

Coded and Uncoded Trace Reconstruction

João Ribeiro

Imperial College London

Deletion channel

0 1 0 0 1 1

Deletion channel

0 1 0 0 1 1



0 1 0 0 1 1

Deletion channel

0 1 0 0 1 1

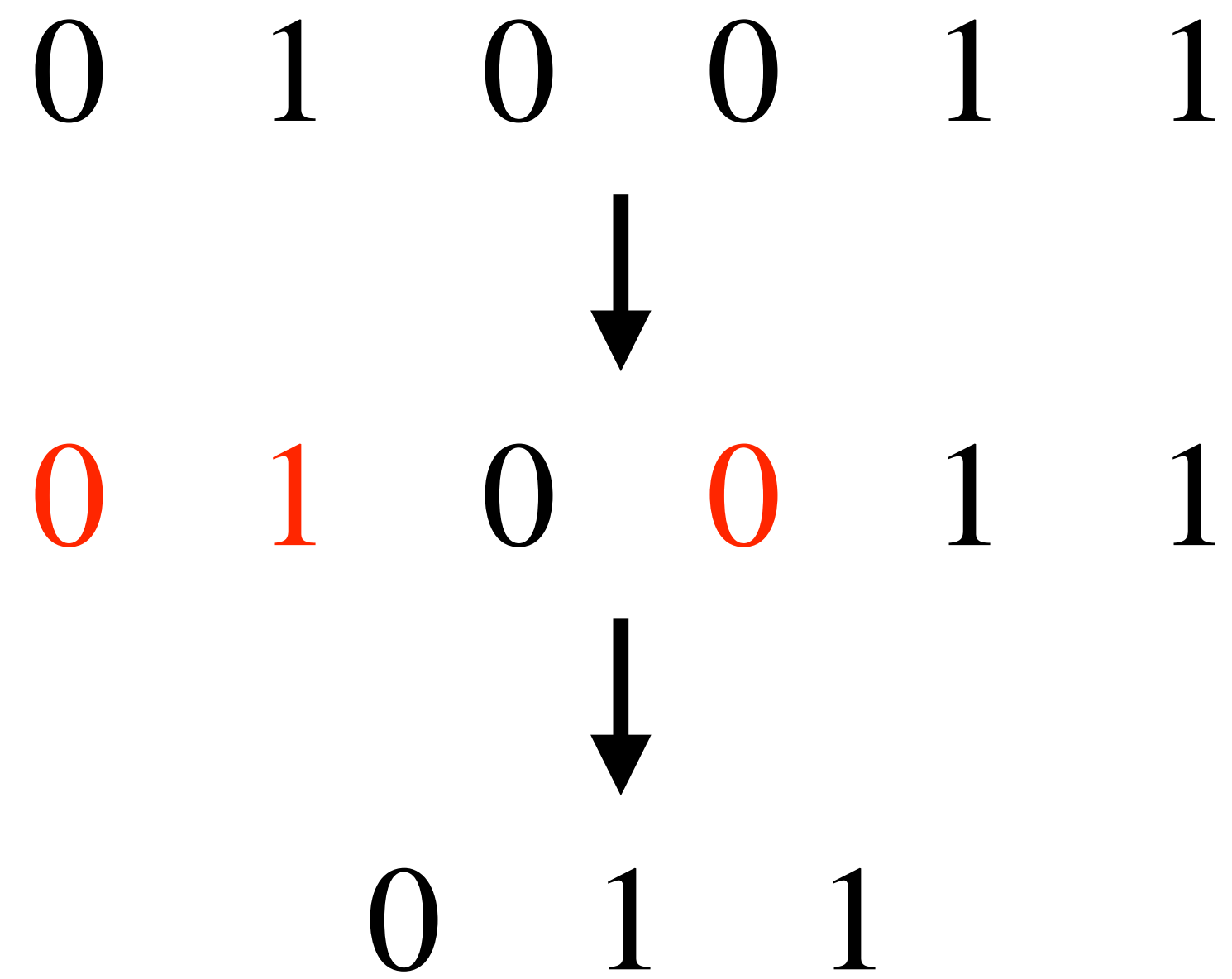


0 1 0 0 1 1

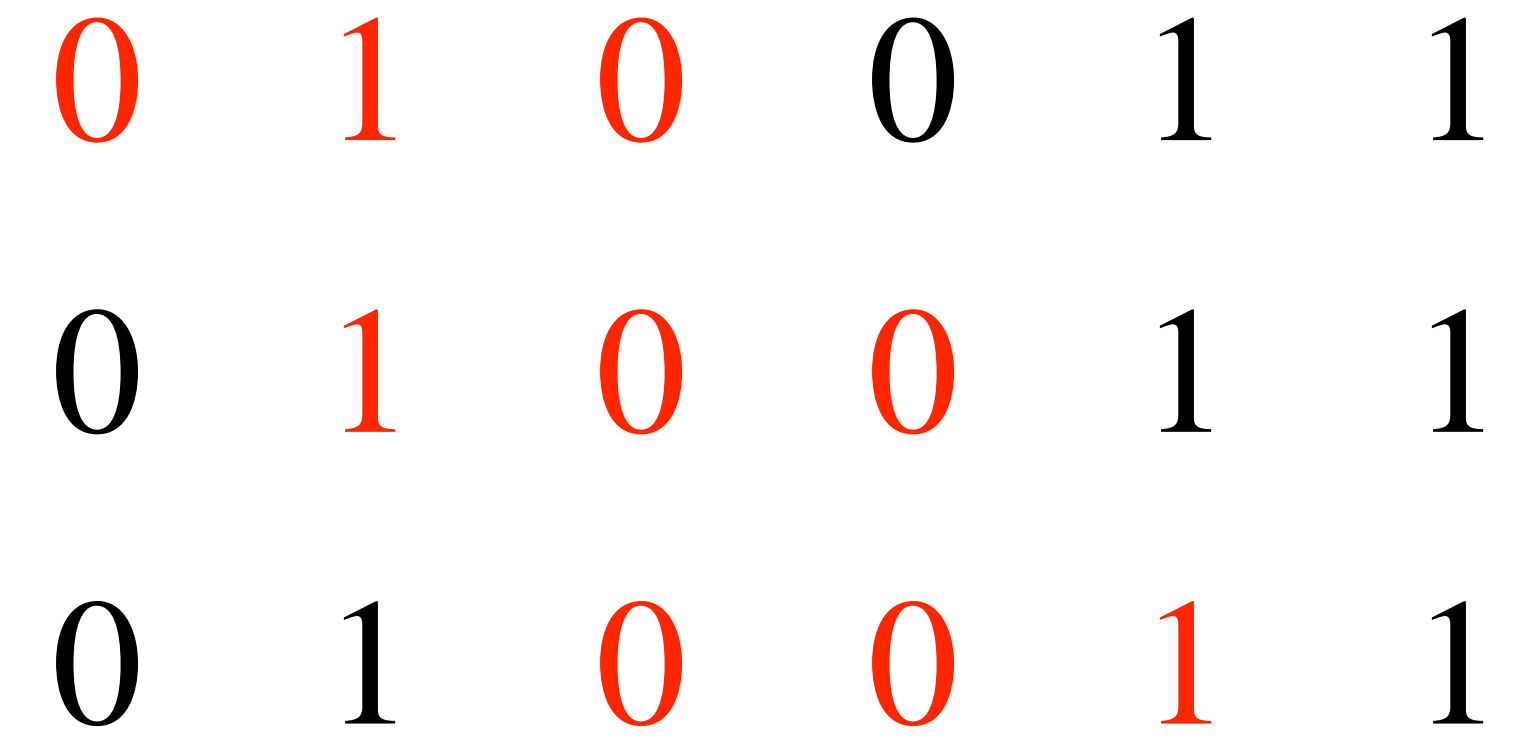


0 1 1

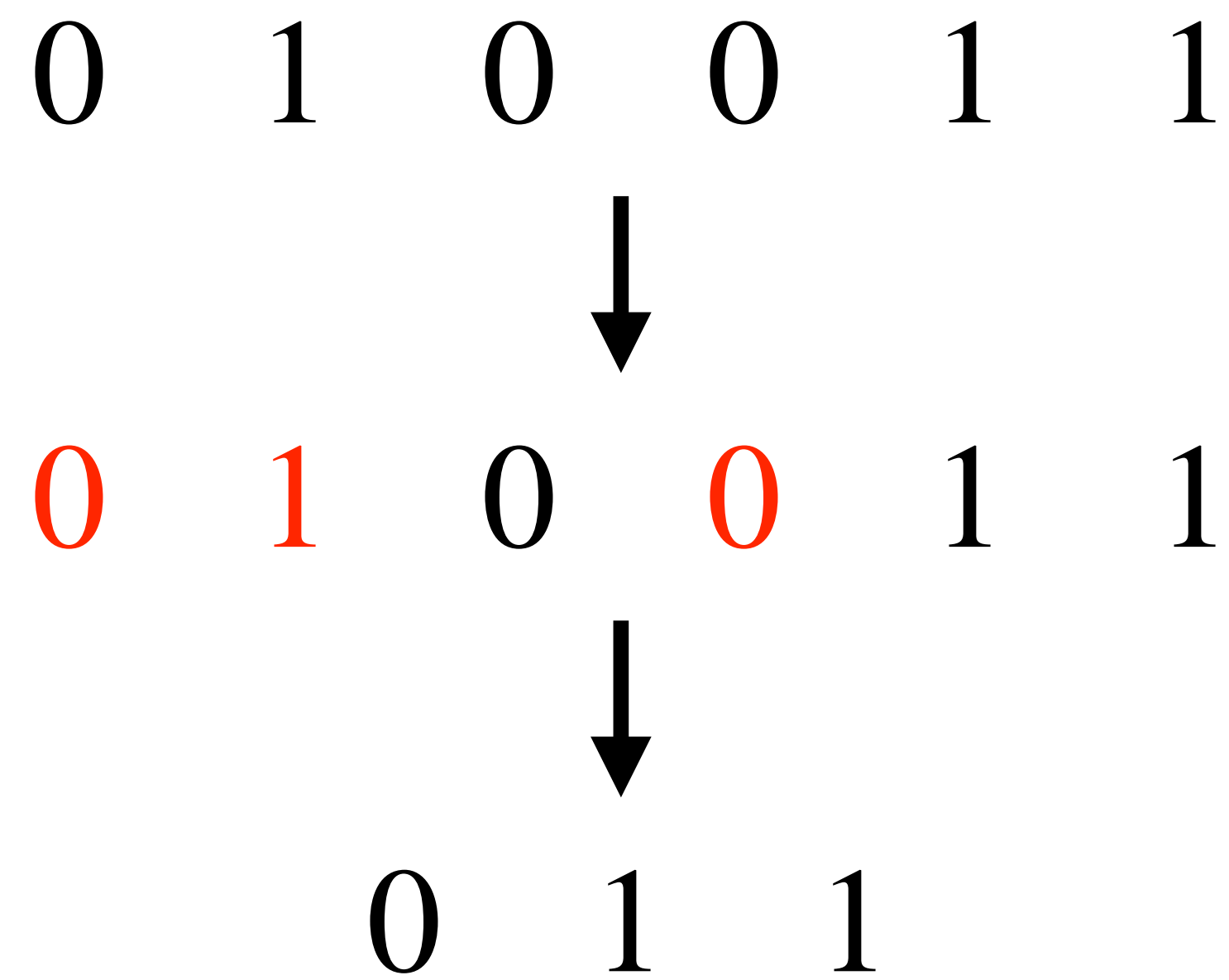
Deletion channel



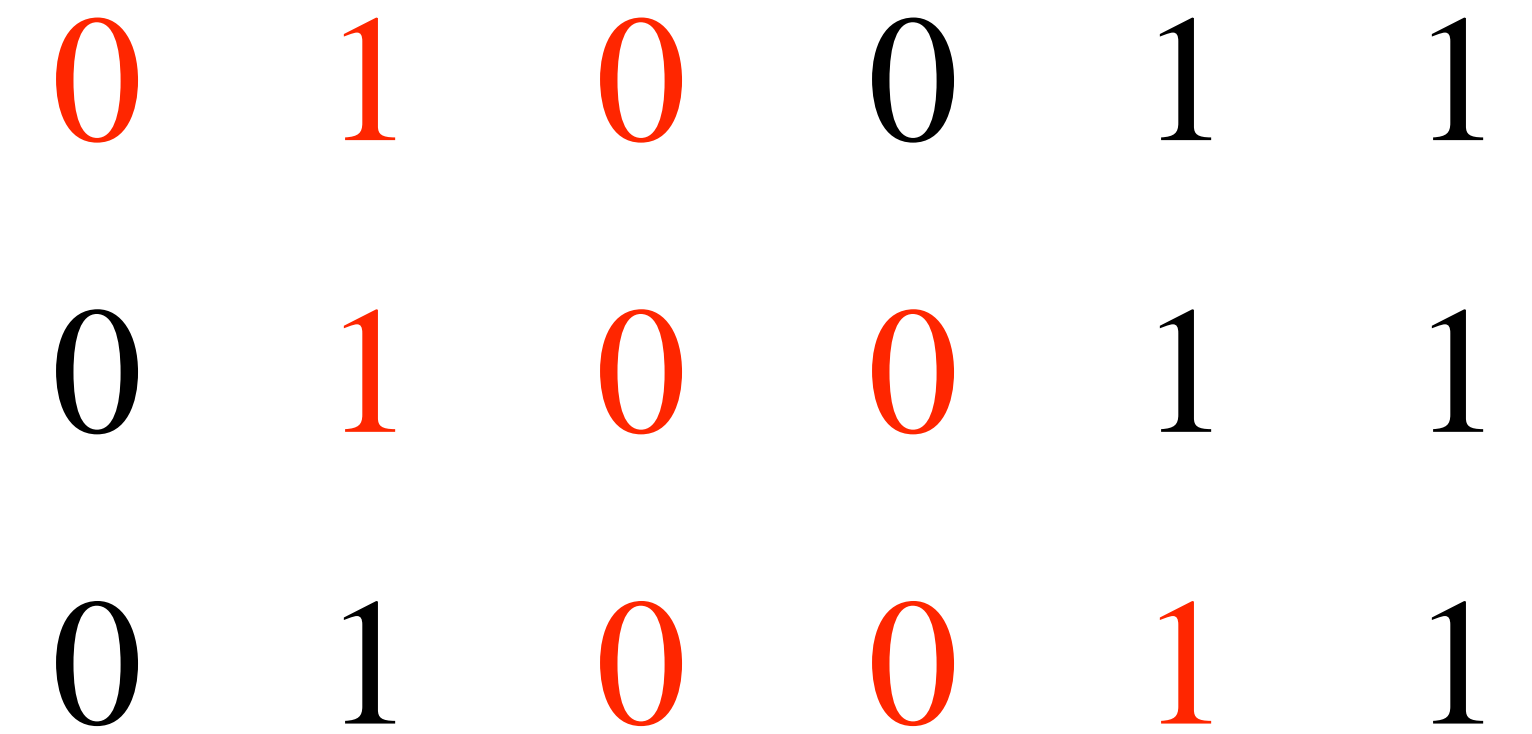
Other deletion patterns that lead to same output



Deletion channel

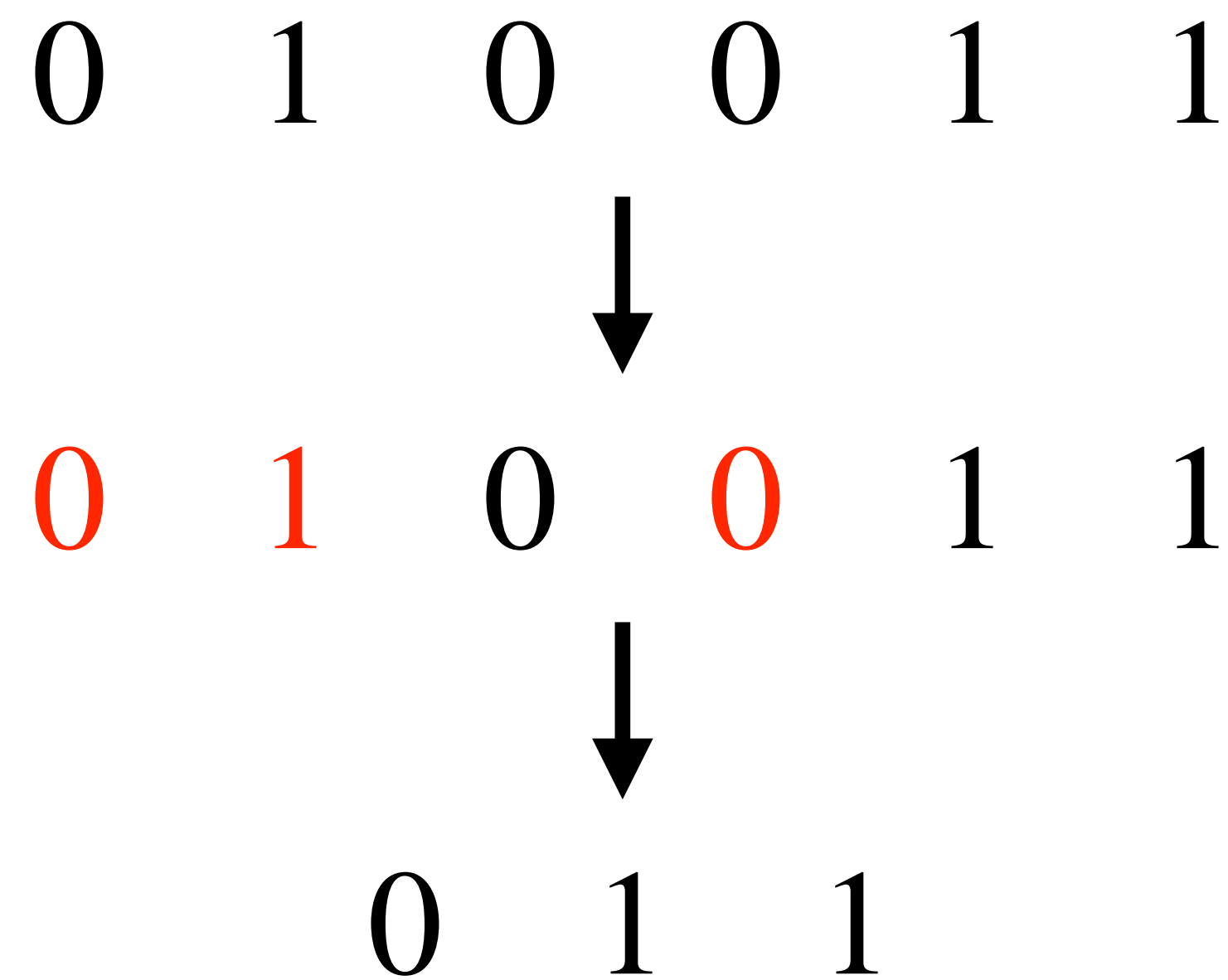


Other deletion patterns that lead to same output

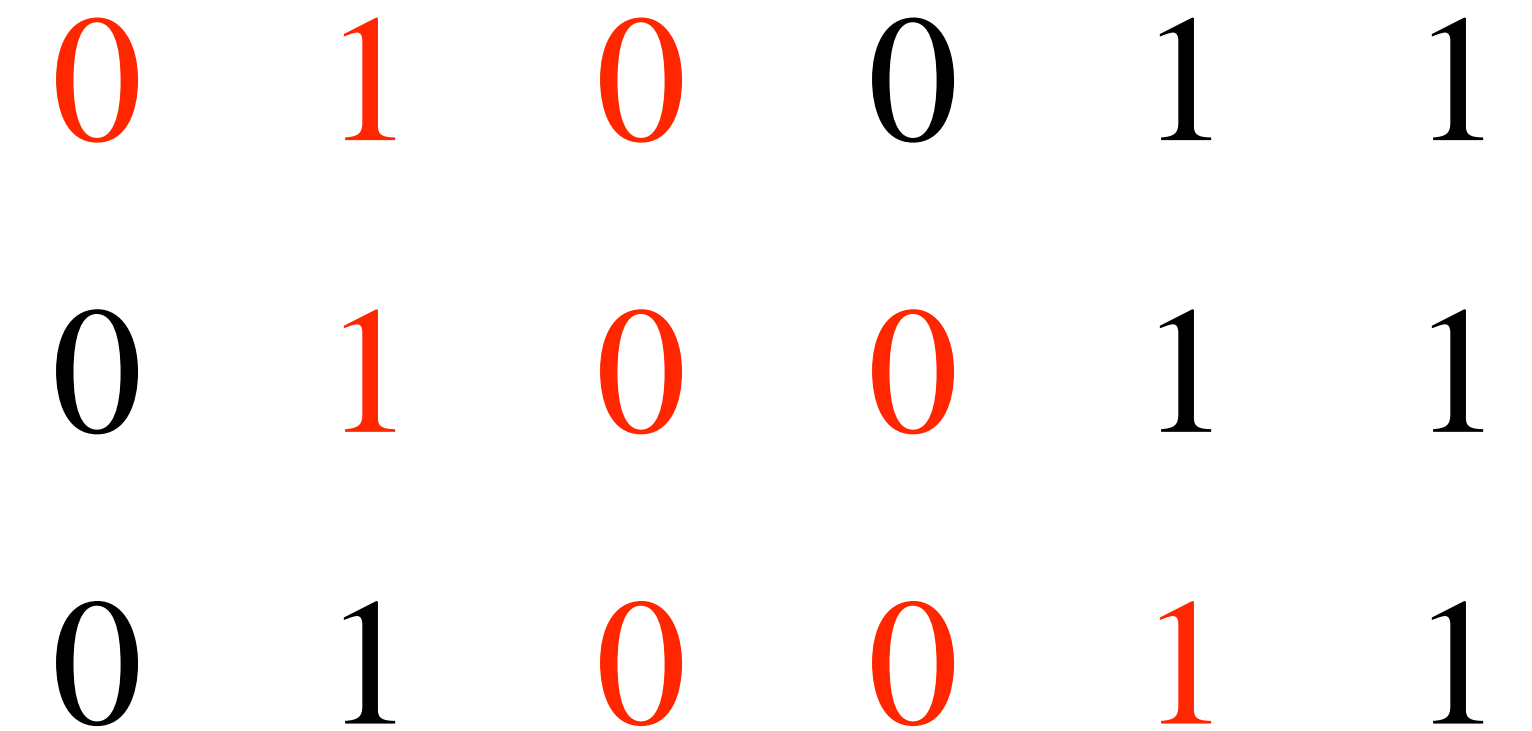


i.i.d. deletions: each bit deleted independently with fixed probability d

Deletion channel



Other deletion patterns that lead to same output

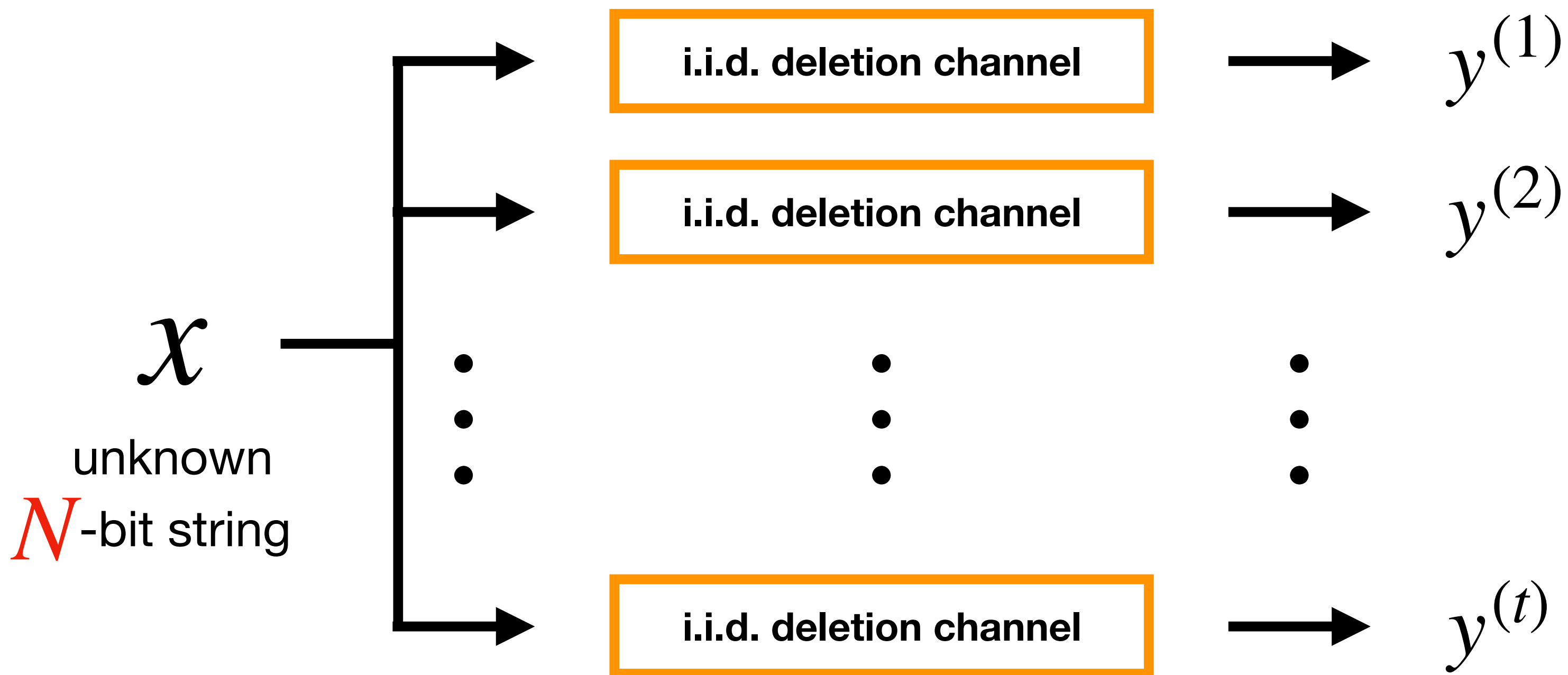


i.i.d. deletions: each bit deleted independently with fixed probability d

Many open problems! E.g., capacity is still unknown...

Trace reconstruction

[Levenshtein01, BatuKannanKhannaMcGregor04]



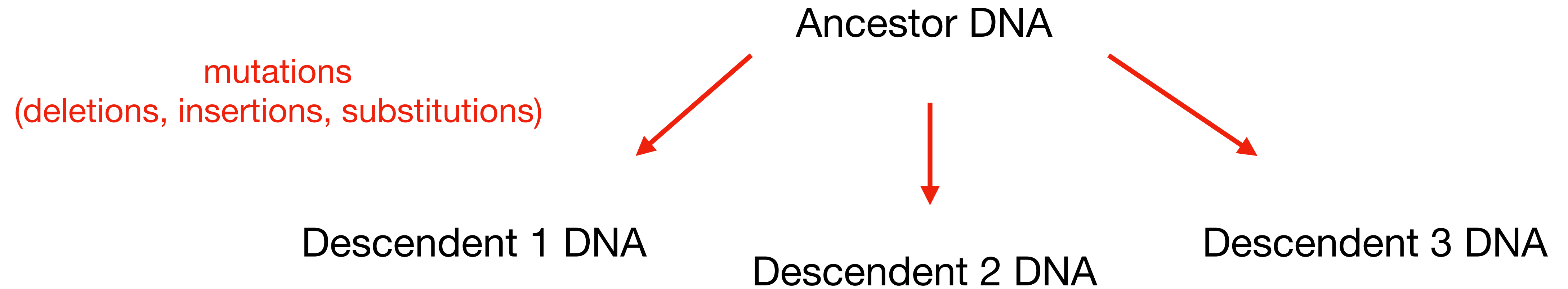
Goal

Reconstruct \mathcal{X} from **traces** $y^{(1)}, \dots, y^{(t)}$ such that:

- We succeed with high probability $1 - O(1/N)$
- We use as few traces as possible

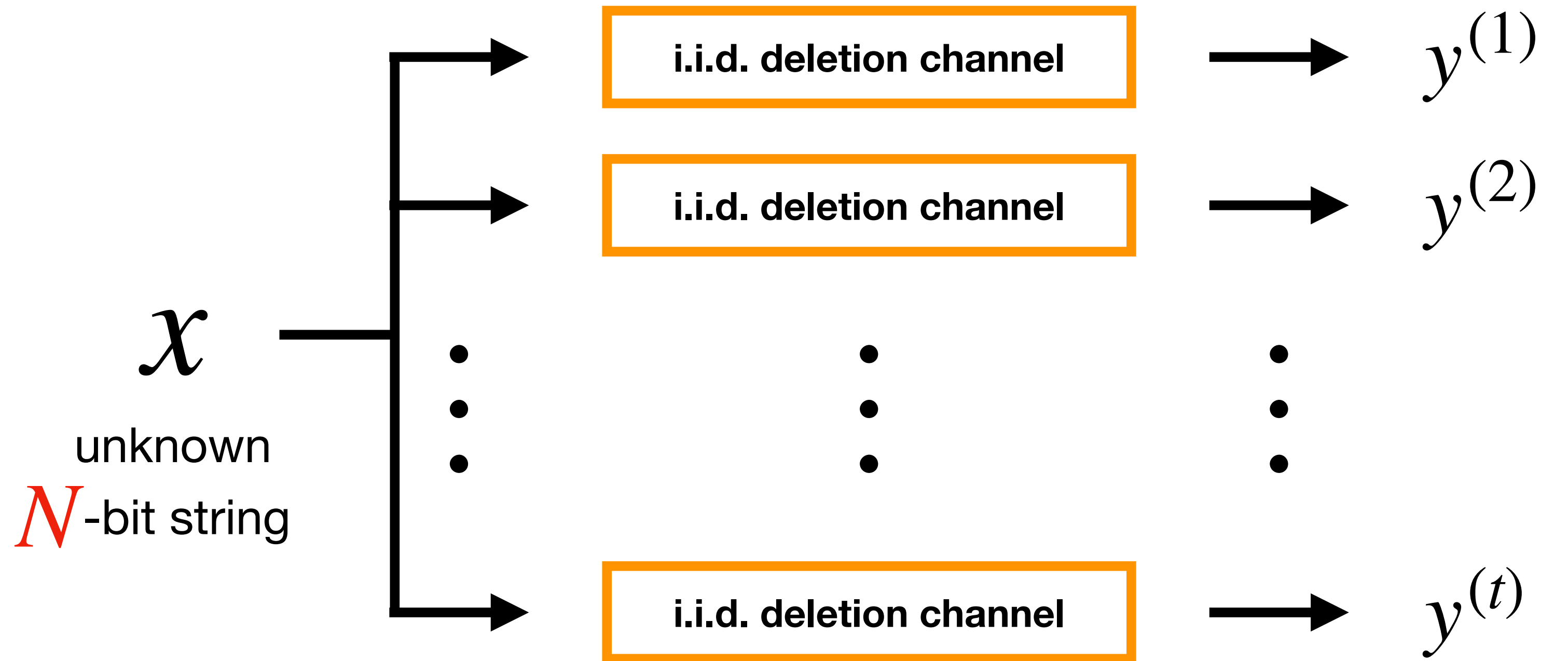
Original motivation for trace reconstruction

[BatuKannanKhannaMcGregor04]



Multiple sequence alignment: Deduce common ancestor DNA from descendants DNA.

Main settings for uncoded trace reconstruction



Worst-case trace reconstruction

*Reconstruction algorithm \mathcal{R} must succeed with high probability simultaneously for **all input strings**.*

$$\forall x : \Pr_{y^{(1)}, \dots, y^{(t)} \leftarrow \text{Del}(x)} [\mathcal{R}(y^{(1)}, \dots, y^{(t)}) = x] \approx 1$$

Average-case trace reconstruction

***Average** error probability of reconstruction algorithm \mathcal{R} over all input strings is small.*

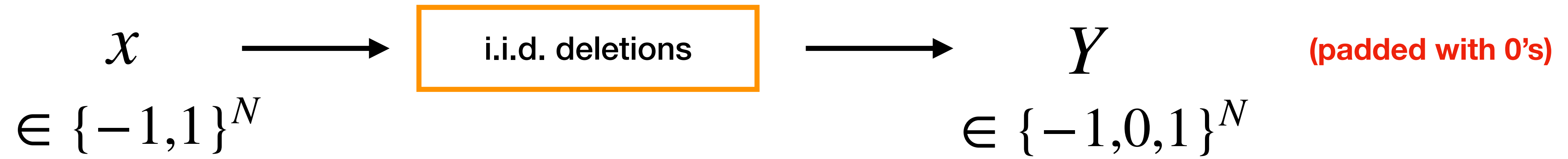
$$\Pr_{x \leftarrow \{0,1\}^N, y^{(1)}, \dots, y^{(t)} \leftarrow \text{Del}(x)} [\mathcal{R}(y^{(1)}, \dots, y^{(t)}) = x] \approx 1$$

Worst-case trace reconstruction

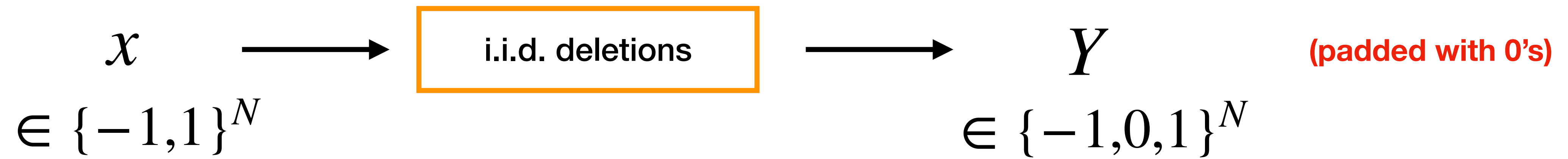
Recall: Reconstruction algorithm must succeed for **all strings** with high probability.

Reconstruction algorithms	#traces	deletion probability
Batu, Kannan, Khanna, McGregor 2004	$N \log N$	$d = 1/N^{1/2+\epsilon}$
Holenstein, Mitzenmacher, Panigrahy, Wieder 2008	$\exp(N^{1/2})$	any constant < 1
De, O'Donnell, Servedio 2017 Nazarov, Peres 2017	$\exp(N^{1/3})$	any constant < 1

Mean-based algorithms & worst-case trace reconstruction

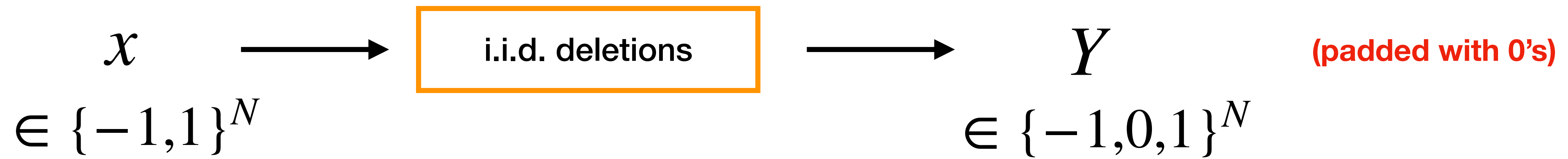


Mean-based algorithms & worst-case trace reconstruction



Mean trace: $\mu(x) = (E[Y_1], E[Y_2], \dots, E[Y_N])$ (real-linear function of x)

Mean-based algorithms & worst-case trace reconstruction



Mean trace: $\mu(x) = (E[Y_1], E[Y_2], \dots, E[Y_N])$ (real-linear function of x)

Mean-based algorithm:

- Compute estimate $\hat{\mu}$ of mean trace $\mu(x)$ from T traces;
- Find x such that $\mu(x)$ is closest to $\hat{\mu}$.

[DeO'DonnellServedio17, NazarovPeres17]:

There exist mean-based algorithms using $T = \exp(N^{1/3})$ traces. Moreover, **this is optimal for mean-based algorithms.**

Mean-based algorithms & complex analysis

Mean trace: $\mu(x) = (E[Y_1], E[Y_2], \dots, E[Y_N])$ (real-linear function of x)

#traces required for accurate estimate $\hat{\mu}$ is dictated by

$$\min_{x \neq x'} ||\mu(x) - \mu(x')||_1 = 2 \times \min_{x \in \{-1, 0, 1\}^N} ||\mu(x)||_1$$

Mean-based algorithms & complex analysis

Mean trace: $\mu(x) = (E[Y_1], E[Y_2], \dots, E[Y_N])$ (real-linear function of x)

#traces required for accurate estimate $\hat{\mu}$ is dictated by

$$\min_{x \neq x'} \|\mu(x) - \mu(x')\|_1 = 2 \times \min_{x \in \{-1, 0, 1\}^N} \|\mu(x)\|_1$$

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Mean-based algorithms & complex analysis

Mean trace: $\mu(x) = (E[Y_1], E[Y_2], \dots, E[Y_N])$ (real-linear function of x)

#traces required for accurate estimate $\hat{\mu}$ is dictated by

$$\min_{x \neq x'} \|\mu(x) - \mu(x')\|_1 = 2 \times \min_{x \in \{-1, 0, 1\}^N} \|\mu(x)\|_1$$

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Bounding $\min \|\mu(x)\|_1 \approx$ **Maximizing special polynomial over arc of complex circle**

(can use powerful complex analytic tools!)

Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

$$P_{d,x}(w) = (1 - d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1 - d)w$$

**Easy to write
in terms of x**

Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1,0,1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

**Easy to write
in terms of x**

Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1,0,1\}$

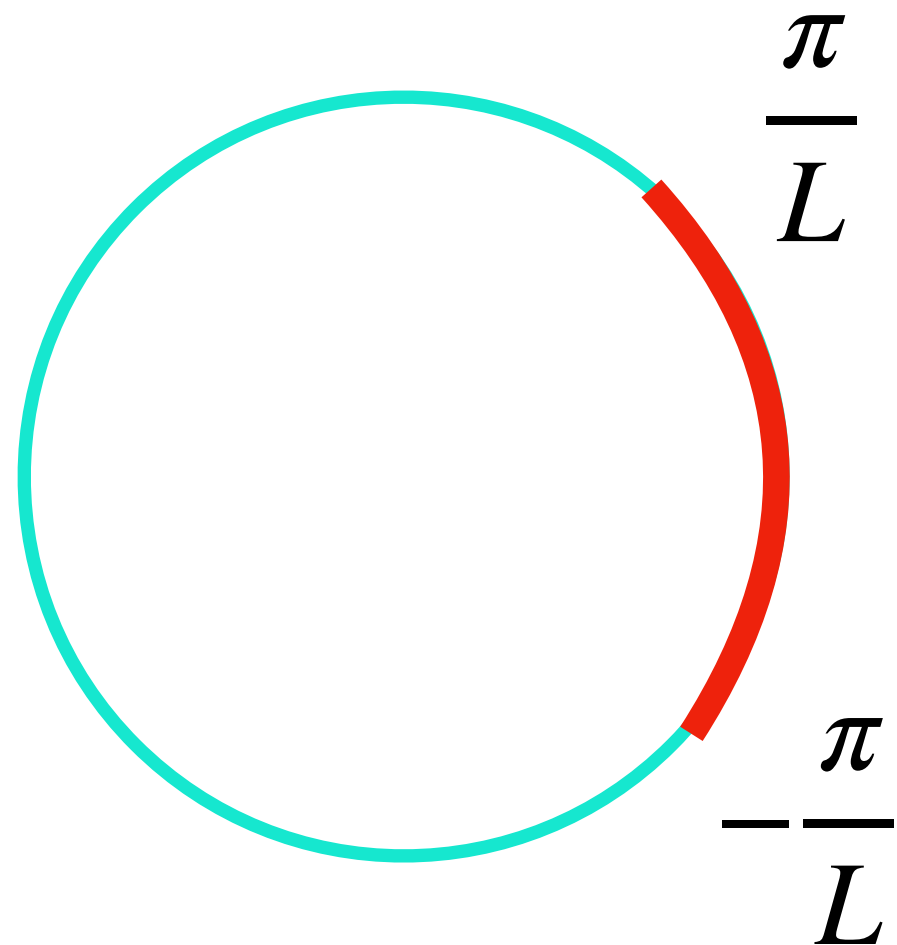
$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

**Easy to write
in terms of x**

[BorweinErdélyi97]

$A(z)$ Littlewood polynomial



Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1,0,1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

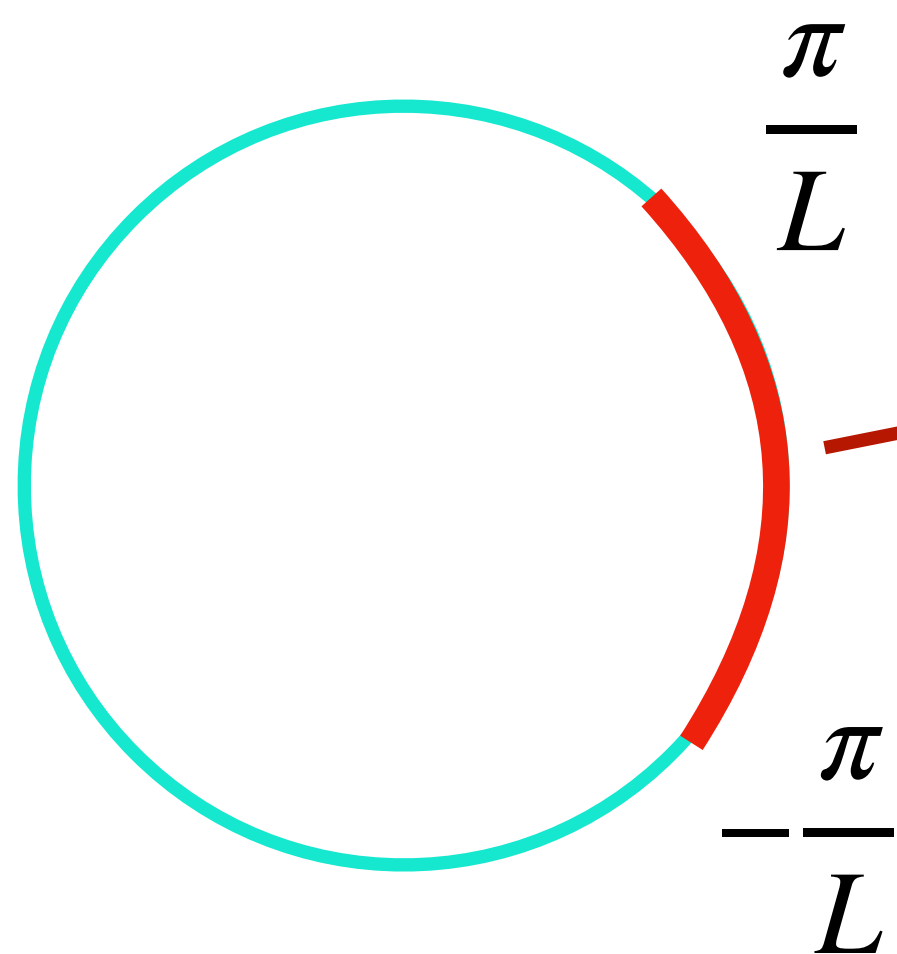
$$z = d + (1-d)w$$

Easy to write
in terms of x

[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$



Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1, 0, 1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

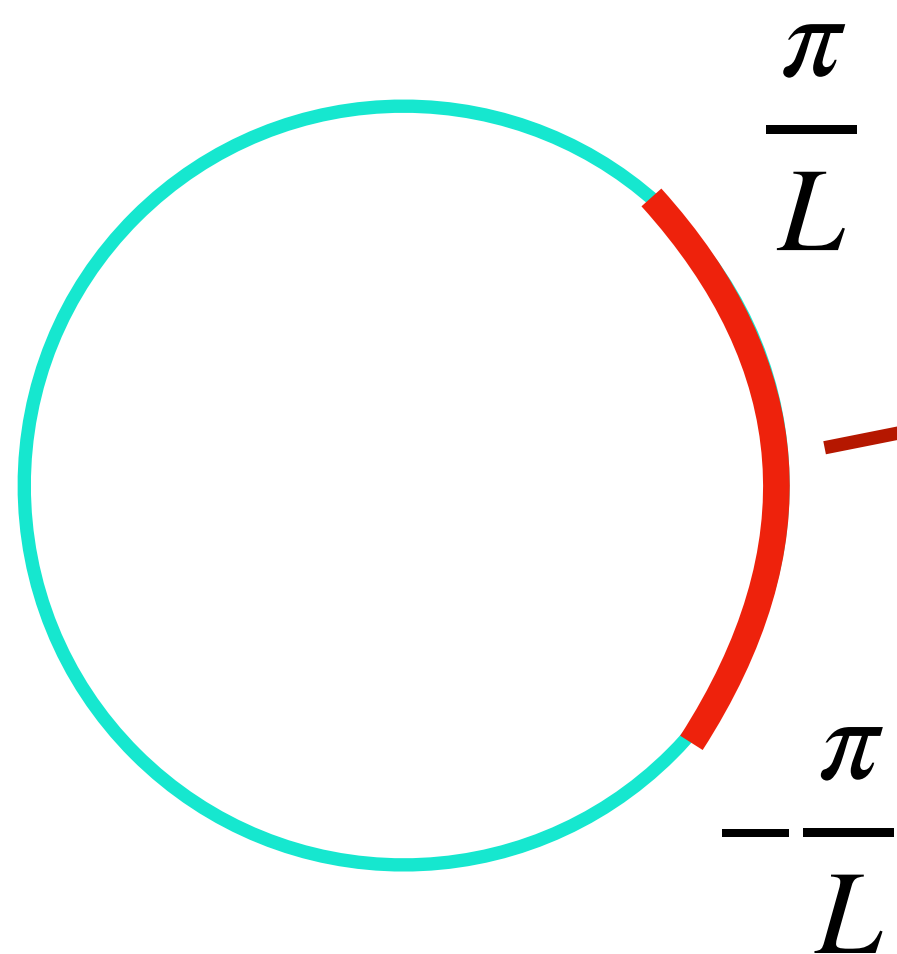
Easy to write
in terms of x

w such that $z = e^{i\theta}, \quad \theta \in [-\pi/L, \pi/L]$

[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$



Mean-based algorithms & Littlewood polynomials

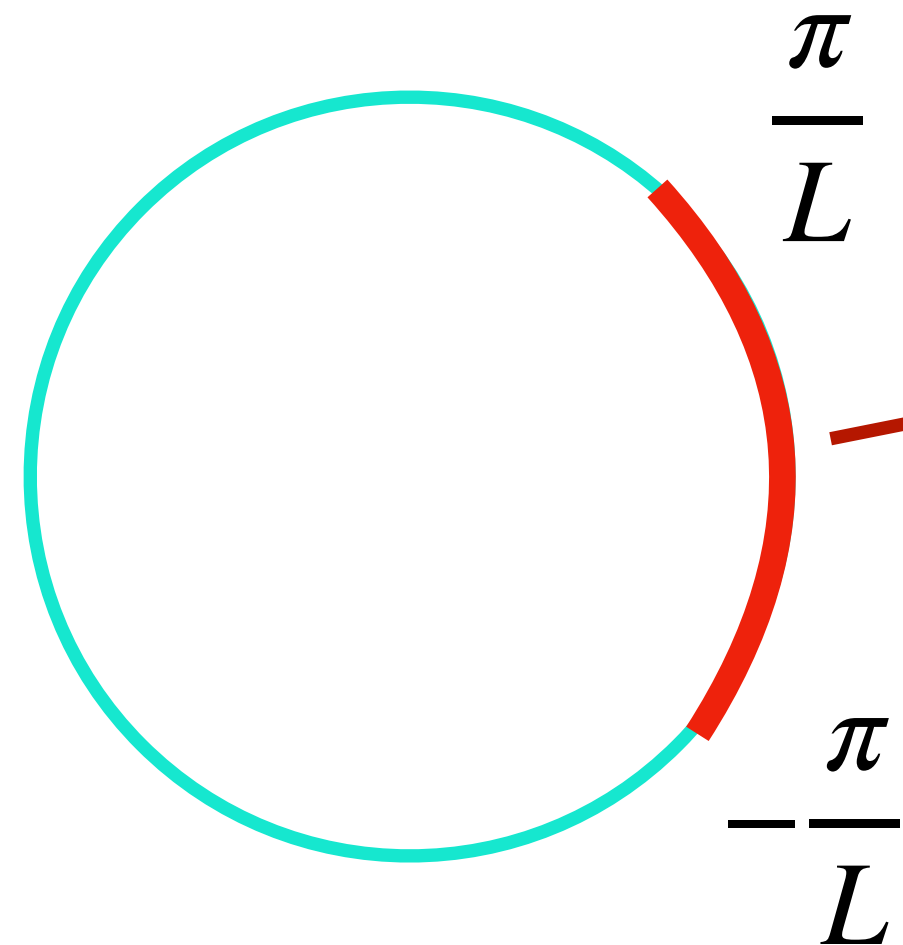
Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1, 0, 1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

Easy to write
in terms of x



[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$

w such that $z = e^{i\theta}, \quad \theta \in [-\pi/L, \pi/L]$

$$\sum_j |\mu_j(x)| \cdot |w|^j \geq (1-d) |A(e^{i\theta})|$$

(Δ -ineq)

Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1, 0, 1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

Easy to write
in terms of x

[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$

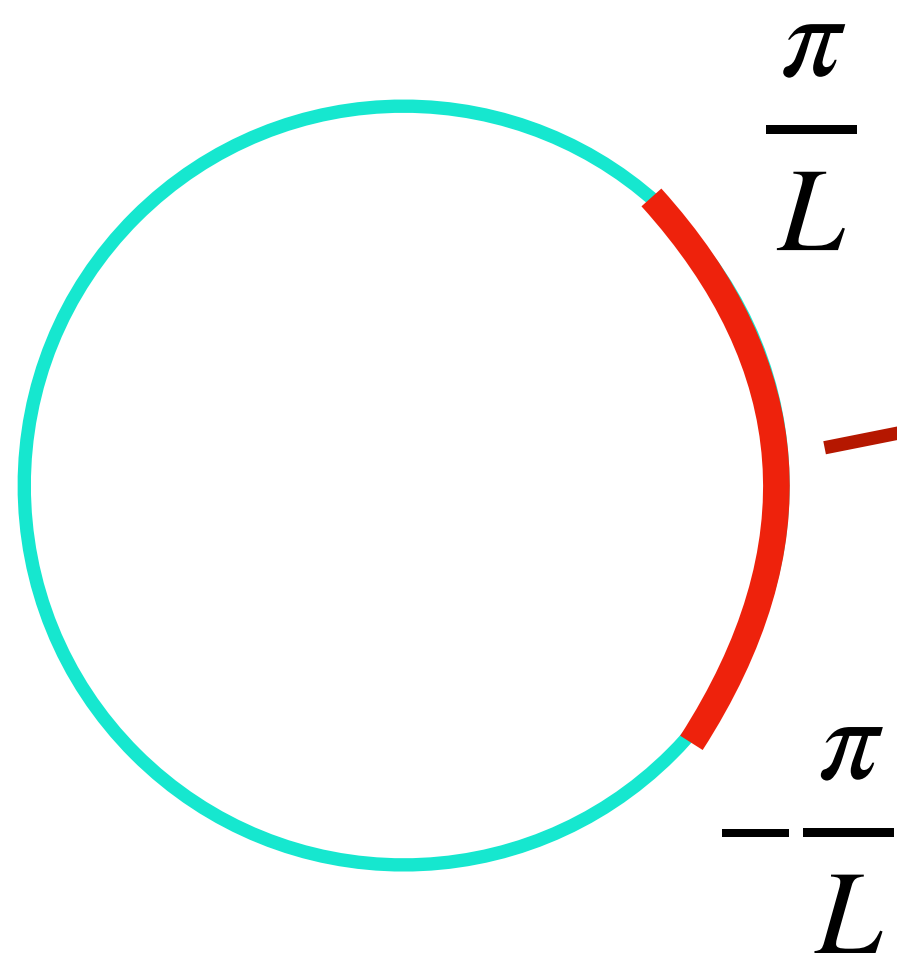
w such that $z = e^{i\theta}, \quad \theta \in [-\pi/L, \pi/L]$

$$\sum_j |\mu_j(x)| \cdot |w|^j \geq (1-d) |A(e^{i\theta})|$$

$$\implies \|\mu(x)\|_1 \geq e^{-\frac{c'N}{L^2}} \cdot |A(e^{i\theta})|$$

(Δ -ineq)

(simple trig)



Mean-based algorithms & Littlewood polynomials

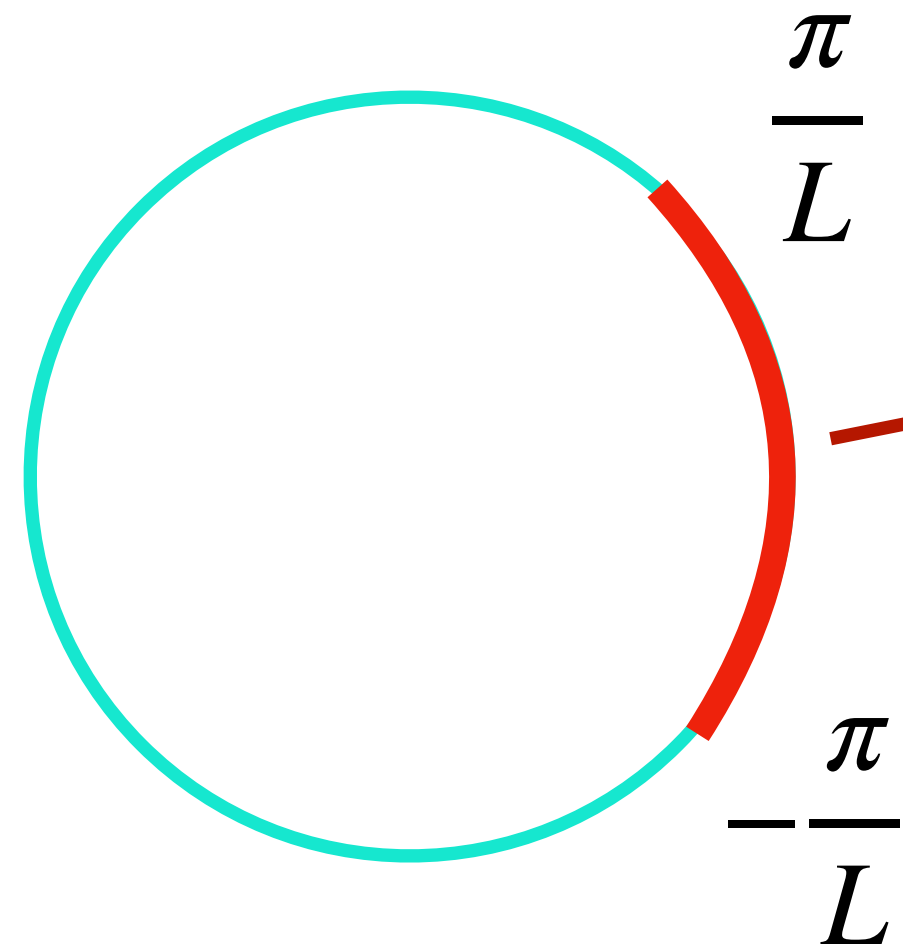
Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1,0,1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

Easy to write
in terms of x



[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$

w such that $z = e^{i\theta}, \quad \theta \in [-\pi/L, \pi/L]$

$$\sum_j |\mu_j(x)| \cdot |w|^j \geq (1-d) |A(e^{i\theta})|$$

(Δ -ineq)

$$\implies \|\mu(x)\|_1 \geq e^{-\frac{cN}{L^2}} \cdot |A(e^{i\theta})|$$

(simple trig)

$$\implies \|\mu(x)\|_1 \geq e^{-\frac{cN}{L^2}} \cdot \max_{\theta \in [-\pi/L, \pi/L]} |A(e^{i\theta})|$$

Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1,0,1\}$

$$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$$

$$z = d + (1-d)w$$

Easy to write
in terms of x

[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$

w such that $z = e^{i\theta}, \quad \theta \in [-\pi/L, \pi/L]$

$$\sum_j |\mu_j(x)| \cdot |w|^j \geq (1-d) |A(e^{i\theta})|$$

(Δ -ineq)

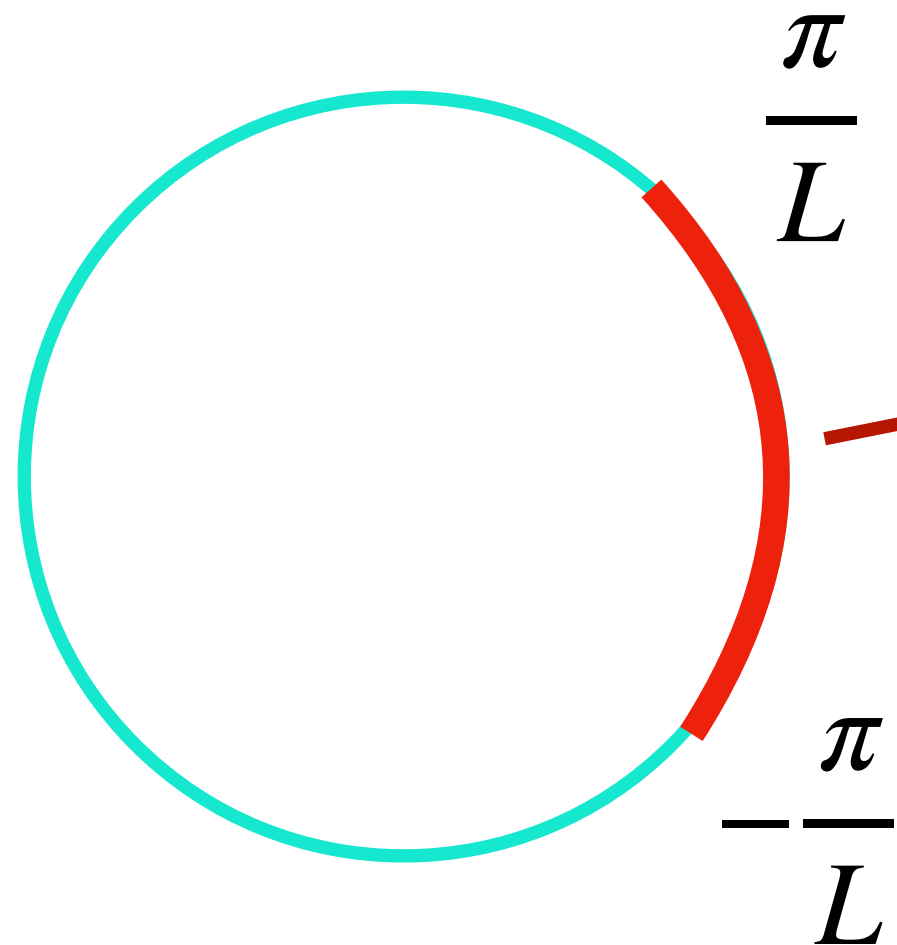
$$\implies \|\mu(x)\|_1 \geq e^{-\frac{cN}{L^2}} \cdot |A(e^{i\theta})|$$

(simple trig)

$$\implies \|\mu(x)\|_1 \geq e^{-\frac{cN}{L^2}} \cdot \max_{\theta \in [-\pi/L, \pi/L]} |A(e^{i\theta})|$$

$$\geq e^{-cN^{1/3}}$$

([BE97] + Max wrt L)

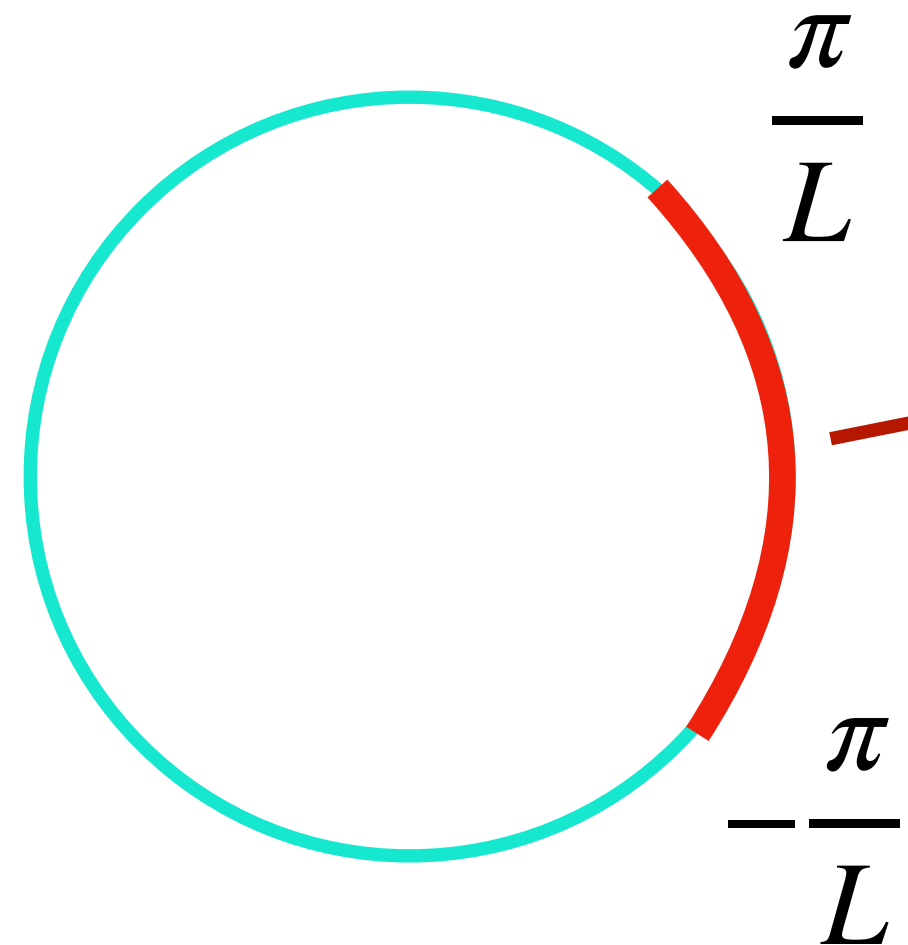


Mean-based algorithms & Littlewood polynomials

Deletion channel polynomial $P_{d,x}(w) = \sum_{j=1}^N \mu(x)_j \cdot w^j, \quad w \in \mathbb{C}$

Littlewood polynomial $A(z), \quad x_j \in \{-1,0,1\}$

$P_{d,x}(w) = (1-d) \sum_{j=1}^N x_j \cdot z^j$	$z = d + (1-d)w$	<p>Easy to write in terms of x</p>
---	------------------	---



[BorweinErdélyi97]

$A(z)$ Littlewood polynomial

$$\max |A(e^{i\theta})| \geq e^{-cL}$$

w such that $z = e^{i\theta}, \quad \theta \in [-\pi/L, \pi/L]$

$$\sum_j |\mu_j(x)| \cdot |w|^j \geq (1-d) |A(e^{i\theta})|$$

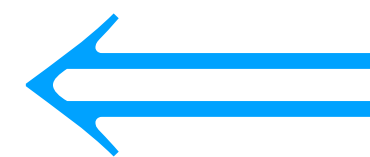
(Δ -ineq)

$$\implies \|\mu(x)\|_1 \geq e^{-\frac{cN}{L^2}} \cdot |A(e^{i\theta})|$$

(simple trig)

$$\implies \|\mu(x)\|_1 \geq e^{-\frac{cN}{L^2}} \cdot \max_{\theta \in [-\pi/L, \pi/L]} |A(e^{i\theta})|$$

$e^{O(N^{1/3})}$ traces suffice to distinguish mean traces



$$\geq e^{-cN^{1/3}}$$

([BE97] + Max wrt L)

Average-case trace reconstruction

Recall: Average error probability of reconstruction algorithm must be $O(1/N)$.

Reconstruction algorithms	#traces	deletion probability
Batu, Kannan, Khanna, McGregor 2004	$\log N$	$d = 1/\log N$
Holenstein, Mitzenmacher, Panigrahy, Wieder 2008	poly (N)	$d < c$ for small absolute constant c
Peres, Zhai 2017	$\exp(\log^{1/2} N)$	$d < 1/2$
Holden, Pemantle, Peres 2018	$\exp(\log^{1/3} N)$	$d < 1$

A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

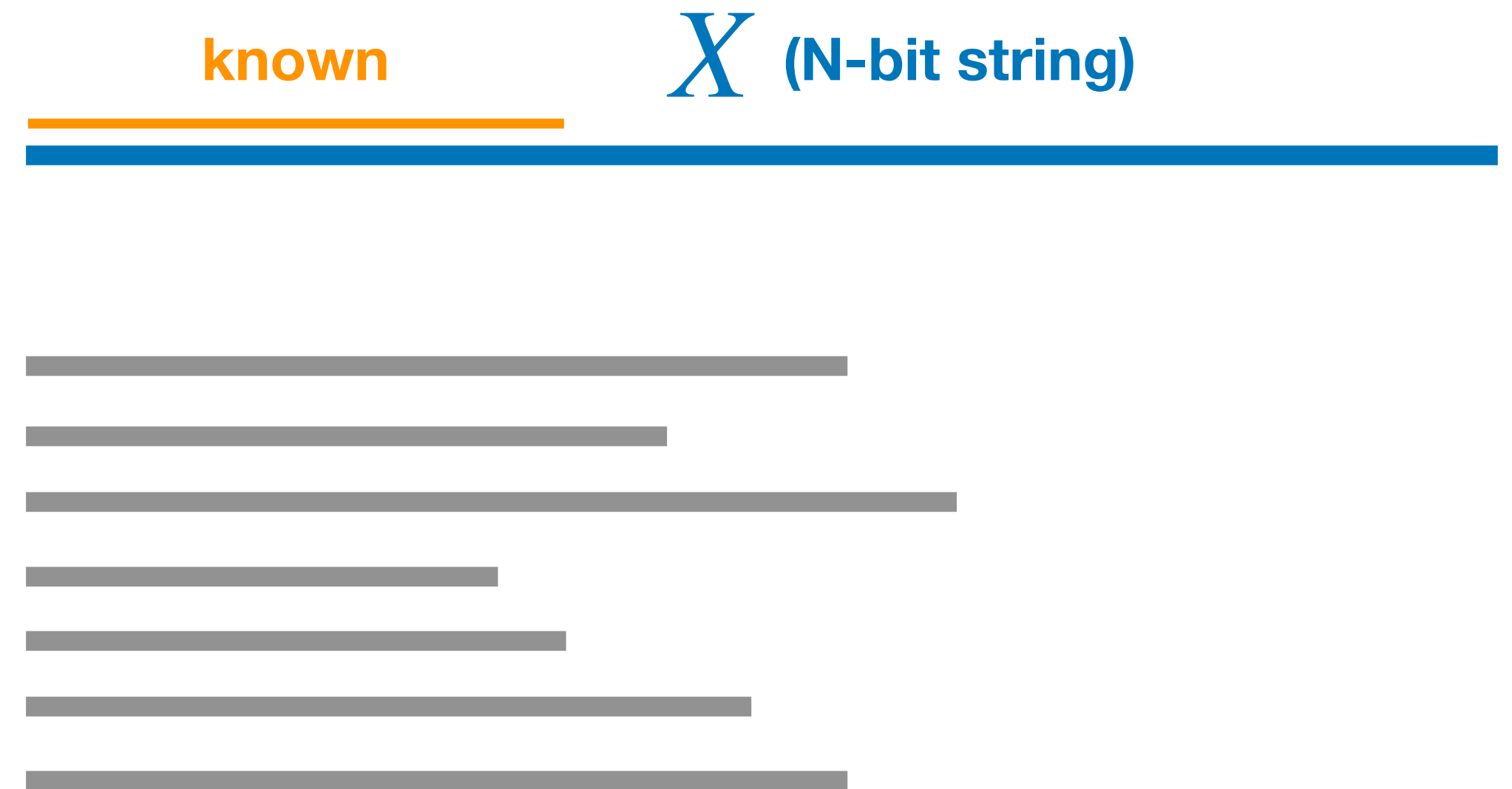
X (N-bit string)



A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

1) **Bootstrapping:** Learn first bits of X “for free”



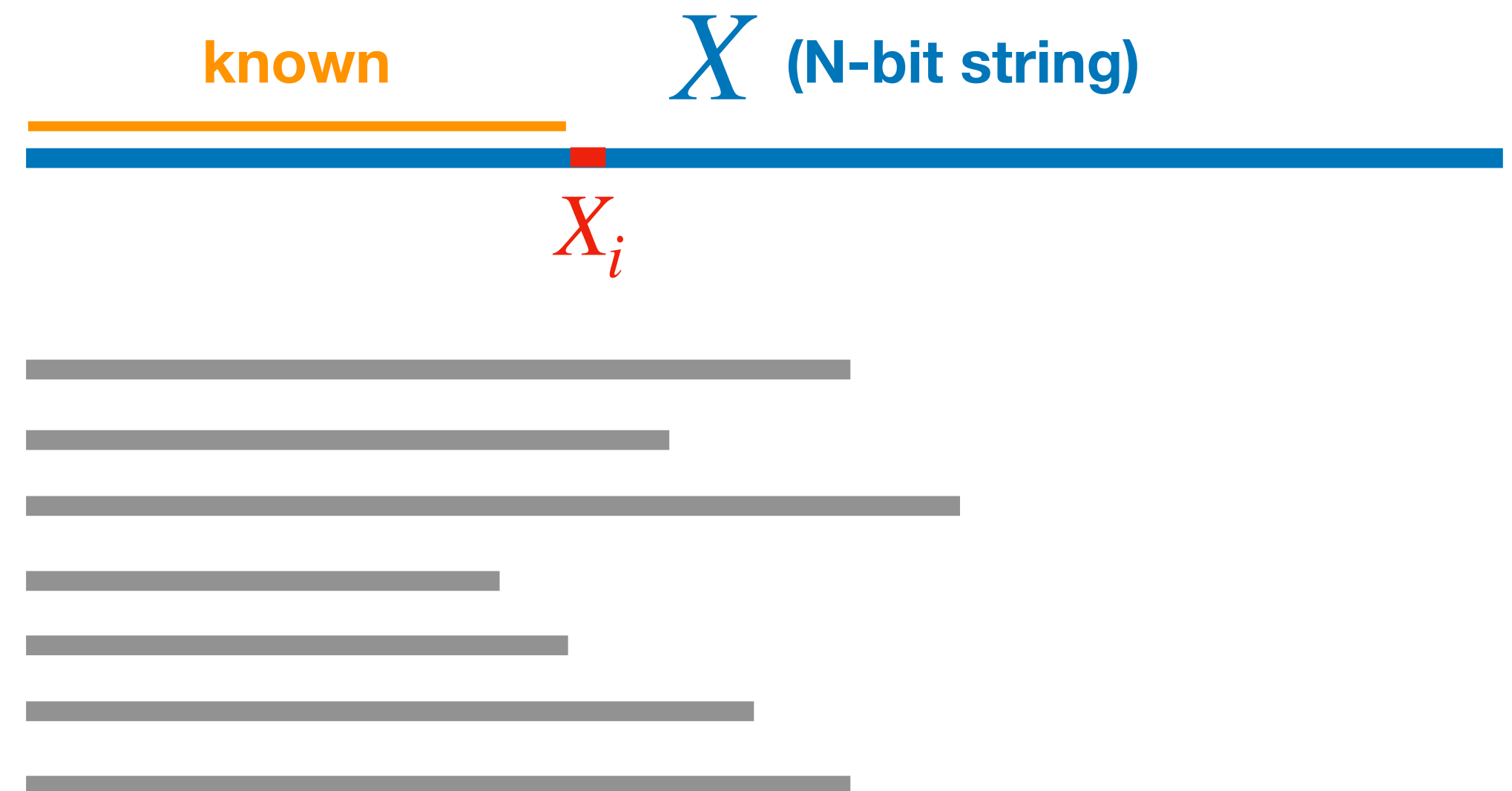
A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

1) Bootstrapping: Learn first bits of X “for free”

Suppose we know X_1, X_2, \dots, X_{i-1}

Goal: Find X_i



A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

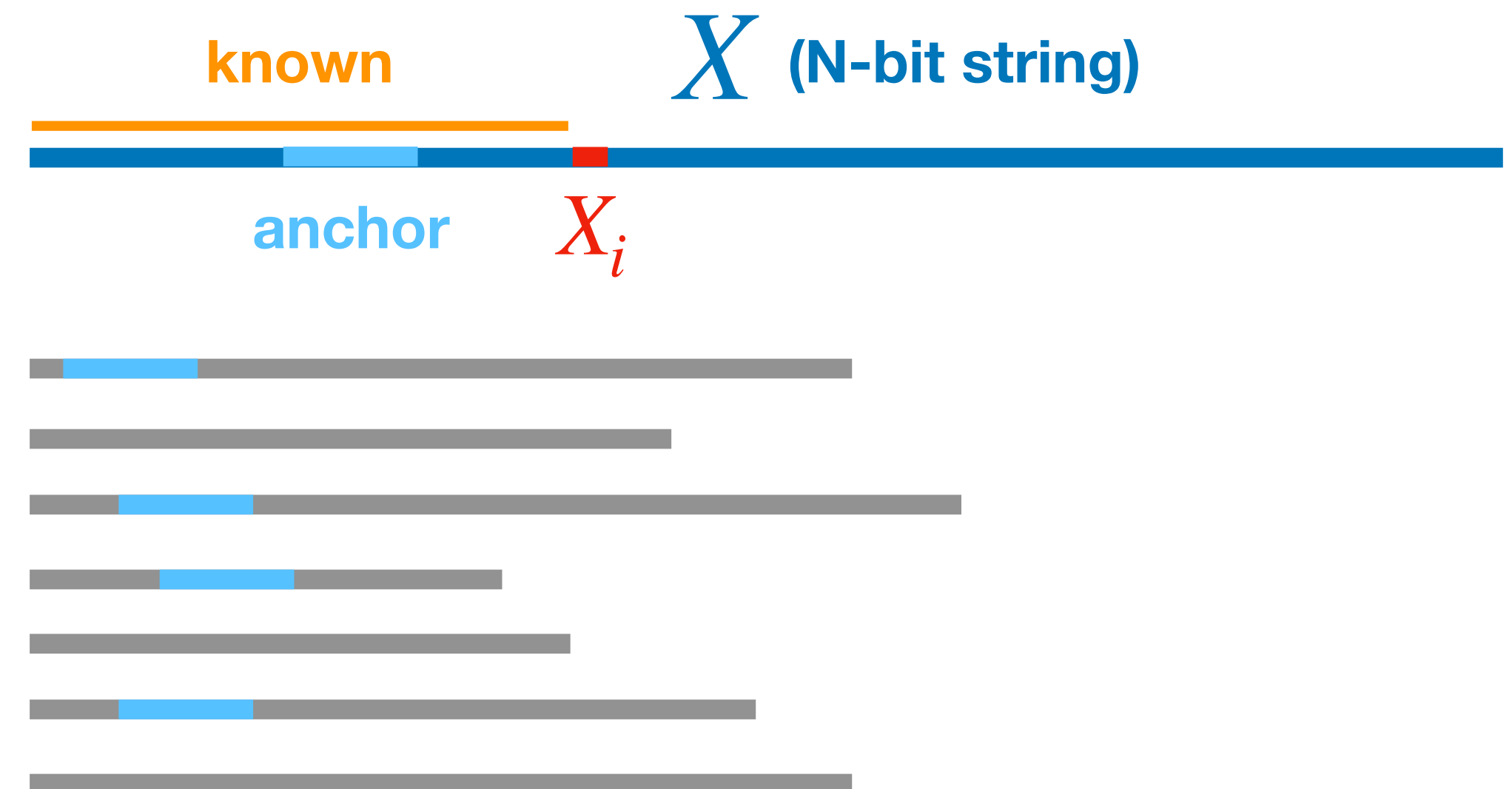
1) Bootstrapping: Learn first bits of X “for free”

Suppose we know X_1, X_2, \dots, X_{i-1}

Goal: Find X_i

2) Trace alignment: Align by **anchor** close to X_i

“If X is random, whp anchor in trace comes from anchor in X .”



A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

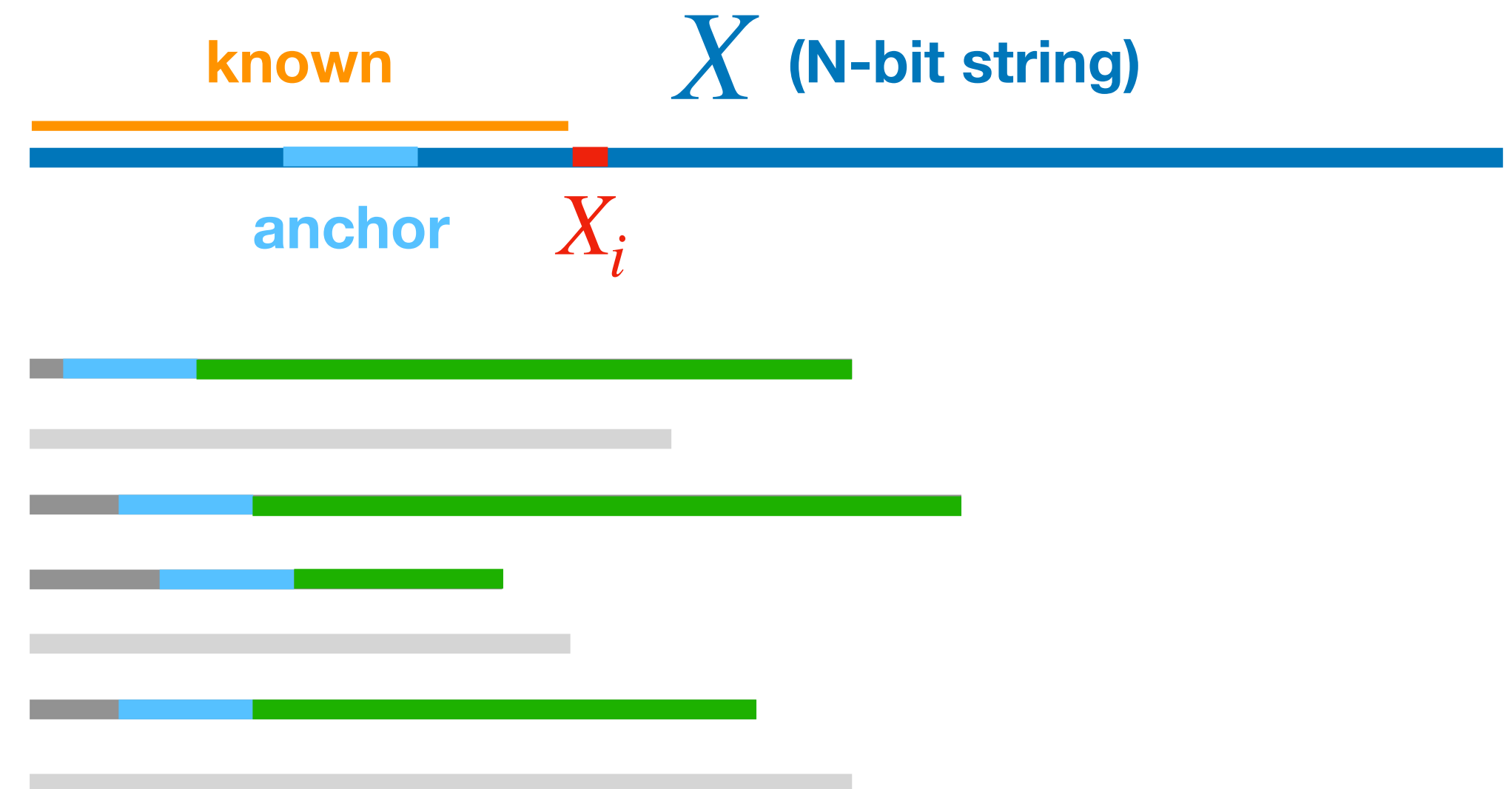
1) Bootstrapping: Learn first bits of X “for free”

Suppose we know X_1, X_2, \dots, X_{i-1}

Goal: Find X_i

2) Trace alignment: Align by **anchor** close to X_i

“If X is random, whp anchor in trace comes from anchor in X .”



A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

1) Bootstrapping: Learn first bits of X “for free”

Suppose we know X_1, X_2, \dots, X_{i-1}

Goal: Find X_i

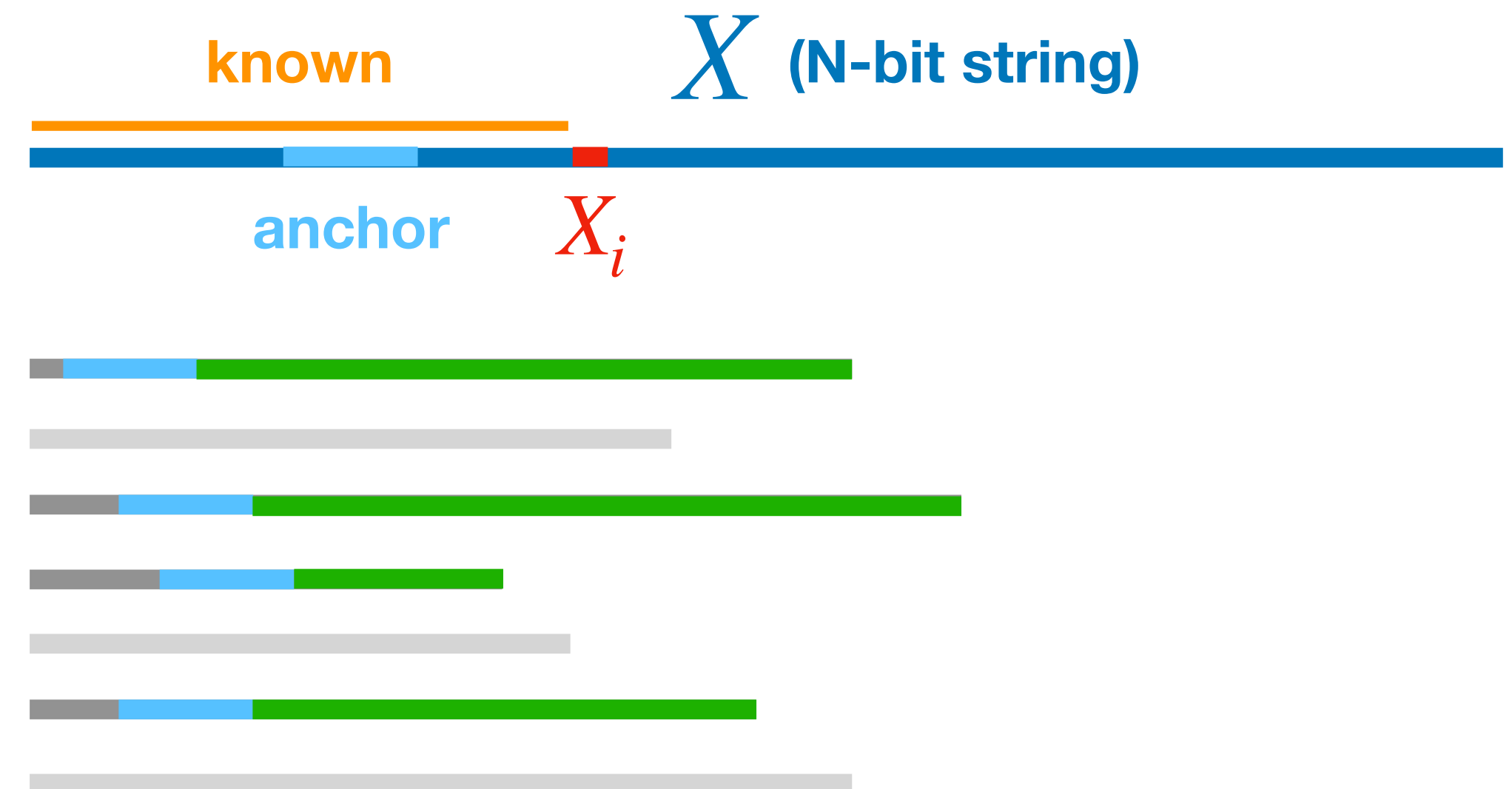
2) Trace alignment: Align by **anchor** close to X_i

“If X is random, whp anchor in trace comes from anchor in X .”

3) Reconstruction: Estimate special bit

$Y =$ distribution of **“trace after anchor”**

“There is special position Y_{j^\star} that is decently influenced by X_i ”



$$|\Pr[Y_{j^\star} = 1 | X_i = \underline{1}] - \Pr[Y_{j^\star} = 1 | X_i = \underline{0}]| \geq \frac{1}{N^C}$$

A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

1) Bootstrapping: Learn first bits of X “for free”

Suppose we know X_1, X_2, \dots, X_{i-1}

Goal: Find X_i

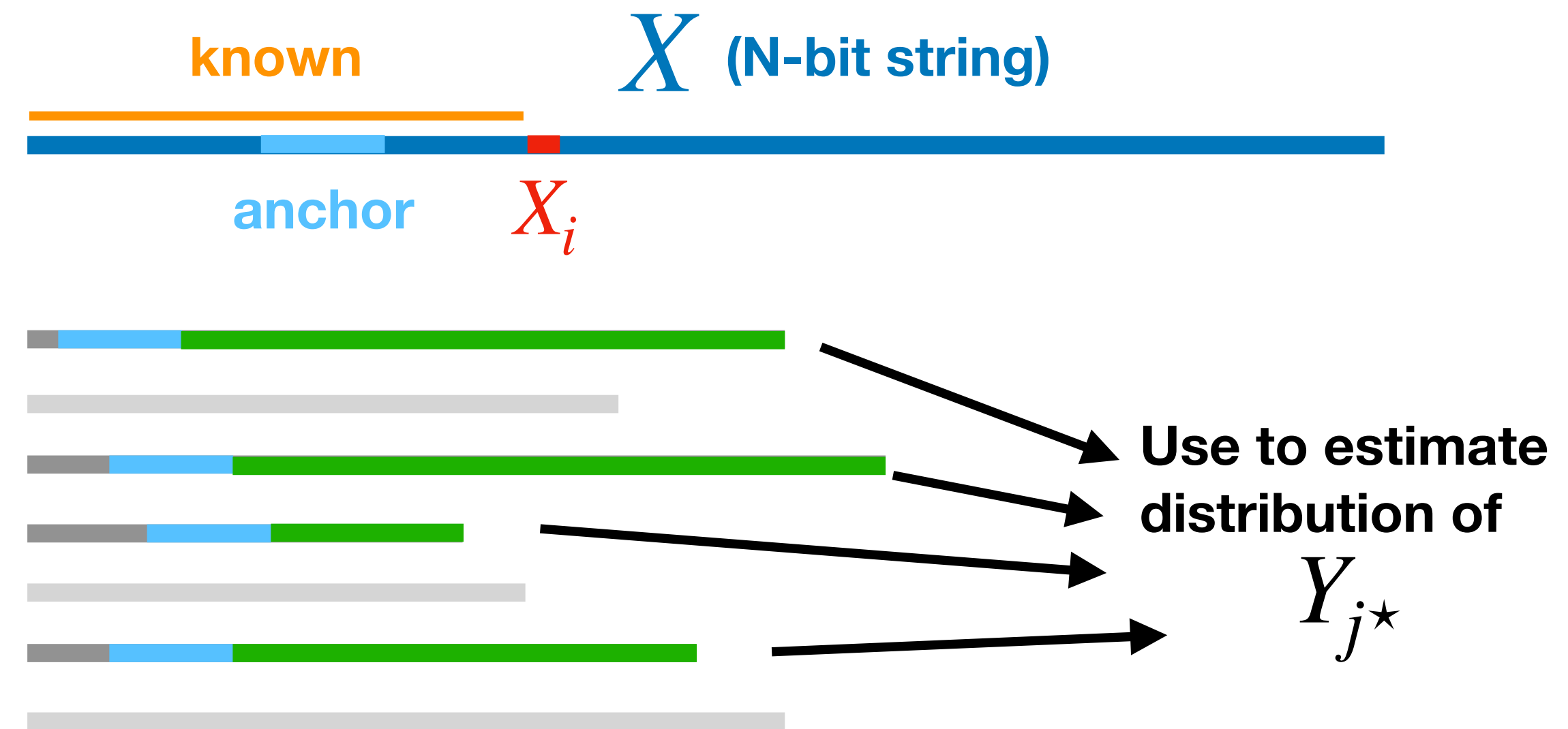
2) Trace alignment: Align by **anchor** close to X_i

“If X is random, whp anchor in trace comes from anchor in X .”

3) Reconstruction: Estimate special bit

$Y =$ distribution of **“trace after anchor”**

“There is special position Y_{j^\star} that is decently influenced by X_i ”



$$|\Pr[Y_{j^\star} = 1 | X_i = \underline{1}] - \Pr[Y_{j^\star} = 1 | X_i = \underline{0}]| \geq \frac{1}{N^C}$$

A recipe for average-case trace reconstruction

[HolensteinMitzenmacherPanigrahyWieder08]

1) Bootstrapping: Learn first bits of X “for free”

Suppose we know X_1, X_2, \dots, X_{i-1}

Goal: Find X_i

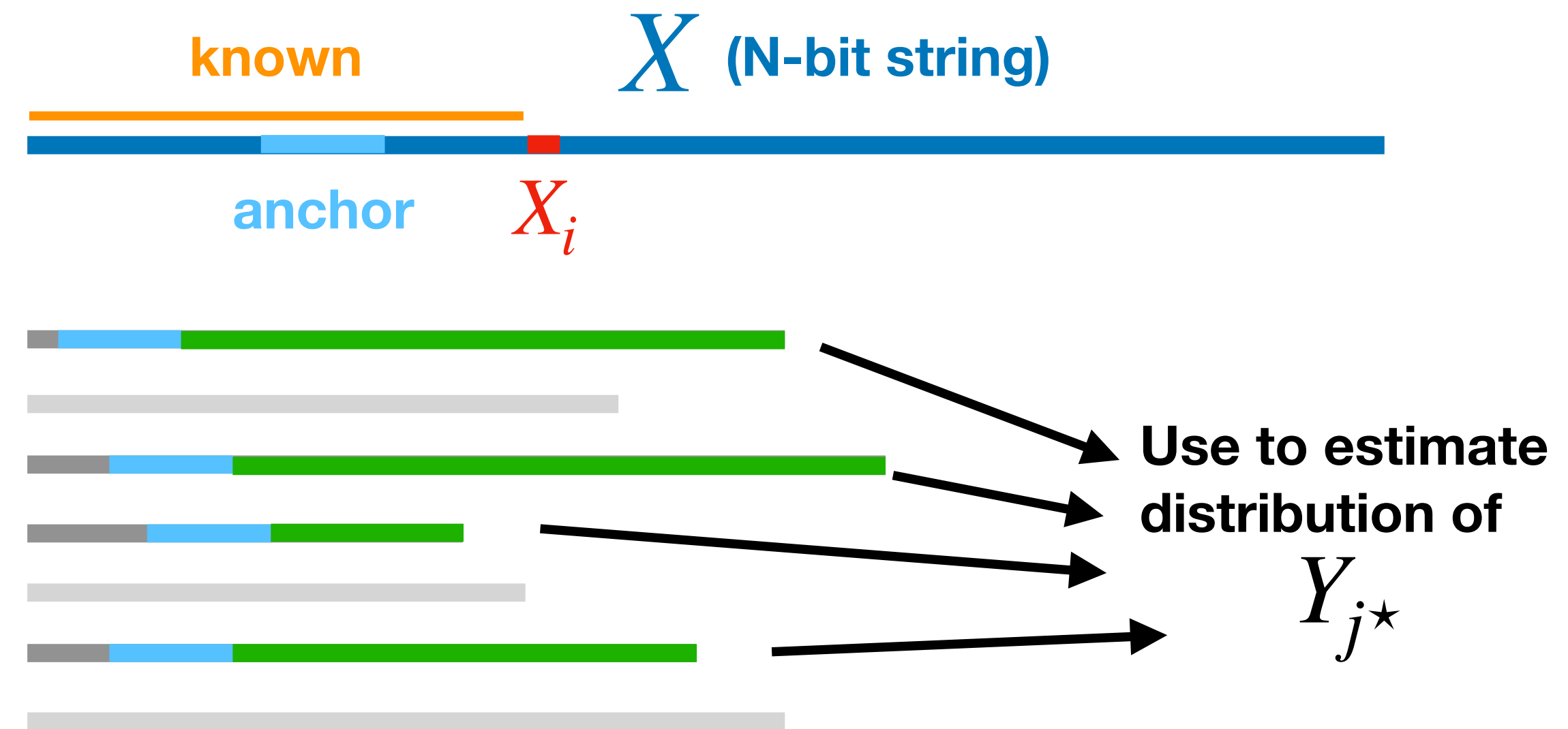
2) Trace alignment: Align by **anchor** close to X_i

“If X is random, whp anchor in trace comes from anchor in X .”

3) Reconstruction: Estimate special bit

$Y =$ distribution of **“trace after anchor”**

“There is special position Y_{j^} that is decently influenced by X_i ”*



$$|\Pr[Y_{j^*} = 1 | X_i = \underline{1}] - \Pr[Y_{j^*} = 1 | X_i = \underline{0}]| \geq \frac{1}{N^C}$$

Recover X_i whp using **poly**(N) traces!

Lower bounds for trace reconstruction

General recipe for lower bounds:

- **Worst-case:** Show it is hard to distinguish between two specific strings with few traces;
- **Average-case:** Worst-case LB + we expect bad string of length $\frac{\log N}{2}$ to show up $\approx \sqrt{N}$ times in random N -bit string.

Lower bounds #traces	Worst-case	Average-case
McGregor, Price, Vorotnikova 2014	N	$\log^2 N$
Holden, Lyons 2018	$N^{1.25}$	$\log^{2.25} N$
Chase 2019	$N^{1.5}$	$\log^{2.5} N$

Lower bounds for trace reconstruction

[McGregorPriceVorotnikova14]

N traces **necessary and sufficient** to distinguish

$$0^{N-1}10^N \quad \text{vs.} \quad 0^N10^{N-1}$$

0000100000
0000010000

[HoldenLyons18, Chase19]

$N^{1.5}$ traces **necessary and sufficient** to distinguish

$$(01)^{N-1}10(01)^N \quad \text{vs.} \quad (01)^N10(01)^{N-1}$$

010110010101
010101100101

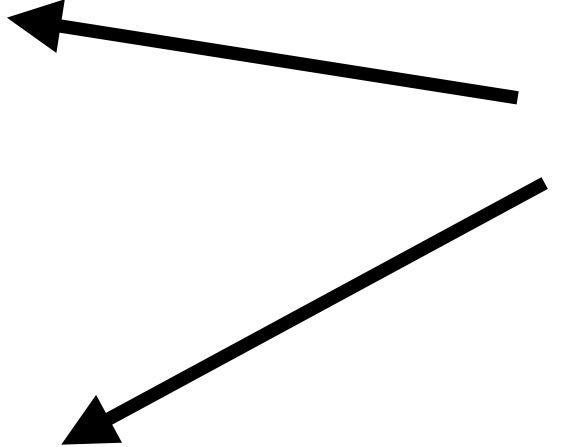
Summing up...

	Worst-case trace reconstruction	Average-case trace reconstruction
Upper bounds #traces	$\exp(O(N^{1/3}))$ [DeO'DonnellServedio17, NazarovPeres17]	$\exp(O(\log^{1/3} N))$ [HoldenPemantlePeres18]
Lower bounds #traces	$\approx N^{3/2}$ [HoldenLyons18, Chase19]	$\approx \log^{5/2} N$ [HoldenLyons18, Chase19]

Summing up...

	Worst-case trace reconstruction	Average-case trace reconstruction
Upper bounds #traces	$\exp(O(N^{1/3}))$ [DeO'DonnellServedio17, NazarovPeres17]	$\exp(O(\log^{1/3} N))$ [HoldenPemantlePeres18]
Lower bounds #traces	$\approx N^{3/2}$ [HoldenLyons18, Chase19]	$\approx \log^{5/2} N$ [HoldenLyons18, Chase19]

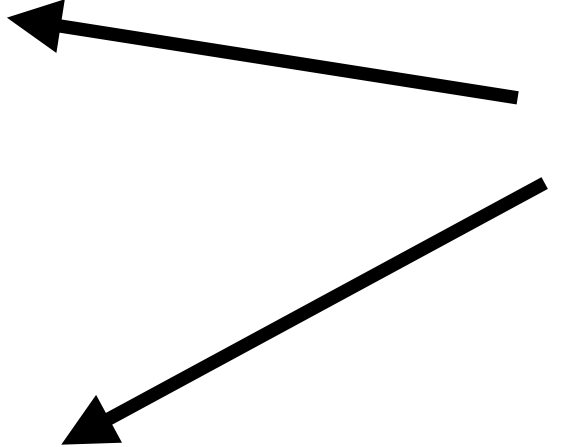
**almost
exponential gap**



Summing up...

	Worst-case trace reconstruction	Average-case trace reconstruction
Upper bounds #traces	$\exp(O(N^{1/3}))$ [DeO'DonnellServedio17, NazarovPeres17]	$\exp(O(\log^{1/3} N))$ [HoldenPemantlePeres18]
Lower bounds #traces	$\approx N^{3/2}$ [HoldenLyons18, Chase19]	$\approx \log^{5/2} N$ [HoldenLyons18, Chase19]

**almost
exponential gap**



***All results over $\{0,1\}^N$,
as it is the hardest setting**

Trace reconstruction and portable DNA-based storage

[YazdiGabrysMilenkovic17, Organick+18]

Write process: Data is encoded into {A,G,T,C} alphabet, then synthesized into DNA strand.

Trace reconstruction and portable DNA-based storage

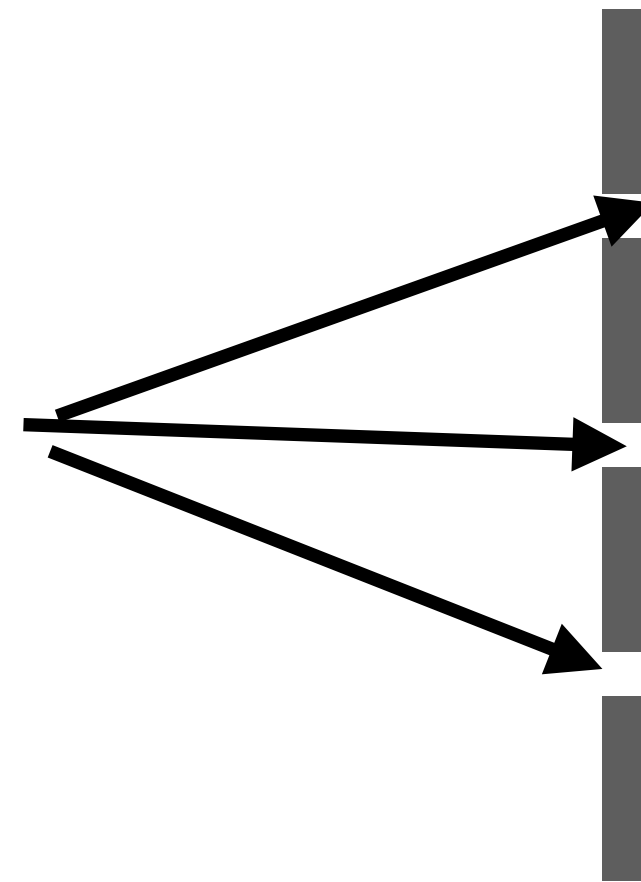
[YazdiGabrysMilenkovic17, Organick+18]

Write process: Data is encoded into {A,G,T,C} alphabet, then synthesized into DNA strand.

**Read from DNA via
nanopore sequencers**

nanopores

ATCGACTCAAGCGTAGAC...



ATCGACTCAAGCGTAGAC...

ATCGACTCAAGCGTAGAC...

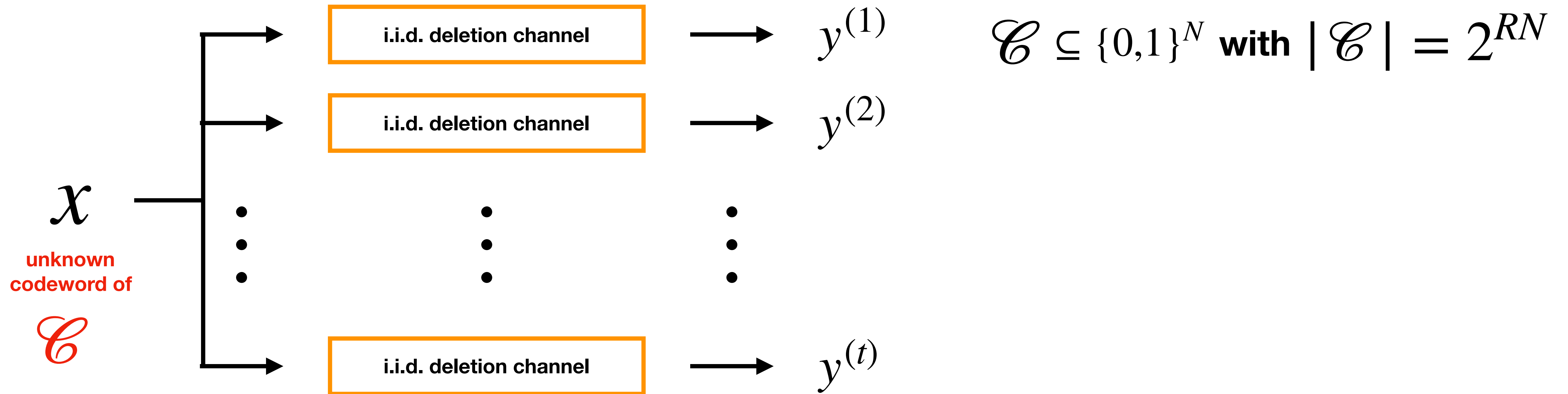
ATCGACTCAAGCGTAGAC...

**noisy
traces**

Trace reconstruction for the coding schemes in [YGM17, OAC+18] is based on heuristics.

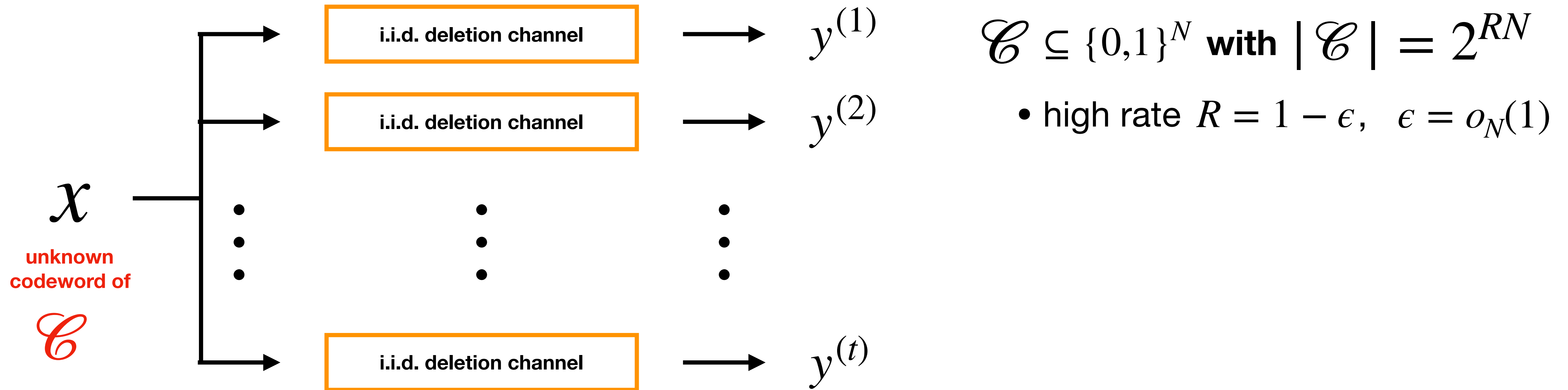
Coded trace reconstruction

[CheraghchiGabrysMilenkovicRibeiro19]



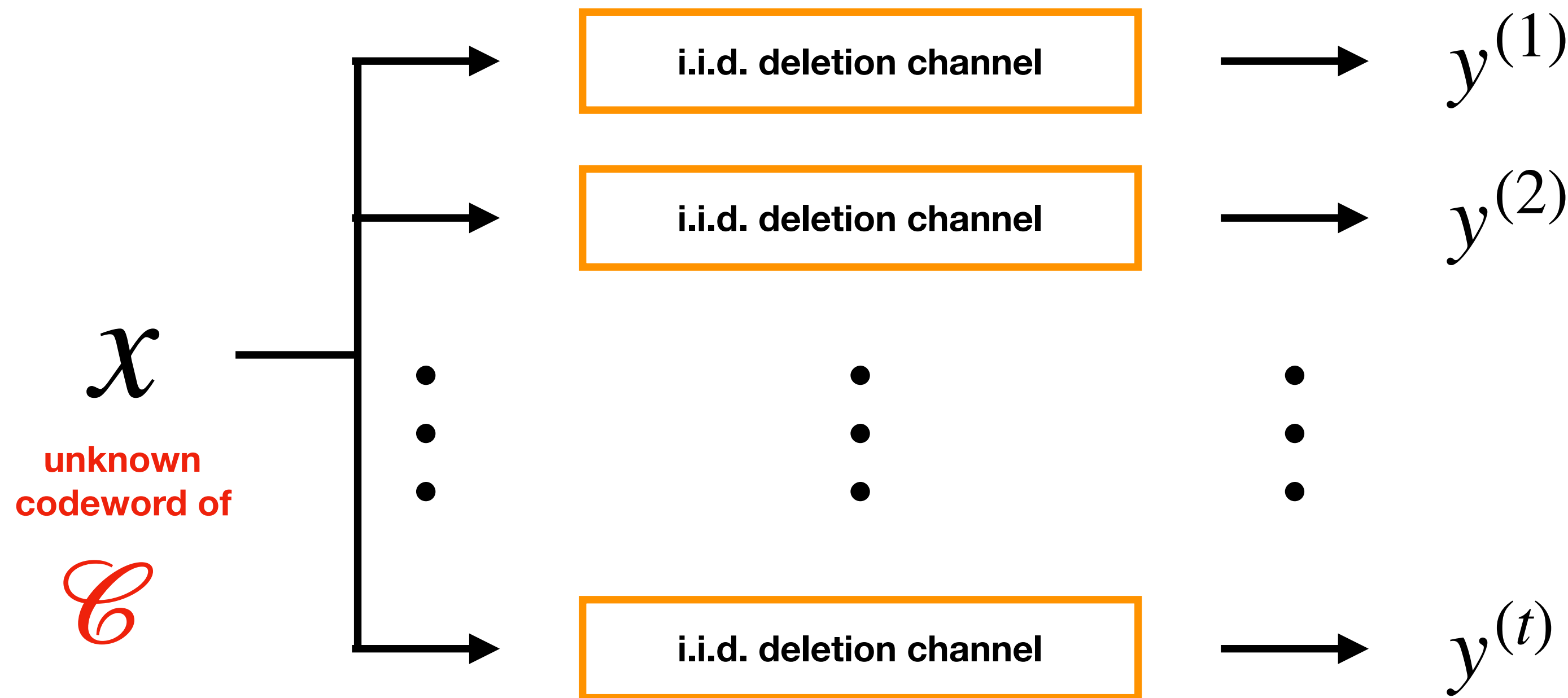
Coded trace reconstruction

[CheraghchiGabrysMilenkovicRibeiro19]



Coded trace reconstruction

[CheraghchiGabrysMilenkovicRibeiro19]

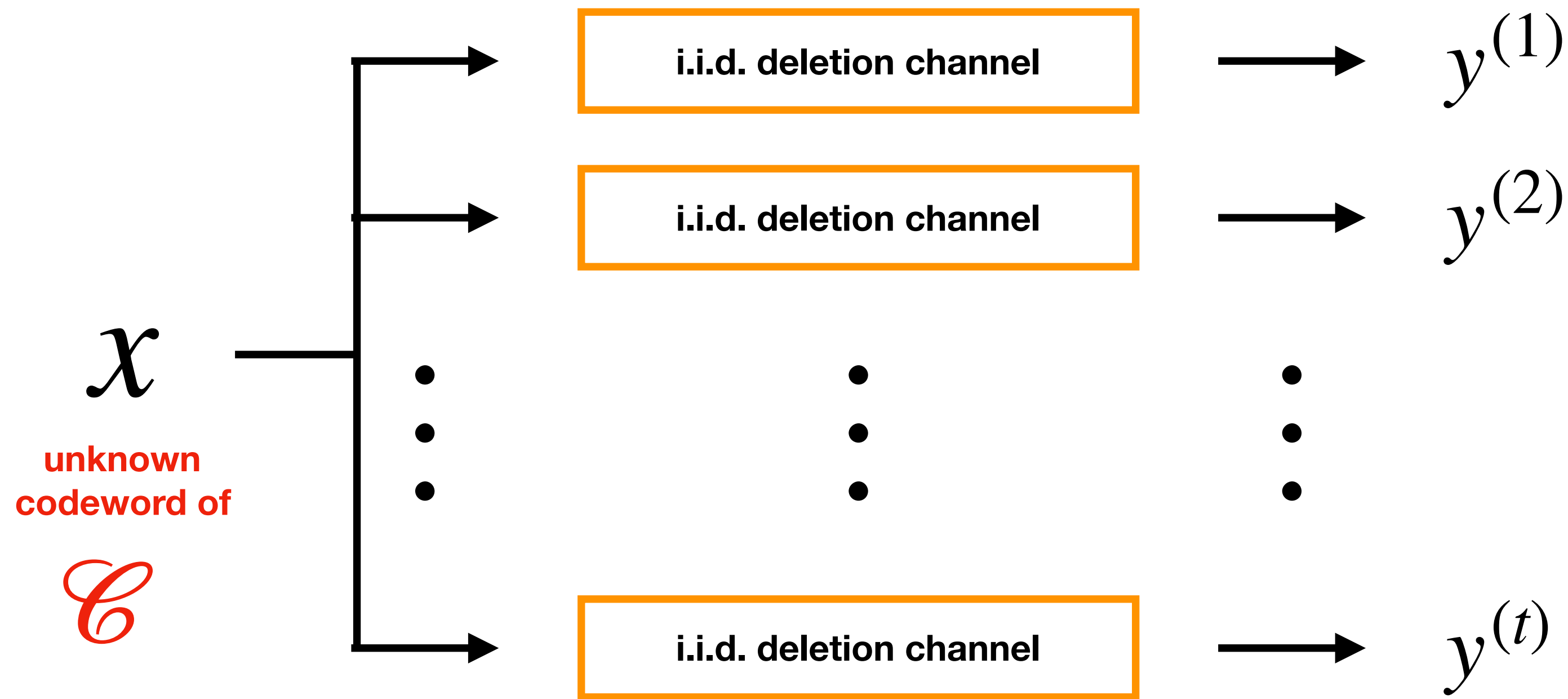


$$\mathcal{C} \subseteq \{0,1\}^N \text{ with } |\mathcal{C}| = 2^{RN}$$

- high rate $R = 1 - \epsilon$, $\epsilon = o_N(1)$
- reconstructed from few traces ($t = t(1/\epsilon)$ slow-growing) with probability $1 - O(1/N)$

Coded trace reconstruction

[CheraghchiGabrysMilenkovicRibeiro19]



$$\mathcal{C} \subseteq \{0,1\}^N \text{ with } |\mathcal{C}| = 2^{RN}$$

- high rate $R = 1 - \epsilon$, $\epsilon = o_N(1)$
- reconstructed from few traces ($t = t(1/\epsilon)$ slow-growing) with probability $1 - O(1/N)$
- **efficiently** encodable and reconstructable (time **poly**(N))

Some related work

Haeupler, Mitzenmacher
2014

First-order capacity for small deletion probability
(constant number of traces only)

Abroshan, Venkataramanan,
Dolecek, Guillén i Fàbregas
2019

Coding for multiple deletion channels with **constant number of deletions only**.
Uses VT codes. Builds up on work about file synchronization.

**Average-case
trace reconstruction**

Implies existence of large codes reconstructable from few traces.
No efficient encoding + we want to use even fewer traces.

PART I:

Markers + worst-case trace reconstruction

Basic construction

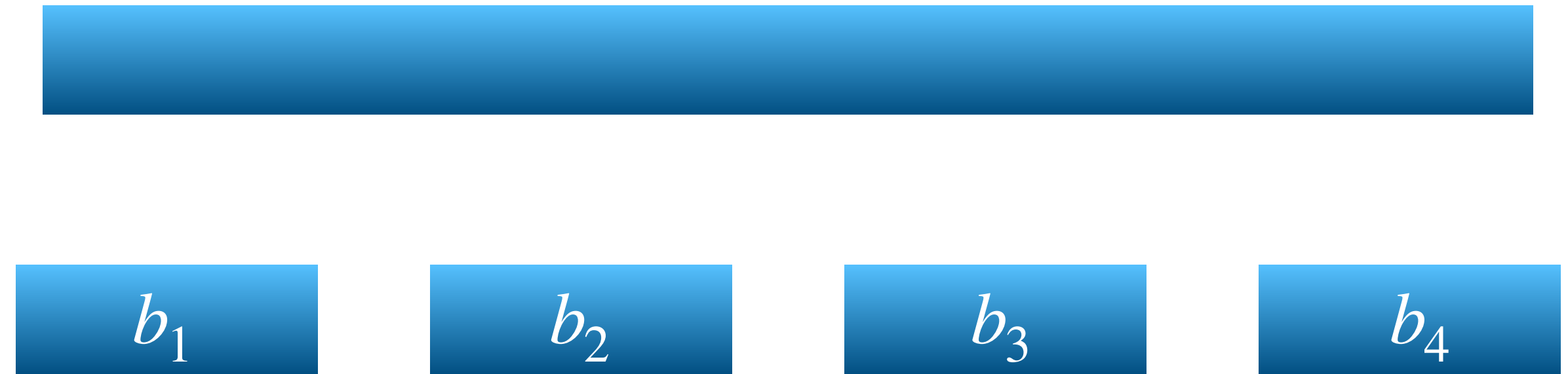
x N -bit string



Basic construction

\mathcal{X} N -bit string

split into blocks with
 $m = O(\log^2 N)$ bits

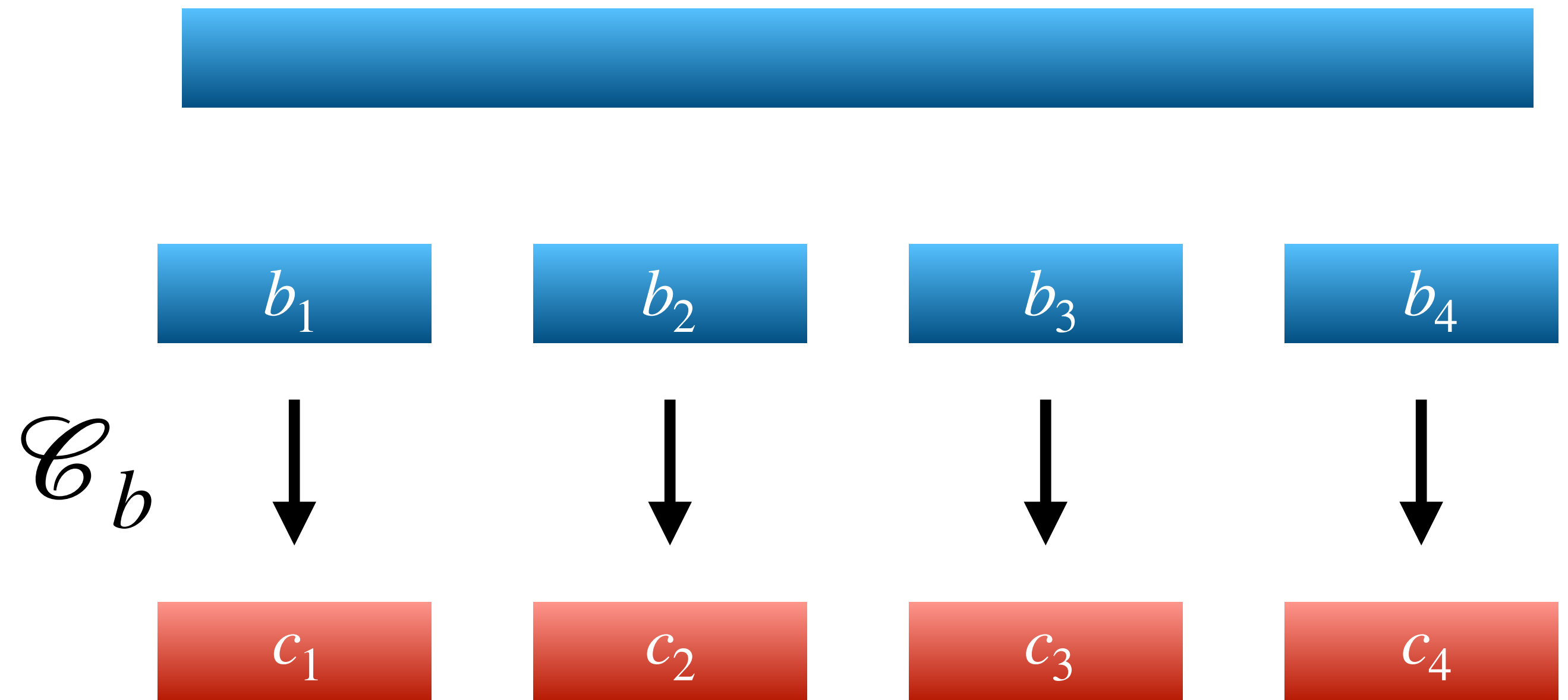


Basic construction

x N -bit string

split into blocks with
 $m = O(\log^2 N)$ bits

encode each block
using appropriate
sub-code \mathcal{C}_b





Basic construction

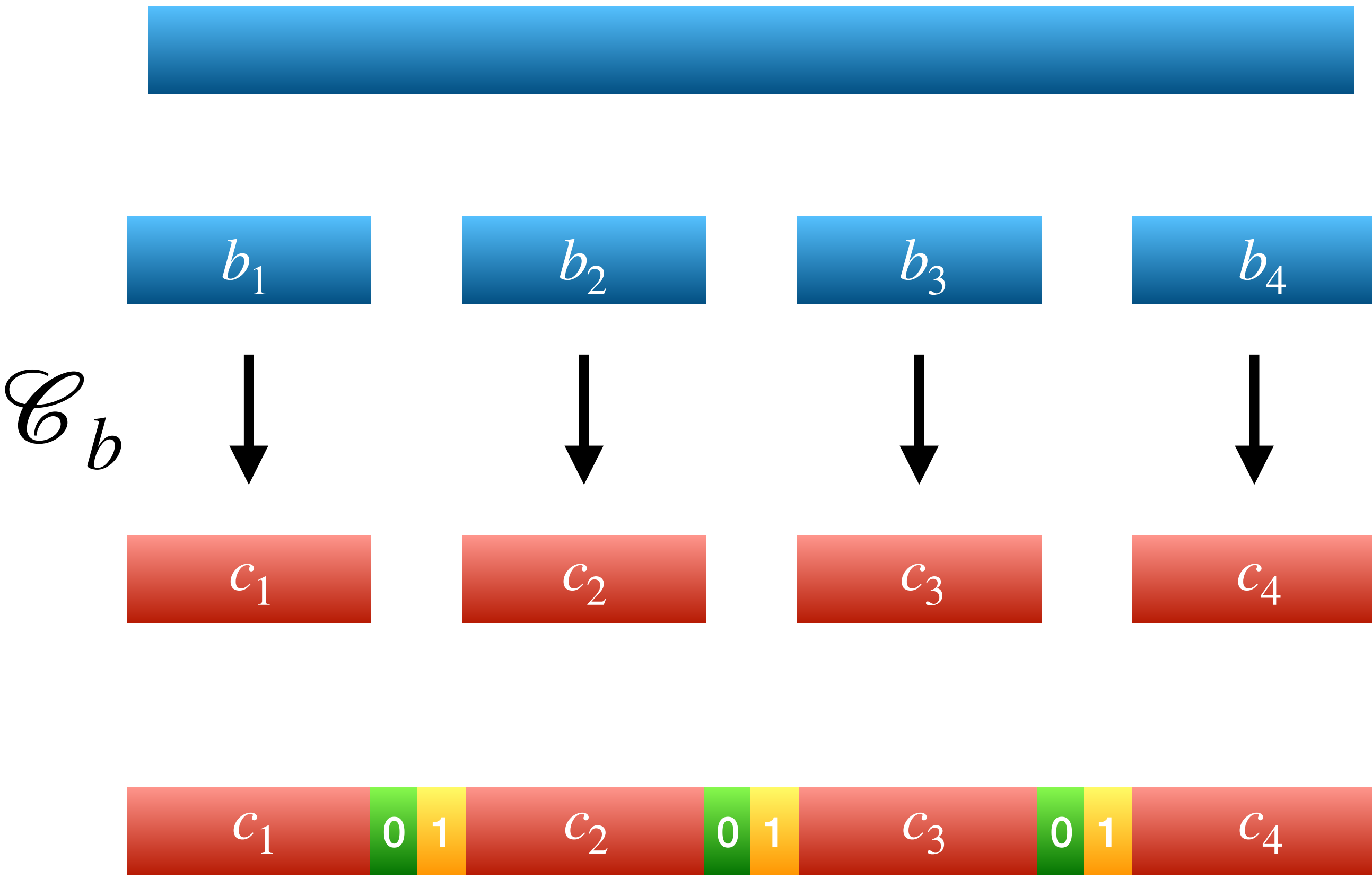
x N -bit string

split into blocks with
 $m = O(\log^2 N)$ bits

encode each block
using appropriate
sub-code \mathcal{C}_b

add markers
between blocks


 = 0000...00
O(log N) bits
 = 1111...11
O(log N) bits



Reconstruction

 = 0000...00
O(log N) bits


 = 1111...11
O(log N) bits


 $\in \mathcal{C}_b$



Reconstruction

 = 0000...00
O(log N) bits

 = 1111...11
O(log N) bits

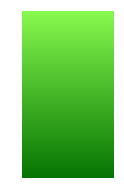
 $\in \mathcal{C}_b$


Trace of codeword


Key observation: part of trace coming from a marker still looks like a marker.



Reconstruction

 = 0000...00
O(log N) bits

 = 1111...11
O(log N) bits

 $\in \mathcal{C}_b$

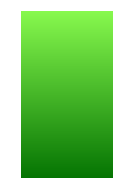
Trace of codeword





Key observation: part of trace coming from a marker still looks like a marker.

If \mathcal{C}_b is chosen appropriately, long runs of 0's in trace **only** come from markers

Reconstruction

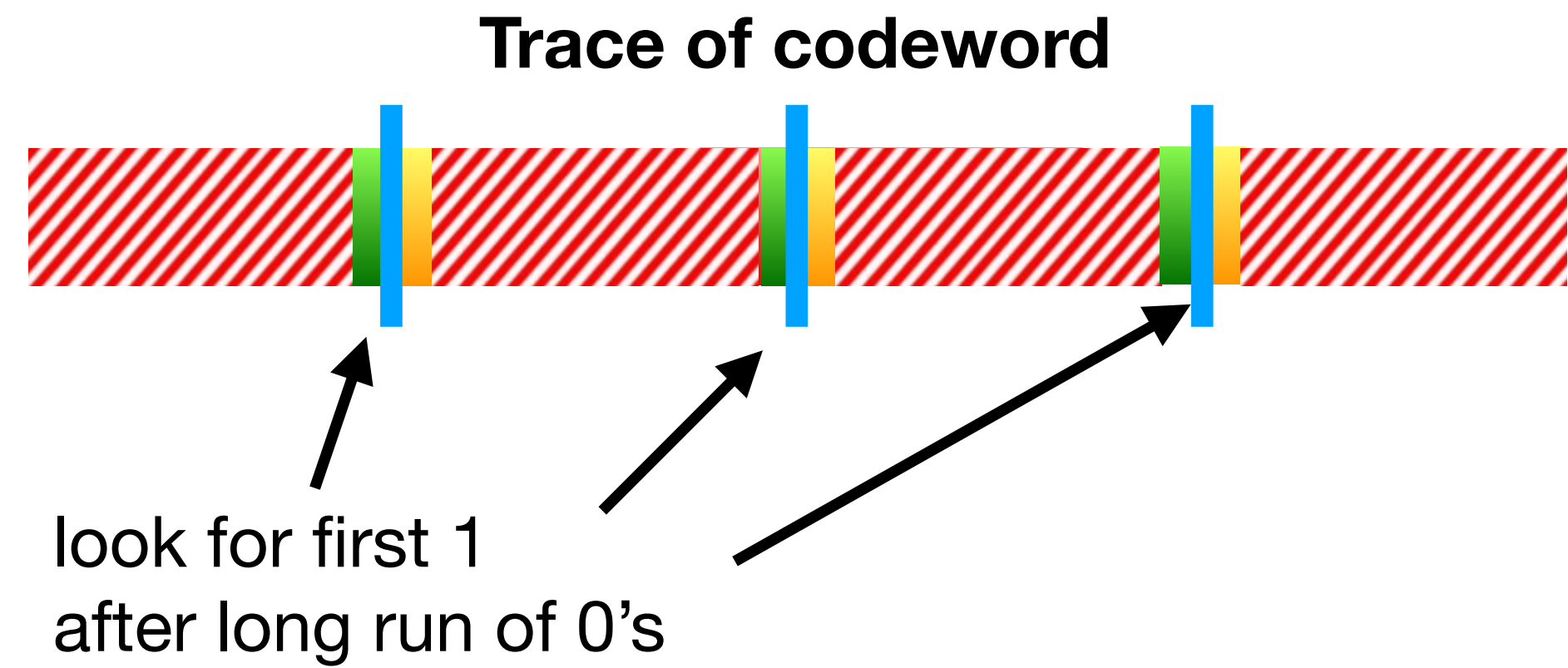
 = 0000...00
O(log N) bits

 = 1111...11
O(log N) bits

 $\in \mathcal{C}_b$


Key observation: part of trace coming from a marker still looks like a marker.


If \mathcal{C}_b is chosen appropriately, long runs of 0's in trace **only** come from markers



Reconstruction

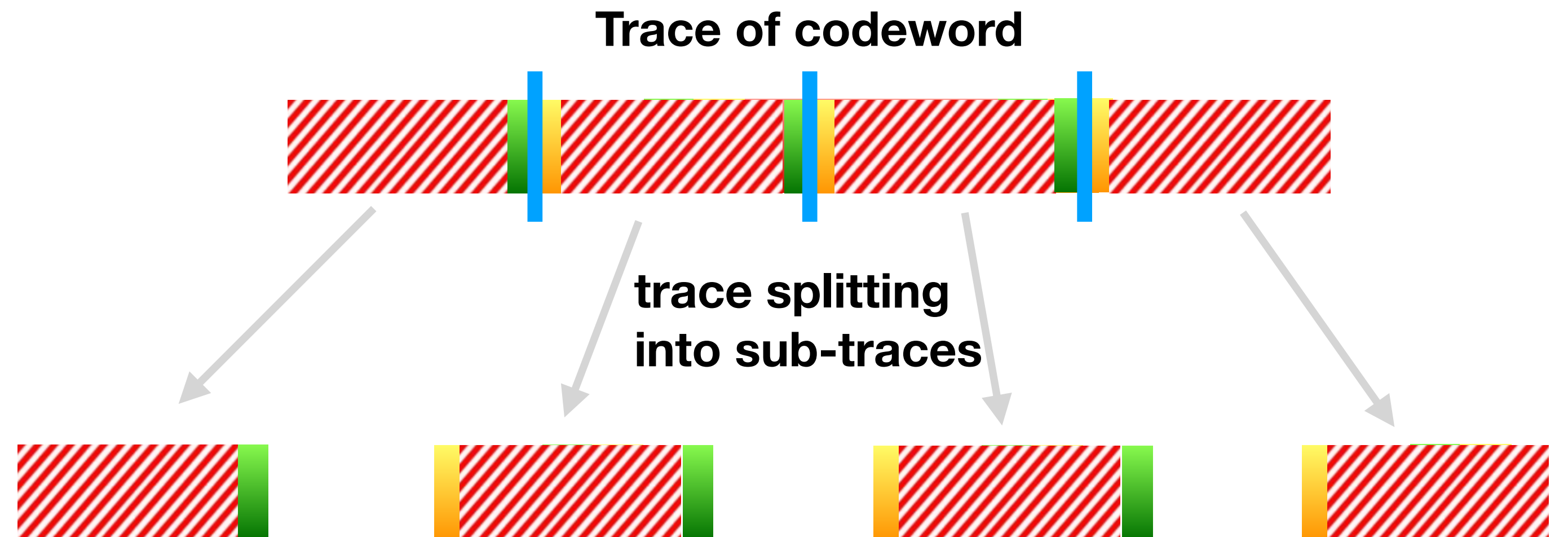
 = 0000...00
O(log N) bits

 = 1111...11
O(log N) bits

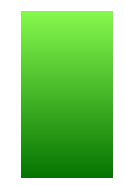
 $\in \mathcal{C}_b$


Key observation: part of trace coming from a marker still looks like a marker.


If \mathcal{C}_b is chosen appropriately, long runs of 0's in trace **only** come from markers



Reconstruction

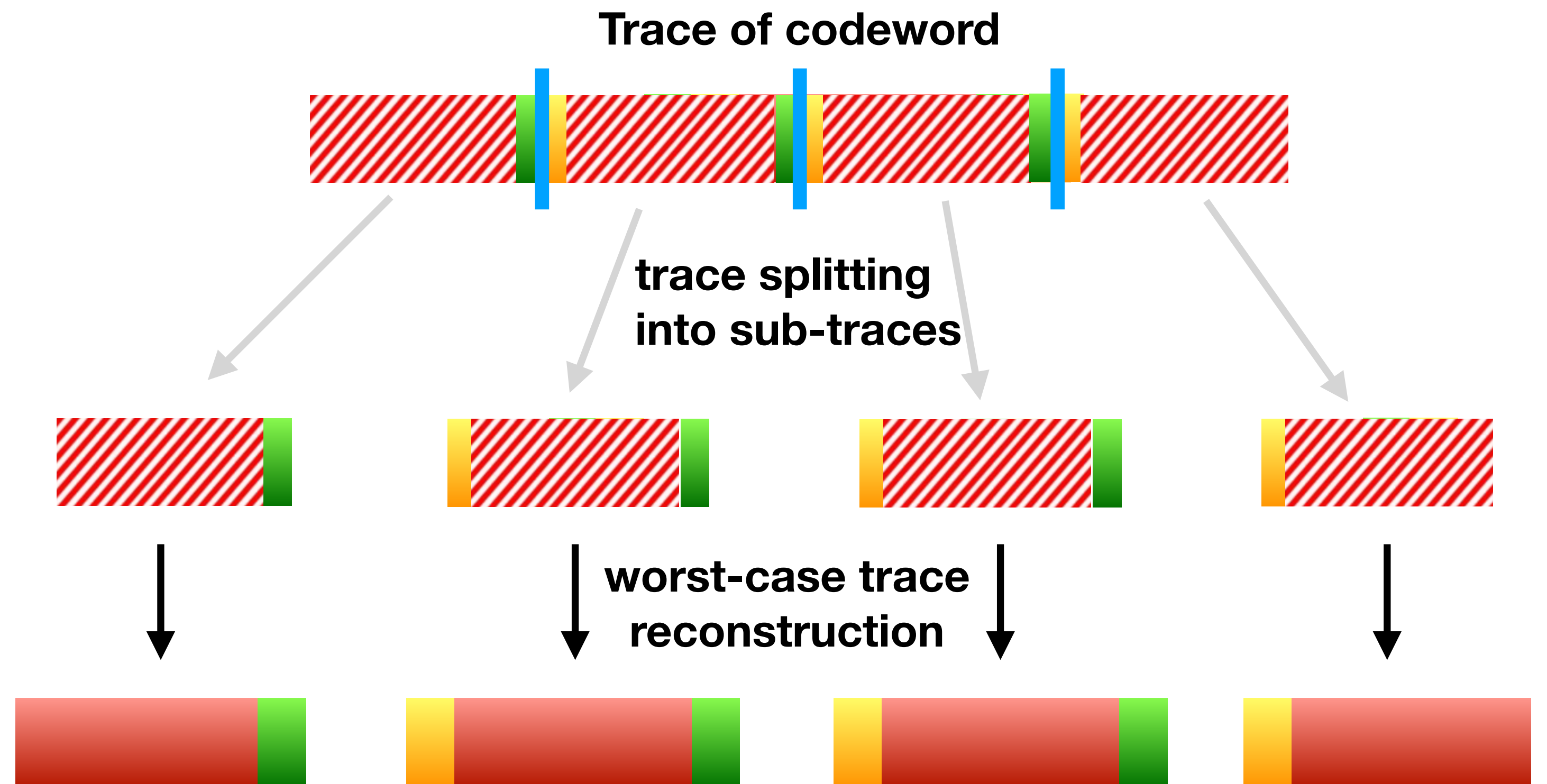
 = 0000...00
O(log N) bits

 = 1111...11
O(log N) bits

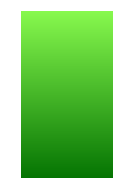
 $\in \mathcal{C}_b$


Key observation: part of trace coming from a marker still looks like a marker.


If \mathcal{C}_b is chosen appropriately, long runs of 0's in trace **only** come from markers



Reconstruction

 = 0000...00
O(log N) bits

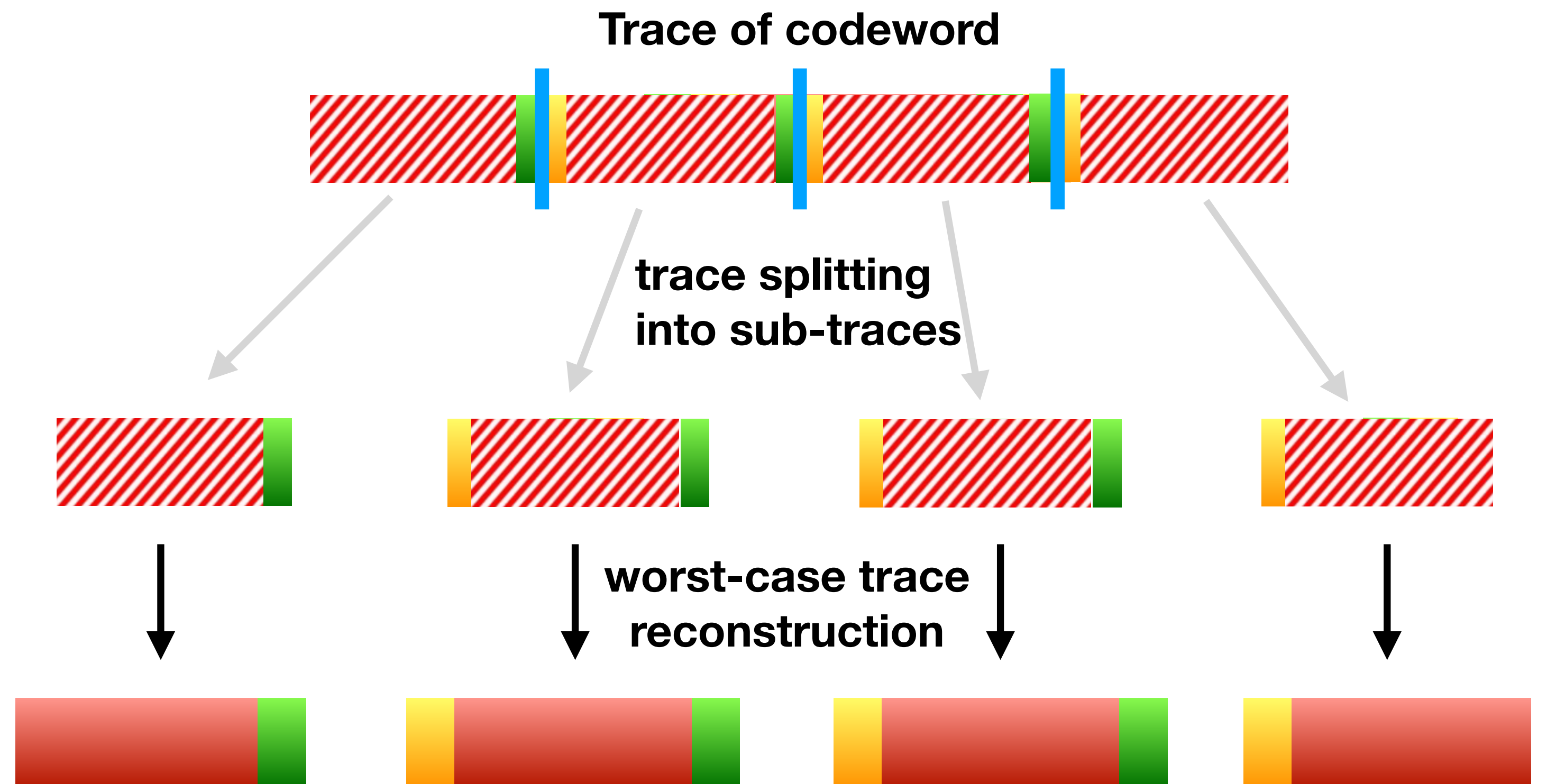
 = 1111...11
O(log N) bits

 $\in \mathcal{C}_b$

Key observation: part of trace coming from a marker still looks like a marker.

If \mathcal{C}_b is chosen appropriately, long runs of 0's in trace **only** come from markers

Efficiently encodable/reconstructable code with rate $1 - O(1/\log N)$ using $\exp(\log^{1/3+\gamma} N)$ traces for any constant deletion prob.



A tool for designing the sub-code: Almost k -wise independent strings

ϵ -almost k -wise independent distribution: $X \in \{0,1\}^m$

$$\forall i_1, \dots, i_k, x_1, \dots, x_k : |\Pr[X_{i_1} = x_1, \dots, X_{i_k} = x_k] - 2^{-k}| \leq \epsilon$$

A tool for designing the sub-code: Almost k -wise independent strings

ϵ -almost k -wise independent distribution: $X \in \{0,1\}^m$

$$\forall i_1, \dots, i_k, x_1, \dots, x_k : |\Pr[X_{i_1} = x_1, \dots, X_{i_k} = x_k] - 2^{-k}| \leq \epsilon$$

[AlonGoldreichHåstadPeralta92]:

“For decent parameters, can efficiently generate such strings from few uniformly random bits”

A tool for designing the sub-code: Almost k -wise independent strings

ϵ -almost k -wise independent distribution: $X \in \{0,1\}^m$

$$\forall i_1, \dots, i_k, x_1, \dots, x_k : |\Pr[X_{i_1} = x_1, \dots, X_{i_k} = x_k] - 2^{-k}| \leq \epsilon$$

[AlonGoldreichHåstadPeralta92]:

“For decent parameters, can efficiently generate such strings from few uniformly random bits”

For every m , $k = O(\log m)$ and $\epsilon = \frac{1}{\text{poly}(m)}$ there is a **poly**(m)-computable function

$$g : \{0,1\}^t \rightarrow \{0,1\}^m$$

with $t = O(\log m)$ such that $g(U_t)$ is ϵ -almost k -wise independent.

Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

$g : \{0,1\}^t \rightarrow \{0,1\}^m$ generator from [AGHP92] with $t = O(\log m)$

$\implies g(U_t)$ satisfies property with high probability

Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

$g : \{0,1\}^t \rightarrow \{0,1\}^m$ generator from [AGHP92] with $t = O(\log m)$

$\implies g(U_t)$ satisfies property with high probability

To encode $x \in \{0,1\}^m$:



x

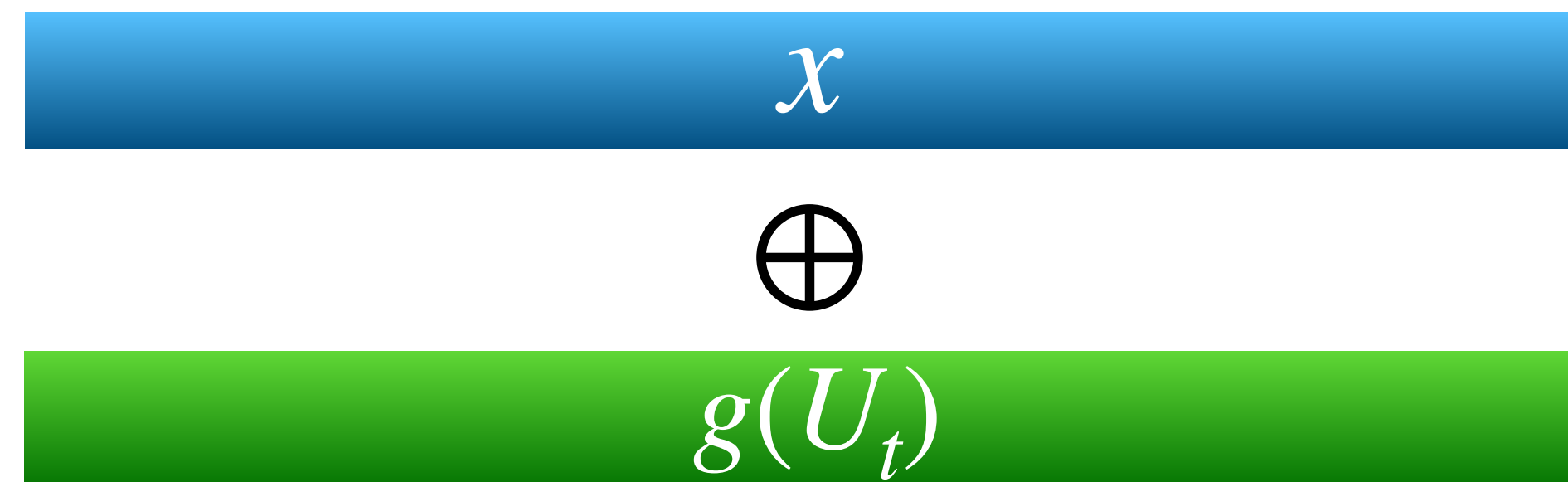
Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

$g : \{0,1\}^t \rightarrow \{0,1\}^m$ generator from [AGHP92] with $t = O(\log m)$

$\implies g(U_t)$ satisfies property with high probability

To encode $x \in \{0,1\}^m$:



Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

$g : \{0,1\}^t \rightarrow \{0,1\}^m$ generator from [AGHP92] with $t = O(\log m)$

$\implies g(U_t)$ satisfies property with high probability

To encode $x \in \{0,1\}^m$:

x

satisfies property whp

\oplus

$g(U_t)$

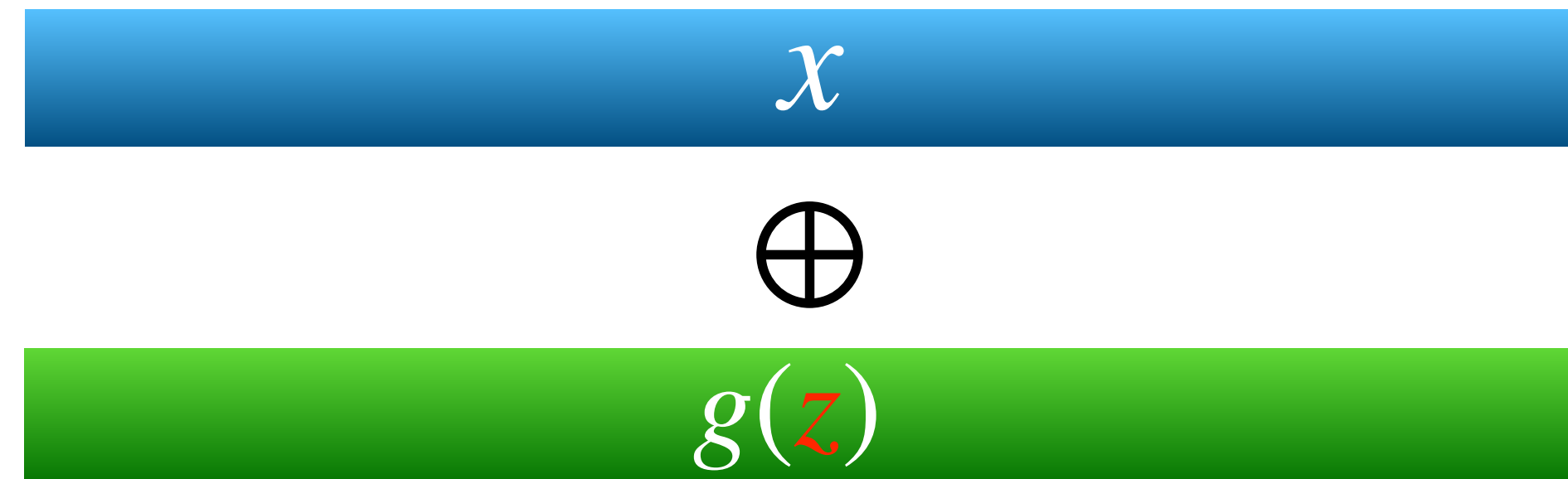
Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

$g : \{0,1\}^t \rightarrow \{0,1\}^m$ generator from [AGHP92] with $t = O(\log m)$

$\implies g(U_t)$ satisfies property with high probability

To encode $x \in \{0,1\}^m$:



satisfies property whp

So there is **good** fixing $U_t = z$ that enforces property

Designing the sub-code

Desired property: No long runs of 0's in trace of $c \in \mathcal{C}_b \rightarrow$ allows for trace splitting

$g : \{0,1\}^t \rightarrow \{0,1\}^m$ generator from [AGHP92] with $t = O(\log m)$

$\implies g(U_t)$ satisfies property with high probability

To encode $x \in \{0,1\}^m$:

$\text{Enc}(x) = \underbrace{z}_{\log m} \oplus x \oplus g(z)$ (z chosen so that property is satisfied)

PART II:
Markers +
modified average-case trace reconstruction

Subsequence-unique strings

[HolensteinMitzenmacherPanigrahyWieder08]

A string \mathcal{X} is w -**subsequence-unique** if no substring of length w can be obtained by deleting bits of another substring of length $1.1w$ (except for trivial containment).

Subsequence-unique strings

[HolensteinMitzenmacherPanigrahyWieder08]

A string \mathcal{X} is w -**subsequence-unique** if no substring of length w can be obtained by deleting bits of another substring of length $1.1w$ (except for trivial containment).

Deletion probability
small enough constant



Efficient trace reconstruction algorithm for all
length m w -subsequence-unique strings
with $w = O(\log m)$ using **poly**(m) traces.

Subsequence-unique strings

[HolensteinMitzenmacherPanigrahyWieder08]

A string \mathcal{X} is w -**subsequence-unique** if no substring of length w can be obtained by deleting bits of another substring of length $1.1w$ (except for trivial containment).

Deletion probability
small enough constant



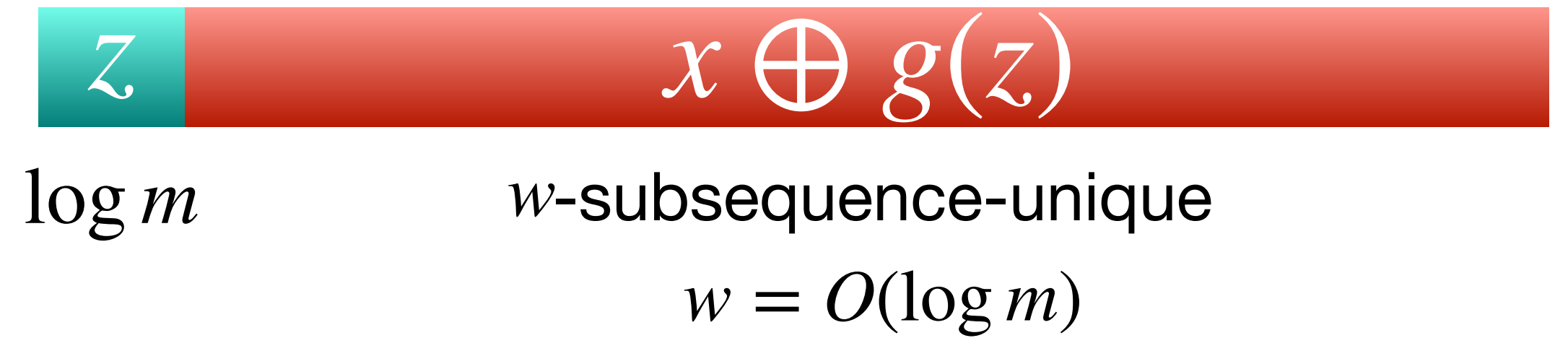
Efficient trace reconstruction algorithm for all
length m w -subsequence-unique strings
with $w = O(\log m)$ using **poly**(m) traces.

Key observation: an ϵ -almost k -wise independent string is w -subsequence-unique with high probability for decent parameters k and ϵ .

Better codes for small deletion probability

Use masking to encode
blocks into “almost”
subsequence-unique strings

x
 $\in \{0,1\}^m$ \longrightarrow encoded as

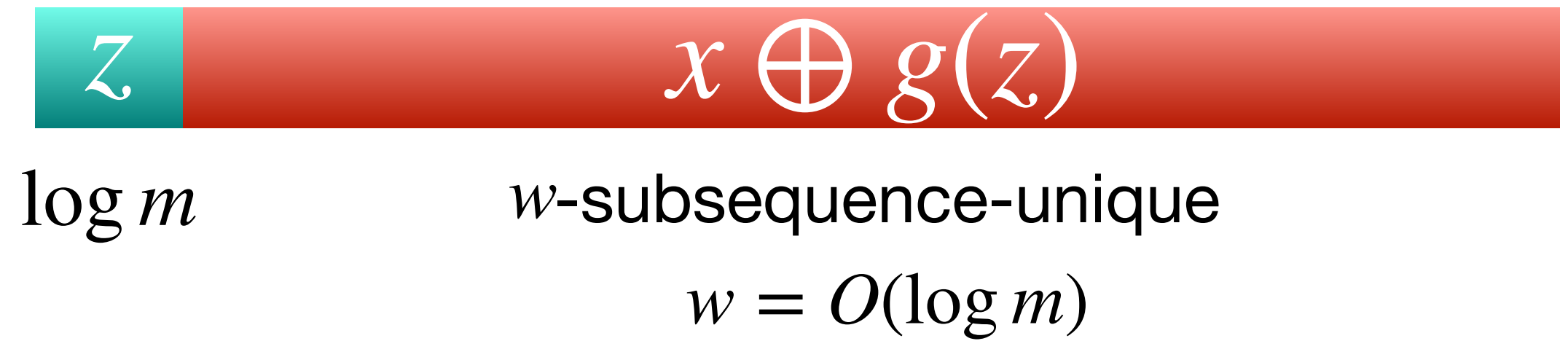


Better codes for small deletion probability

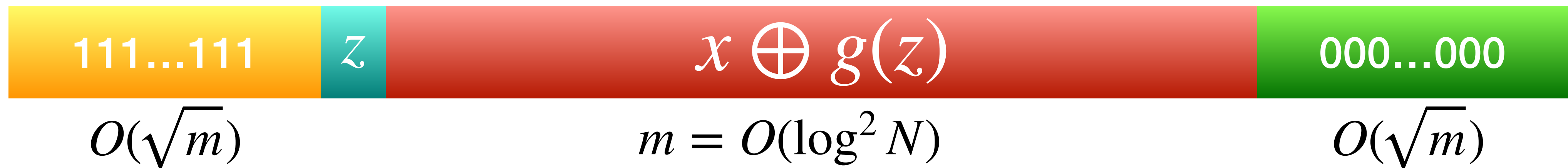
Use masking to encode blocks into “almost” subsequence-unique strings

$$x \in \{0,1\}^m$$

→
encoded as



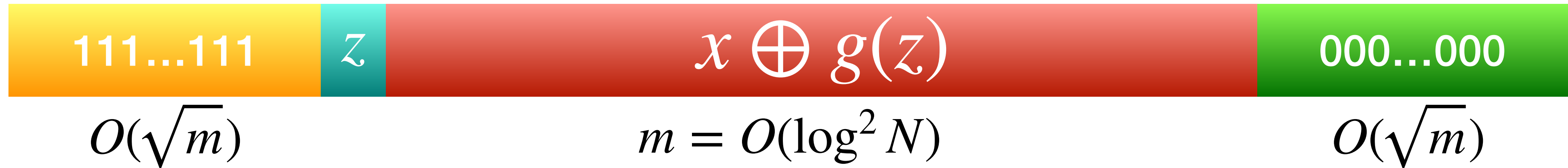
$$O(\log m)$$



HMPW trace reconstruction algorithm doesn't work for strings like this...

Better codes for small deletion probability

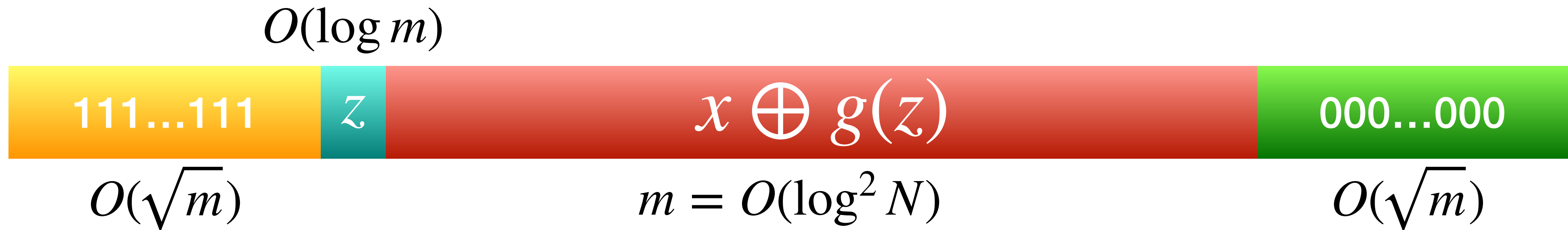
$O(\log m)$



First barrier: HMPW algorithm needs to know first $\log m$ bits of $x \oplus g(z)$ “for free”

Original bootstrapping procedure requires $2^{\sqrt{m}} = \mathbf{poly}(N)$ traces in this case...

Better codes for small deletion probability



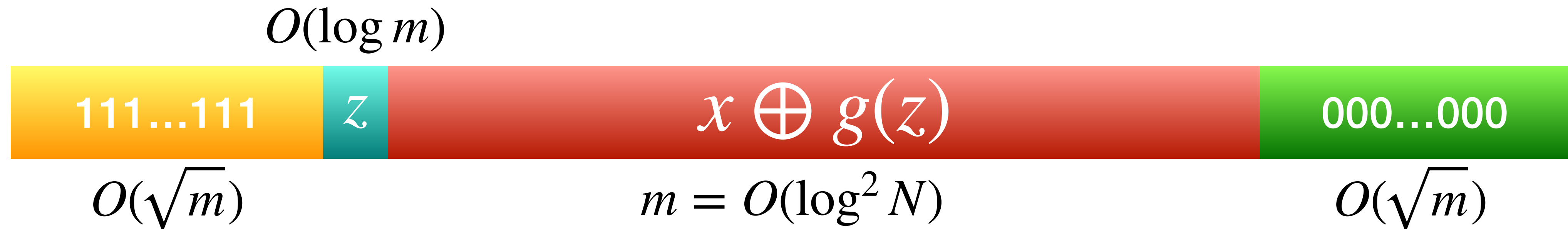
First barrier: HMPW algorithm needs to know first $\log m$ bits of $x \oplus g(z)$ “for free”

Original bootstrapping procedure requires $2^{\sqrt{m}} = \mathbf{poly}(N)$ traces in this case...

Exploit properties of masking:

$$x' = 00\dots00 \parallel x \xrightarrow{\text{encoded as}} z \parallel (x' \oplus g(z))$$

Better codes for small deletion probability



First barrier: HMPW algorithm needs to know first $\log m$ bits of $x \oplus g(z)$ “for free”

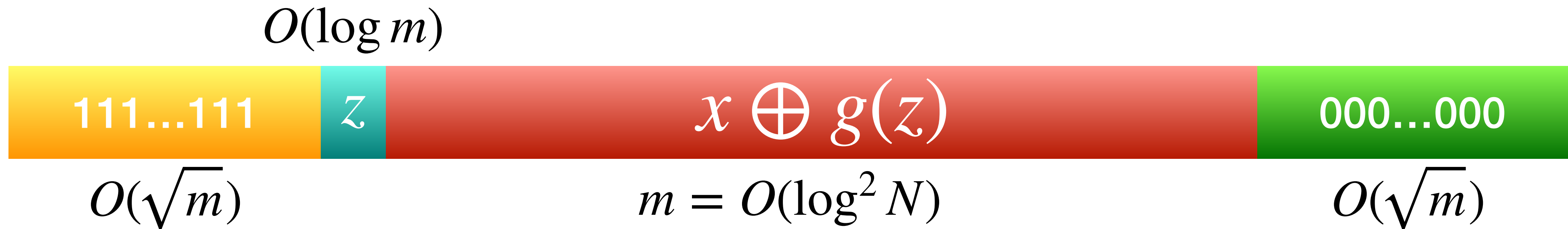
Original bootstrapping procedure requires $2^{\sqrt{m}} = \mathbf{poly}(N)$ traces in this case...

Exploit properties of masking:

$$x' = 00\dots00 \parallel x \xrightarrow{\text{encoded as}} z \parallel (x' \oplus g(z))$$

z' = systematic encoding of $0 \parallel z$ robust against constant fraction of deletions+insertions

Better codes for small deletion probability



First barrier: HMPW algorithm needs to know first $\log m$ bits of $x \oplus g(z)$ “for free”

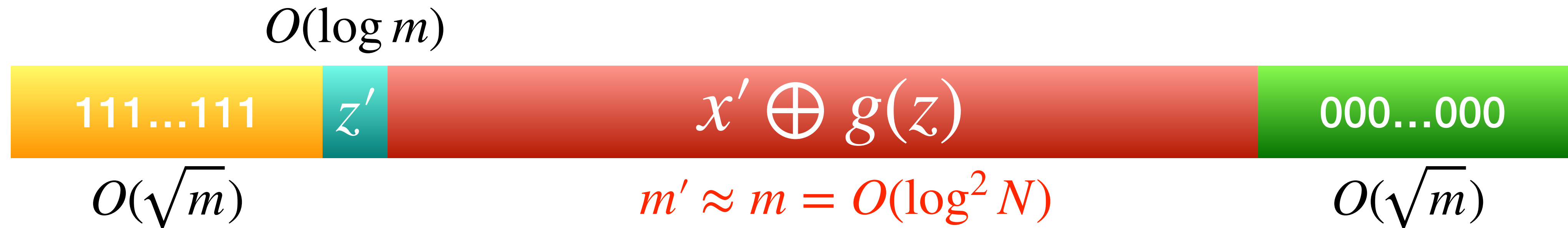
Original bootstrapping procedure requires $2^{\sqrt{m}} = \mathbf{poly}(N)$ traces in this case...

Exploit properties of masking:

$$x' = 00\dots00 \parallel x \xrightarrow{\text{encoded as}} z' \parallel (x' \oplus g(z))$$

$z' =$ systematic encoding of $0 \parallel z$ robust against constant fraction of deletions+insertions

Better codes for small deletion probability



First barrier: HMPW algorithm needs to know first $\log m$ bits of $x \oplus g(z)$ “for free”

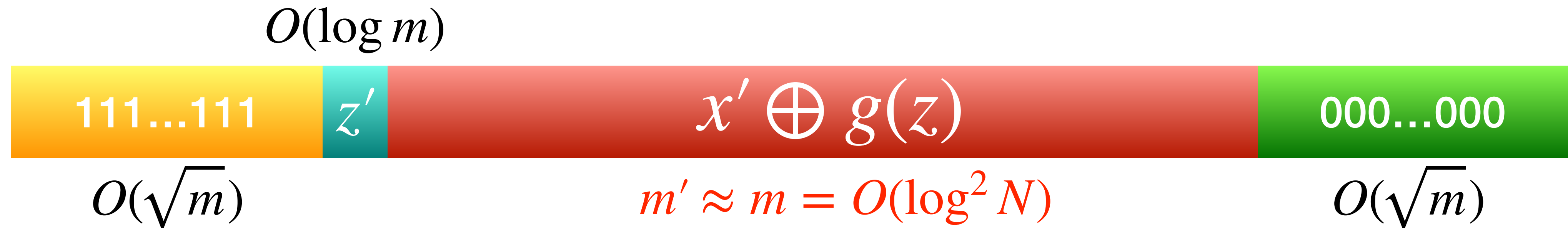
Original bootstrapping procedure requires $2^{\sqrt{m}} = \mathbf{poly}(N)$ traces in this case...

Exploit properties of masking:

$$x' = 00\dots00 \parallel x \xrightarrow{\text{encoded as}} z' \parallel (x' \oplus g(z))$$

$z' =$ systematic encoding of $0 \parallel z$ robust against constant fraction of deletions+insertions

Better codes for small deletion probability



First barrier: HMPW algorithm needs to know first $\log m$ bits of $x \oplus g(z)$ “for free”

Original bootstrapping procedure requires $2^{\sqrt{m}} = \mathbf{poly}(N)$ traces in this case...

Exploit properties of masking:

$$x' = 00\dots00 \parallel x \xrightarrow{\text{encoded as}} z' \parallel (x' \oplus g(z))$$

$z' =$ systematic encoding of $0 \parallel z$ robust against constant fraction of deletions+insertions

Bootstrapping only requires $O(m)$ traces now!

Better codes for small deletion probability

A few more things:

- Actually need a stronger property than subsequence-uniqueness because of markers;
- Need to make sure trace splitting works.

Almost k -wise independent string satisfies both of them with high probability.

There is an absolute constant $d^* \in (0,1)$ such that for all $d \leq d^*$ there exists an efficiently encodable/reconstructable code with rate $1 - O(1/\log N)$ using **polylog**(N) traces against i.i.d. deletions with probability d .

Coded trace reconstruction over large alphabets

[BrakensiekLiSpang19]

There exist efficiently encodable/reconstructable codes with rate $1 - \epsilon$ over an alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces.

Moreover, **this is tight.**

Coded trace reconstruction over large alphabets

[BrakensiekLiSpang19]

There exist efficiently encodable/reconstructable codes with rate $1 - \epsilon$ over an alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces.

Moreover, **this is tight.**

Lower bound:

Capacity of T -use deletion channel \leq Capacity of T -use erasure channel $= 1 - d^T$

Coded trace reconstruction over large alphabets

[BrakensiekLiSpang19]

There exist efficiently encodable/reconstructable codes with rate $1 - \epsilon$ over an alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces.

Moreover, **this is tight.**

Lower bound:

Capacity of T -use deletion channel \leq Capacity of T -use erasure channel $= 1 - d^T$

$$\implies 1 - \epsilon \leq 1 - d^T \implies T \geq \log_{1/d}(1/\epsilon)$$

A tool for the upper bound: Synchronization strings

[HaeuplerShahrasbi17]

A string S is a τ -**synchronization string** if for all $i < j < k$

$$ED(S_{[i,j]}, S_{[j,k]}) > (1 - \tau)(k - i)$$

A tool for the upper bound: Synchronization strings

[HaeuplerShahrasbi17]

A string S is a τ -**synchronization string** if for all $i < j < k$

$$ED(S_{[i,j]}, S_{[j,k]}) > (1 - \tau)(k - i)$$

 “close” to maximal edit distance everywhere

A tool for the upper bound: Synchronization strings

[HaeuplerShahrasbi17]

A string S is a τ -**synchronization string** if for all $i < j < k$

$$ED(S_{[i,j]}, S_{[j,k]}) > (1 - \tau)(k - i)$$

“close” to maximal edit distance everywhere

For every τ , can efficiently construct a τ -synch-string of length N over an alphabet of size **poly**($1/\tau$).

Why are synchronization strings useful?

$$c = (c_1, c_2, \dots, c_N) \in \Sigma^N \xrightarrow{\text{naive indexing}} ((c_1, 1), (c_2, 2), \dots, (c_N, N)) \in (\Sigma \times [N])^N$$

Why are synchronization strings useful?

$$c = (c_1, c_2, \dots, c_N) \in \Sigma^N \xrightarrow{\text{naive indexing}} ((c_1, 1), (c_2, 2), \dots, (c_N, N)) \in (\Sigma \times [N])^N$$

$(c_1, 1)$ $(c_2, 2)$ $(c_3, 3)$ $(c_4, 4)$



$(c_1, 1)$ $(c_2, 2)$ $(c_3, 3)$ $(c_4, 4)$



deletions

$(c_1, 1)$ $(c_3, 3)$

Why are synchronization strings useful?

$$c = (c_1, c_2, \dots, c_N) \in \Sigma^N \xrightarrow{\text{naive indexing}} ((c_1, 1), (c_2, 2), \dots, (c_N, N)) \in (\Sigma \times [N])^N$$

$(c_1, 1)$ $(c_2, 2)$ $(c_3, 3)$ $(c_4, 4)$

c_1 c_2 c_3 c_4



$(c_1, 1)$ $(c_2, 2)$ $(c_3, 3)$ $(c_4, 4)$

equivalent
to

c_1 c_2 c_3 c_4



deletions



erasures

$(c_1, 1)$ $(c_3, 3)$

c_1 $?$ c_3 $?$

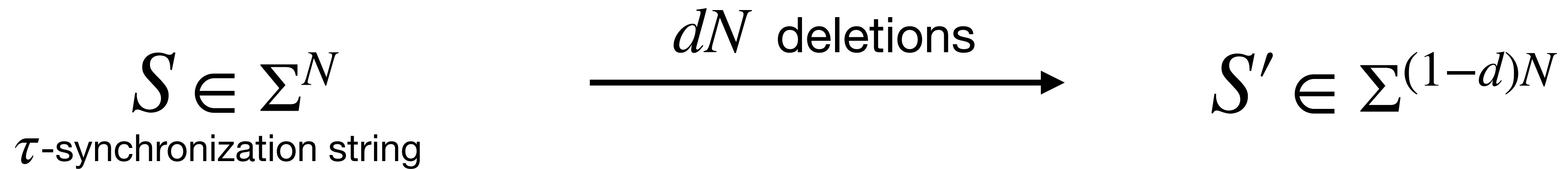
Why are synchronization strings useful?

[HaeuplerShahrasbi17]: “Can efficiently transform deletions into (worst-case) erasures by indexing with synchronization string, with **constant** alphabet size blowup”

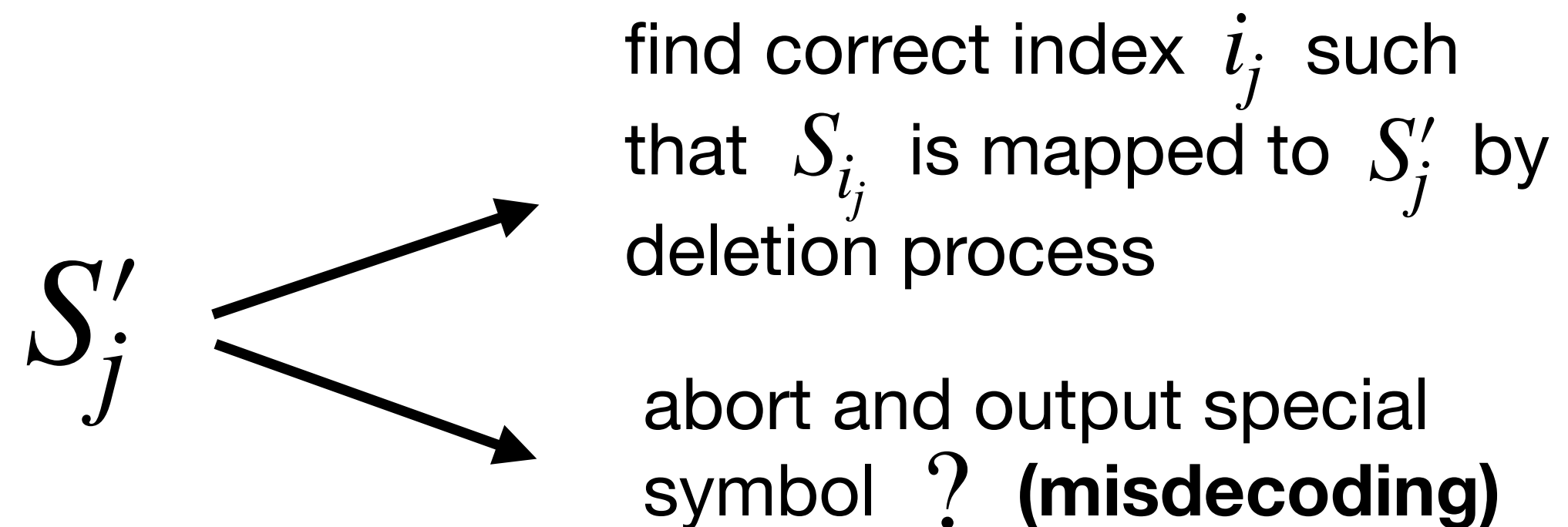
$$\begin{array}{ccc} S \in \Sigma^N & \xrightarrow{dN \text{ deletions}} & S' \in \Sigma^{(1-d)N} \\ \tau\text{-synchronization string} & & \end{array}$$

Why are synchronization strings useful?

[HaeuplerShahrasbi17]: “Can efficiently transform deletions into (worst-case) erasures by indexing with synchronization string, with **constant** alphabet size blowup”



“**Error-free**” indexing algorithm:

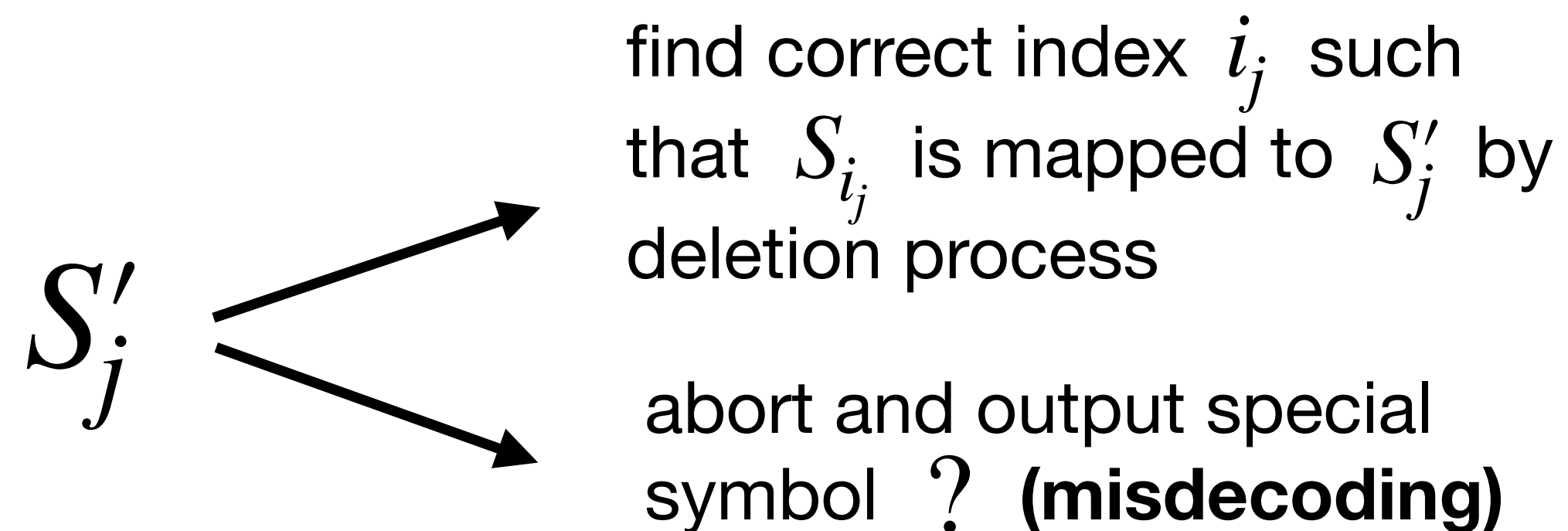


Why are synchronization strings useful?

[HaeuplerShahrasbi17]: “Can efficiently transform deletions into (worst-case) erasures by indexing with synchronization string, with **constant** alphabet size blowup”

$$\begin{array}{ccc} S \in \Sigma^N & \xrightarrow{dN \text{ deletions}} & S' \in \Sigma^{(1-d)N} \\ \tau\text{-synchronization string} & & \end{array}$$

“**Error-free**” indexing algorithm:



[HaeuplerShahrasbi17]:

There is efficient error-free indexing algorithm with at most

$$\frac{\tau dN}{1 - \tau}$$

misdecodings.

Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

S $(\tau = \epsilon^C)$ -synchronization string over alphabet of size $O_\epsilon(1)$

\mathcal{C} rate $\approx 1 - \epsilon$ code robust against $\approx \epsilon^3$ fraction of erasures

Construction: *Index codewords of \mathcal{C} with S*

Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

S $(\tau = \epsilon^C)$ -synchronization string over alphabet of size $O_\epsilon(1)$

\mathcal{C} rate $\approx 1 - \epsilon$ code robust against $\approx \epsilon^3$ fraction of erasures

Construction: *Index codewords of \mathcal{C} with S*

$(c_1, \dots, c_N) \in \mathcal{C} \longrightarrow \boxed{(c_1, S_1)} \quad \boxed{(c_2, S_2)} \quad \dots \quad \boxed{(c_N, S_N)}$

Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

Reconstruction:

Trace of codeword

(x_{11}, x_{12})

(x_{21}, x_{22})

(x_{31}, x_{32})

(x_{41}, x_{42})

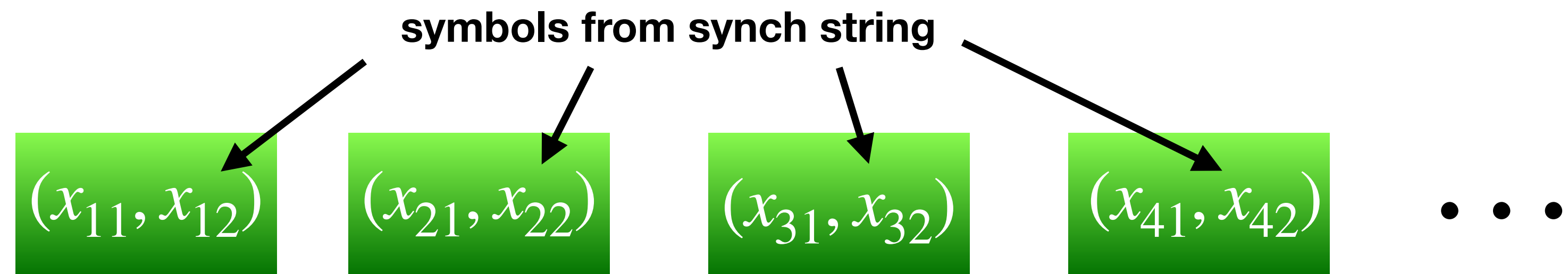
• • •

Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

Reconstruction:

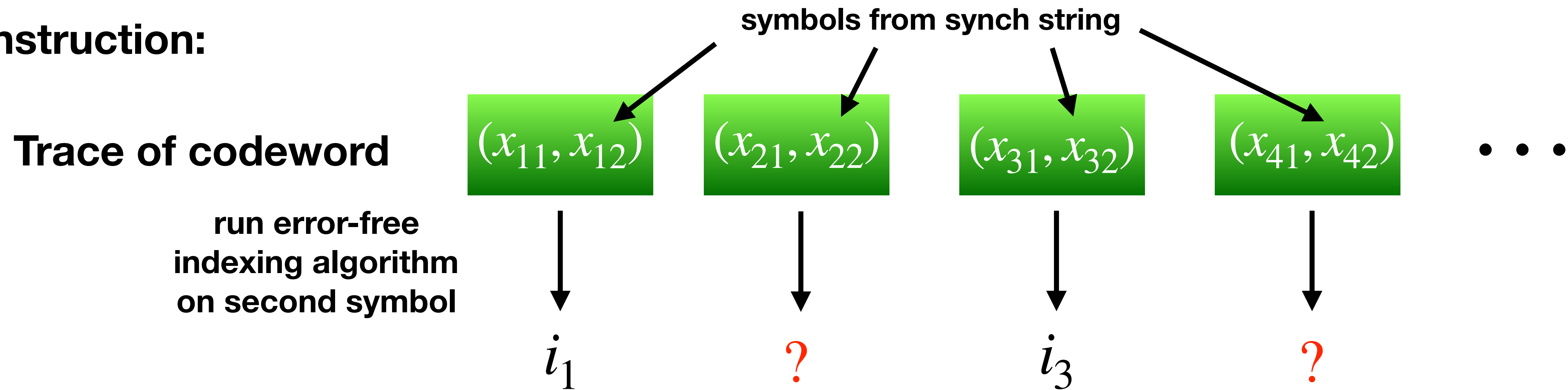
Trace of codeword



Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

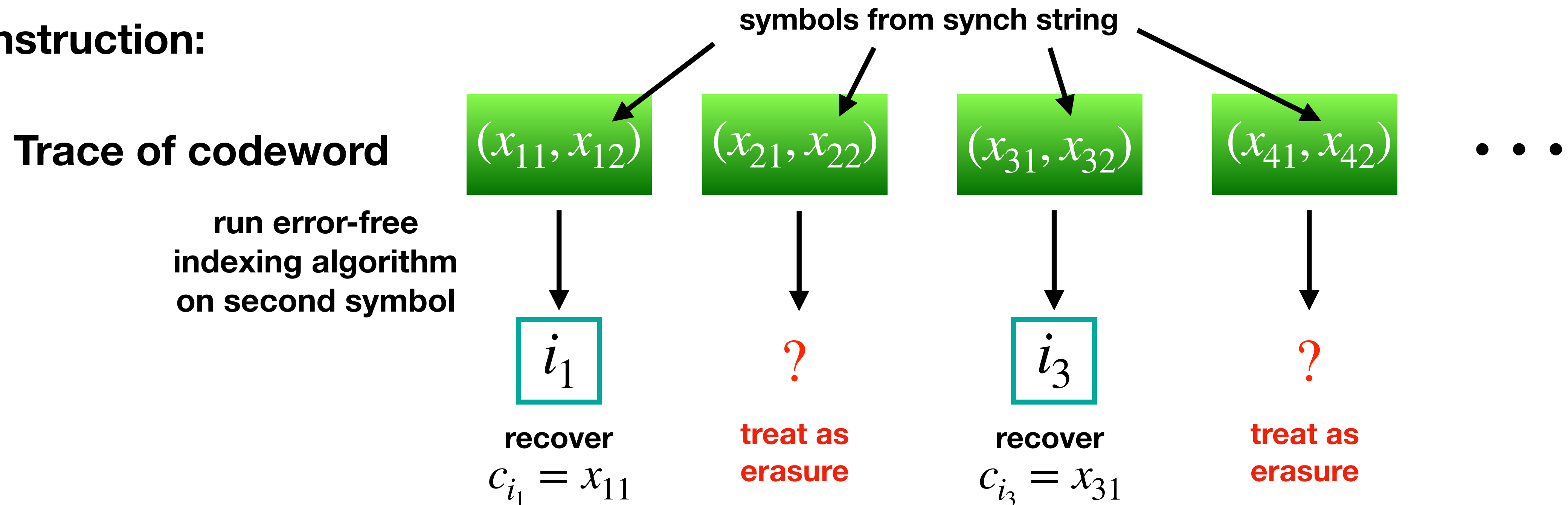
Reconstruction:



Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

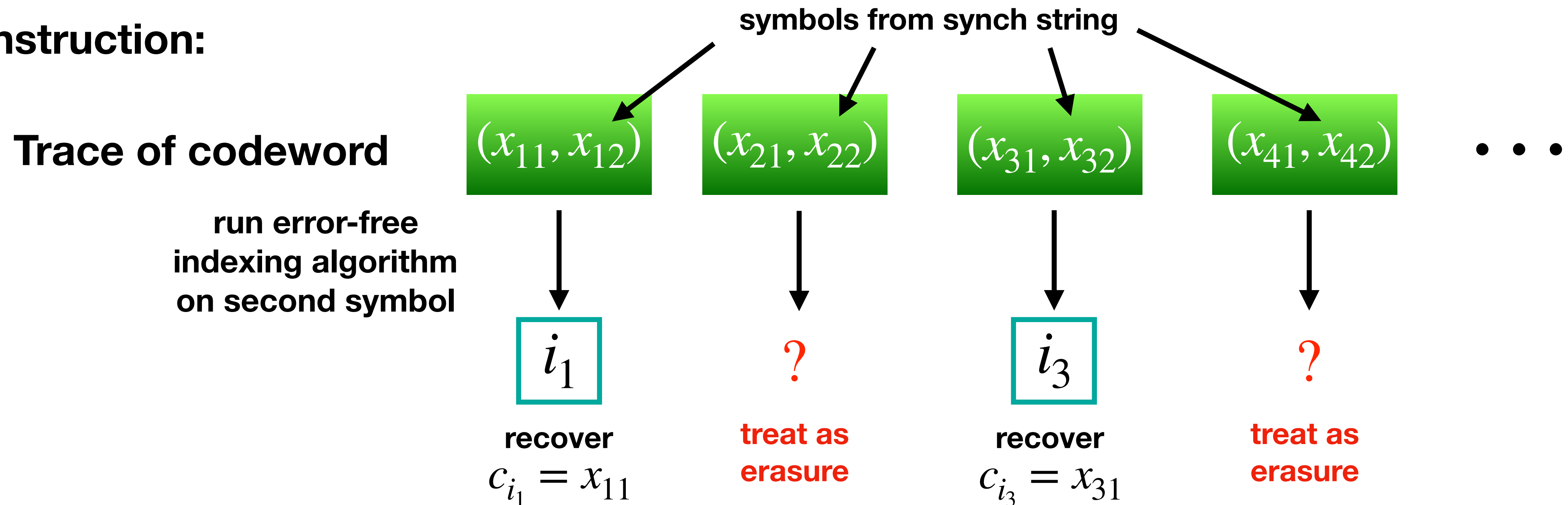
Reconstruction:



Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

Reconstruction:

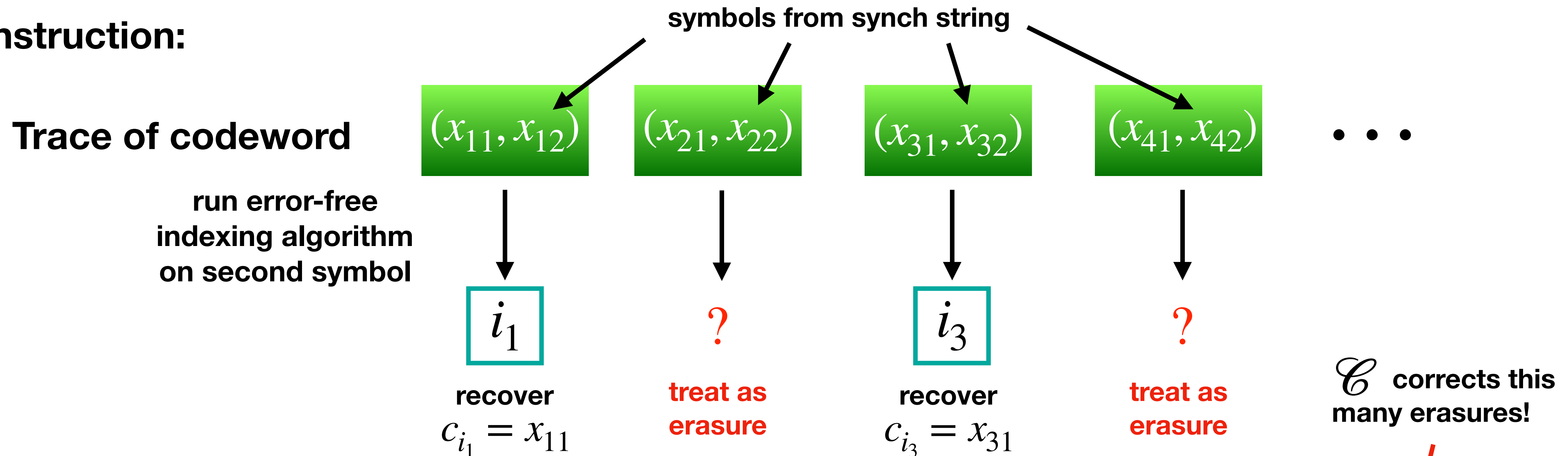


$$\# \text{erasures} \leq \# \text{misdecodings} + \#(\text{symbols deleted in every trace}) \leq \frac{\epsilon^3 N}{2} + d^T N \leq \epsilon^3 N$$

Coded trace reconstruction from synchronization strings

[BrakensiekLiSpang19]: Efficient code over alphabet of size $O_\epsilon(1)$ using $O(\log_{1/d}(1/\epsilon))$ traces

Reconstruction:



$$\# \text{erasures} \leq \# \text{misdecodings} + \#(\text{symbols deleted in every trace}) \leq \frac{\epsilon^3 N}{2} + d^T N \leq \epsilon^3 N$$

From large alphabets to binary codes

Construction: *Concatenate good code over large alphabet with marker-based inner code*

x N -bit string



From large alphabets to binary codes

Construction: *Concatenate good code over large alphabet with marker-based inner code*

x N -bit string



\mathcal{C}_{out} large alphabet code correcting many substitutions



From large alphabets to binary codes

Construction: *Concatenate good code over large alphabet with marker-based inner code*

x N -bit string



\mathcal{C}_{out} large alphabet code correcting many substitutions



S 1/3-synchronization string over alphabet of size $O(1)$



From large alphabets to binary codes

Construction: Concatenate good code over large alphabet with marker-based inner code

x N -bit string



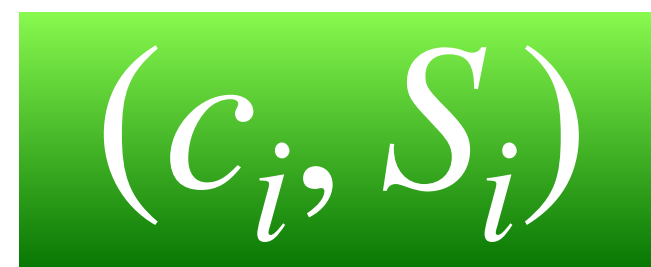
\mathcal{C}_{out} large alphabet code correcting many substitutions



S 1/3-synchronization string over alphabet of size $O(1)$



concatenate with binary code



From large alphabets to binary codes

Construction: Concatenate good code over large alphabet with marker-based inner code

x N -bit string



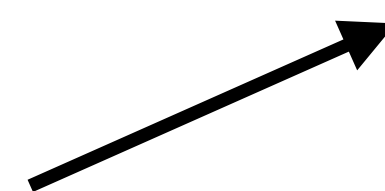
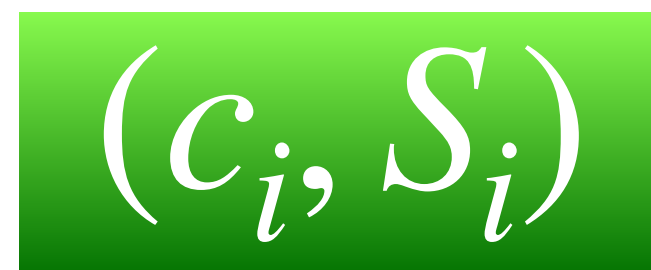
\mathcal{C}_{out} large alphabet code correcting many substitutions



S 1/3-synchronization string over alphabet of size $O(1)$



concatenate with binary code



$$\text{Enc}_R(c_i) = \begin{array}{|c|c|c|} \hline 0 & \bar{c}_i & 1 \\ \hline \log m & m & \log m \\ \hline \end{array}$$

From large alphabets to binary codes

Construction: Concatenate good code over large alphabet with marker-based inner code

x N -bit string



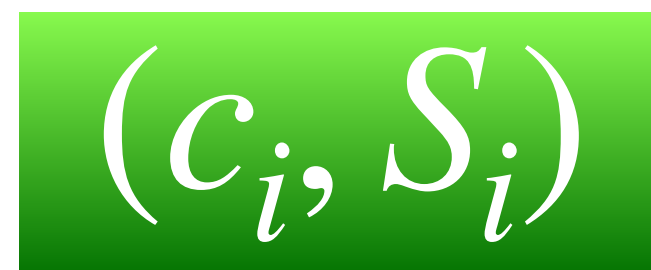
\mathcal{C}_{out} large alphabet code correcting many substitutions



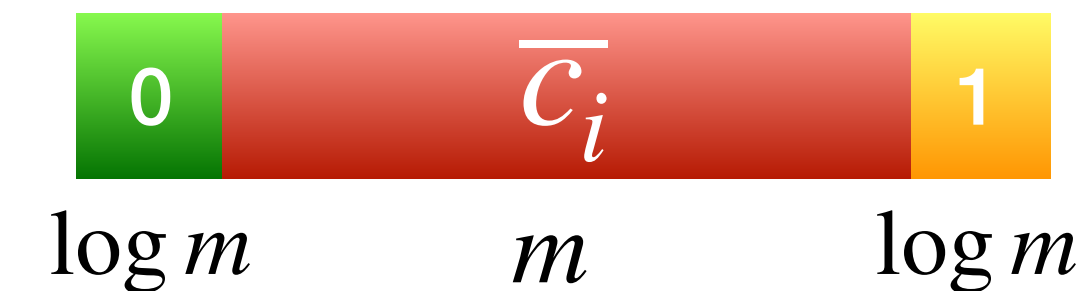
S 1/3-synchronization string over alphabet of size $O(1)$



concatenate with binary code



$\text{Enc}_R(c_i) =$



\bar{c}_i "Dense" codeword TBD

From large alphabets to binary codes

Construction: Concatenate good code over large alphabet with marker-based inner code

x N -bit string



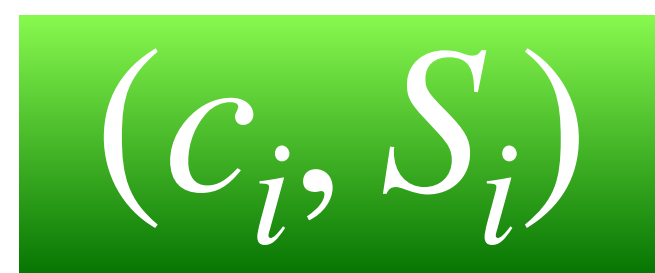
\mathcal{C}_{out} large alphabet code correcting many substitutions



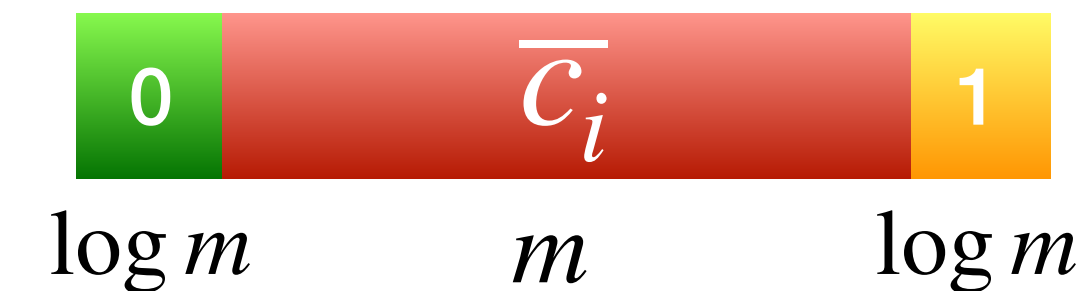
S 1/3-synchronization string over alphabet of size $O(1)$



concatenate with binary code



$$\text{Enc}_R(c_i) =$$



\bar{c}_i "Dense" codeword TBD

$$\text{Enc}_S(S_i) =$$



From large alphabets to binary codes

Construction: Concatenate good code over large alphabet with marker-based inner code

x N -bit string



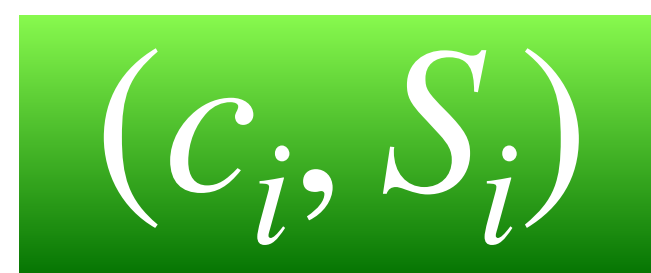
\mathcal{C}_{out} large alphabet code correcting many substitutions



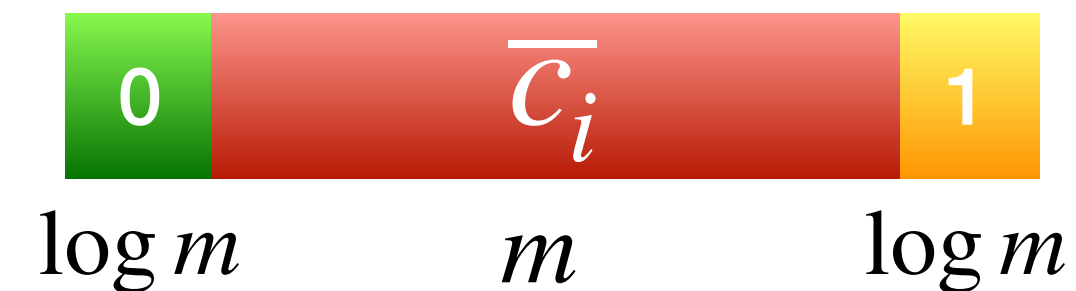
S 1/3-synchronization string over alphabet of size $O(1)$



concatenate with binary code

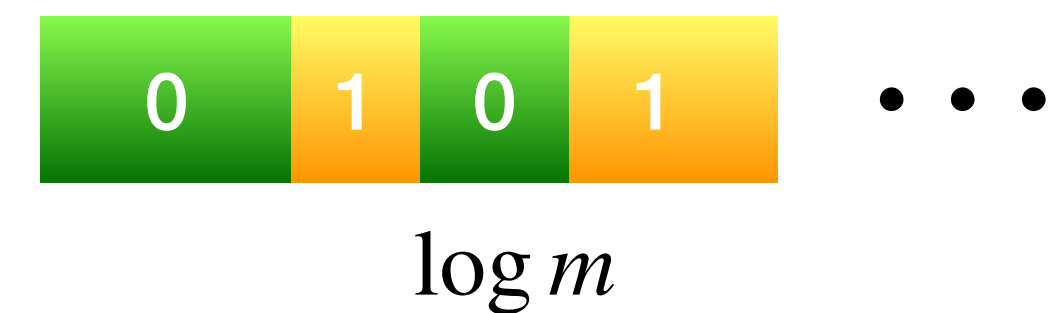


$$\text{Enc}_R(c_i) =$$



\bar{c}_i "Dense" codeword TBD

$$\text{Enc}_S(S_i) =$$



Robust against deletions, with **two** possible runlengths, same number of runs

Reconstruction



Reconstruction



Reconstruction



1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols

Reconstruction



1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols



Reconstruction

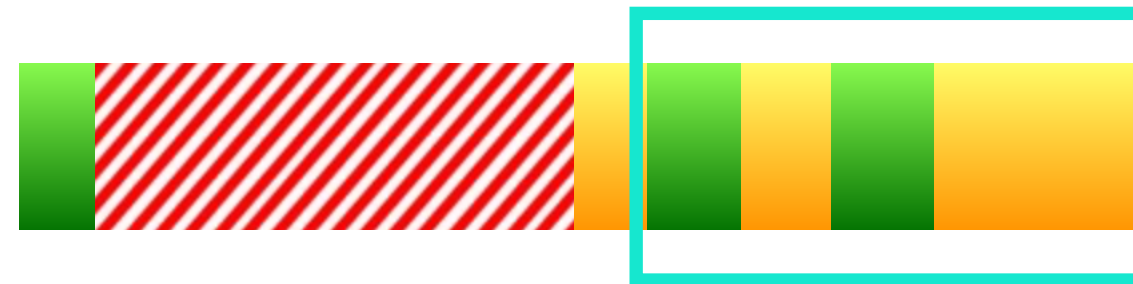


1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols



recover \tilde{S}

Reconstruction

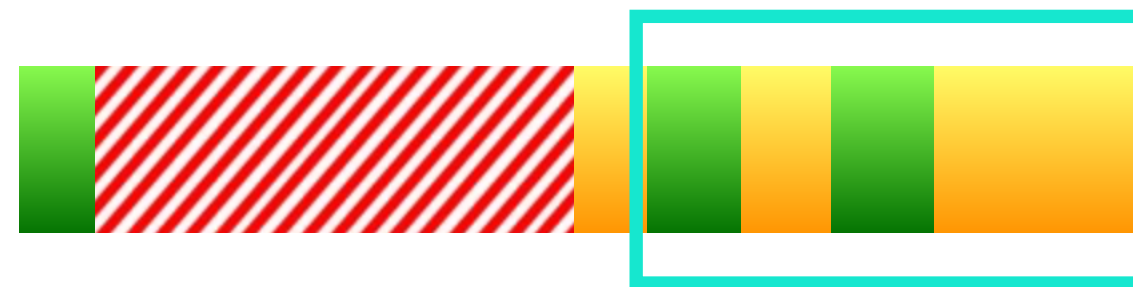


1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols



recover \tilde{S}



run error-free
indexing algo

Reconstruction

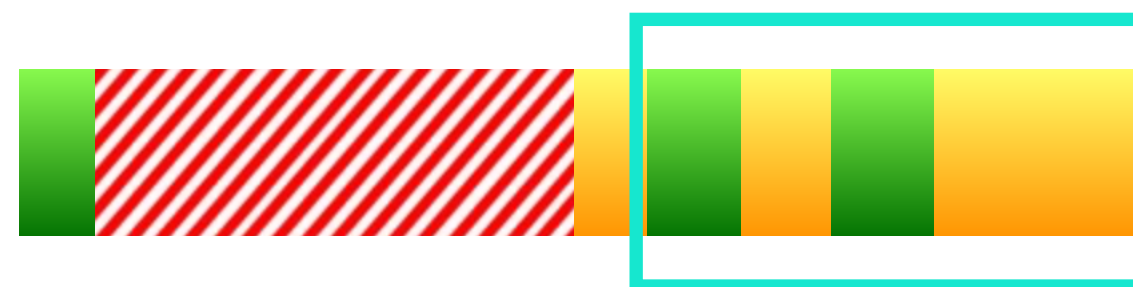


1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols



recover \tilde{S}



run error-free
indexing algo



Recover index i

Reconstruction

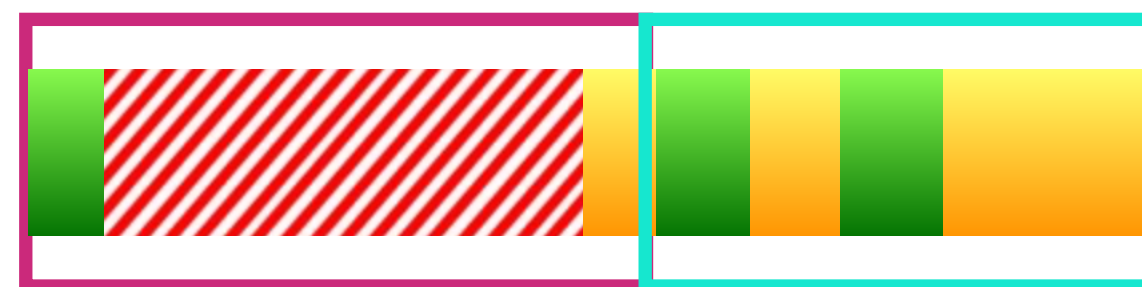


1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols



trace of C_i
whp

recover \tilde{S}



run error-free
indexing algo



Recover index i

Reconstruction

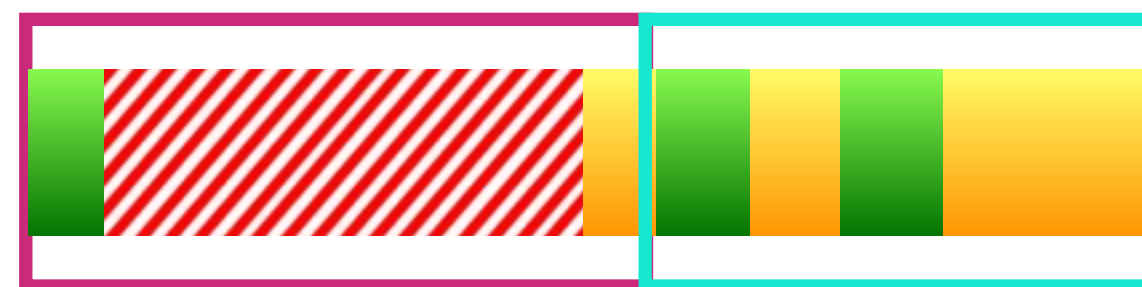


1) Trace alignment:

Long runs **only come from markers**



Find traces of (most) synch symbols



trace of C_i
whp

recover \tilde{S}



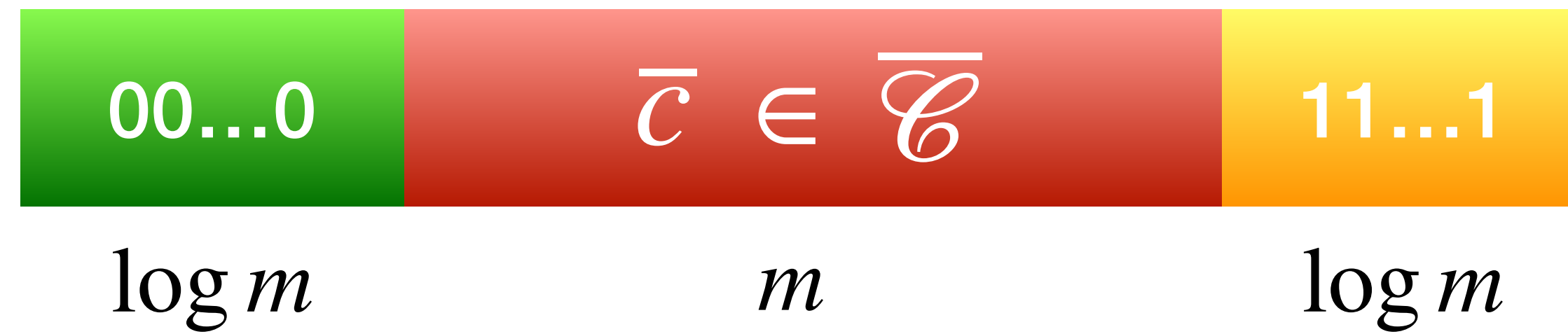
run error-free
indexing algo



Recover index i

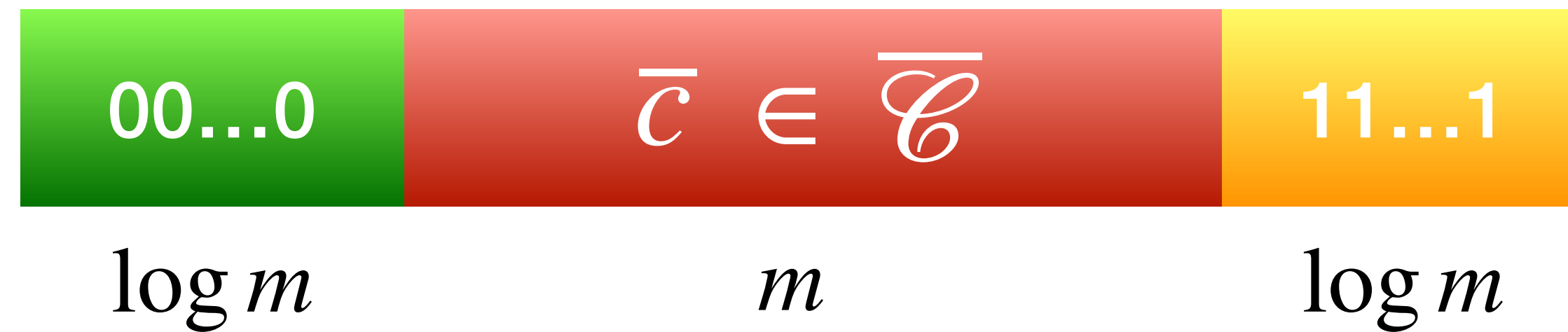
2) Use special trace rec properties of $\overline{\mathcal{C}}$ + error-correction properties of \mathcal{C}_{out}

Designing the (existential) inner code



Key observation: Since markers are short, **nearly all strings of this type are good for any average-case trace rec algorithm!**

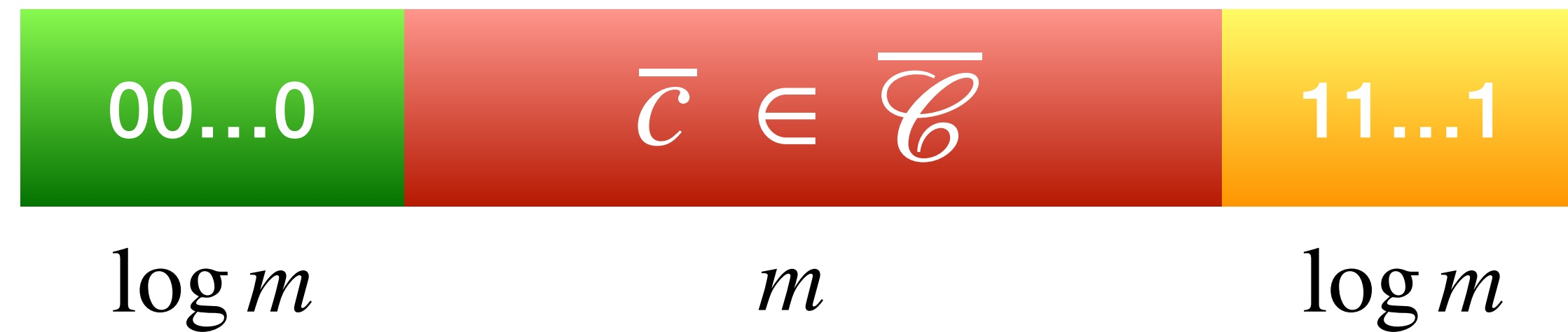
Designing the (existential) inner code



Key observation: Since markers are short, **nearly all strings of this type are good for any average-case trace rec algorithm!**

⇒ If m is **small**, can **brute force** high-rate code $\bar{\mathcal{C}}$ to require **few traces + be dense!**

Designing the (existential) inner code



Key observation: Since markers are short, **nearly all strings of this type are good for any average-case trace rec algorithm!**

⇒ If m is **small**, can **brute force** high-rate code $\bar{\mathcal{C}}$ to require **few traces + be dense!**

⇒ Using [HPP18], reconstruct $\bar{\mathcal{C}}$ with $\exp(\log^{1/3} m)$ traces.

Combining everything

Set inner code length to $m = \frac{1}{\epsilon} \log \left(\frac{1}{\epsilon} \right)$ to ensure rate $\geq 1 - \epsilon$

[BrakensiekLiSpang19]

For any constant deletion probability and every ϵ , there exists a code of rate $1 - \epsilon$ that can be reconstructed from $\exp(\log^{1/3}(1/\epsilon))$ traces.

Combining everything

Set inner code length to $m = \frac{1}{\epsilon} \log \left(\frac{1}{\epsilon} \right)$ to ensure rate $\geq 1 - \epsilon$

[BrakensiekLiSpang19]

For any constant deletion probability and every ϵ , there exists a code of rate $1 - \epsilon$ that can be reconstructed from $\exp(\log^{1/3}(1/\epsilon))$ traces.

Caveat: Inner code construction takes time $\exp(m)$.



Overall code construction only efficient for $\epsilon \geq \frac{\log \log N}{\log N}$

Summing up...

We can exploit worst-case and average-case trace reconstruction to design efficient high-rate codes requiring significantly fewer traces, or satisfying other nice properties.

	rate = $1 - \epsilon$	#traces	efficient encoding & reconstruction	observations
Markers + modified avg-case trace rec [CGMR19]	$\epsilon = \frac{1}{\log N}$	$\log^c N$	✓	Deletion probability smaller than absolute constant
Synch string + Markers + existential trace rec [BLS19]	$\epsilon = \frac{1}{\log N}$	$\exp((\log \log N)^{1/3})$	✓	Code construction is not efficient
Synch string + Markers + existential trace rec [BLS19]	$\epsilon \geq \frac{\log \log N}{\log N}$	$\exp(\log^{1/3}(1/\epsilon))$	✓	Arbitrary constant deletion probability

Summing up...

We can exploit worst-case and average-case trace reconstruction to design efficient high-rate codes requiring significantly fewer traces, or satisfying other nice properties.

	rate = $1 - \epsilon$	#traces	efficient encoding & reconstruction	observations
Markers + modified avg-case trace rec [CGMR19]	$\epsilon = \frac{1}{\log N}$	$\log^c N$	✓	Deletion probability smaller than absolute constant
Synch string + Markers + existential trace rec [BLS19]	$\epsilon = \frac{1}{\log N}$	$\exp((\log \log N)^{1/3})$	✓	Code construction is not efficient
Synch string + Markers + existential trace rec [BLS19]	$\epsilon \geq \frac{\log \log N}{\log N}$	$\exp(\log^{1/3}(1/\epsilon))$	✓	Arbitrary constant deletion probability

Lower bound [BLS19]: Arbitrary code of rate $1 - \epsilon \implies$ requires $\approx \log^{5/2} \left(\frac{1}{\epsilon} \right)$ traces

Future work

- Efficient high-rate codes using even fewer traces;
- Bridge gap between bounds for (coded and uncoded) trace reconstruction;
- High-rate codes that handle deletions **and** random insertions with few traces;

Future work

- Efficient high-rate codes using even fewer traces;
- Bridge gap between bounds for (coded and uncoded) trace reconstruction;
- High-rate codes that handle deletions **and** random insertions with few traces;

Thanks!