

55th Annual Allerton Conference on Communication, Control, and Computing

Correcting Localized Deletions Using Guess & Check Codes

Salim El Rouayheb

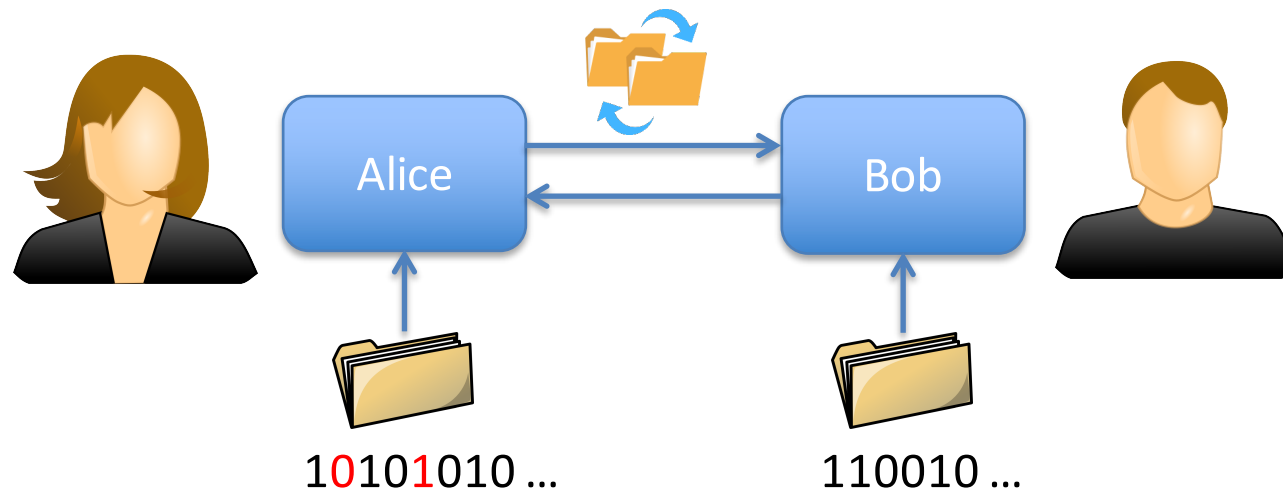
Rutgers University

Joint work with

Serge Kas Hanna and Hieu Nguyen

Motivation

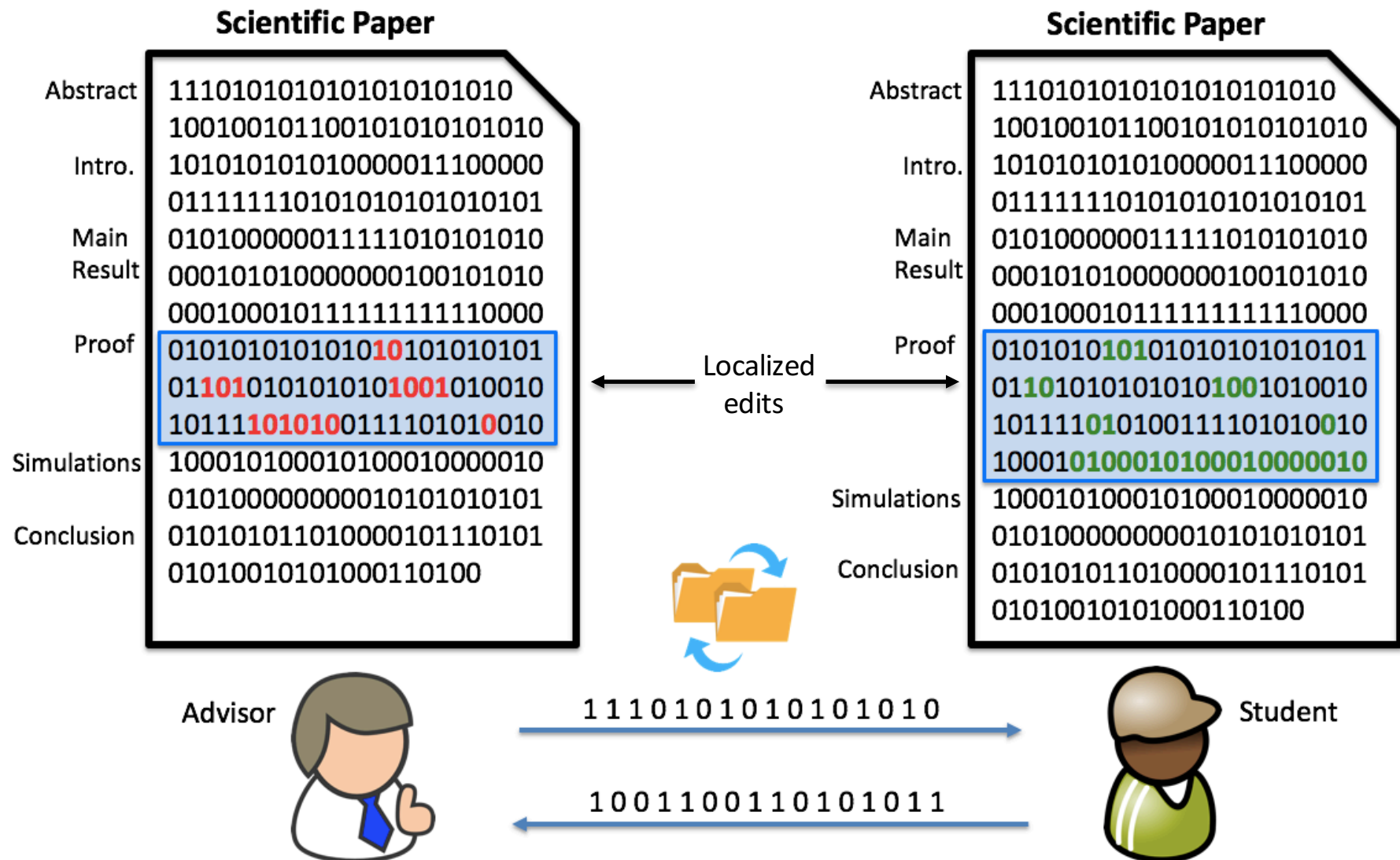
- Deletions: $\overset{\text{Transmitted}}{10101010} \longrightarrow \overset{\text{Received}}{110010}$
- Deletions were first studied by Varshamov-Tenengolts ('65) and Levenshtein ('66)
- Our motivation: file synchronization, E.g. Dropbox



- Recent application: DNA-based storage

Localized Deletions

- Motivation: file synchronization, E.g. Dropbox



Previous Work on Deletions

➤ Unrestricted deletions

- *Information theoretic approach*: [Gallager '61], [Dobrushin '67]; lower and upper bounds on the capacity: [Mitzenmacher and Drinea '06], [Diggavi et al. '07], [Kanoria and Montanari '13], [Venkataramanan et al. '13] ...
- *Recent file synchronization algorithms*: [Yazdi and Dolecek '14], [Venkataramanan et al. '15], [Sala et al. '17] ...
- *Code constructions and fundamental limits*: [Varshamov and Tenenglots '65], [Levenshtein '66], [Schulman and Zuckerman '99], [Helberg and Ferreira '02], [Cullina and Kiyavash '14], [Gabrys et al. '16], [Brankensiek et al. '16], [Thomas et al. '17] ...

➤ Bursty deletions

- File synchronization: [Ma et al. '11]
- Code constructions: [Levenshtein '67], [Cheng et al. 14], [Schoeny et al '17]

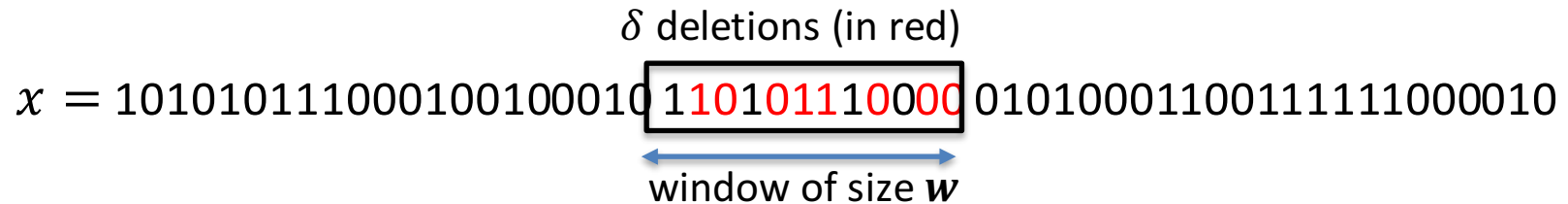
Existence of codes for localized model $w=3,4$

Model and Contribution

δ deletions (in red)

$x = 10101011100010010001011101011100000101010001100111111000010$

window of size w

The diagram shows a binary string x. A subsequence of the string is enclosed in a black box, representing a window of size w. Within this window, several bits are colored red, representing deletions. A blue double-headed arrow below the window indicates its size w.

- $\delta \leq w$ deletions localized in a window of size w
- Hard problem for $w = n$
- [Schoeny et al. '17]: existence of codes for $w = 3,4$
- Our assumptions: 1) positions of the deletions are independent of the codeword; 2) information message is uniform iid
- Contribution: Explicit codes with deterministic polynomial time encoding and decoding that can correct localized deletions whp

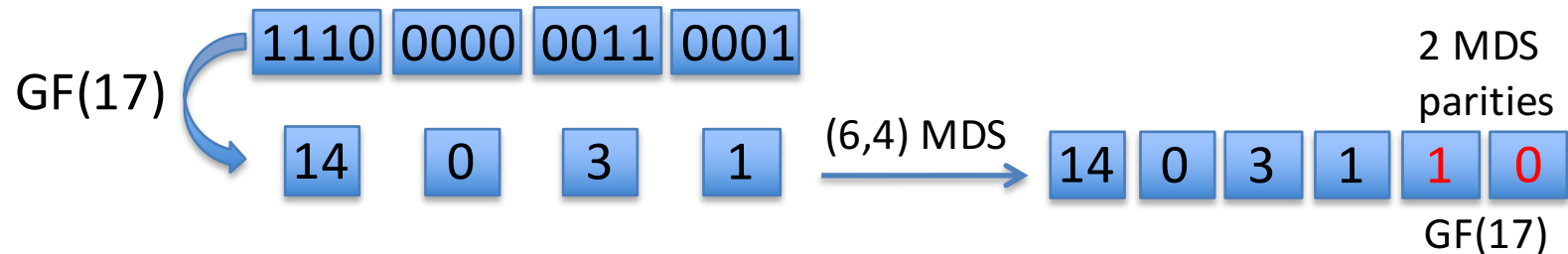
Guess &
Check (GC)
Codes

- Logarithmic redundancy: $n - k = c \log k + w + 1$
- Polynomial time encoding and decoding
- Asymptotically vanishing probability of decoding failure
- Can be generalized to multiple windows

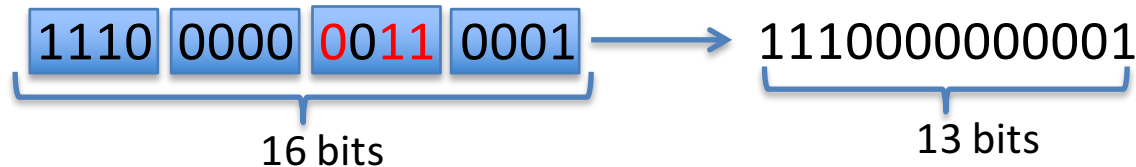
GC Codes Example

Deletions occur in one of these windows

- Encoding the message of length $k=16$: 1 1 1 00 0 0 00 0 1 10 0 0 1



- Assume that the deletions (in red) affect only one systematic block

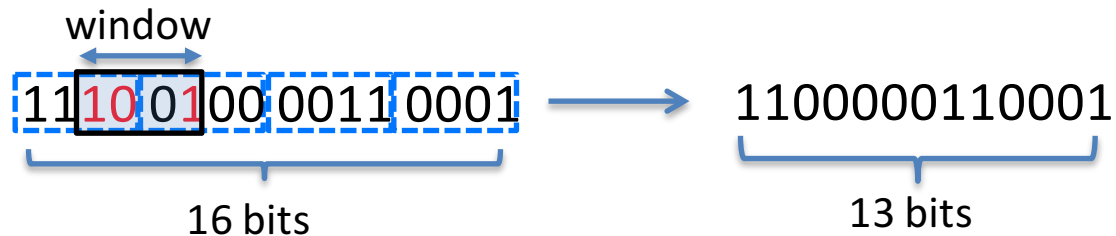


- Decoding

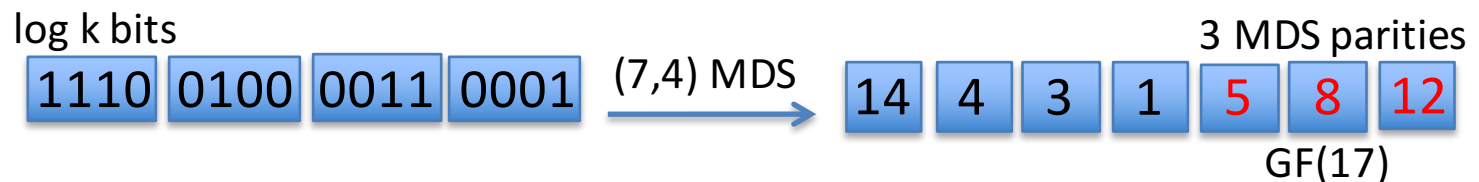
	?	12	0	1	Decoded using 1 st parity					
➤ Guess 1:	1	1100	0000	0001	5	12	0	1	✗	Check with 2 nd parity 0 [Kas Hanna and El Rouayheb ISIT 17']
➤ Guess 2:	1110	0	0000	0001	14	3	0	1	✗	
➤ Guess 3:	1110	0000	0	0001	14	0	3	1	✓	
➤ Guess 4:	1110	0000	0000	1	14	0	0	4	✗	

Generalizing to Any Window Position

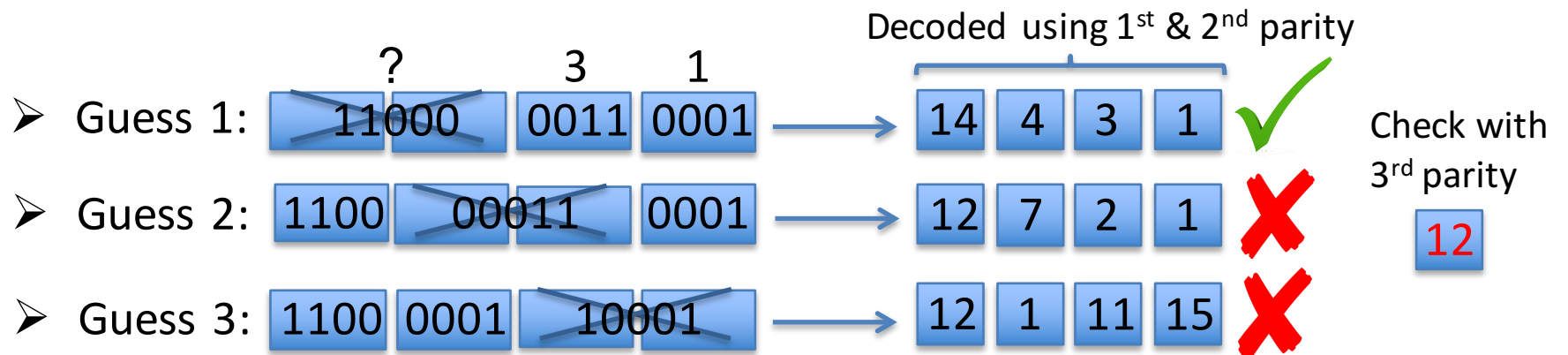
- Assume that 3 deletions (in red) affect systematic bits, $w = \log k = 4$ bits



- Same encoding with one extra parity

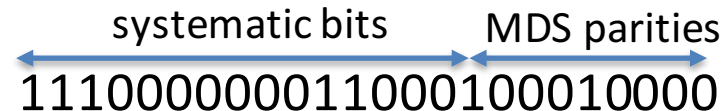


- Decoding, window of size $\log k$ can affect at most 2 **consecutive** blocks

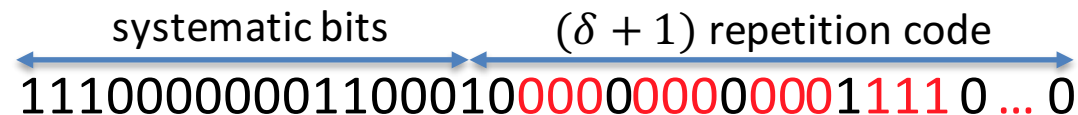


Recovering the MDS parities

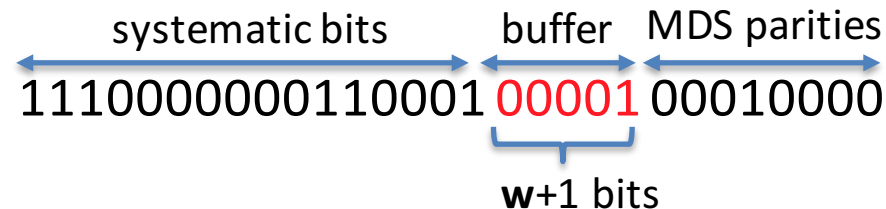
- How to recover the MDS parity symbols at the decoder?



- Trivial solution: repeat the parity bits



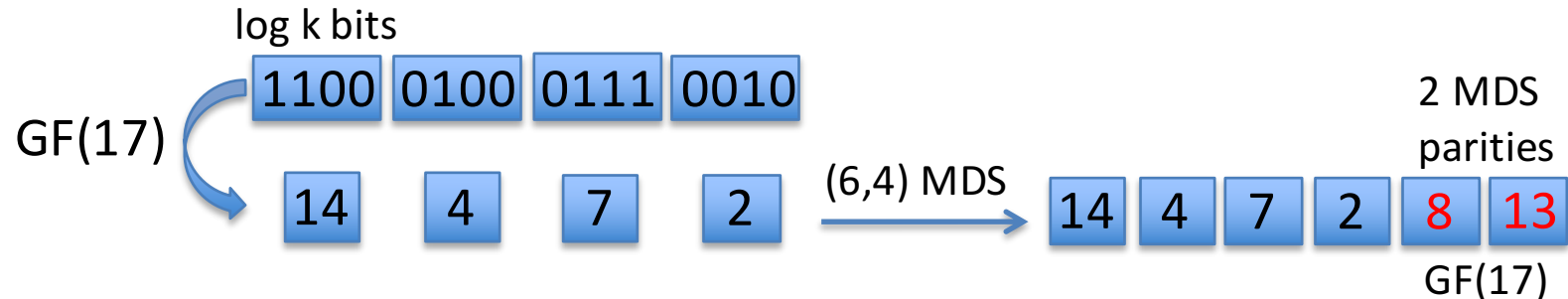
- Better solution: insert a buffer between systematic and MDS parity bits



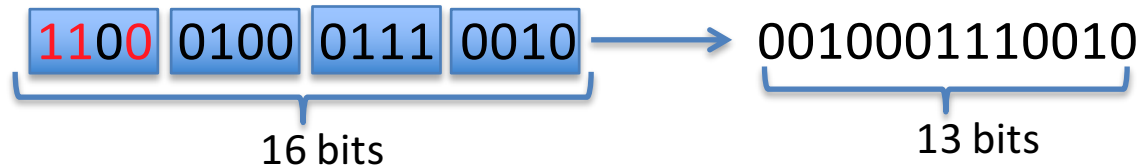
- Buffer: w zeros + a single one
- Deletions cannot affect both systematic and parity bits simultaneously
- If parity bits get affected \rightarrow simply output the first k bits. Else apply Guess & Check decoding

When does decoding fail?

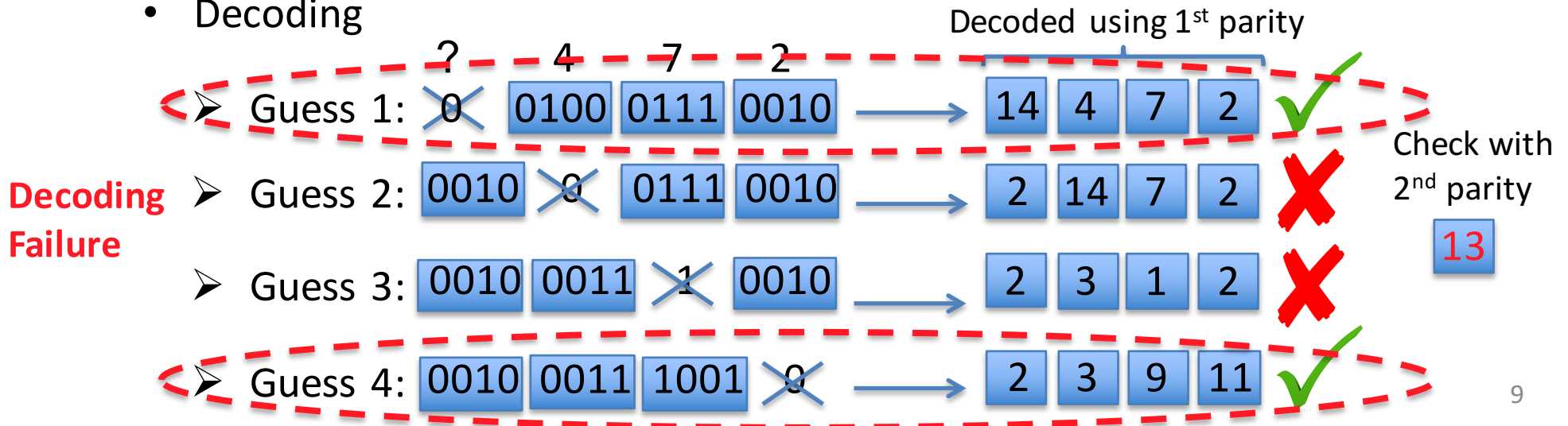
- Encoding the message of length $k=16$: 1 1 0 0 0 1 0 0 0 1 1 1 0 0 1 0



- Assume that the deletions (in red) affect only one systematic block

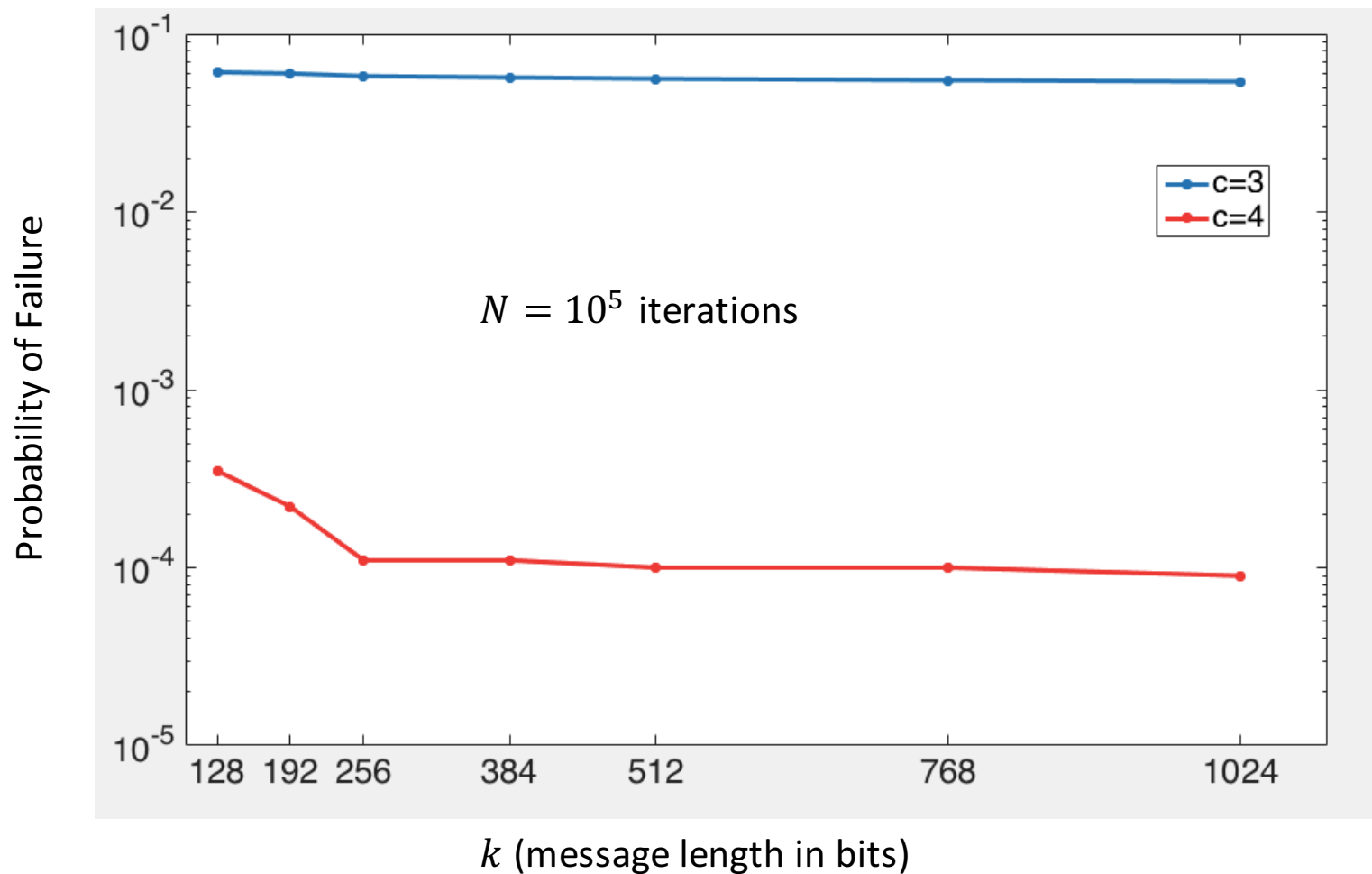


- Decoding



Simulations – Decoding Failure

- Simulation results for: $w = \log k$, $\delta = \log k - 1$, $c = 3, 4$ MDS parities



Main Results

Theorem 1 (One window): Guess & Check (GC) codes can correct in polynomial time up to $\delta \leq w = O(\log k)$ localized deletions, where $m \log k < w < (m + 1) \log k$ for some integer $m \geq 0$. Let $c > m + 2$ be a constant integer.

- Redundancy: $c \log k + w + 1$
- Encoding complexity is $\mathcal{O}(k \log k)$, Decoding complexity is $\mathcal{O}(k^3 / \log k)$
- Probability of decoding failure: $\Pr(F) \leq k^{m+4-c} / \log k$

Sketch of Theorem 2 ($z > 1$ windows):

- Redundancy: $c(zw + 1) \log k$
- Encoding complexity is $\mathcal{O}(k \log k)$, Decoding complexity is $\mathcal{O}(k^{z+2})$
- Probability of decoding failure: $\Pr(F) \leq k^{z(m+4)-c}$

Sketch of Theorem 3 [ISIT '17] (Unrestricted deletions):

- Redundancy: $c(\delta + 1) \log k$
- Encoding complexity is $\mathcal{O}(k \log k)$, Decoding complexity is $\mathcal{O}(k^{\delta+2} / \log^\delta k)$
- Probability of decoding failure: $\Pr(F) = \mathcal{O}(k^{2\delta-c} / \log^\delta k)$

Test GC Codes Online



- Test the codes online using the Jupyter notebook

Go to : <https://try.jupyter.org/>

Upload the notebook files from <http://eceweb1.rutgers.edu/csi/software.html>

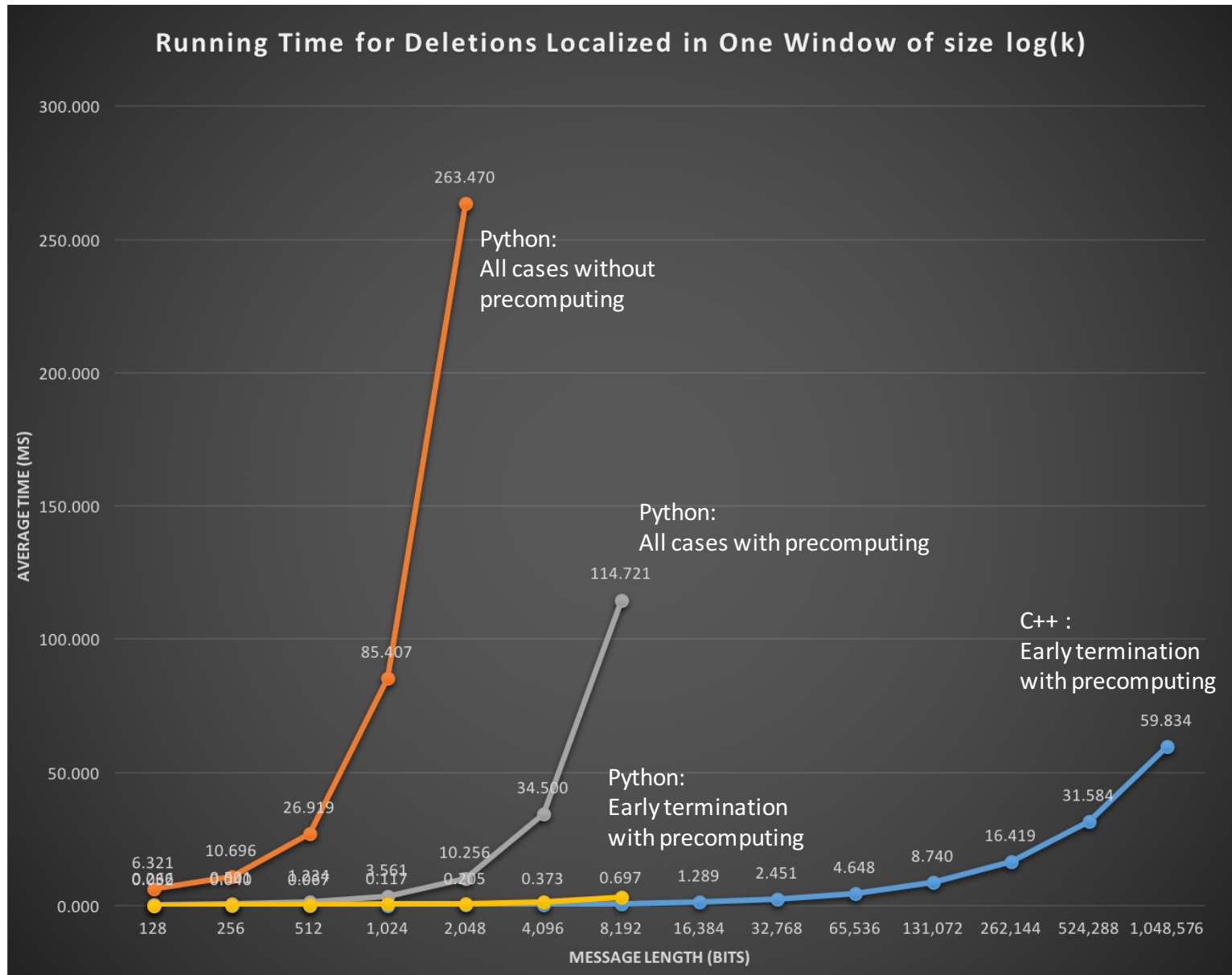
- C++ & Python codes are available on GitHub



GitHub repository: <https://github.com/serge-k-hanna/GC>

- For more details: <http://eceweb1.rutgers.edu/csi/software.html>

Simulations – Running Time



Decoding Failures: What Happened.

- Decoding failure: more than one possible guess, different decoded strings
- Example for one deletion: 16-bit message 0000100011110110

➤ (6,4) MDS encoding over GF(17): 0 8 15 6 12 5 2 parities

➤ Suppose 14th bit gets deleted, decoding:

❖ Guess 1: 8 4 7 10 ✓

❖ Guess 2: 8 4 7 10 ✗

❖ Guess 3: 8 4 7 10 ✗

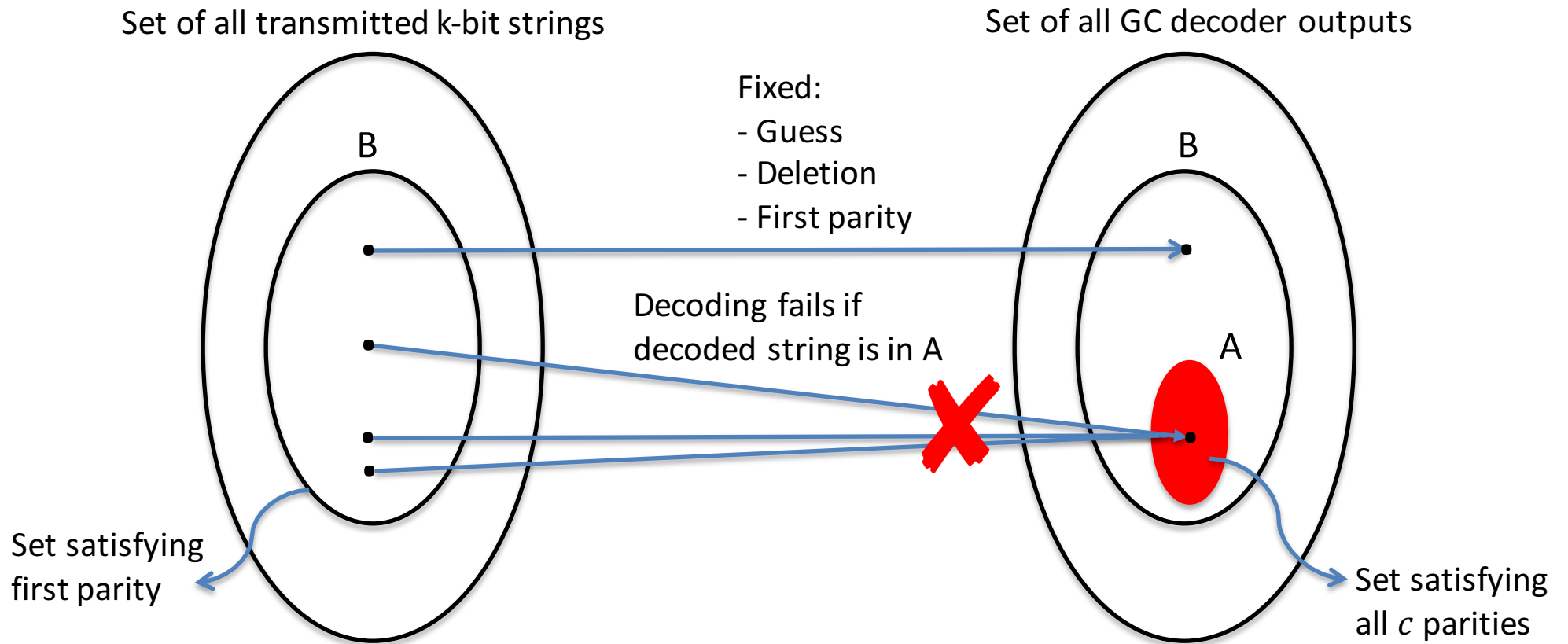
❖ Guess 4: 0 8 15 6 ✓

Guesses 1 & 4 satisfy
the 2 parities

- Probability of decoding failure for a given string: combinatorial problem that depends on the string and deletion position
- Proof approach: assume message is uniform iid, average over all possible messages

Decoding Failure – 1 Deletion

Assume WLOG that Guess 1 is **correct**, observe the output of decoder at **wrong** Guess $i \neq 1$



$$Pr(\text{decoding failure in guess } i \neq 1) = Pr(\text{decoded string is in } A)$$

Lemma: at most 2 different transmitted sequences can lead to the same decoded string in any Guess $i \neq 1$

Proof of Pr(F) for One Deletion

$$\Pr(F) \leq \Pr \left(\bigcup_{i=2}^{k/\log k} \{\mathcal{Y}_i \in A, \mathcal{Y}_i \neq \mathcal{Y}_1\} \right) \quad (1)$$

Union bound

$$\leq \sum_{i=2}^{k/\log k} \Pr(\mathcal{Y}_i \in A, \mathcal{Y}_i \neq \mathcal{Y}_1) \quad (2)$$

$$\leq \sum_{i=2}^{k/\log k} \Pr(\mathcal{Y}_i \in A) \quad (3)$$

$$= \sum_{i=2}^{k/\log k} \sum_{Y \in A} \Pr(\mathcal{Y}_i = Y) \quad (4)$$

Lemma

$$\leq \sum_{i=2}^{k/\log k} \sum_{Y \in A} \frac{2}{|B|} \quad (5)$$

$$= 2 \left(\frac{k}{\log k} - 1 \right) \frac{|A|}{|B|} \quad (6)$$

$$= 2 \left(\frac{k}{\log k} - 1 \right) \frac{q^{k/\log k - c}}{q^{k/\log k - 1}} \quad (7)$$

$$< \frac{2}{k^{c-2} \log k} \quad (8)$$

k : length of message

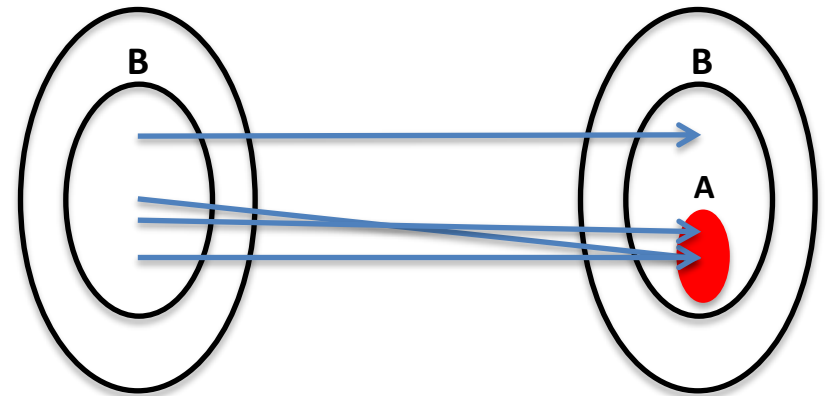
\mathcal{Y}_i : string decoded in Guess i

c : number of parities

A : set satisfying all c parities

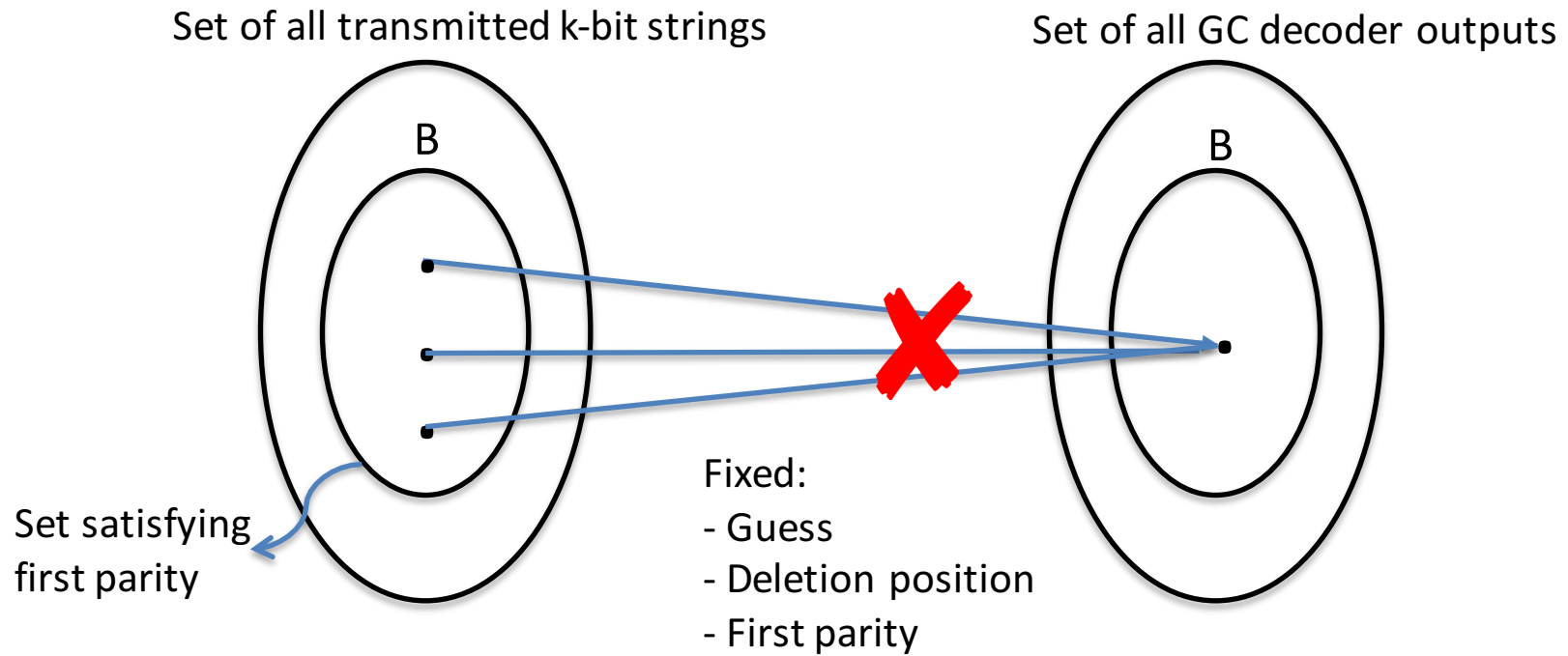
B : set satisfying first parity

q : field size



Claim

- Claim 1 (one deletion): at most 2 different transmitted strings can lead to the same decoded string in any wrong guess

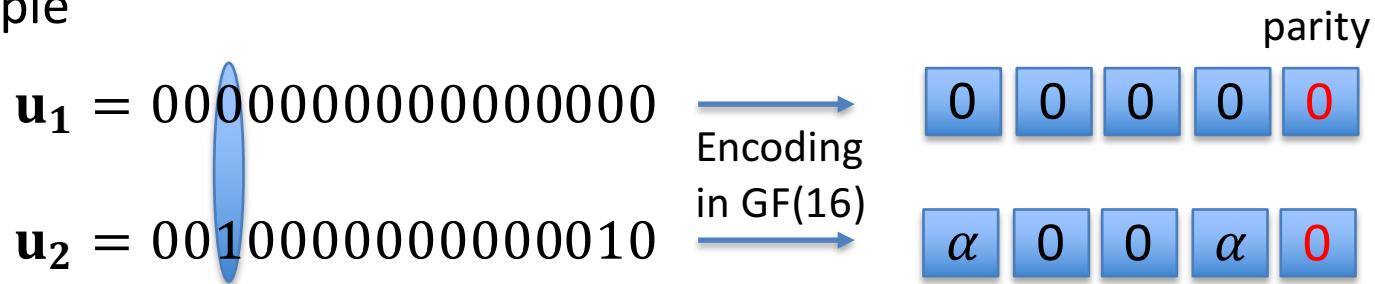


- Claim 2 (δ deletions): a **constant** number of different transmitted strings can lead to the same decoded string in any wrong guess

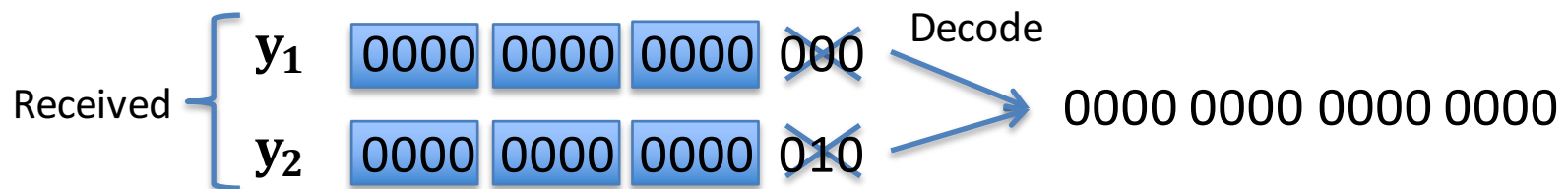
Claim - Example

- Claim 1 (one deletion): at most 2 different transmitted strings can lead to the same decoded string in any wrong guess

- Example



- 3rd bit deleted; Guess: deletion occurred in 4th block



- $\mathbf{u}_3 = 1111111111111111$ and $\mathbf{u}_4 = 0010000000000000$ ✗

- Two conditions: (1) Symmetry constraint; (2) Algebraic linear constraint

Claim

- Suppose 3rd bit is deleted, guess : deletion occurred in 4th block

$b_1 b_2 b_4 b_5$	$b_6 b_7 b_8 b_9$	$b_{10} b_{11} b_{12} b_{13}$	$b_{14} b_{15} b_{16}$
GF(17) $8b_1 + 4b_2 + 2b_4 + b_5$	$8b_6 + 4b_7 + 2b_8 + b_9$	$8b_{10} + 4b_{11} + 2b_{12} + b_{13}$?
Symbol 1	Symbol 2	Symbol 3	Erasure

- How many different messages can lead to same decoded string?
- **Symmetry constraint:** same decoded string \Rightarrow same bit values at positions of symbols 1, 2 and 3
- Bits which can be different: b_{14}, b_{15}, b_{16} and b_3 (deleted bit)
- Algebraic constraint: erasure is decoded using first parity

$$4b_{14} + 2(b_3 + b_{15}) + b_{16} = p_1$$

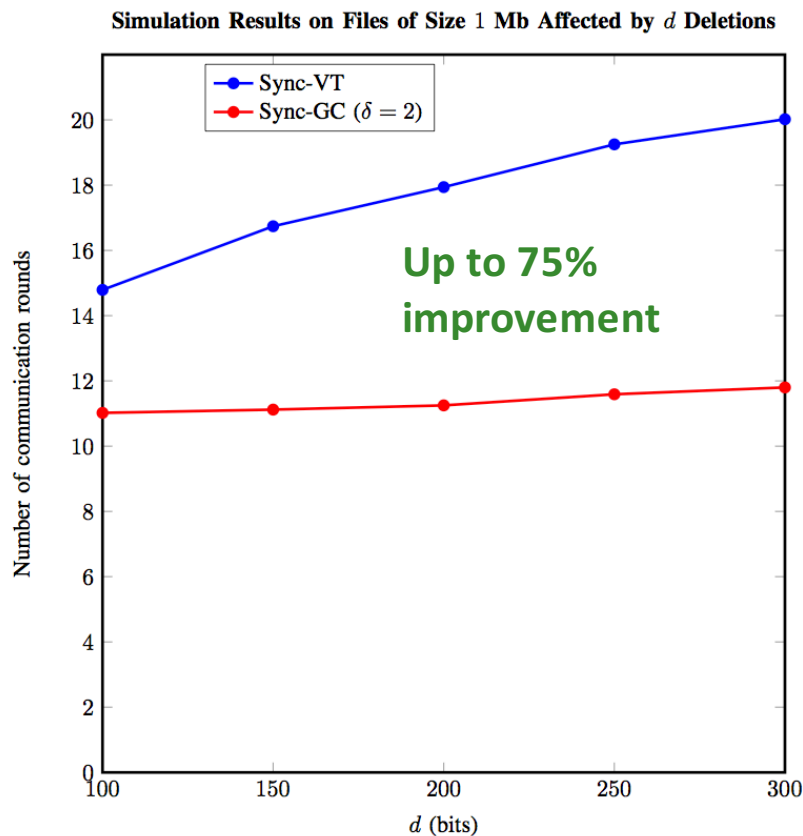
$$b_{14}, b_{15}, b_{16}, b_3 \in GF(2)$$

$$p_1 \in GF(17)$$

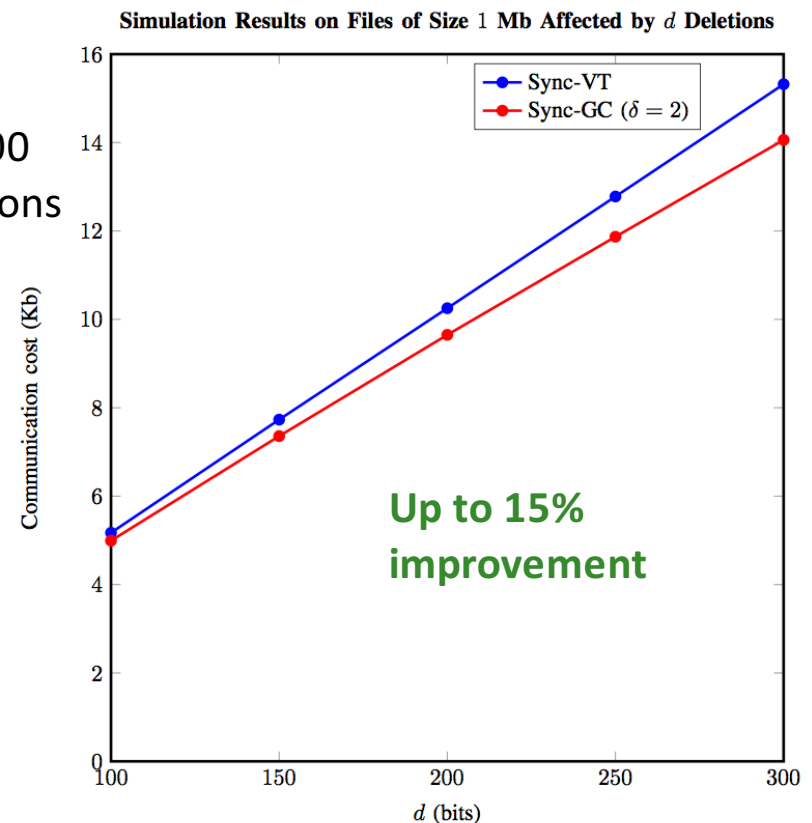
Equation has at most 2 solutions

Application to File Synchronization

- Interactive synchronization algorithm by [Venkataramanan et al. '15]
 - Isolate single deletions, use VT codes
 - Modification: isolate δ or fewer deletions, use GC codes
- Gain: (1) less communication rounds, (2) lower communication cost



N=1000
iterations



Summary

- Guess & Check Codes for localized deletions
 - For single or multiple windows
 - Explicit code construction with logarithmic redundancy
 - Deterministic polynomial time encoding and decoding
 - Asymptotically vanishing probability of decoding failure
- Open problems
- Capacity of deletion channel with localized deletions?
- Codes for adversarial localized deletions
- And of course for “unrestricted” deletion capacity and codes are still open problems

A word cloud of various languages expressing gratitude, including: THANK YOU, Merci, Obrigado, Takk, Danke, Kiitos, Efharisto, Grazie, and Tak. The words are arranged in a dense, overlapping pattern, with 'THANK YOU' and 'Merci' being the most prominent.

Questions?