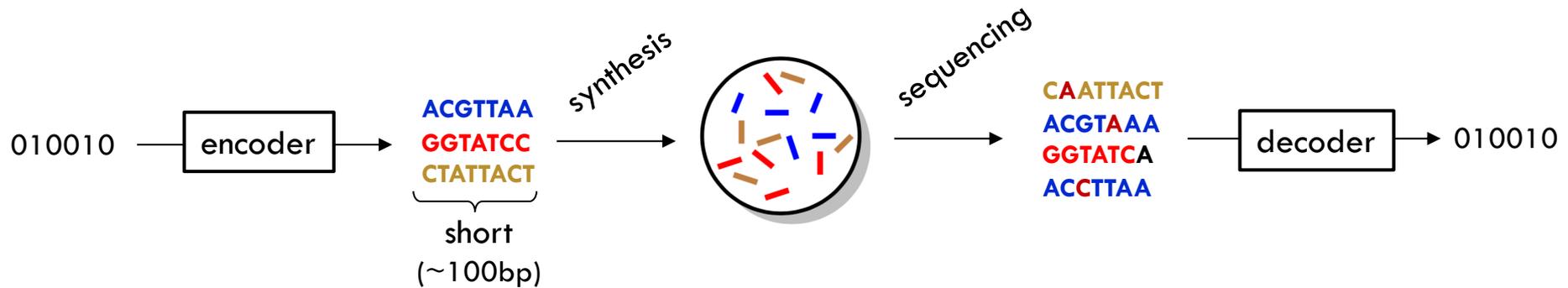


Why store data in DNA?



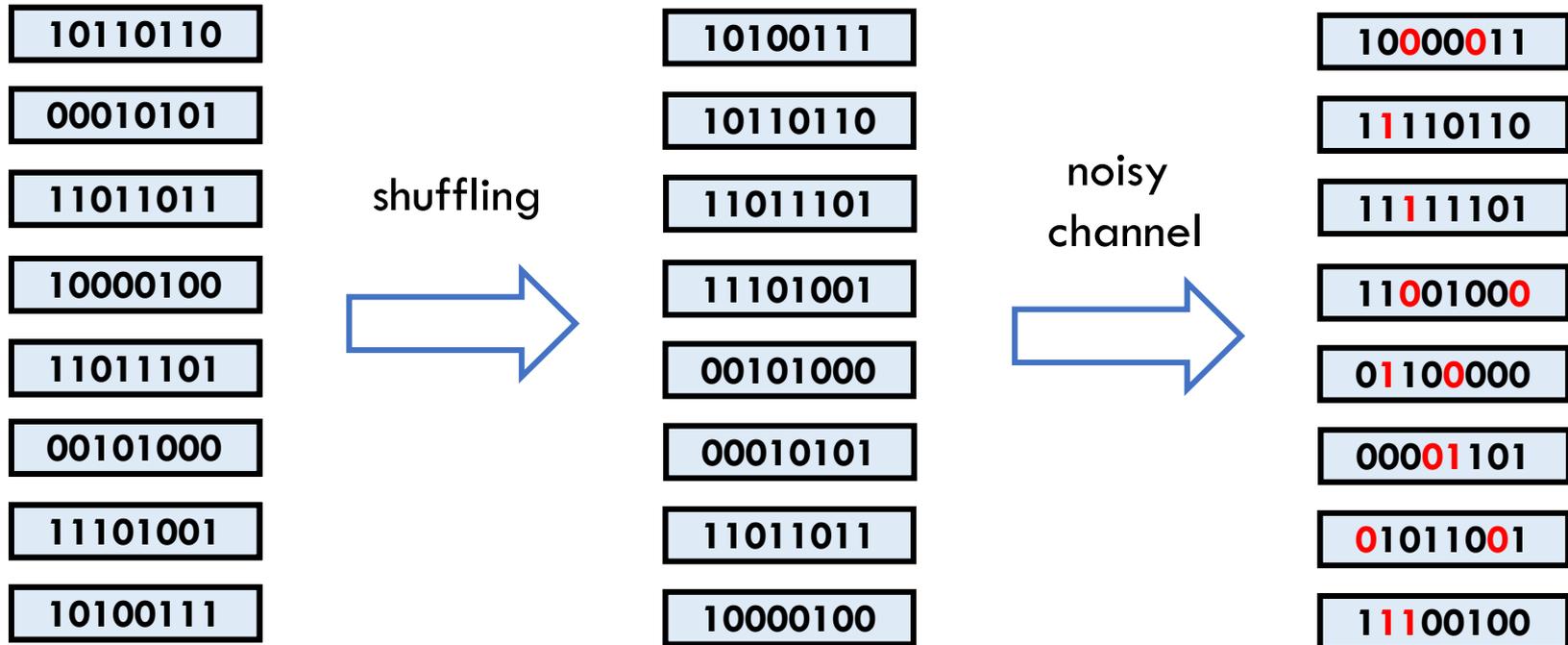
- Good candidate for archival storage

Storing data on DNA

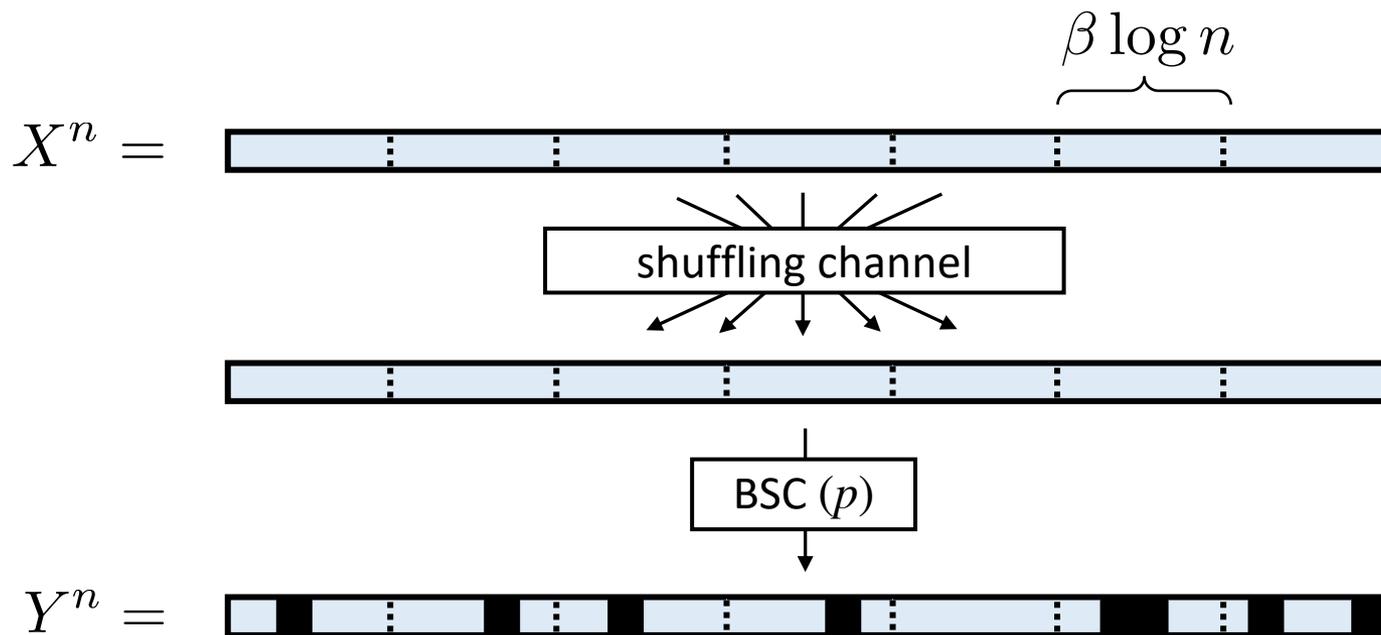


- Can we study the capacity of this channel?
- Distinctive property: out-of-order (shuffled) observations

Noisy Shuffling Channel



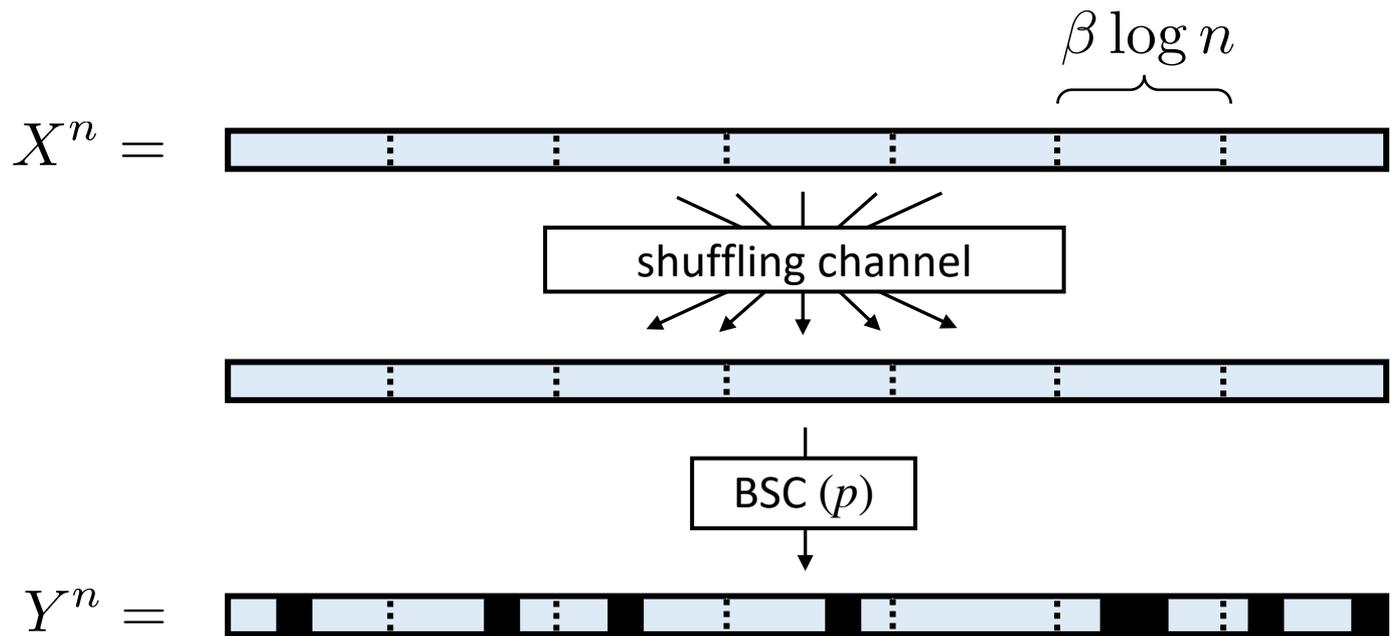
Example: BSC Shuffling Channel



- What is the capacity of this channel?

$$\left(\text{Maximum rate } R = \frac{\log |\mathcal{C}|}{n} \text{ with } P_e \rightarrow 0 \text{ as } n \rightarrow \infty \right)$$

Example: BSC Shuffling Channel

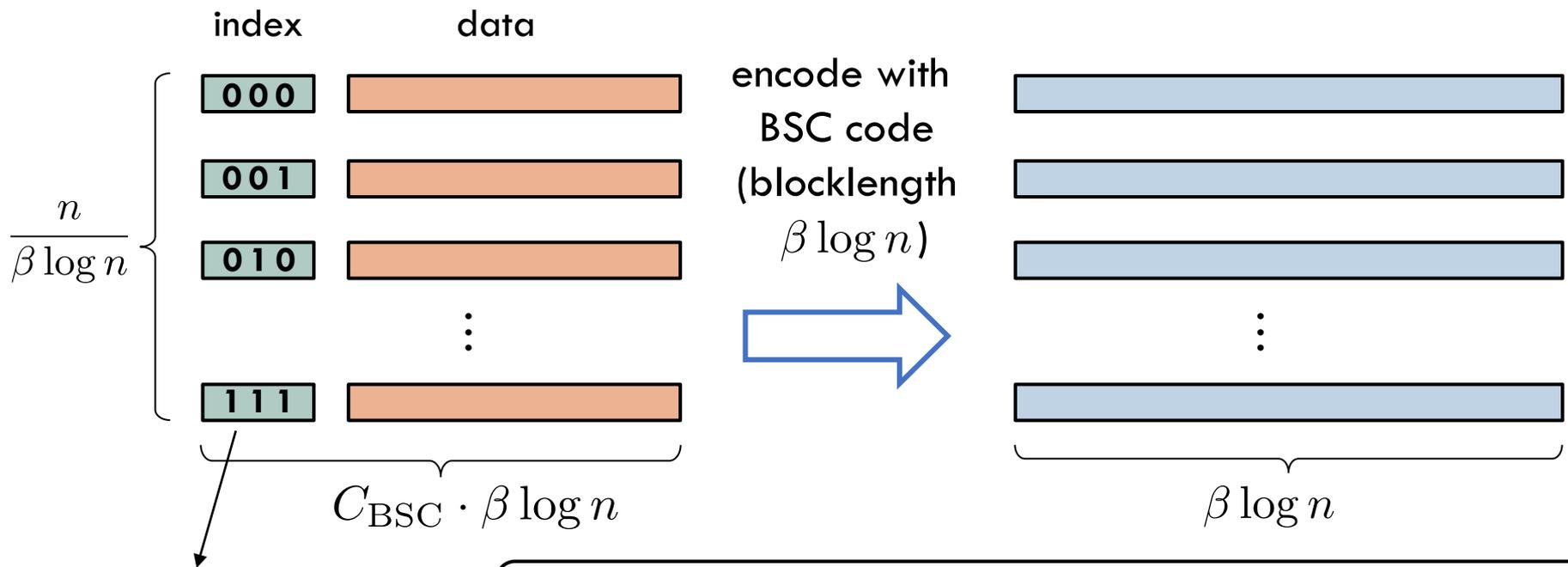


$$C = 1 - H(p) - \frac{1}{\beta}$$

cost of lack of ordering

Simple scheme: index + individual encoding

- Encode **unique index** into each block
- Individually encode block with an optimal BSC code



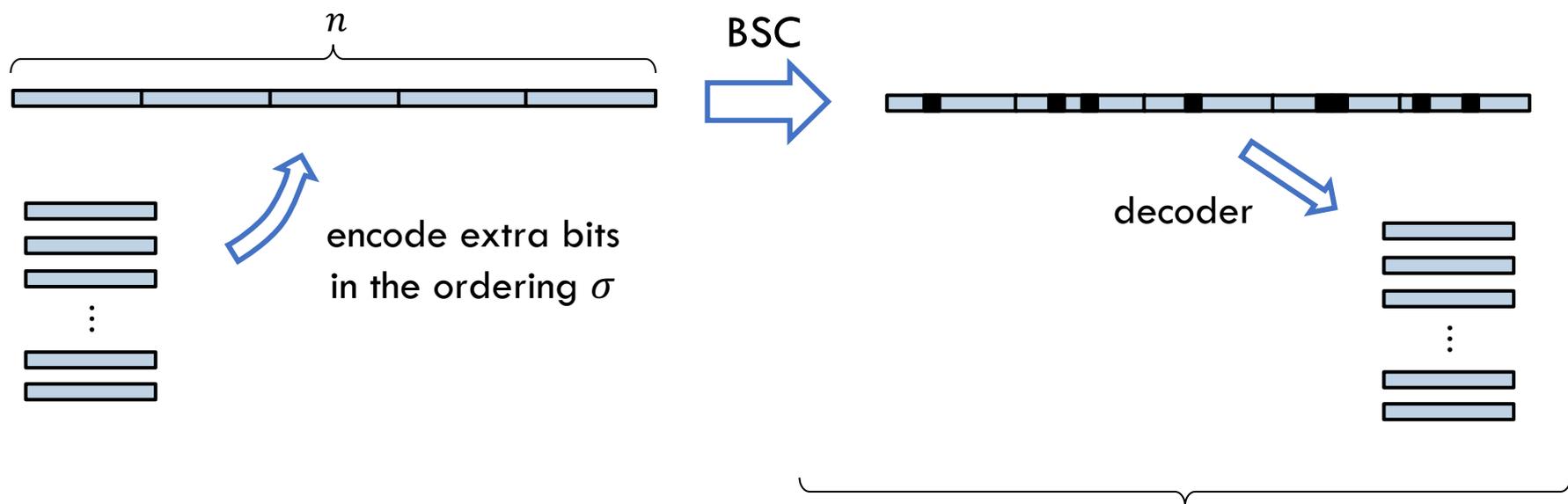
index length:

$$\log \left(\frac{n}{\beta \log n} \right) \approx \log n$$

Rate:
$$\frac{\frac{n}{\beta \log n} (C_{\text{BSC}} \cdot \beta \log n - \log n)}{n}$$

Converse idea

- Consider using code for BSC shuffling channel on a regular BSC

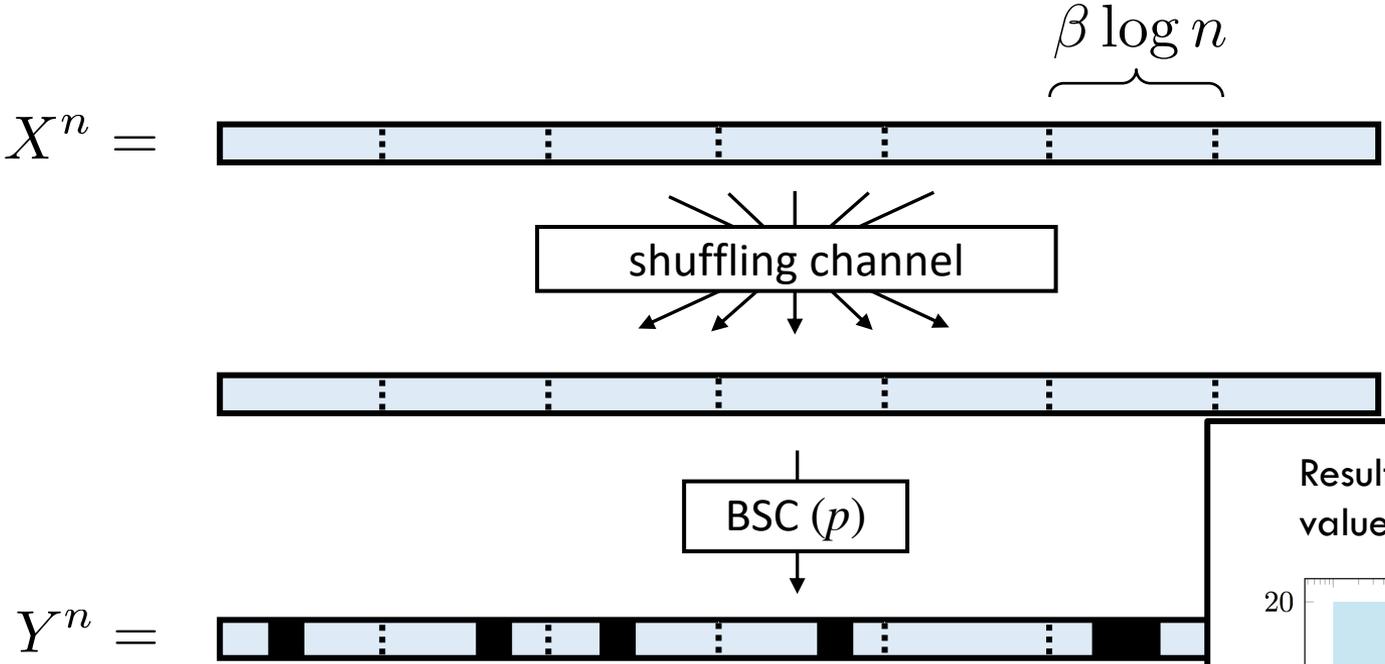


$$\Rightarrow R + \frac{1}{\beta} \leq 1 - H(p)$$

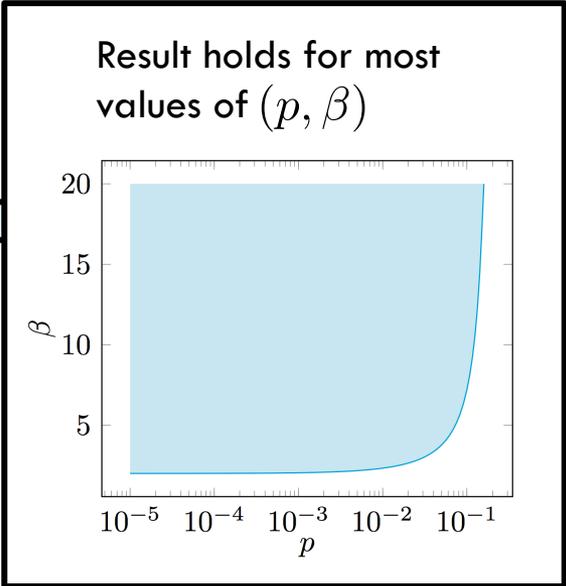
How much of σ can we decode?

$$R_{extra} = \frac{I(\sigma; Y^n)}{n}$$

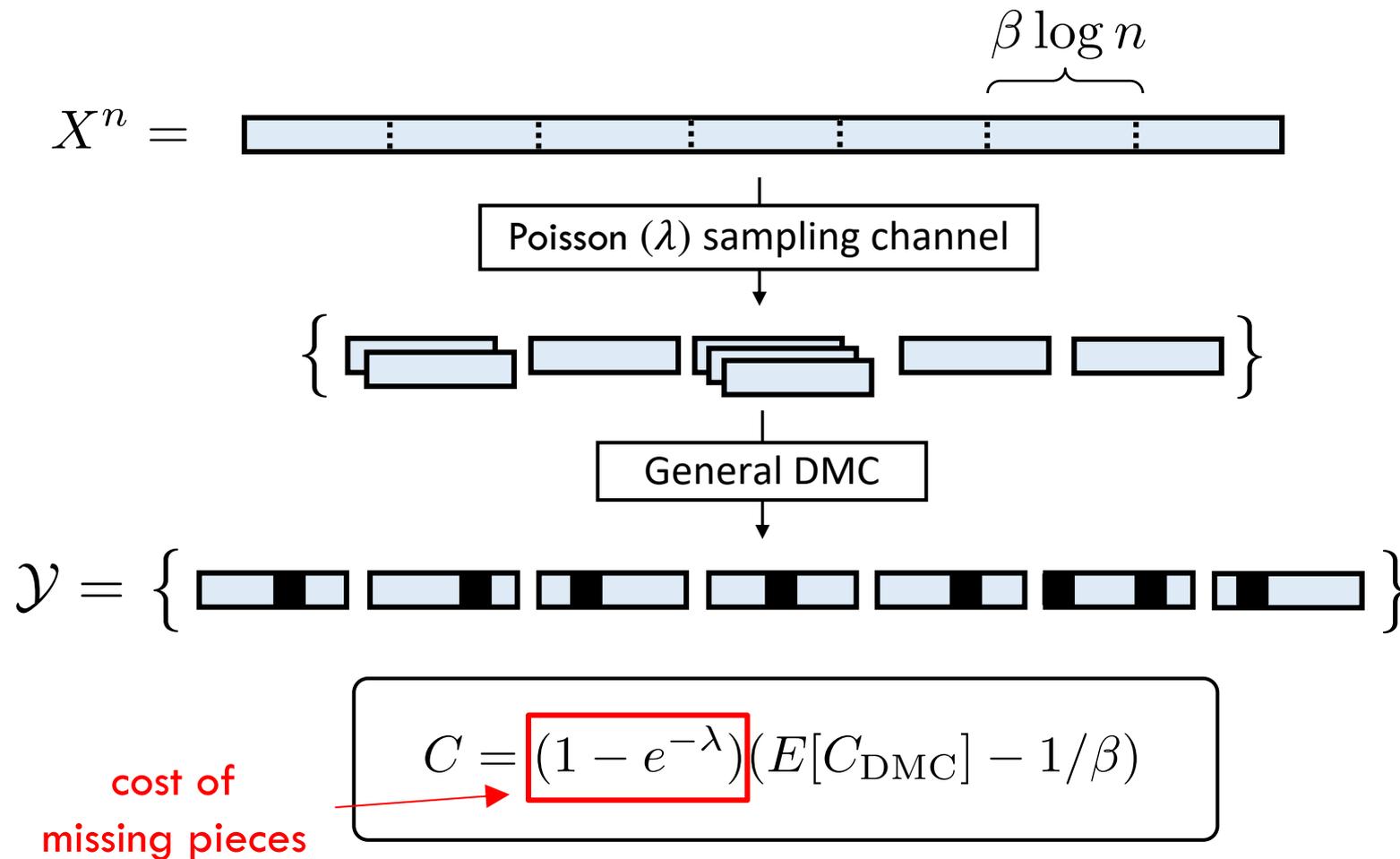
Example: BSC Shuffling Channel



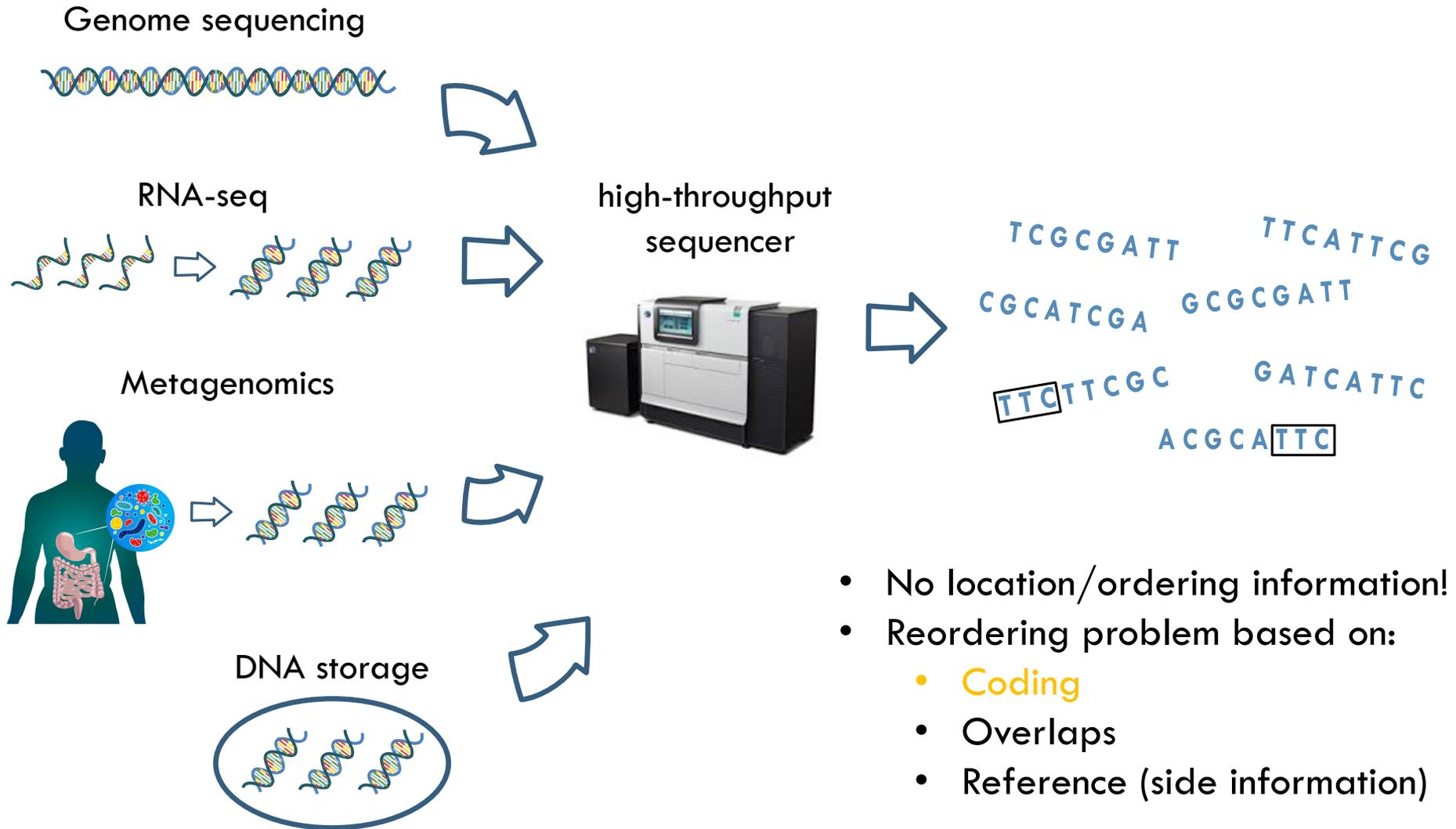
$$C = 1 - H(p) - \frac{1}{\beta}$$



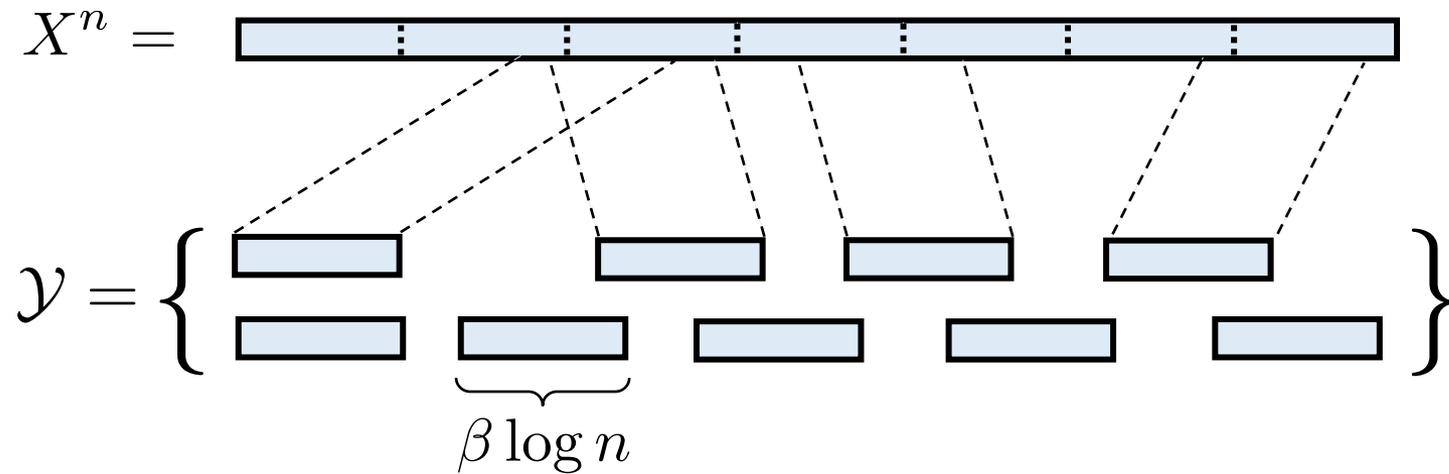
General DNA Storage Channel



The wonders of high-throughput sequencing



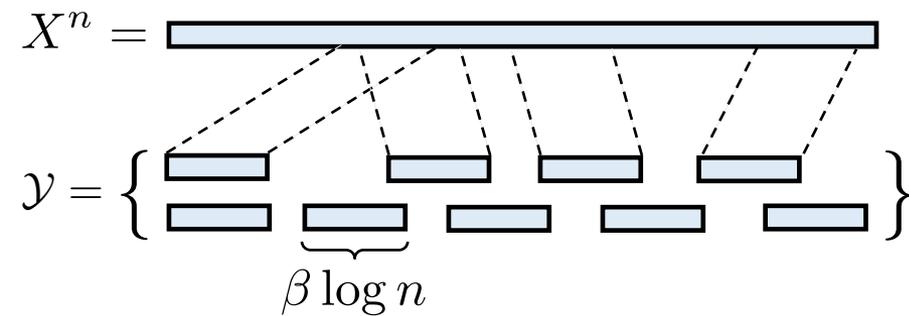
Shotgun Sequencing Channel



- What is the capacity of this channel?

Shotgun Sequencing Channel

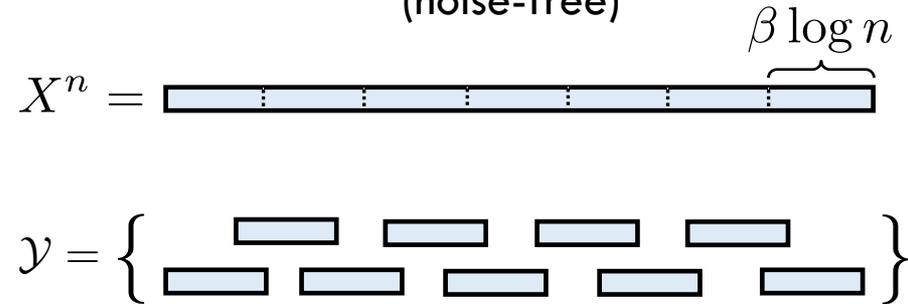
Shotgun Sequencing Channel



- Cannot place unique identifiers
- Overlaps can help reordering

$$C_{\text{SSC}} = 1 - e^{-\lambda(1-1/\beta)}$$

Shuffling Channel
(noise-free)



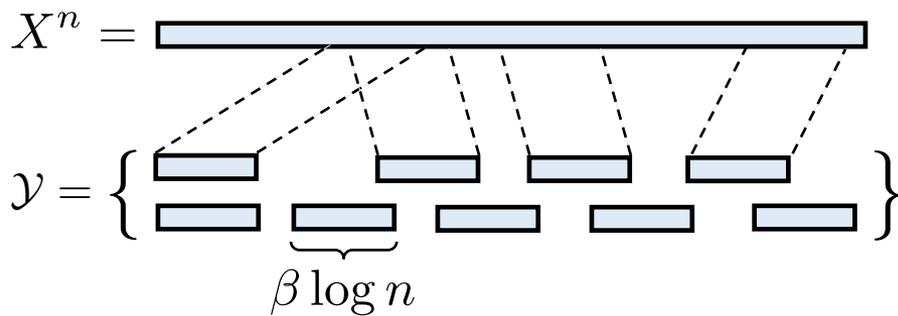
$$C_{\text{shuf}} = (1 - e^{-\lambda})(1 - 1/\beta)$$

sequencing depth: $\lambda = \frac{|\mathcal{Y}| \cdot \beta \log n}{n}$

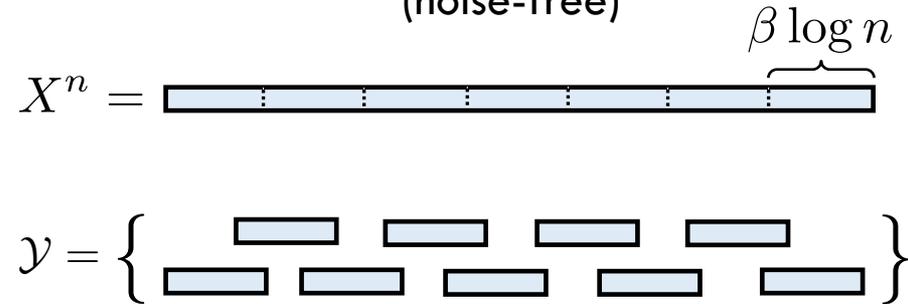


Shotgun Sequencing Channel

Shotgun Sequencing Channel



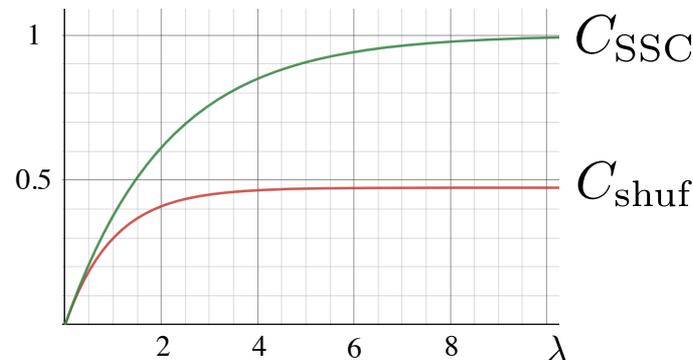
Shuffling Channel
(noise-free)



$$C_{\text{SSC}} = 1 - e^{-\lambda(1-1/\beta)}$$

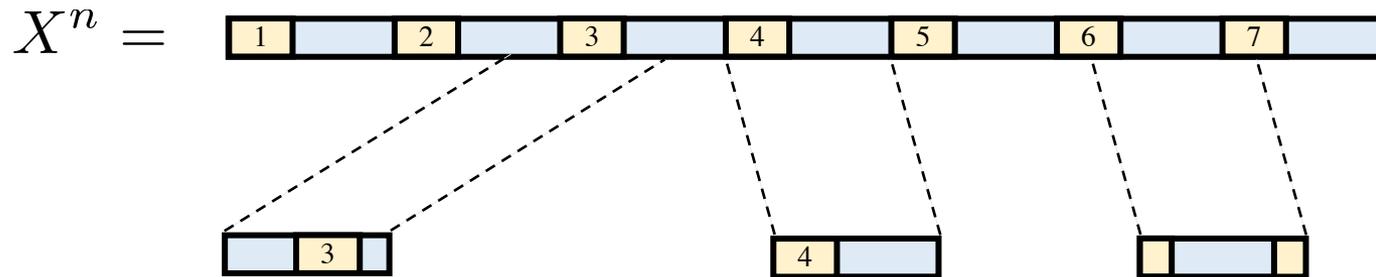
$$C_{\text{shuf}} = (1 - e^{-\lambda})(1 - 1/\beta)$$

overlaps neutralize
lack of ordering
information

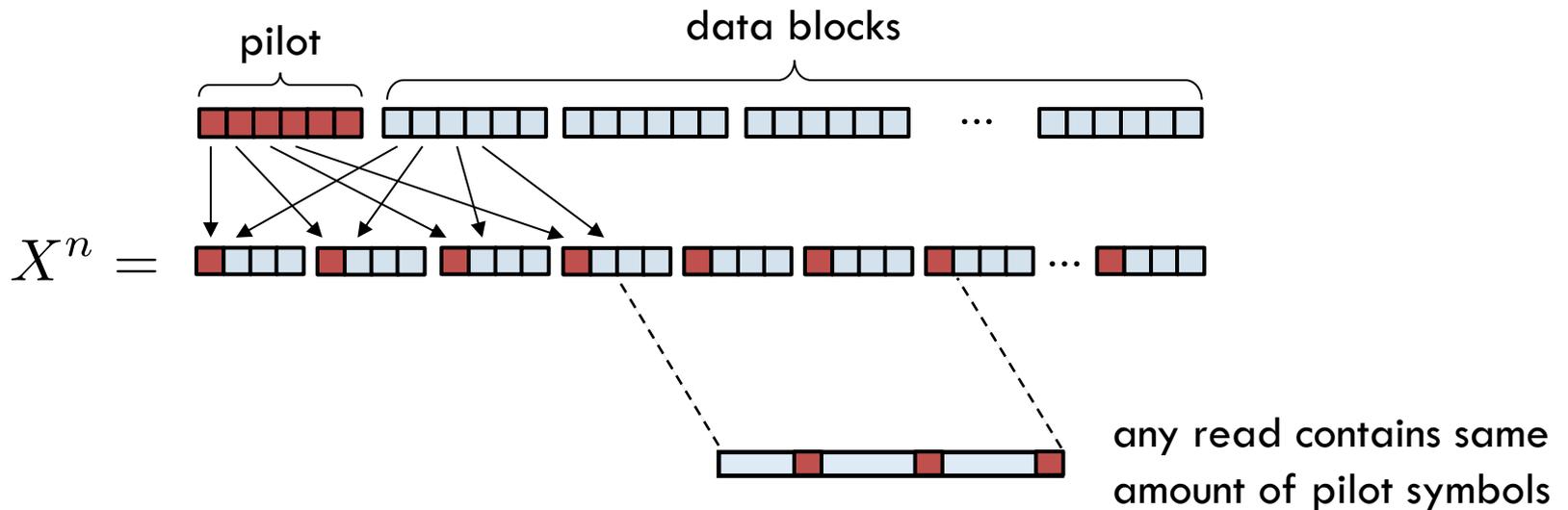


$$\lambda = \frac{|\mathcal{Y}| \cdot \beta \log n}{n}$$

Coding for the Shotgun Sequencing Channel

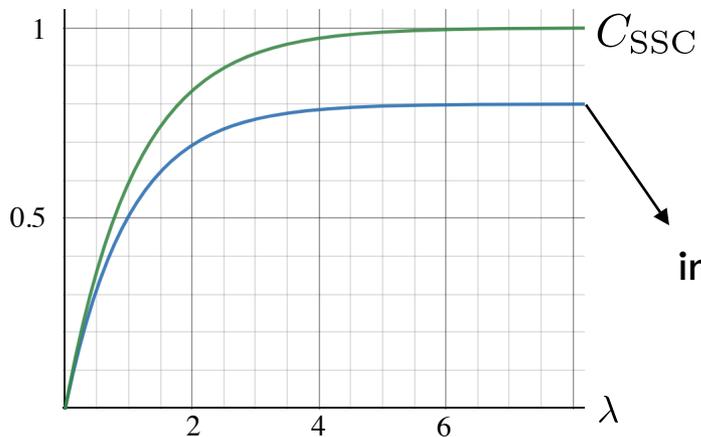
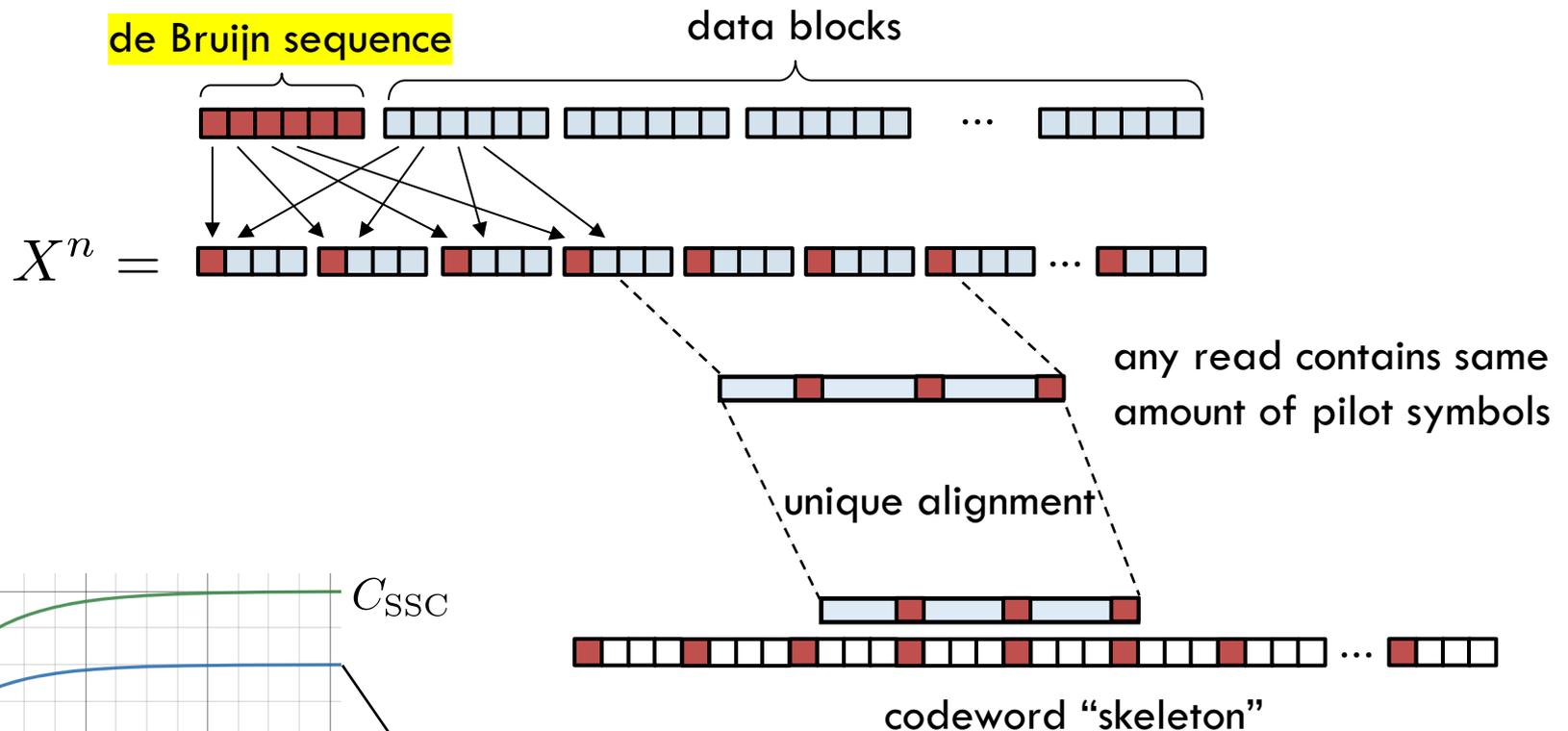


- Better approach: interleave a pilot sequence in X^n



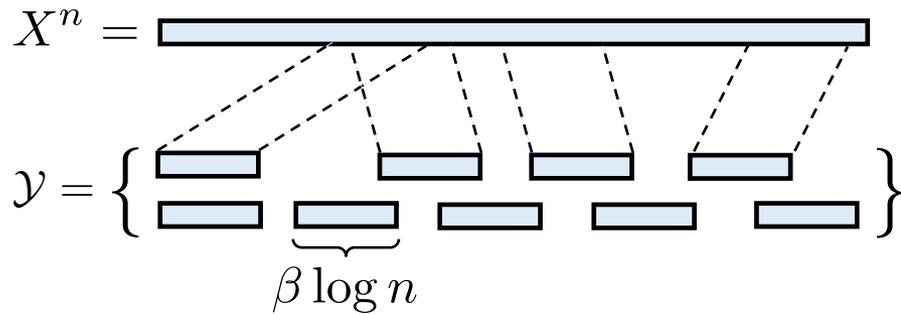
Coding for the Shotgun Sequencing Channel

- Better approach: interleave a pilot sequence in X^n



interleaved pilot: **doesn't take advantage of overlaps**

Coding for the Shotgun Sequencing Channel

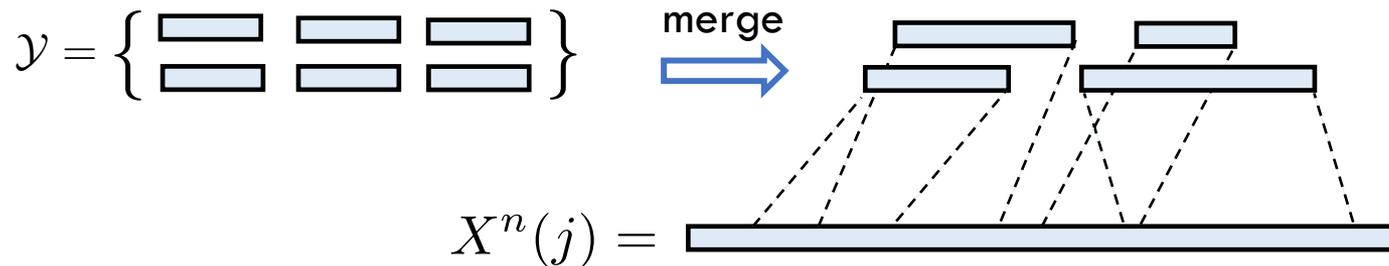


$$C_{\text{SSC}} = 1 - e^{-\lambda(1-1/\beta)}$$

$$\lambda = \frac{|\mathcal{Y}| \cdot \beta \log n}{n}$$

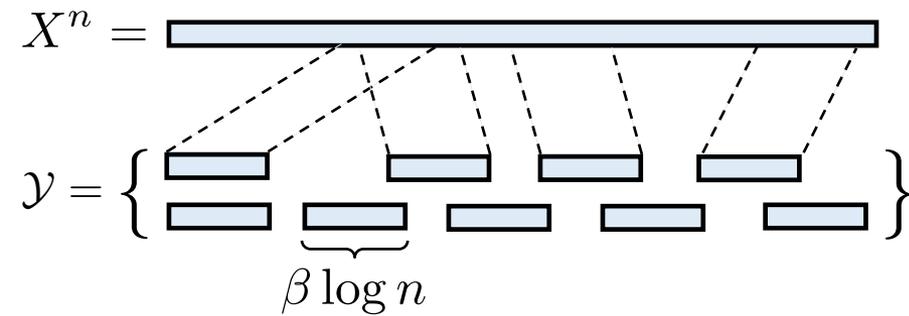
□ Achievability

- ▣ Random code: i.i.d. $\text{Ber}(1/2)$ codewords
- ▣ Decoding: merge reads into “islands” + align to codewords

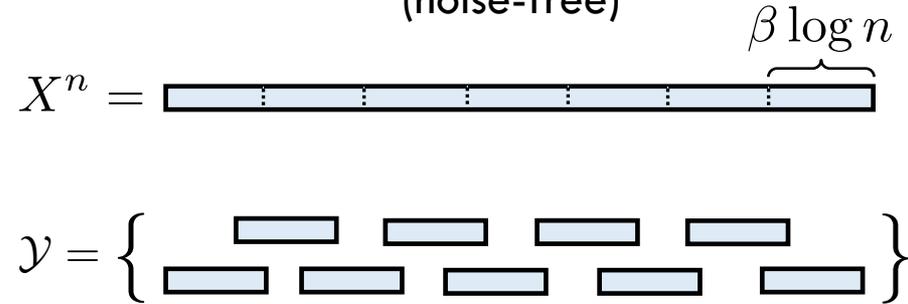


Shotgun Sequencing Channel

Shotgun Sequencing Channel

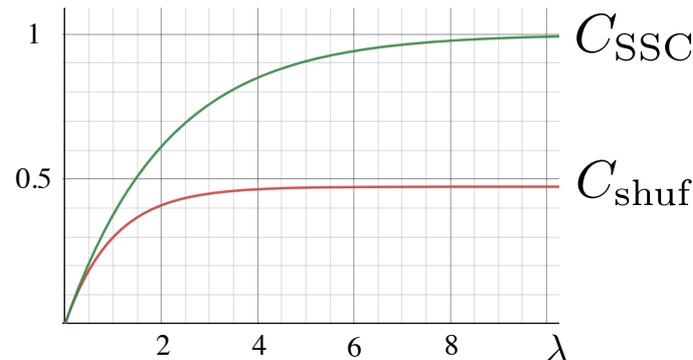


Shuffling Channel
(noise-free)



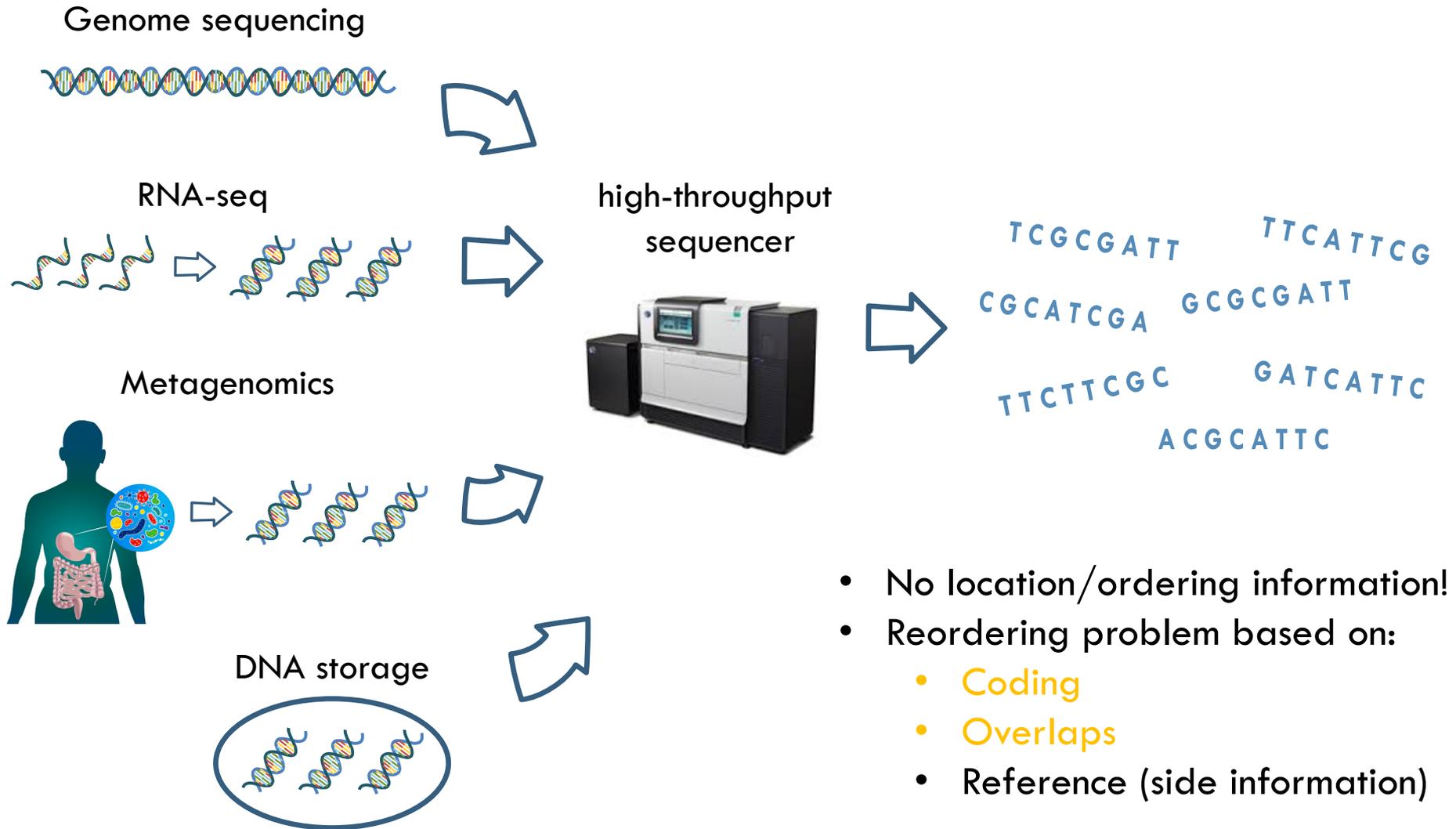
$$C_{\text{SSC}} = 1 - e^{-\lambda(1-1/\beta)}$$

$$C_{\text{shuf}} = (1 - e^{-\lambda})(1 - 1/\beta)$$

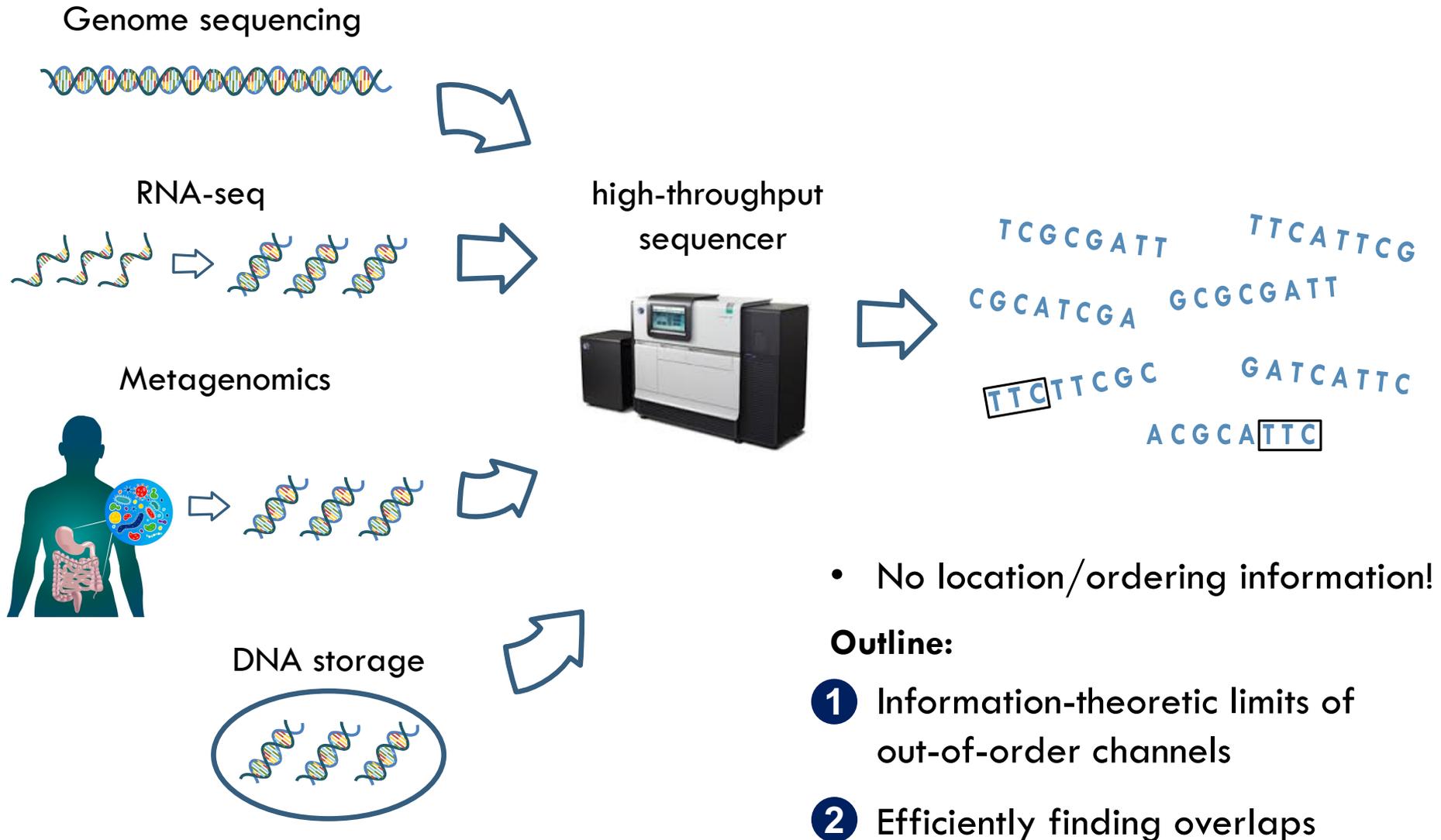


$$\lambda = \frac{|\mathcal{Y}| \cdot \beta \log n}{n}$$

The wonders of high-throughput sequencing



The wonders of high-throughput sequencing



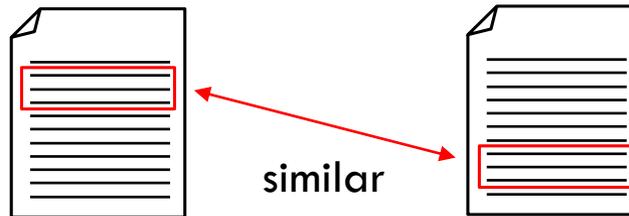
Finding pairwise overlaps

high-throughput
sequencer

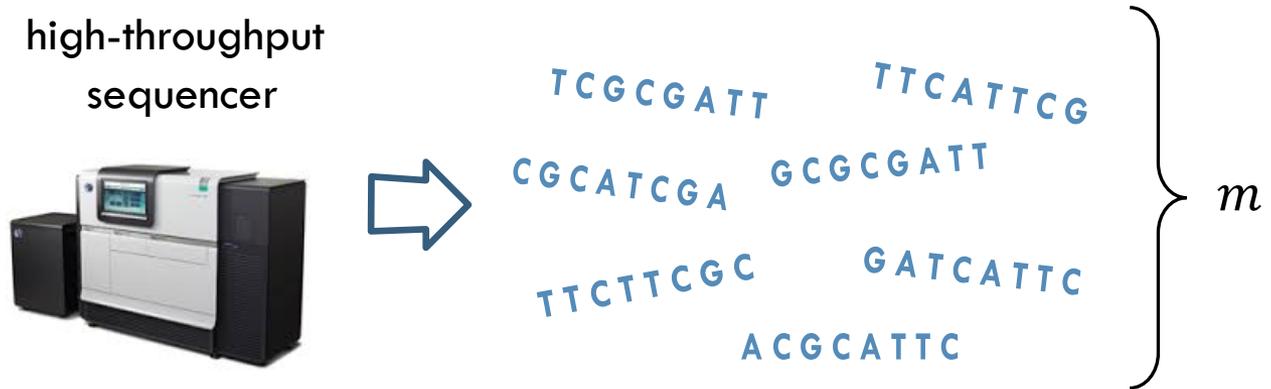


TCGCGATT TTCATTTCG
CGCATCGA GCGCGATT
TTCTTCGC GATCATTC
 ACGCATTC

- Computational bottleneck in bioinformatics
- Other applications: document comparison (plagiarism detection)



Finding pairwise overlaps



- Given two noisy sequences, find best overlap between them

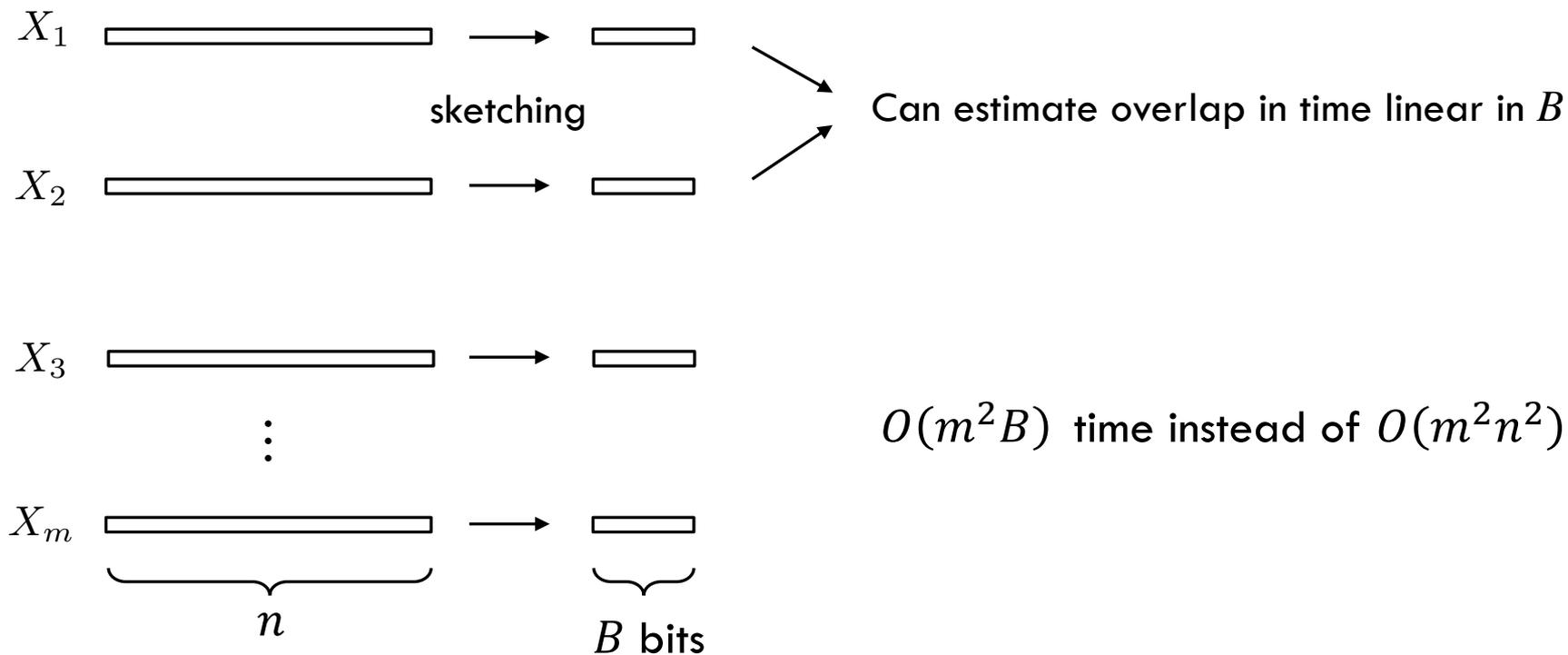
GCCGTAGGCGATTTGCCG-CG
CGGTTTGTCCGACGGCGGCTAG

n

The diagram shows two DNA sequences. The first sequence is GCCGTAGGCGATTTGCCG-CG. The second sequence is CGGTTTGTCCGACGGCGGCTAG. A curly brace under the second sequence is labeled with the variable n .

- Dynamic programming solution: $O(m^2n^2)$ time

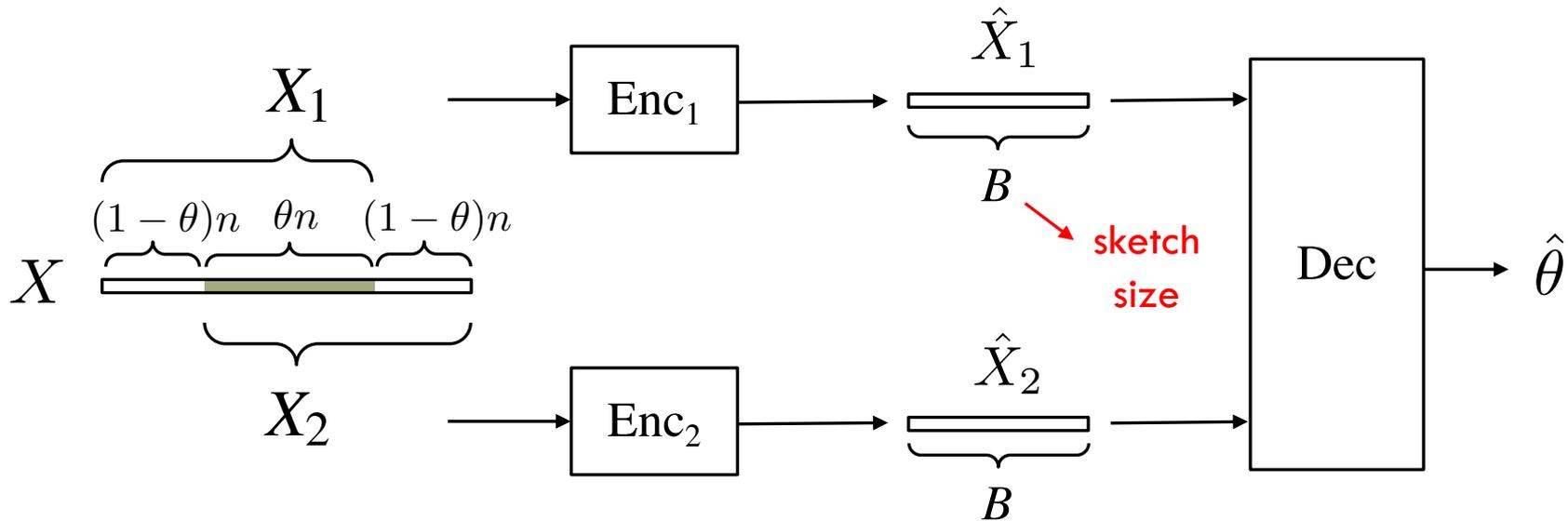
Finding overlaps in practice



- **Fundamental tradeoff:** Sketch size and overlapping accuracy

Distributed source coding formulation

- We focus on estimating the **overlap** size $\theta \in (0,1)$



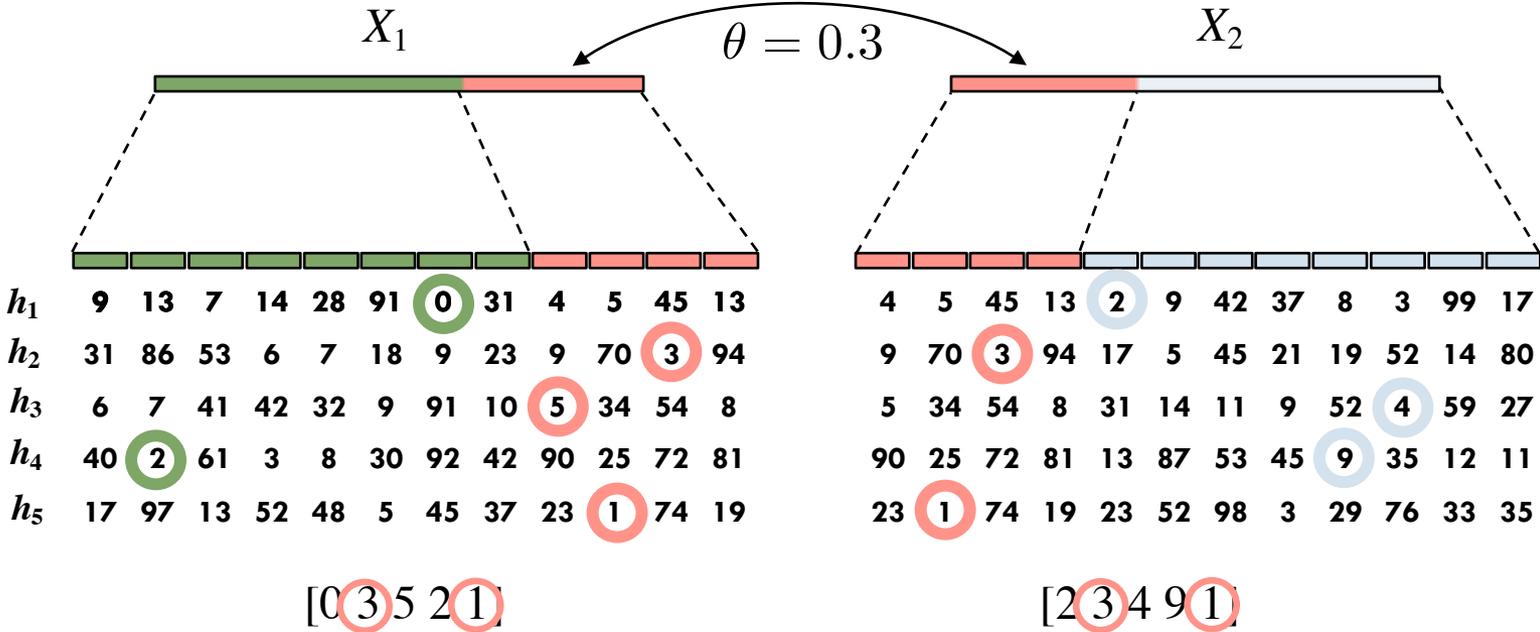
Distortion:

$$E[(\theta - \hat{\theta})^2] \leq D$$

- $B(D)$: Minimum B required to achieve distortion D (as $n \rightarrow \infty$)
- Related to literature on distributed parameter estimation (Amari et al., "Statistical inference under multiterminal data compression," 1998)

Standard approach to sketching

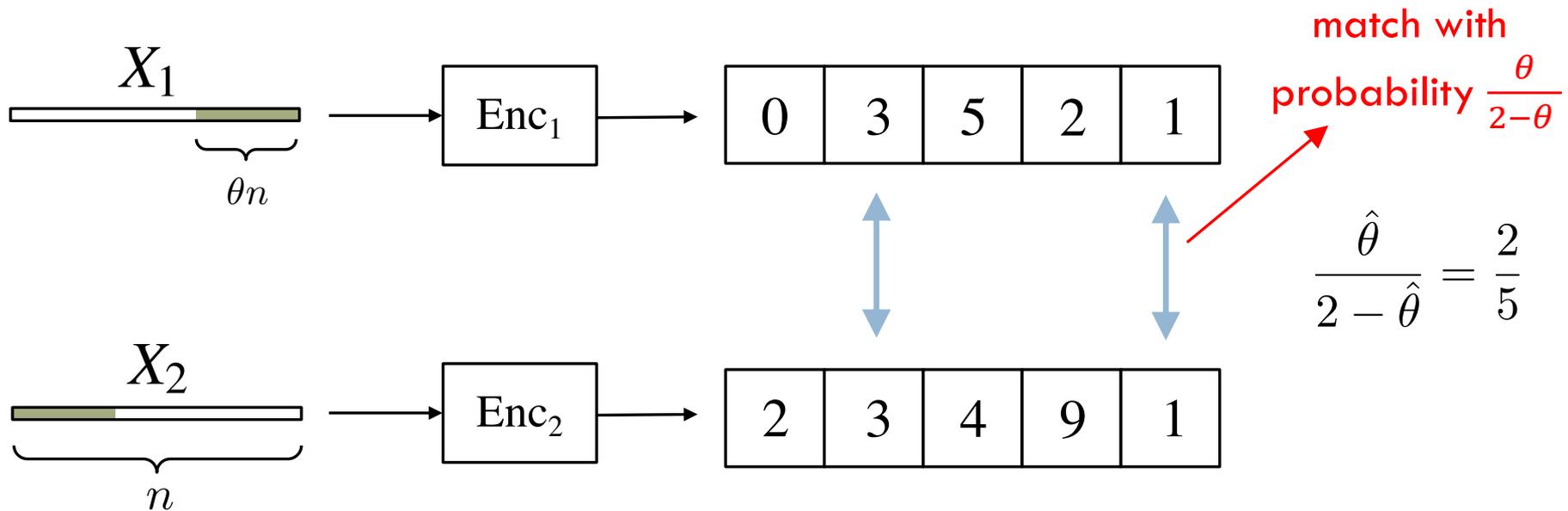
- Sketches based on locality-sensitive hashing (min-hash)



overlap estimate = $2/5$

Standard approach to sketching

- Sketches based on locality-sensitive hashing (min-hash)



- Achievable sketch size: $B \leq \frac{3 \log n}{D}$

Fundamental Limits of Sketching

- Previous approach based on min-hash achieves

$$B(D) \leq \frac{3 \log n}{D}$$

minimum sketch size \rightarrow $B(D)$ \rightarrow distortion

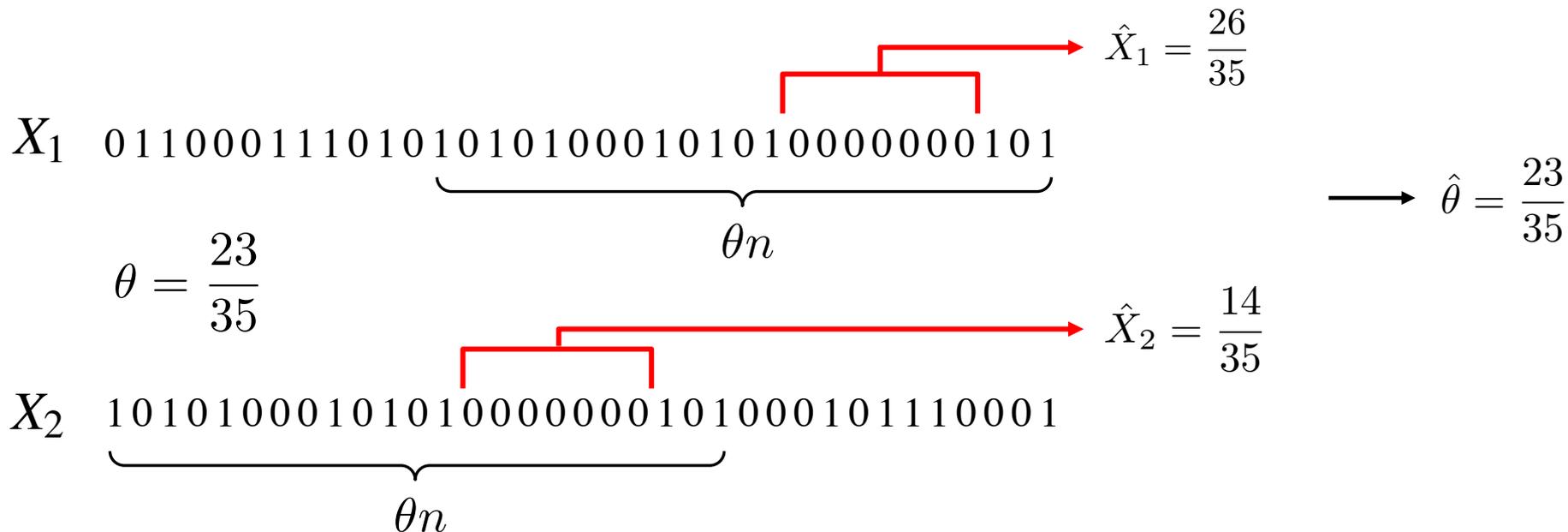
Theorem: The function $B(D)$ satisfies

$$K_1 \cdot \log \left(\frac{1}{D} \right) \leq B(D) \leq K_2 \cdot \log^2 \left(\frac{1}{D} \right)$$

New sketching idea:
Lexicographic Hashing



New idea: Lexicographic Hashing



- Encode location of longest run of zeros
- Overlap estimate: $\hat{\theta} = 1 - (\hat{X}_1 - \hat{X}_2)$

New idea 2: Multiple lexicographic orderings

- Encode location of lexicographically first suffix

different lexicographic ordering

$X_1 \longrightarrow \hat{X}_1$

$\frac{87}{256}$	$\frac{33}{256}$	$\frac{201}{256}$	$\frac{231}{256}$	$\frac{101}{256}$	$\frac{17}{256}$	$\frac{240}{256}$	$\frac{150}{256}$
------------------	------------------	-------------------	-------------------	-------------------	------------------	-------------------	-------------------

$\theta = 0.428$

$X_2 \longrightarrow \hat{X}_2$

$\frac{80}{256}$	$\frac{197}{256}$	$\frac{54}{256}$	$\frac{84}{256}$	$\frac{12}{256}$	$\frac{65}{256}$	$\frac{93}{256}$	$\frac{213}{256}$
------------------	-------------------	------------------	------------------	------------------	------------------	------------------	-------------------

$$\frac{7}{256} - \frac{164}{256} \quad \left(\frac{147}{256} \right) \quad \left(\frac{147}{256} \right) \quad \frac{89}{256} - \frac{48}{256} \quad \left(\frac{147}{256} \right) - \frac{63}{256}$$

$$\hat{\theta} = 1 - \frac{147}{256} = \frac{109}{256} = 0.426$$

Fundamental Limits of Sketching

- Previous approach based on min-hash achieves

$$B(D) \leq \frac{3 \log n}{D}$$

minimum sketch size \rightarrow $B(D)$ \rightarrow distortion

Theorem: The function $B(D)$ satisfies

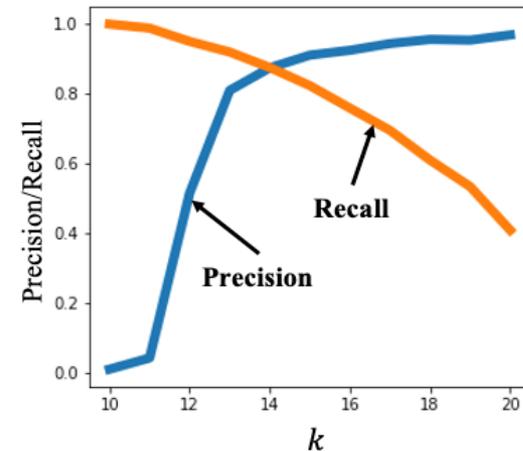
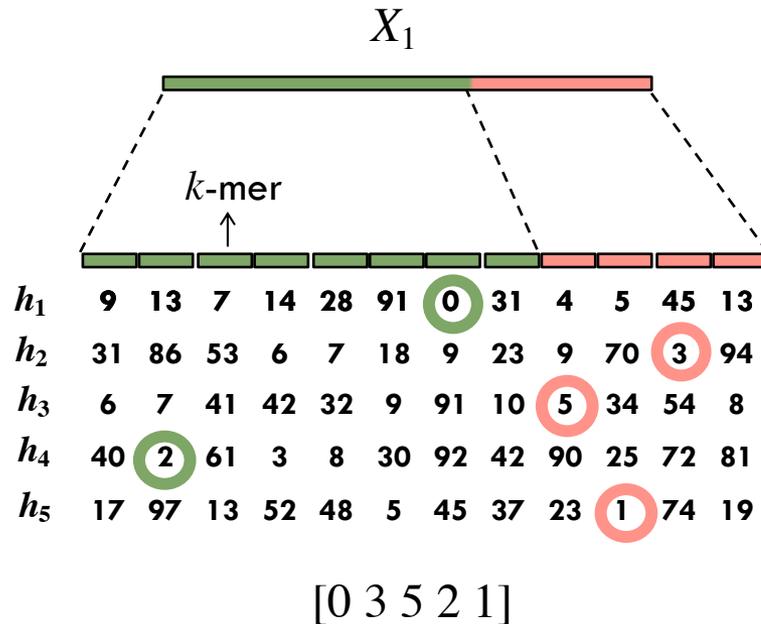
$$K_1 \cdot \log \left(\frac{1}{D} \right) \leq B(D) \leq K_2 \cdot \log^2 \left(\frac{1}{D} \right)$$

New sketching idea:
Lexicographic Hashing



Practical overlap finding: LexicHash

- Sequence alignment tool based on lexicographic hashing
- More accurate overlap estimates
- Avoids need to select parameter k from min-hash

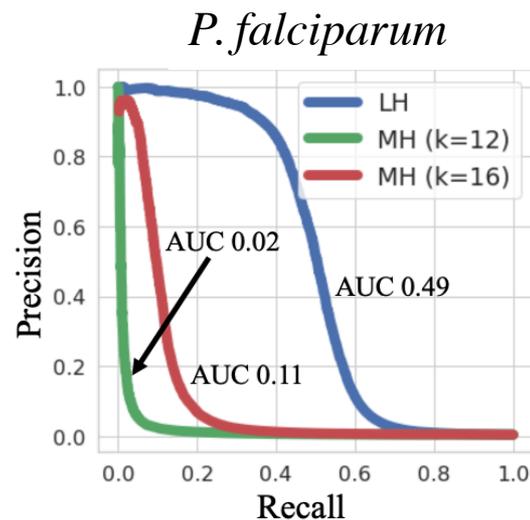
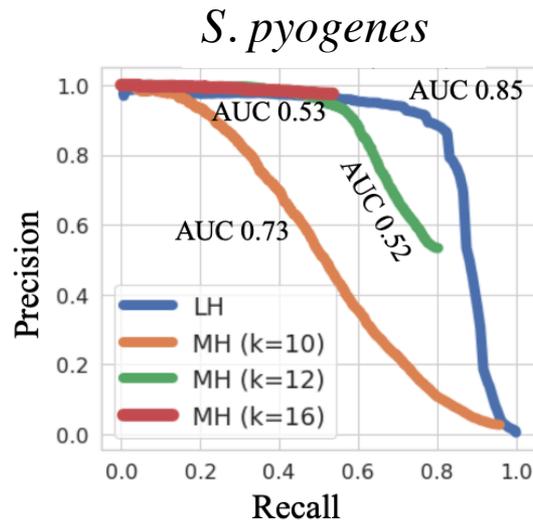


G. Greenberg, A. Ravi, I. Shomorony, LexicHash: Sequence Similarity Estimation via Lexicographic Comparison of Hashes (submitted)



Practical overlap finding: LexicHash

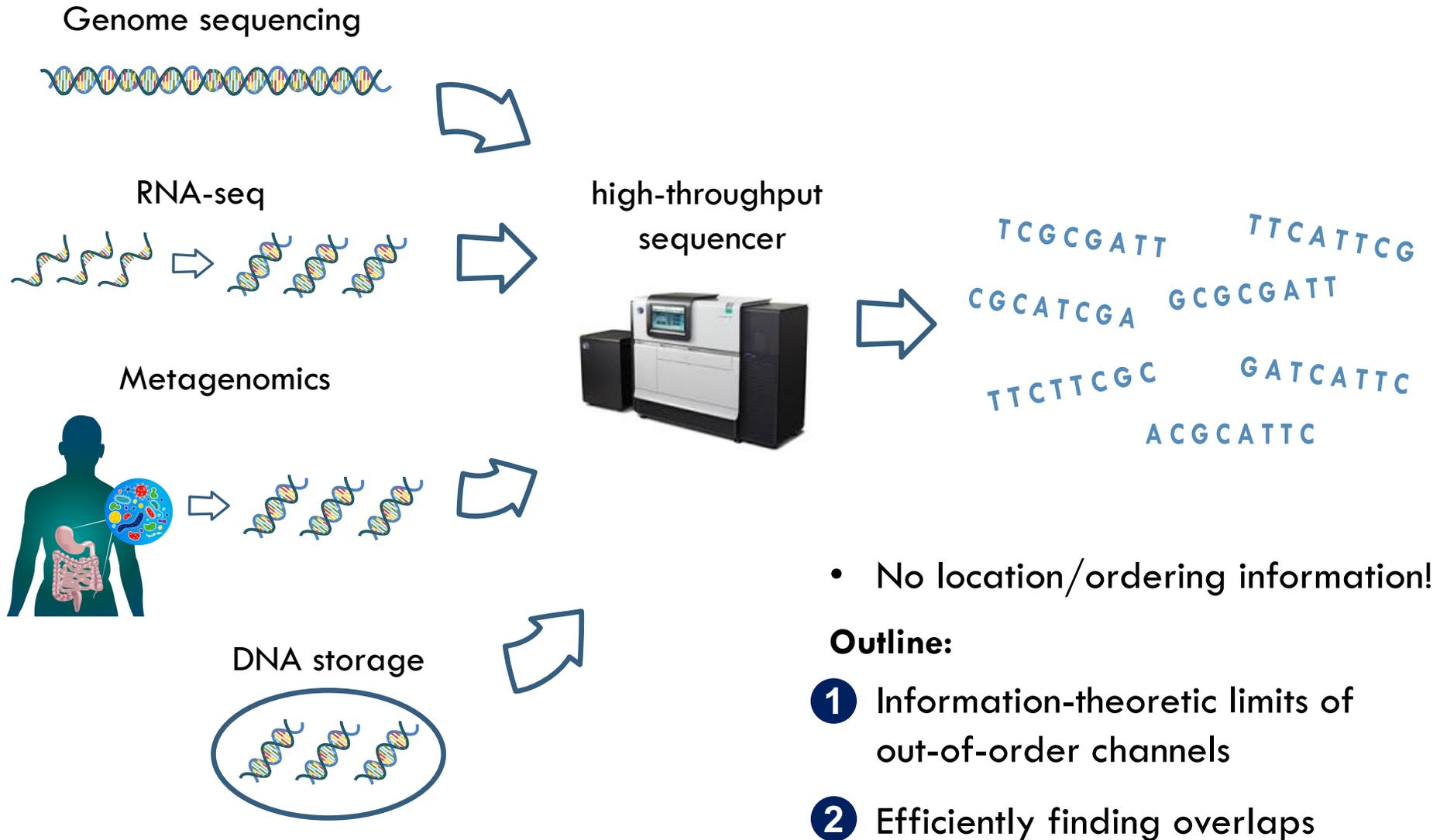
- Sequence alignment tool based on lexicographic hashing
- More accurate overlap estimates
- Avoids need to select parameter k from min-hash
- Significant improvements in overlap classification:



G. Greenberg, A. Ravi, I. Shomorony, LexicHash: Sequence Similarity Estimation via Lexicographic Comparison of Hashes (submitted)



The wonders of high-throughput sequencing



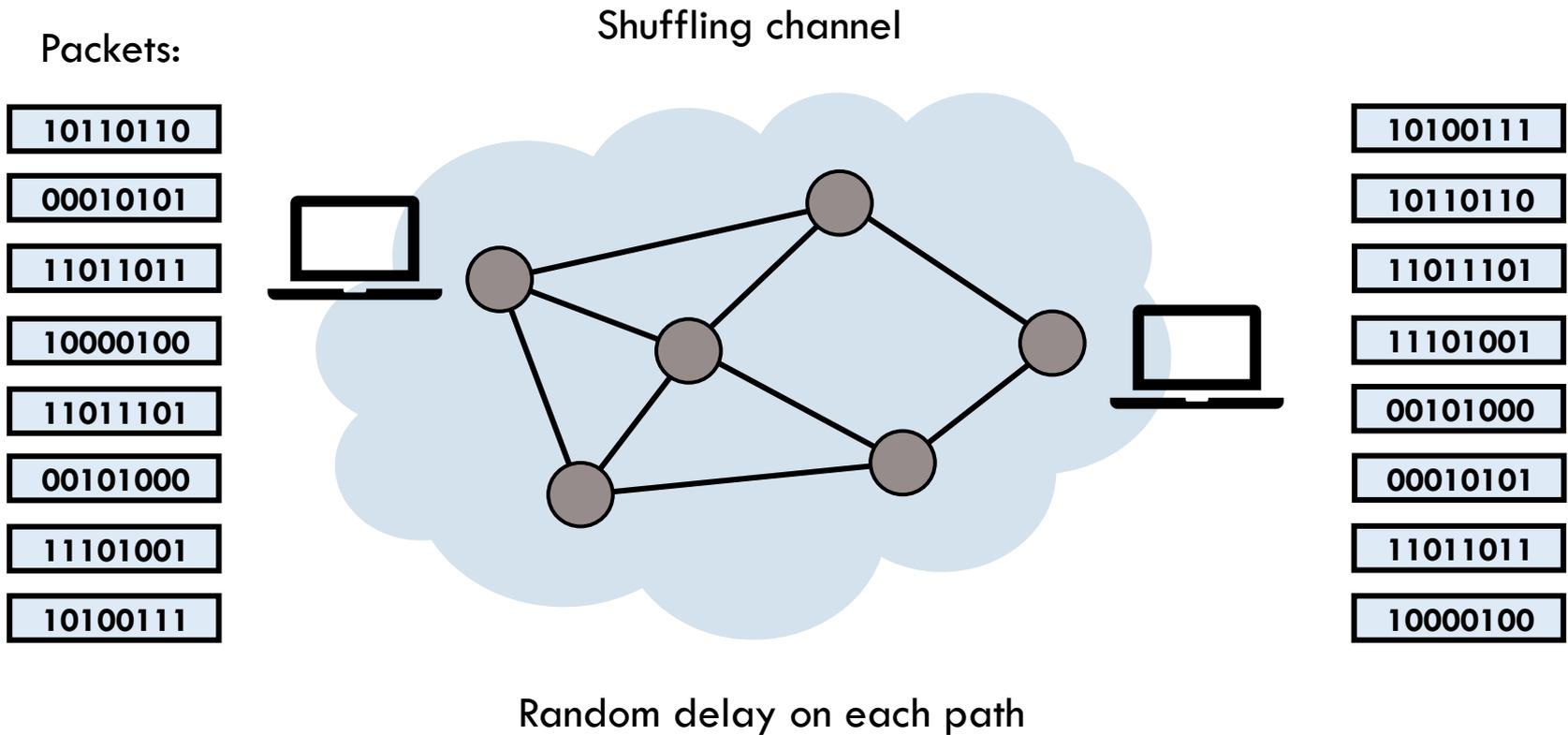
- No location/ordering information!

Outline:

- 1 Information-theoretic limits of out-of-order channels
- 2 Efficiently finding overlaps

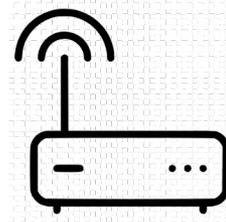
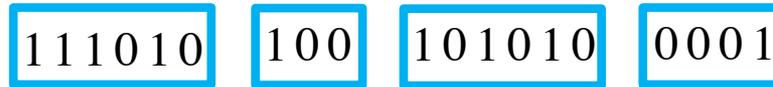
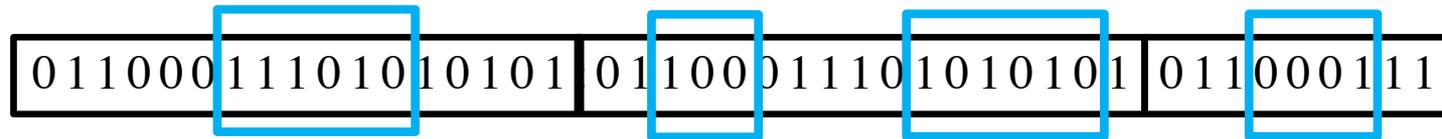
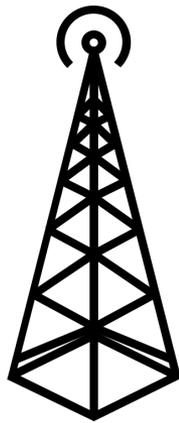
Out-of-order information?

□ Packet networks



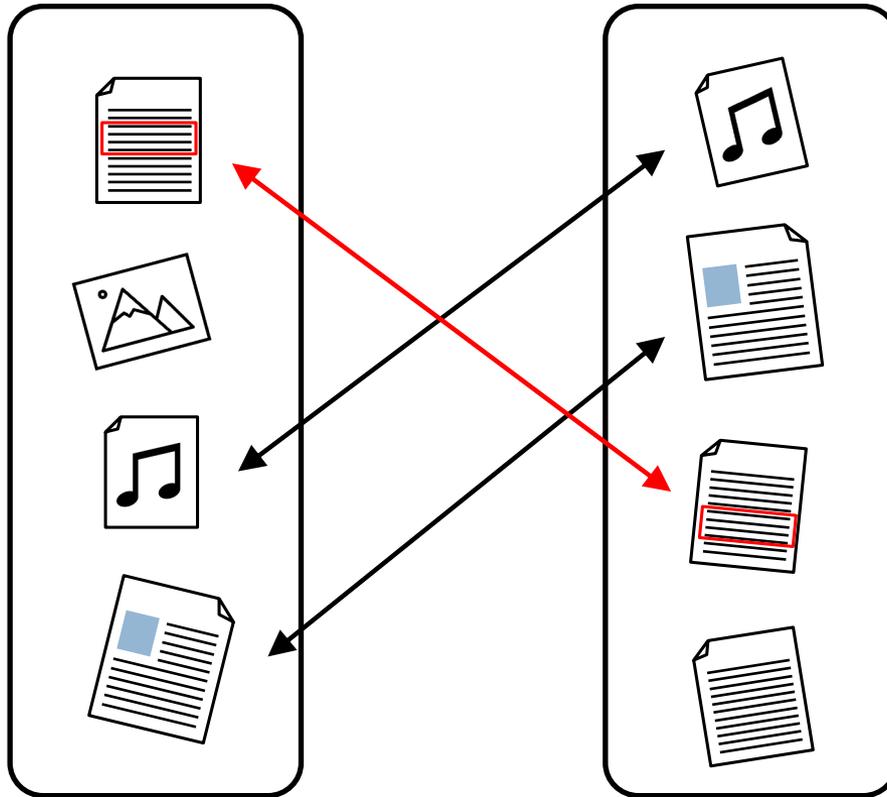
Out-of-order information?

- Intermittent broadcast communication?

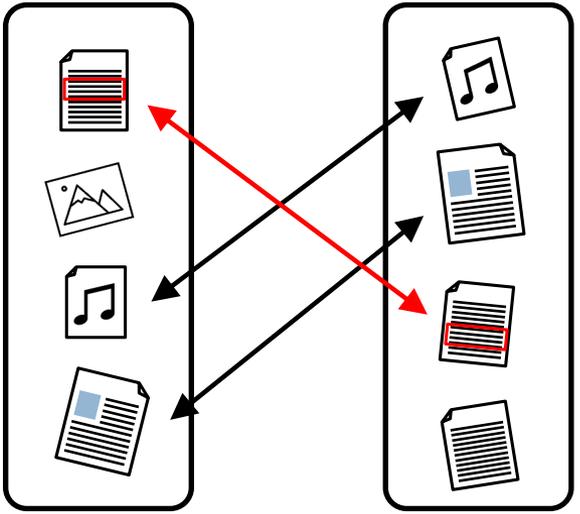
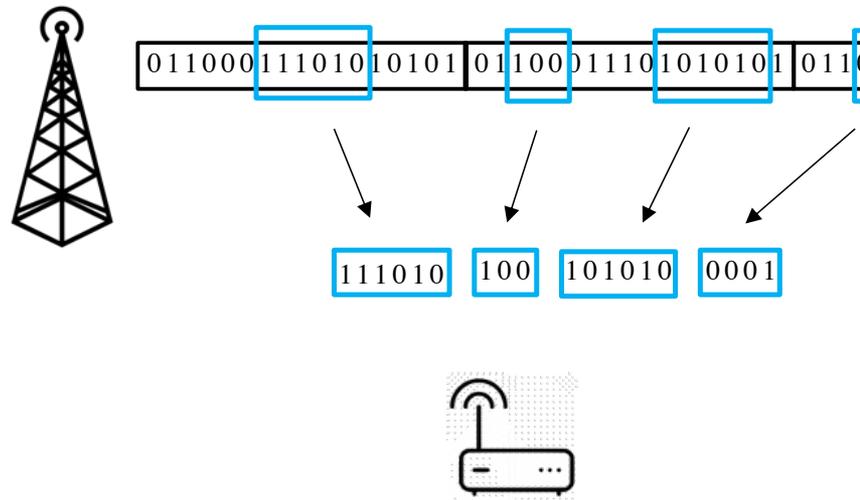
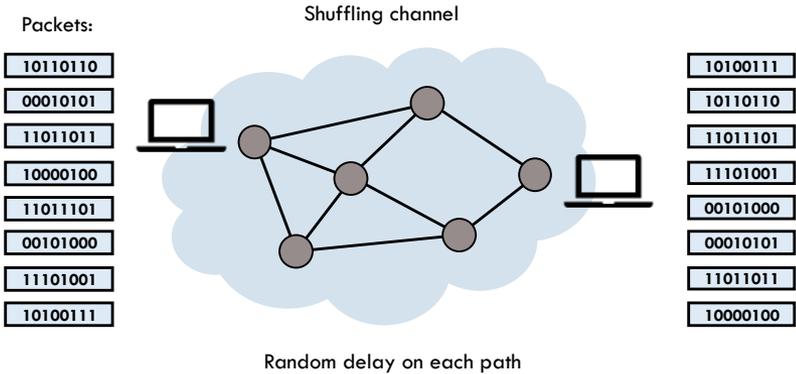


Out-of-order information?

- Dataset alignment?



Out-of-Order Information Theory?



TCGCGATT TTCATTTCG
 CGCATCGA GCGCGATT
 TTCTTCGC GATCATTTC
 ACGCATTTC

Thank you!