## Channel Coding for storage on DNA molecules

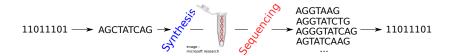
Elsa Dupraz (IMT Atlantique/Lab-STICC)

June 24th, 2025



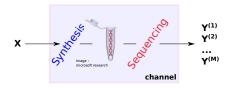


# The DNA data storage workflow



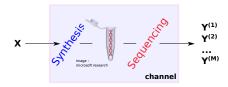
- Synthesis (chemical) is currently very expensive
- Sequencing (nanopore) introduces a large amount of errors

# The DNA data storage workflow



► The (synthesis + sequencing) part can be seen as a **channel** which introduces errors in the read sequences

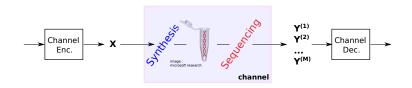
# The DNA data storage workflow



- ► The (synthesis + sequencing) part can be seen as a channel which introduces errors in the read sequences
- This channel introduces three types of errors ex: AGCTATCAG

Substitutions: AGCTATTAG
 Deletions: AGCTATCAG
 Insertions: AGCTATCGAG

# Coding scheme for DNA data storage



### In this presentation

- ▶ How can we **model** errors introduced by the DNA storage process?
- ▶ How can we **correct** errors introduced by the DNA storage channel?
- ▶ How can we use **multiple reads** to improve error correction?

## Outline

Introduction

#### Channel model

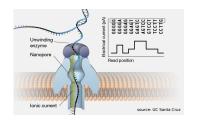
#### Error-correction codes

Pure coding solution (CC codes)

Mixed consensus-coding solution

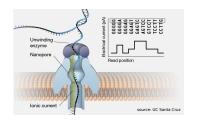
Perspectives

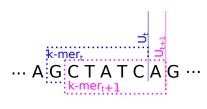
### Nanopore sequencing:



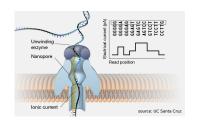


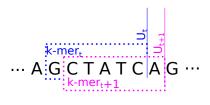
### Nanopore sequencing:





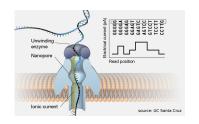
### Nanopore sequencing:

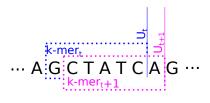




Basecaller: currents  $(..U_{t-1}, U_t, U_{t+1}...) \rightarrow \text{bases} (..Y_{t-1}, Y_t, Y_{t+1}...)$ 

### Nanopore sequencing:





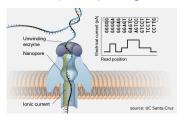
Basecaller: currents  $(..U_{t-1}, U_t, U_{t+1}...) \rightarrow \text{bases } (..Y_{t-1}, Y_t, Y_{t+1}...)$ 

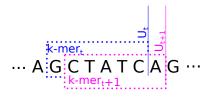
### Model assumptions:

- Error probability depends on the underlying k-mer
- Error probability depends on the previous error event

# Proposed channel model with memory<sup>1</sup>

### Nanopore sequencing:





#### Notation:

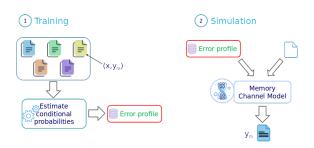
- kmer<sub>t</sub>: K-mer at position t
- ▶  $E_t \in \{\text{Ins, Del, Sub, Match}\}$ : error event at position t

#### Probabilities:

- ▶  $\mathbb{P}(E_t | \mathsf{kmer}_t, E_{t-1})$ : Probability of error event  $E_t$
- $ightharpoonup \mathbb{P}(L|\mathsf{kmer}_t, E_t = \mathsf{Ins})$ : Probability of insertion length L
- ▶  $\mathbb{P}(B|\mathsf{kmer}_t, E_t = \mathsf{Sub})$ : Probability of substitution by base B

<sup>&</sup>lt;sup>1</sup>Belaid Hamoum, Elsa Dupraz, Laura Conde-Canencia, Dominique Lavenier, Channel Model with Memory for DNA Data Storage with Nanopore Sequencing, ISTC 2021

## Training and simulations



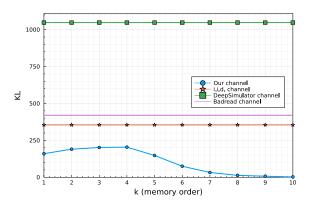
### Training over different datasets:

Experimental data (2020): Thermo Fisher synthesis, **error rate**: 10% Genomics data (2021): Streptococcus thermophilus bacteria, **error rate**: 3% Experimental data (2021): Twist synthesis, Kodak Images, **Error rate**: 4% Experimental data (2024): IDT synthesis, **Error rate**: 5%

Simulator available on GitHub: https://github.com/BHam-1/DNArSim

# Kulback Leibler divergence<sup>1</sup>

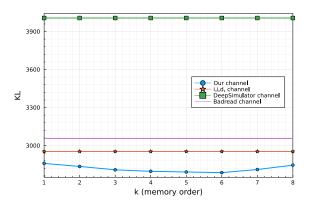
► Model trained on experimental data



<sup>&</sup>lt;sup>1</sup>Belaid Hamoum, Elsa Dupraz, Laura Conde-Canencia, Dominique Lavenier, Channel Model with Memory for DNA Data Storage with Nanopore Sequencing, ISTC 2021

# Kulback Leibler divergence<sup>1</sup>

► Model trained on genomics data



<sup>&</sup>lt;sup>1</sup>Belaid Hamoum, Elsa Dupraz, Laura Conde-Canencia, Dominique Lavenier, Channel Model with Memory for DNA Data Storage with Nanopore Sequencing, ISTC 2021

## Outline

Introduction

Channel model

### Error-correction codes

Introduction
Pure coding solution (CC codes)
Mixed consensus-coding solution

Perspectives

# Channel coding



### Channel coding:

- (coding) Structure redundancy in the data
- ▶ (decoding) Algorithm that leverages redundancy to correct errors

### Coding rate:

$$R=\frac{k}{n}$$

## Channel codes for DNA data storage

#### Adversarial models:

- Varshamov-Tenengolts (VT) codes can correct one deletion [Varshamov65]
- VT codes have been extended to the correction of a few insertions and/or deletions [Abroshan19,Xue20,KasHanna22,etc.]

Explicit bounds on the redundancy:

$$R \ge 1 - \lceil \log(n_b + 1) \rceil / n_b$$

## Channel codes for DNA data storage

#### Adversarial models:

- Varshamov-Tenengolts (VT) codes can correct one deletion [Varshamov65]
- VT codes have been extended to the correction of a few insertions and/or deletions [Abroshan19,Xue20,KasHanna22,etc.]

Explicit bounds on the redundancy:

$$R \geq 1 - \lceil \log(n_b + 1) \rceil / n_b$$

#### Statistical models:

- Pure coding:
  - ► Convolutional codes [Lenz21,etc.]
  - ► LDPC codes [Wang11,Shibata19,etc.]
- Combination of consensus and coding

## Outline

Introduction

Channel mode

### Error-correction codes

Introduction
Pure coding solution (CC codes)
Mixed consensus-coding solution

Perspectives

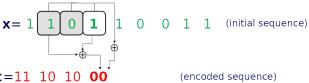
### Standard convolutional codes

#### Convolutional codes:

- Were initially developed to correct substitution errors
- ► Encode blocks of arbitrary length

### Example:

$$(k_c=1,n_c=2,\nu=2)$$
 convolutional code with  $\Delta=[\Delta^2+1,\Delta^2+\Delta+1].$ 



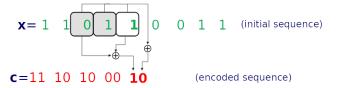
### Standard convolutional codes

#### Convolutional codes:

- Were initially developed to correct substitution errors
- ▶ Encode blocks of arbitrary length

### Example:

$$(k_c=1,n_c=2,\nu=2)$$
 convolutional code with  $\Delta=[\Delta^2+1,\Delta^2+\Delta+1]$ .



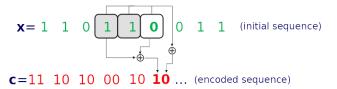
### Standard convolutional codes

#### Convolutional codes:

- Were initially developped to correct substitution errors
- ▶ Encode blocks of arbitrary length

### Example:

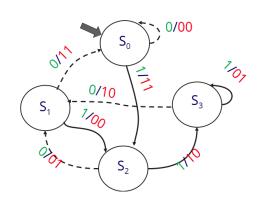
$$(k_c=1,n_c=2,\nu=2)$$
 convolutional code with  $\Delta=[\Delta^2+1,\Delta^2+\Delta+1].$ 



# Standard CC decoding

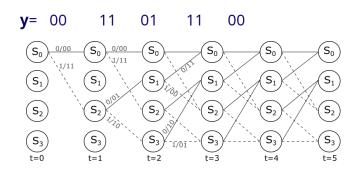
- ▶ A CC is often represented as a state diagram
- Example:

$$\mathbf{x} = 1 \ 1 \ 0 \ 0..$$
  
 $\mathbf{c} = 11 \ 10 \ 10 \ 00..$ 



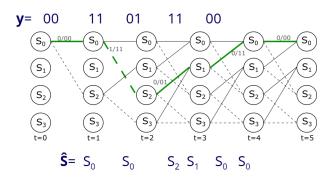
## Standard CC decoding

- **Decoder task**:  $\max_{\mathbf{x}^n} \mathbb{P}(\mathbf{x}^n | \mathbf{y}^n)$
- (Markov property) Each symbol (2 bits)  $y_t$  depends only on the states  $S_t$  and  $S_{t+1}$
- ► Decoding from a trellis:



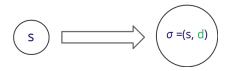
## Standard CC decoding

- **Decoder task**:  $\max_{\mathbf{y}^n} \mathbb{P}(\mathbf{x}^n | \mathbf{y}^n)$
- (Markov property) Each symbol (2 bits)  $y_t$  depends only on the states  $S_t$  and  $S_{t+1}$
- ▶ Decoding from a trellis:



## CC decoding with insertions and deletions

▶ To consider synchronization errors, another variable, called Drift is added to the decoder states  $(S_0, S_1, S_2, S_3)$ 

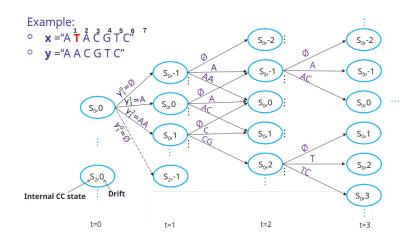


▶  $d_t = \text{Nb}(\text{Ins})_t - \text{Nb}(\text{Del})_t$  is the drift value before transmitting the symbol  $x_t$ .

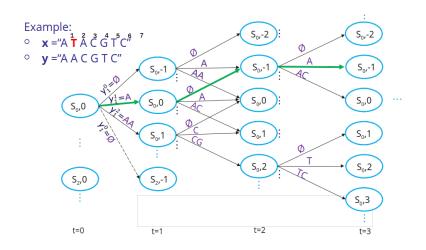
### Example:

$$x = ACAAT\_AGA$$
  
 $y = A\_AATTAGA$ 

# CC decoding with insertions and deletions



# CC decoding with insertions and deletions



## Full solution: Convolutionnal codes + LDPC codes<sup>1</sup>

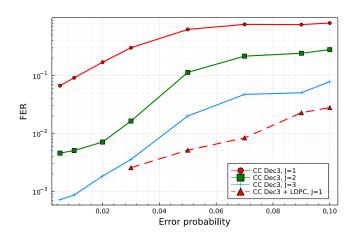


- Solution initially proposed in [Lenz21] for i.i.d. channel models
- It uses a modified BCJR algorithm for CC decoding to correct most synchronization errors
- We updated the BCJR algorithm to include the knowledge of our channel model

<sup>&</sup>lt;sup>1</sup>Belaid Hamoum, Elsa Dupraz, Channel Model and Decoder with Memory for DNA Data Storage with Nanopore Sequencing, IEEE Access, pp. 52075-52087, May 2023

## Results with the dynamic channel model

Comparison with the concatenated construction, CC+LDPC, R=1/4, N=204



## Outline

Introduction

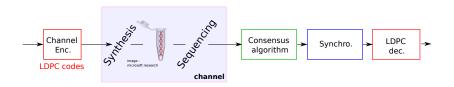
Channel mode

### Error-correction codes

Introduction
Pure coding solution (CC codes)
Mixed consensus-coding solution

Perspectives

# Mixed solution: Consensus algorithm + LDPC codes<sup>1</sup>



- ► The consensus algorithm [Lavenier21] takes **m sequences** as inputs
- ► A few errors remain after consensus
- ► The proposed synchronization method uses the LDPC code structure to correct (ins+sub) or (del+sub)

<sup>&</sup>lt;sup>1</sup>Belaid Hamoum, Aref Ezzeddine, Elsa Dupraz, Synchronization algorithms from high-rate LDPC codes for DNA data storage, DSP 2023

LDPC codes: linear block codes

**x** (Init. information)



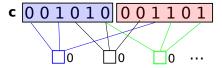
LDPC codes: linear block codes

**Codeword**:  $\mathbf{c}^n = G\mathbf{x}^k$  such that  $H\mathbf{c}^n = \mathbf{0}^m$  (HG = 0)

LDPC codes: linear block codes

- x (Init. information) p (parity bits)
  c 0 0 1 0 1 0 0 0 1 1 0 1
- **Codeword**:  $\mathbf{c}^n = G\mathbf{x}^k$  such that  $H\mathbf{c}^n = \mathbf{0}^m$  (HG = 0)

LDPC decoding by using the code parity check equations

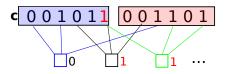


LDPC codes: linear block codes

**Codeword**: 
$$\mathbf{c}^n = G\mathbf{x}^k$$
 such that  $H\mathbf{c}^n = \mathbf{0}^m$   $(HG = 0)$ 

LDPC decoding by using the code parity check equations

#### Substitution error:



### LDPC codes

LDPC codes: linear block codes

**Codeword**: 
$$\mathbf{c}^n = G\mathbf{x}^k$$
 such that  $H\mathbf{c}^n = \mathbf{0}^m$   $(HG = 0)$ 

LDPC decoding by using the code parity check equations

**Deletion error**:

Synchronization by testing all possible positions for deletions

#### Original message:

0100101111111000011

#### Synchronization method:

- ▶ Block length: L<sub>s</sub>
- Try inserting two bits in each block so as to maximize the score
- Score: number of satisfied parity check equations
- **Example**:  $(L_s = 4)$

 $0 \; 1 \; 0 \; 0 \; | \; 1 \; 0 \; 1 \; 1 \; | \; 1 \; 1 \; 0 \; 0 \; | \; 0 \; 0 \; 1 \; 1$ 

#### Original message:

01001011**11**111000011

#### Synchronization method:

- ▶ Block length: L<sub>s</sub>
- Try inserting two bits in each block so as to maximize the score
- Score: number of satisfied parity check equations
- **Example**:  $(L_s = 4)$

$$\begin{smallmatrix} 0 & 0 & 0 & 1 & 0 & 0 & | & 1 & 0 & 1 & 1 & | & 1 & 1 & 0 & 0 & | & 0 & 0 & 1 & 1 \end{smallmatrix}$$

score: 3

#### Original message:

01001011**11**111000011

#### Synchronization method:

- ▶ Block length: L<sub>s</sub>
- Try inserting two bits in each block so as to maximize the score
- Score: number of satisfied parity check equations
- **Example**:  $(L_s = 4)$

score: 4

#### Original message:

01001011**11**111000011

#### Synchronization method:

- ▶ Block length: L<sub>s</sub>
- Try inserting two bits in each block so as to maximize the score
- Score: number of satisfied parity check equations
- **Example**:  $(L_s = 4)$

0 1 0 0 | 1 0 1 1 | 0 0 1 1 0 0 | 0 0 1 1

score: 8

#### Original message:

01001011**11**111000011

#### Synchronization method:

- ▶ Block length: L<sub>s</sub>
- Try inserting two bits in each block so as to maximize the score
- Score: number of satisfied parity check equations
- **Example**:  $(L_s = 4)$

0 1 0 0 | 1 0 1 1 | 1 1 0 0 | 0 0 0 0 1 1

score: 4

Greedy approach: insert pairs of bit one after each other until t pairs have been inserted

$$0\;1\;0\;0\;|\;1\;0\;1\;1\;|\;1\;1\;0\;0\;|\;0\;0\;1\;1$$

**Exhaustive approach**: insert the *t* pairs of bits at once

$$0\;1\;0\;0\;|\;1\;0\;1\;1\;|\;1\;1\;0\;0\;|\;0\;0\;1\;1$$

Greedy approach: insert pairs of bit one after each other until t pairs have been inserted

score: 8

**Exhaustive approach**: insert the *t* pairs of bits at once

$$0\;1\;0\;0\;|\;1\;0\;1\;1\;|\;1\;1\;0\;0\;|\;0\;0\;1\;1$$

Greedy approach: insert pairs of bit one after each other until t pairs have been inserted

$$0\ 1\ 0\ 0\ |\ 1\ 0\ 1\ 1\ |0\ 0\ 1\ 1\ 0\ 0\ |0\ 0\ 0\ 0\ 1\ 1$$

score: 14

**Exhaustive approach**: insert the t pairs of bits at once

$$0\;1\;0\;0\;|\;1\;0\;1\;1\;|\;1\;1\;0\;0\;|\;0\;0\;1\;1$$

Greedy approach: insert pairs of bit one after each other until t pairs have been inserted

$$0\ 1\ 0\ 0\ |\ 1\ 0\ 1\ 1\ |0\ 0\ 1\ 1\ 0\ 0\ |0\ 0\ 0\ 0\ 1\ 1$$

score: 14

**Exhaustive approach**: insert the *t* pairs of bits at once

 $\begin{smallmatrix} 0 & 0 & 0 & 1 & 0 & 0 & | 0 & 0 & 1 & 0 & 1 & 1 & | & 1 & 1 & 0 & 0 & | & 0 & 0 & 1 & 1 \end{smallmatrix}$ 

score: 6

Greedy approach: insert pairs of bit one after each other until t pairs have been inserted

$$0\ 1\ 0\ 0\ |\ 1\ 0\ 1\ 1\ |0\ 0\ 1\ 1\ 0\ 0\ |0\ 0\ 0\ 0\ 1\ 1$$

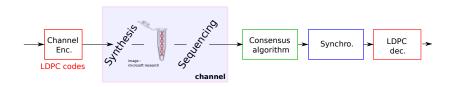
score: 14

**Exhaustive approach**: insert the *t* pairs of bits at once

$$\begin{smallmatrix} 0 & 0 & 0 & 1 & 0 & 0 & | & 1 & 0 & 1 & 1 & | 0 & 0 & 1 & 1 & 0 & 0 & | & 0 & 0 & 1 & 1 \\ \end{smallmatrix}$$

score: 12

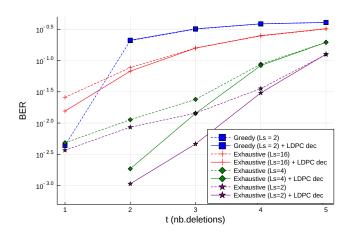
# Mixed solution: Consensus algorithm + LDPC codes<sup>1</sup>



- ► The consensus algorithm [Lavenier21] takes **m sequences** as inputs
- ► A few errors remain after consensus
- ► The proposed synchronization method uses the LDPC code structure to correct (ins+sub) or (del+sub)

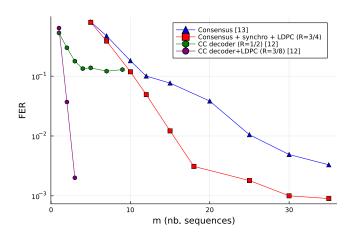
<sup>&</sup>lt;sup>1</sup>Belaid Hamoum, Aref Ezzeddine, Elsa Dupraz, Synchronization algorithms from high-rate LDPC codes for DNA data storage, DSP 2023

## Results for synchronization+LDPC



- ▶ Regular LDPC code, code rate R = 3/4, N = 256
- ► Conclusion: the exhaustive approach is more efficient but more complex

### Results for the full solution



- ▶ Regular (3,4) LDPC codes with  $N = 256 \rightarrow 192$  bits of information
- ▶ Concatenated construction of [Lenz21]:  $\rightarrow$  96 bits of information

### Outline

Introduction

Channel mode

Error-correction codes

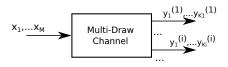
Introduction

Pure coding solution (CC codes)

Mixed consensus-coding solution

### Perspectives

### Multi-Draw channel model



#### Model

- ► Consider M sequences  $\mathbf{x}_1, \dots, \mathbf{x}_M$
- ▶ Multi-draw channel: for each  $\mathbf{x}_i$ ,  $K_i$  outputs  $\mathbf{y}_i^{(k)}$  ( $K_i$  is random)

#### Coding may help to:

- ightharpoonup Cluster the reads  $\mathbf{y}_{i}^{(k)}$
- Decode the M sequences x<sub>i</sub>
- ▶ Re-order them
- Retrieve non-observed inputs  $(K_i = 0)$
- Address the coupon collector problem