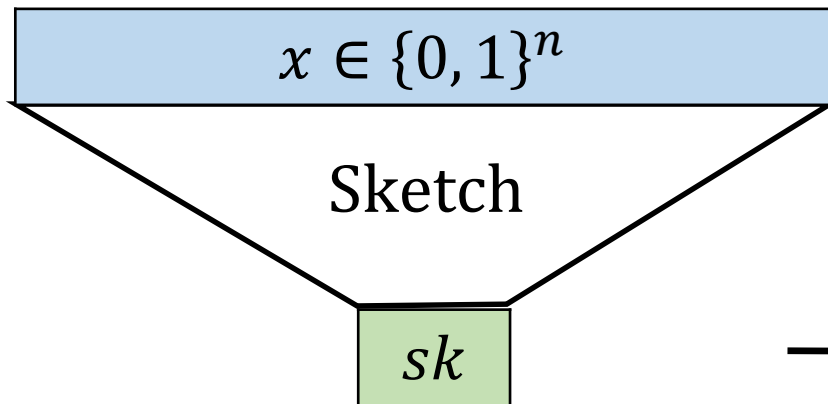


Deterministic Document Exchange Protocols, and Almost Optimal Binary Codes for Edit Errors

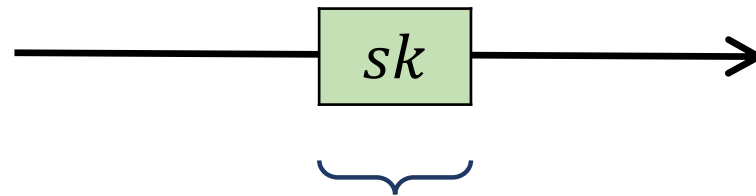
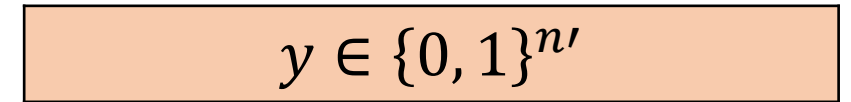
Kuan Cheng, **Zhengzhong Jin**, Xin Li, Ke Wu



Document Exchange Protocol



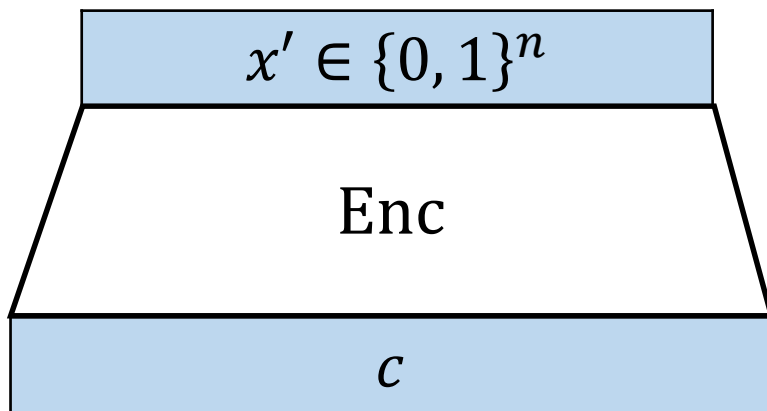
$$ED(x, y) \leq k$$



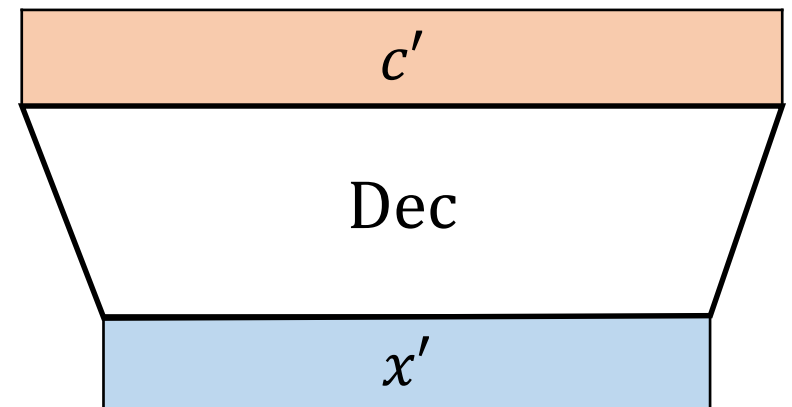
$$x = \text{Recover}(sk, y)$$

Minimize Sketch Size

Error Correcting Code



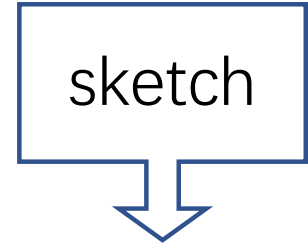
$\leq k$ Edit Errors



- Minimize Redundancy = $|c| - |x'|$
- Maximize Information Rate = $|x'|/|c|$

Previous Results For Hamming Errors

- Document Exchange from Systematic Error Correcting Code



	Reed-Solomon Code	Algebraic Geometry Code (Optimal)
Redundancy (sketch size)	$k \log n$	$\Theta\left(k \log \frac{n}{k}\right)$
Information Rate	$1 - \Theta(\epsilon \log n)$	$1 - \Theta\left(\epsilon \log \frac{1}{\epsilon}\right)$

- $\epsilon = k/n$

Results

	Previous	Our Results	Bounds
Error Correcting Code for Edit Distance:			
Redundancy	$O(k^2 \log k \log n)$ for $k = O(1)$ [BGZ17]	$O(k \log n)$	$\Theta\left(k \log \frac{n}{k}\right)$
	$O(k^2 + k \log^2 n)$ [Bel15]	$O\left(k \log^2 \frac{n}{k}\right)$	
Information Rate	$1 - \tilde{O}(\sqrt{\epsilon})$ [GL16, GW17]	$1 - \tilde{O}(\epsilon)$	$1 - \Theta\left(\epsilon \log\left(\frac{1}{\epsilon}\right)\right)$
Document Exchange:			
Sketch Size	$O(k^2 + k \log^2 n)$ [D] [Bel15] $O\left(k \log^2 \frac{n}{k}\right)$ [R] [IMS05]	$O\left(k \log^2 \frac{n}{k}\right)$ [D]	$\Theta\left(k \log \frac{n}{k}\right)$

- $\epsilon = k/n$
- $O(k \log n) = O\left(k \log \frac{n}{k}\right)$, when $k = n^{1-a}$, for any $a > 0$.

Overview: Almost optimal ECC for Edit Distance

Document exchange for a uniform random string



Derandomize the random string
using Pseudorandom Generator (PRG)



Error Correcting Code

Document Exchange for a Uniform Random String

Stage I:

- Alice partitions her string into blocks of size $\text{poly}(\log n)$.
- Send $O(k \log n)$ bits to help Bob learn her partition.

Stage II:

- $O(1)$ levels.
- In each level, divide each block into $O(\log^{0.4} n)$ smaller blocks.
- In each level, Alice sends $O(k \log n)$ bits to help Bob learn most blocks, except $O(k)$ of them are **incorrect/unfilled**.

Stage I (Alice)

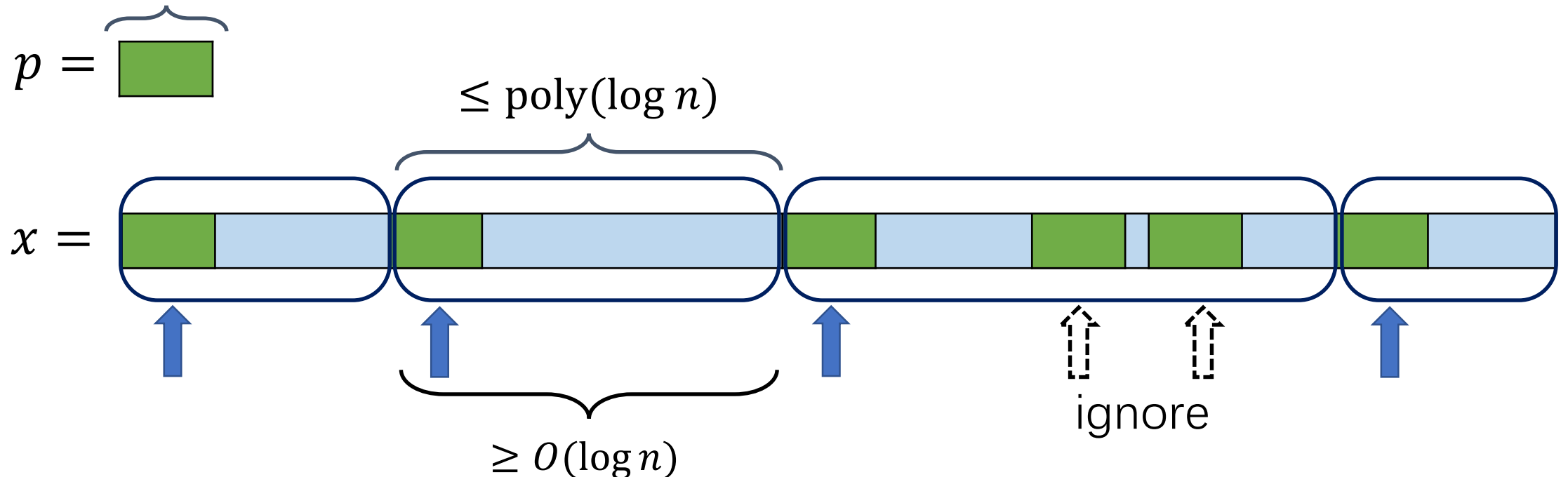
- Suppose x is uniformly random, use a fixed string p to partition x .

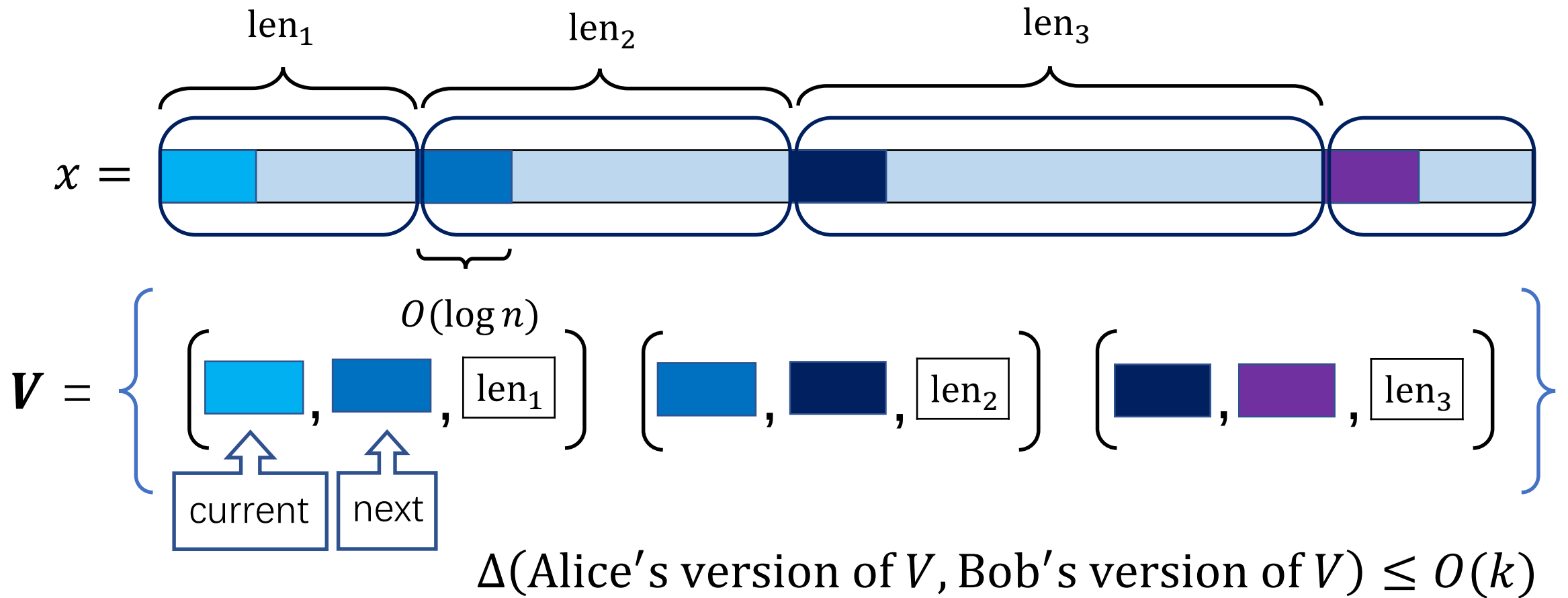
With probability $1 - 1/\text{poly}(n)$,

- Any two substrings of length $O(\log n)$ are distinct.

$\log \log n + O(1)$

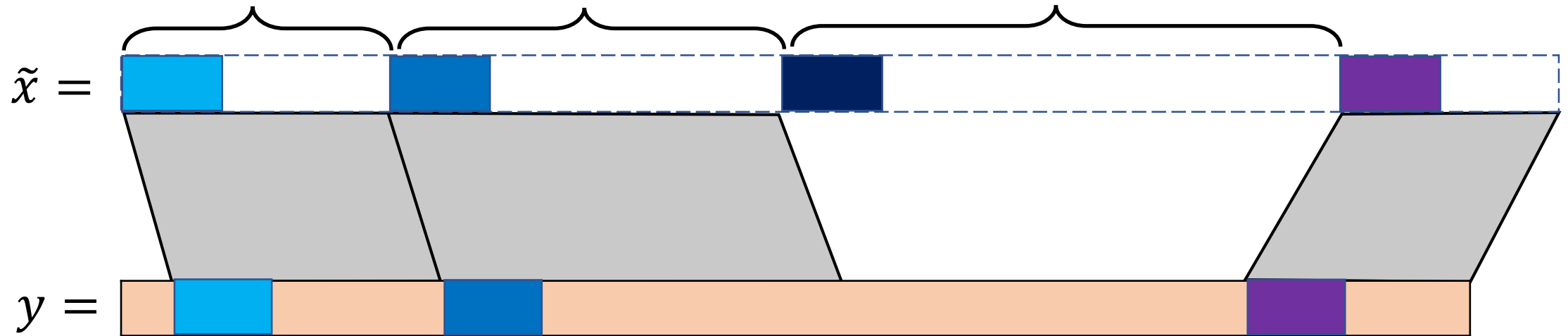
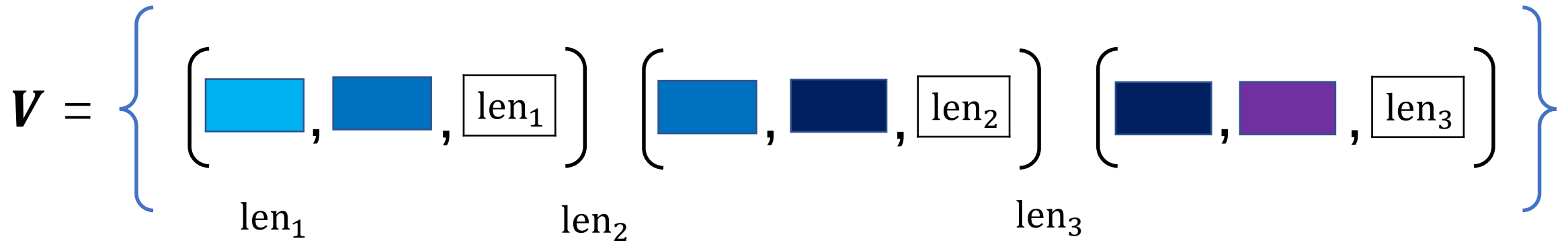
- Block length $\leq \text{poly}(\log n)$





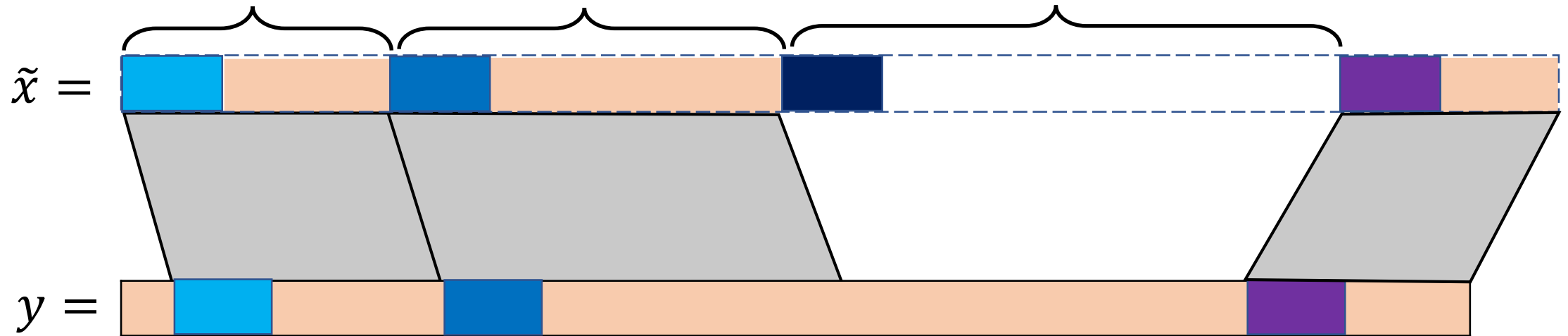
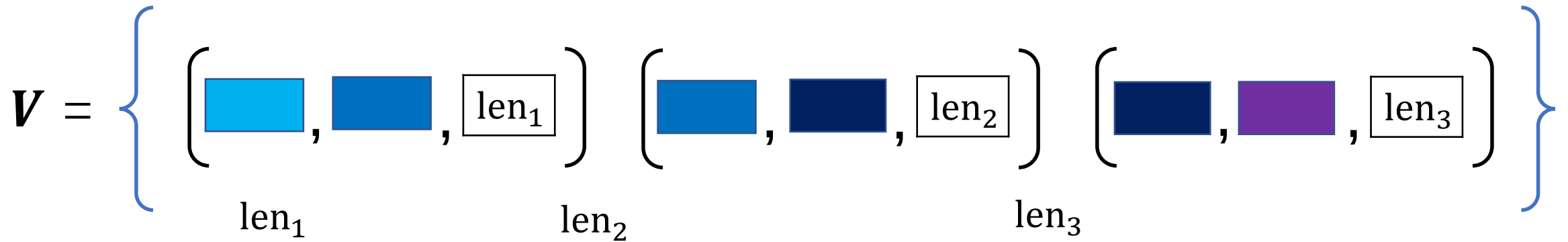
- Using Reed-Solomon Code, Alice computes a redundancy of $O(k \log n)$ bits to help Bob correct V from $O(k)$ Hamming Errors.

Stage I (Bob)



- Match the string y using the $O(\log n)$ -prefixes, then fill \tilde{x} using y .

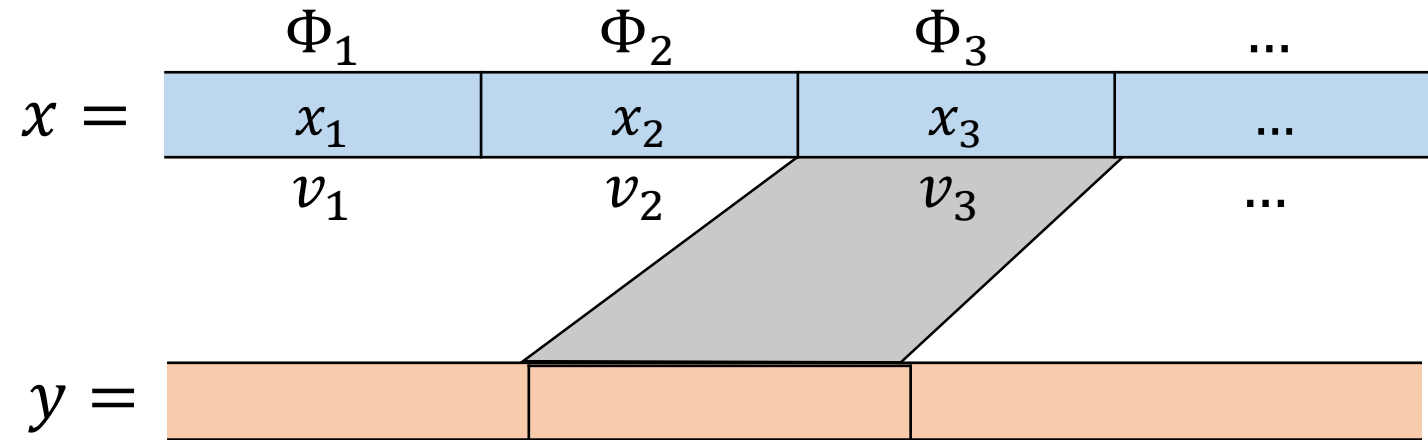
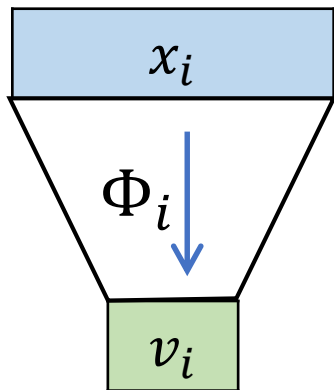
Stage I (Bob)



- Match the string y using the $O(\log n)$ -prefixes, then fill \tilde{x} using y .

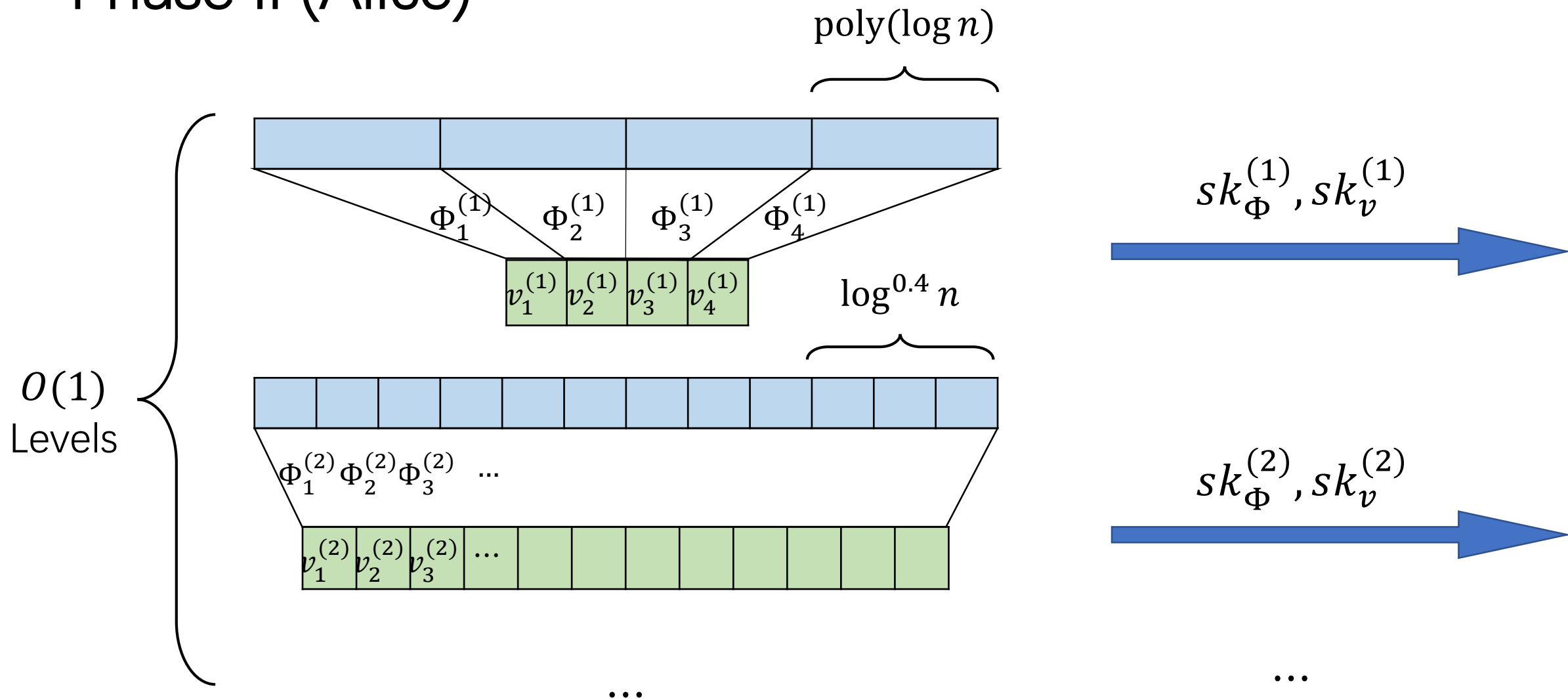
Matching under Hash Functions

- Match between x and y under hash functions
- Block Size = T , Φ_i are hash functions.

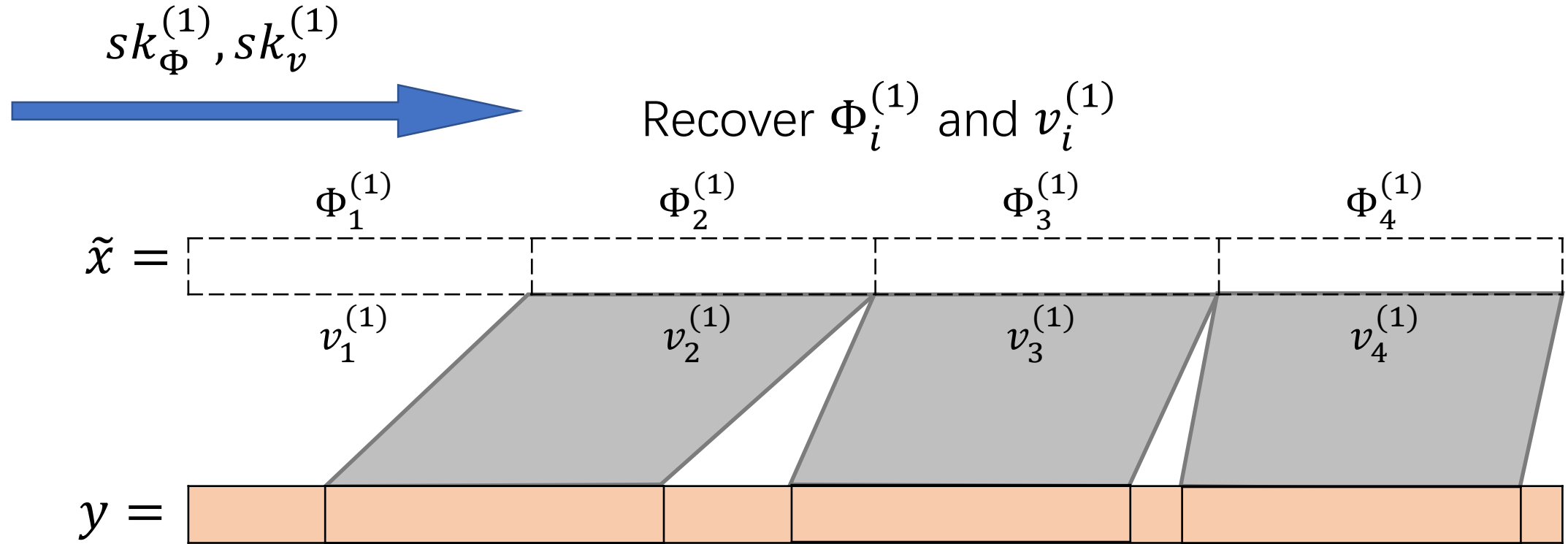


$$\text{Matched} \iff \Phi_i(\text{orange box}) = \text{green box } v_i$$

Phase II (Alice)



Phase II (Bob)



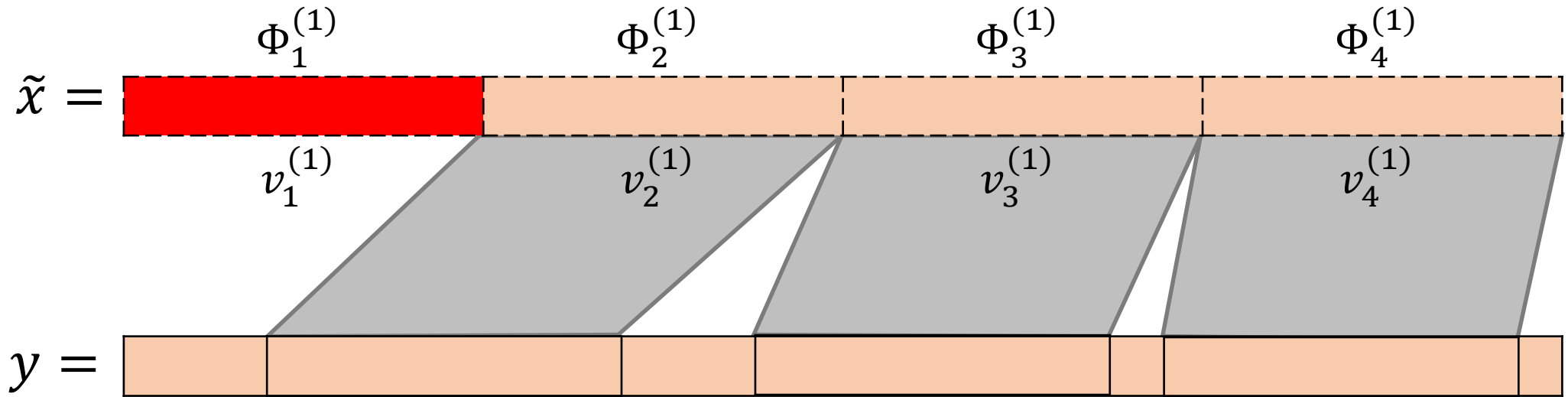
- Find the maximum monotone matching between x and y under Φ_i .
- Fill \tilde{x} using the matching. Objective: **#incorrect / unmatched blocks** $\leq O(k)$

Phase II (Bob)

$sk_{\Phi}^{(1)}, sk_v^{(1)}$

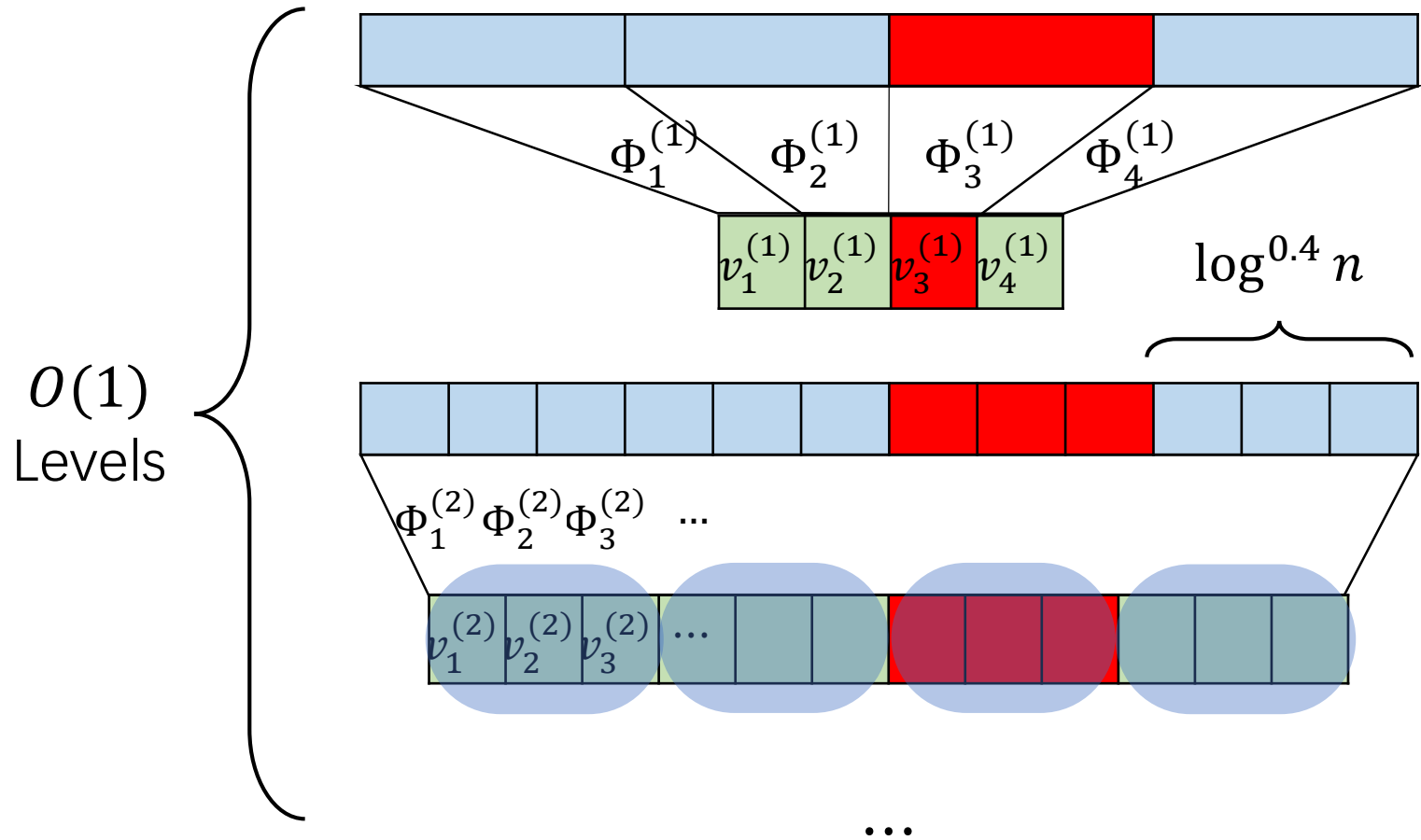


Recover $\Phi_i^{(1)}$ and $v_i^{(1)}$



- Find the maximum monotone matching between x and y under Φ_i .
- Fill \tilde{x} using the matching. Objective: **#incorrect / unmatched blocks** $\leq O(k)$

Phase II : Packing



$O(k)$ incorrect/unfilled blocks



~~$O(k \log^{0.4} n)$
incorrect/unfilled blocks~~

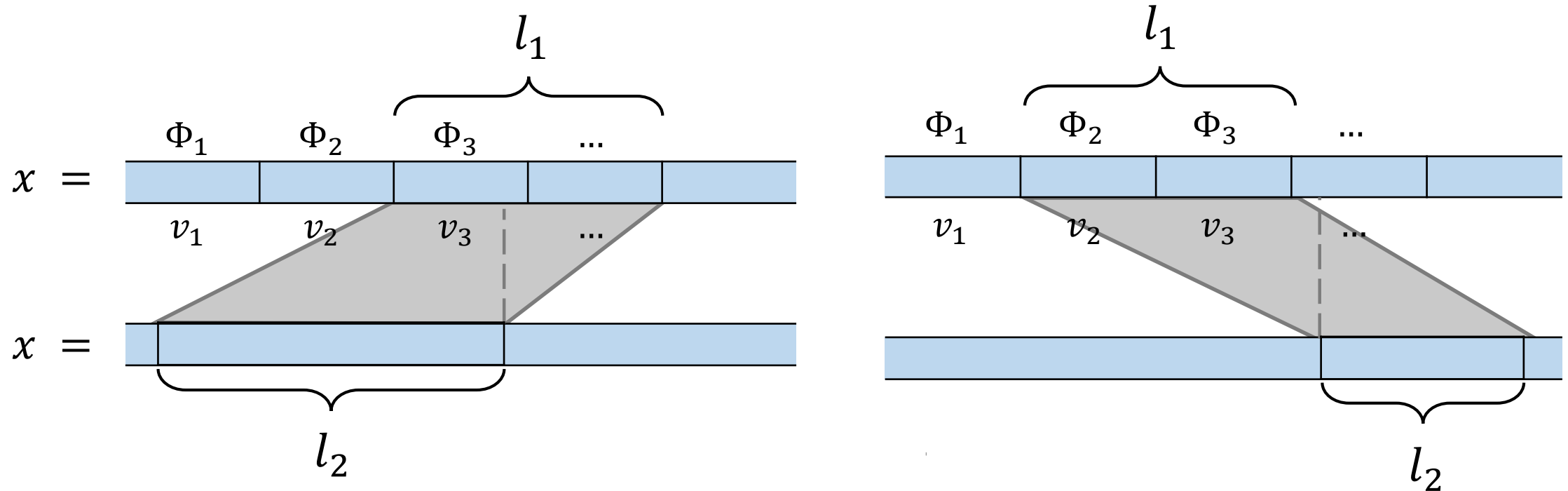
$O(k)$ incorrect/unfilled packages

- Pack every $O(\log^{0.4} n)$ successive hash values.

- Minimize the package size (or hash value size).

ϵ -Synchronization Hash Functions: Definition

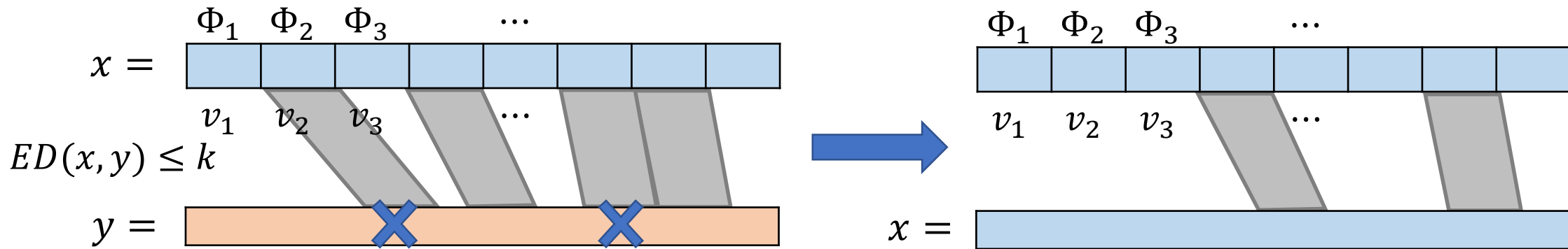
- ‘Skew-Shape’: Two adjacent intervals.



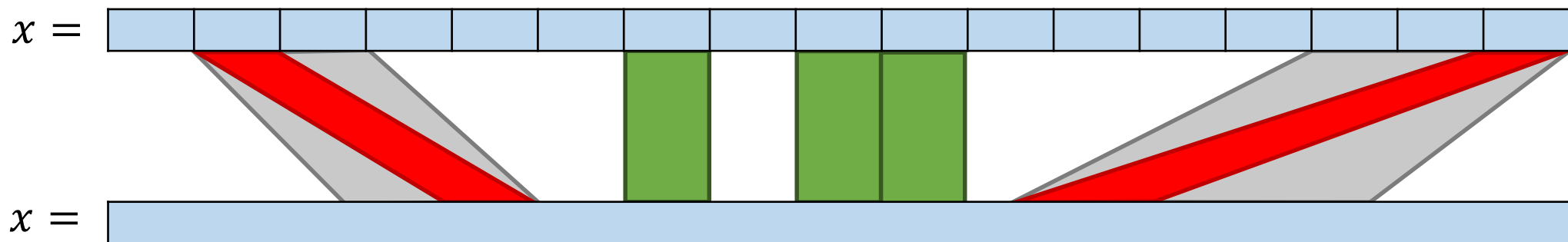
- \forall ‘Skew-Shape’, $\#MATCHING \leq \epsilon \left(l_1 + \frac{l_2}{T} \right)$, ($T =$ block size).

ϵ -Synchronization Hash Functions: Properties

- Matching between x and y induces a 'self-matching' on x itself.



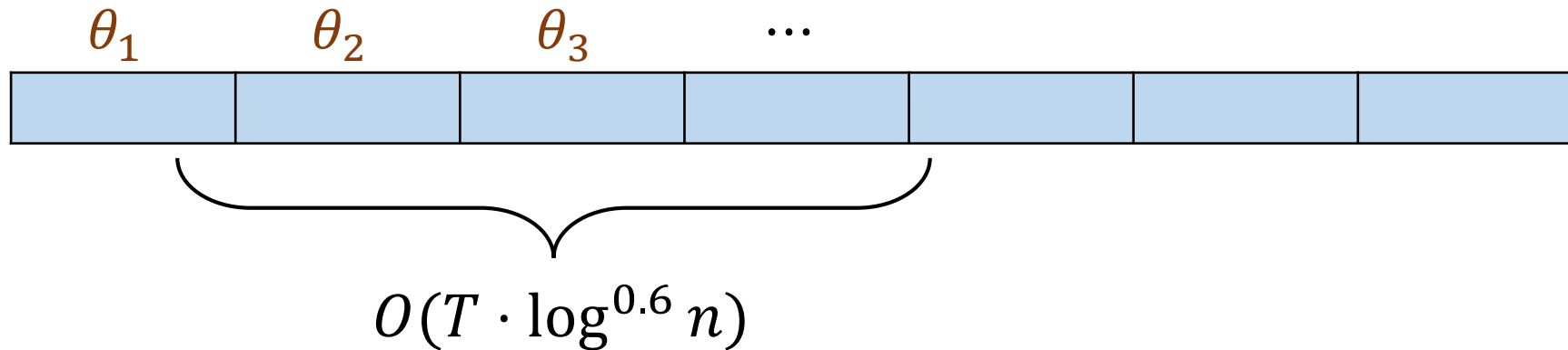
- The **wrong self-matchings** can be grouped into 'Z-Shapes'



- The two properties keeps **#incorrect/unfilled blocks** $\leq O(k)$, while minimize the hash value size.

ϵ -Synchronization Hash Functions: Construction

- For large intervals ($l_1 + l_2/T \geq \log^{0.6} n$): use uniform random functions ϕ_i (Probability Method).
- For small intervals: Use 'locally injective' hash functions θ_i .



- $\Phi_i(\text{[]}) = (\phi_i(\text{[]}), \theta_i(\text{[]})), \forall i$

Pseudorandom Generator (PRG)

- A generator $g: \{0, 1\}^r \rightarrow \{0, 1\}^n$ is a PRG against a function $f: \{0, 1\}^n \rightarrow \{0, 1\}$ with error ϵ , if

$$|\Pr[f(U_n) = 1] - \Pr[f(g(U_r)) = 1]| \leq \epsilon$$

- ϵ -almost κ -wise independence generator: An explicit construction of $g: \{0, 1\}^d \rightarrow \{0, 1\}^n$, $(X_1, X_2, \dots, X_n) = g(U_r)$, s.t.

$$\forall i_1, i_2, \dots, i_\kappa \in [n], \quad |\Pr[X_{i_1}, X_{i_2}, \dots, X_{i_\kappa} = x] - 2^{-\kappa}| \leq \epsilon$$

ϵ -Synchronization Hash Functions: Derandomization

Use almost κ -wise independence generator:

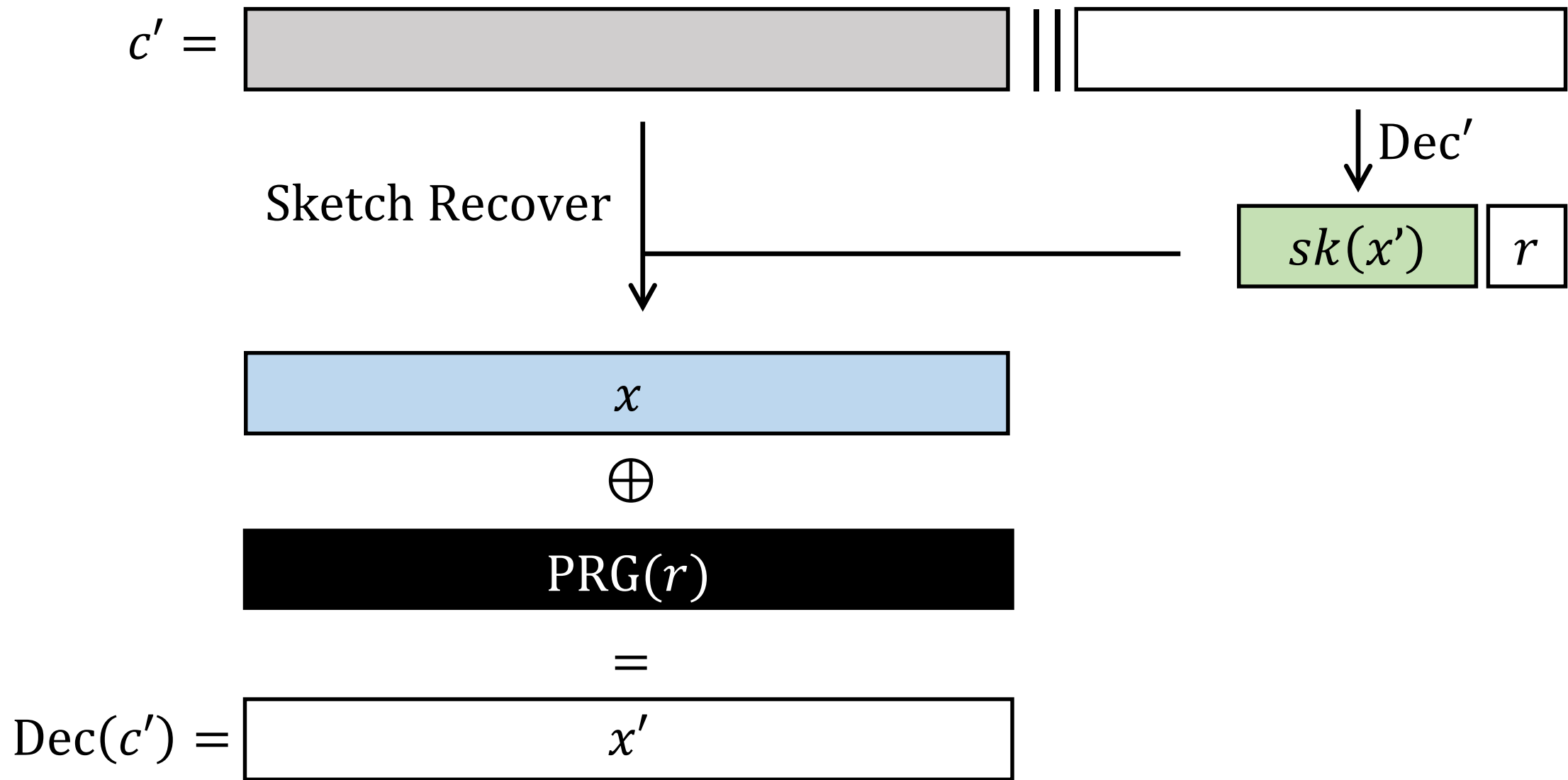
- A large fraction of random seed satisfies the properties we want.
- Use exhaustive search to find the appropriate random seed.
- For ϕ_i , random seed size = $O(\log n)$. Send it directly.
- For θ_i , random seed size = $O(\log \log n)$ bits, packing them and send a redundancy in $O(k \log n)$ bits.

Error Correcting Codes from Sketch

$$\begin{array}{c} \boxed{x'} \\ \oplus \\ \boxed{\text{PRG}(r)} \\ = \\ \text{Enc}(x) \stackrel{\text{def}}{=} \boxed{x} \quad || \quad \text{Enc}'(\boxed{sk(x)} \boxed{r}) \end{array}$$

With probability $1 - 1/\text{poly}(n)$,

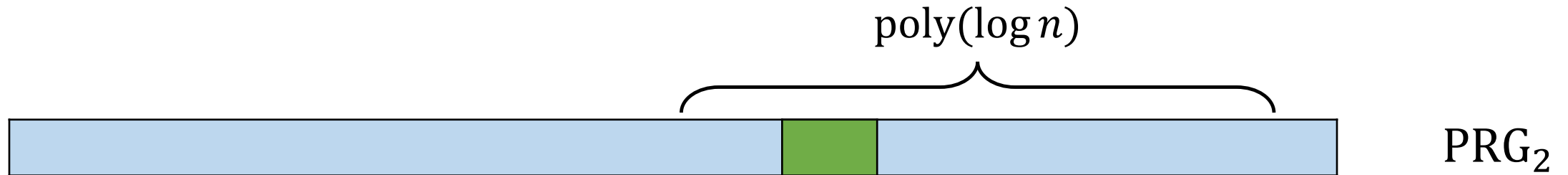
- Any two substrings of length $O(\log n)$ are distinct.
- Block length $\leq \text{poly}(\log n)$



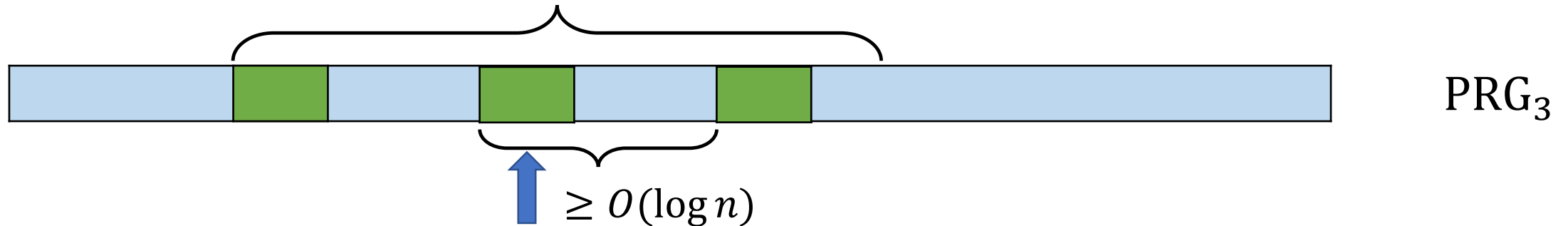
Derandomize the Uniform Random String

With probability $1 - 1/\text{poly}(n)$:

1. Any two substrings of length $O(\log n)$ are distinct. PRG₁
2. Any interval of length $\text{poly}(\log n)$ contains p .



3. Any interval of length $\text{poly}(\log n)$ starting from p contains a chosen p .



$$\text{PRG}(r) = \text{PRG}_1(r) \oplus \text{PRG}_2(r) \oplus \text{PRG}_3(r), \quad |r| = O(\log n)$$

Open Questions

- Optimal Error Correcting Code and Deterministic document exchange for all k .

Thank You!