

Some recent results on high rate local codes

Shubhangi Saraf

Rutgers

Joint works with

Sivakanth Gopi, Swastik Kopparty, Or Meir, Rafael Oliveira, Noga Ron-Zewi, Mary Wootters

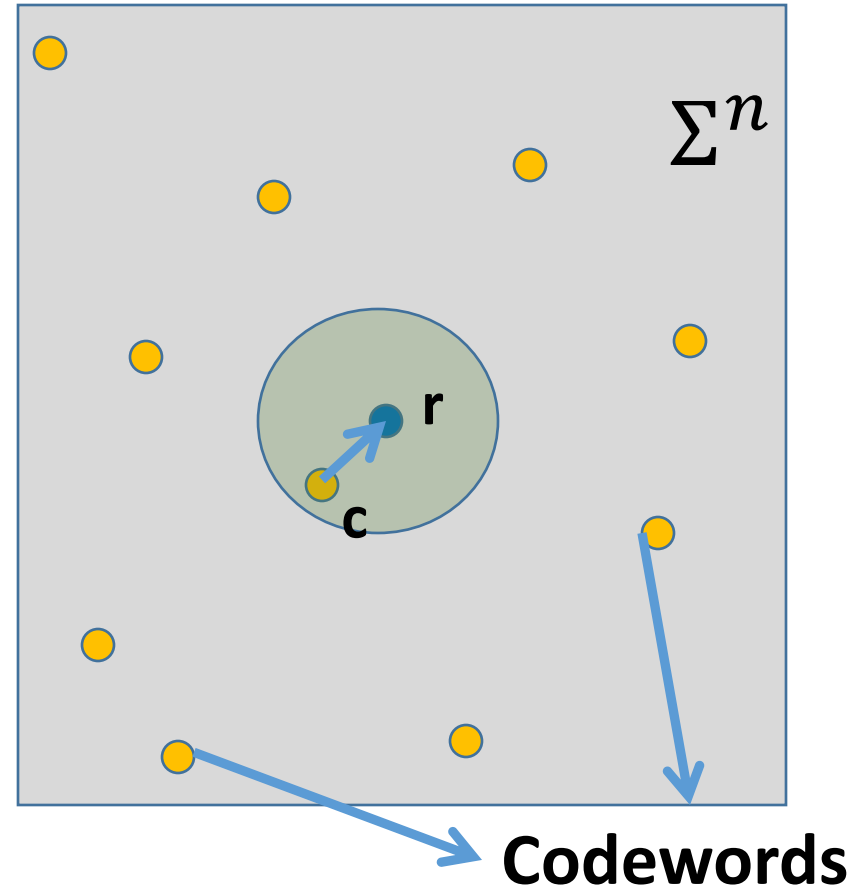
This talk

- Error-correcting codes with:
 - low redundancy
 - robust to large fraction of errors
 - *sublinear time* error-detection and error-correction algorithms

Error-correcting codes

- Alphabet Σ (often $\{0,1\}$)
- Encoding:
 - $E: \Sigma^k \rightarrow \Sigma^n$
 - Maps data to “codeword”
- Code $C = \text{Image}(E)$
 - ▶ Rate = k/n
 - ▶ (Hamming) Distance δ :
Any 2 codewords differ on at least δ fraction coordinates,

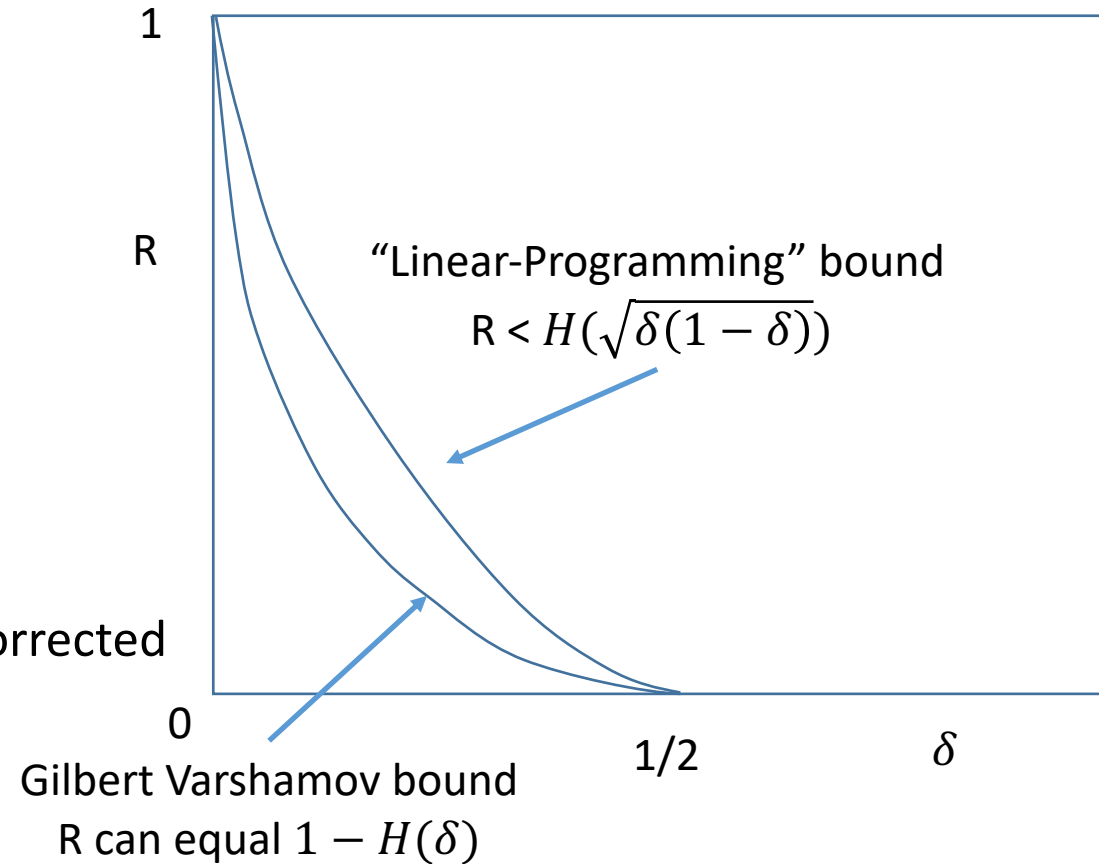
$\frac{\delta}{2}$ *fraction errors can be corrected*



Binary Error-correcting codes

- $C \subseteq \{0,1\}^n$ (with Hamming metric)
- Rate R :
 - $|C| = 2^{Rn}$
- Distance δ :
 - $\Delta(x, y) \geq \delta n$
for distinct $x, y \in C$
 - Implies $\delta/2$ -fraction errors can be corrected

- **Rate vs. Distance?**
 - **OPEN**



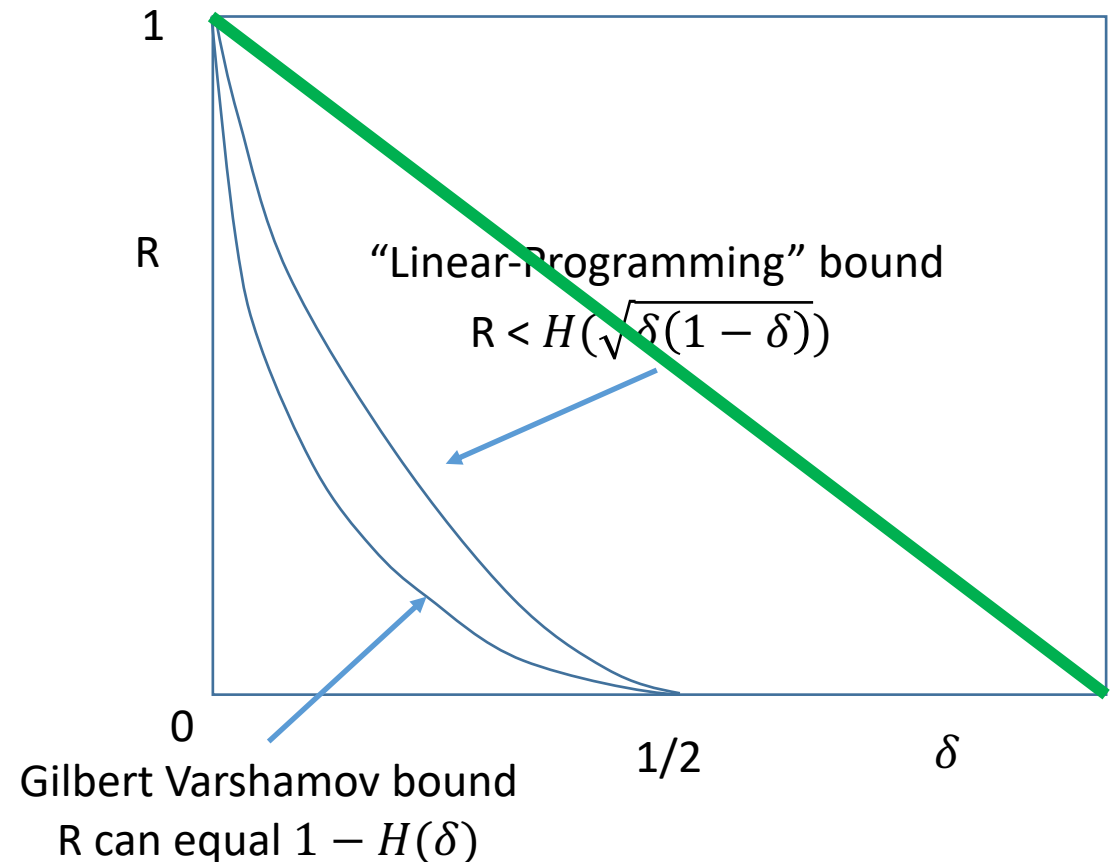
Gilbert Varshamov bound

- GV Bound: There exist codes with $R \geq 1 - H(\delta)$

Over large alphabets

$R = 1 - \delta$ is the optimal tradeoff
(a.k.a. SINGLETON BOUND)
Achieved explicitly

- Great open questions:
 - Is the GV bound tight?
 - Do there exist explicit codes meeting the GV bound?



Goals of classical coding theory

- Basic algorithmic tasks:
 - Encoding
 - Testing (error detection)
 - Decoding (error correction)
- Today we know codes with:
 - good rate-distance tradeoff
 - efficient encoding, testing, decoding
 - Linear/near-linear time

Local Codes

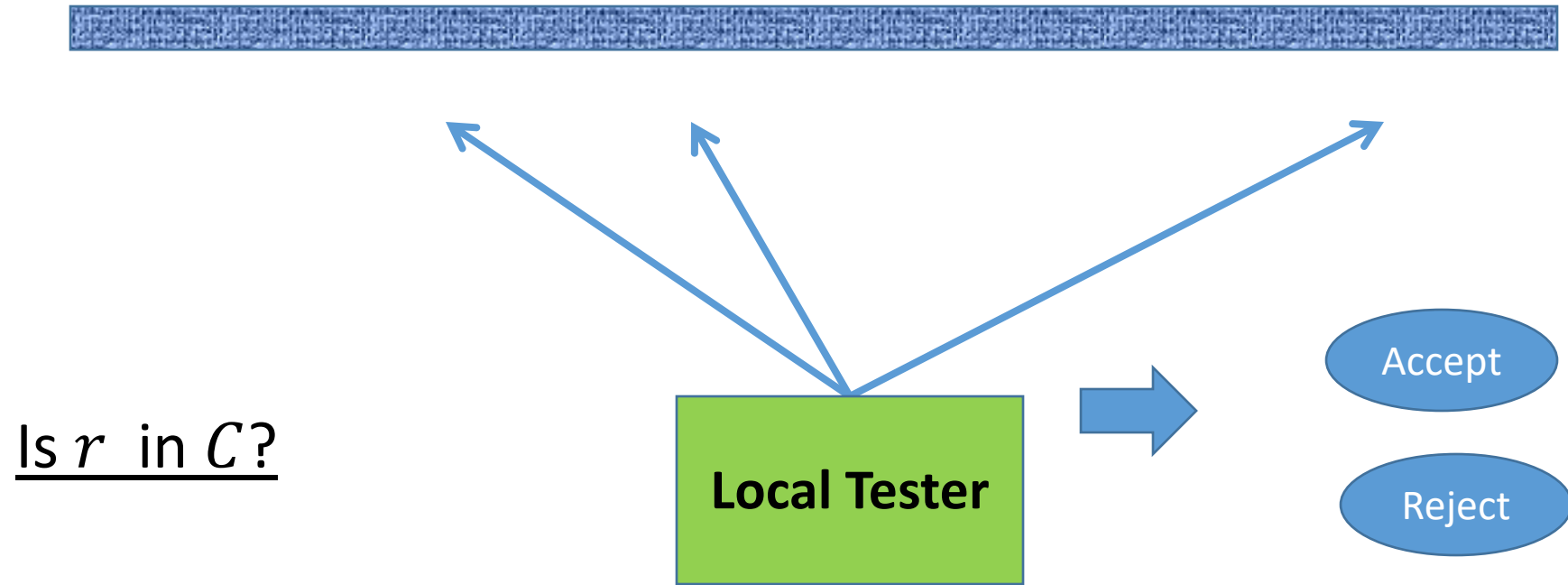
- Meanwhile, in early 90s complexity theory:
 - answers to questions that had never been asked
- Can we work with codes in sublinear time?
- In particular, what can we do with sublinear # queries?

Algorithmic Tasks associated with Error Correction

- **Error Detection:** Given $r \in \Sigma^n$, determine if $r \in \mathcal{C}$
 - Given $r \in \Sigma^n$, with sublinear queries to r , distinguish between $r \in \mathcal{C}$ and $\Delta(r, \mathcal{C}) > \epsilon n$
- **Error Correction:** Given $r \in \Sigma^n$, if $\exists m$ such that $\Delta(r, E(m)) < \epsilon n$, find m
 - Given $r \in \Sigma^n$ and $i \in [k]$ if $\exists m$ such that $\Delta(r, E(m)) < \epsilon n$, with sublinear queries to r find m_i

Locally Testable Code

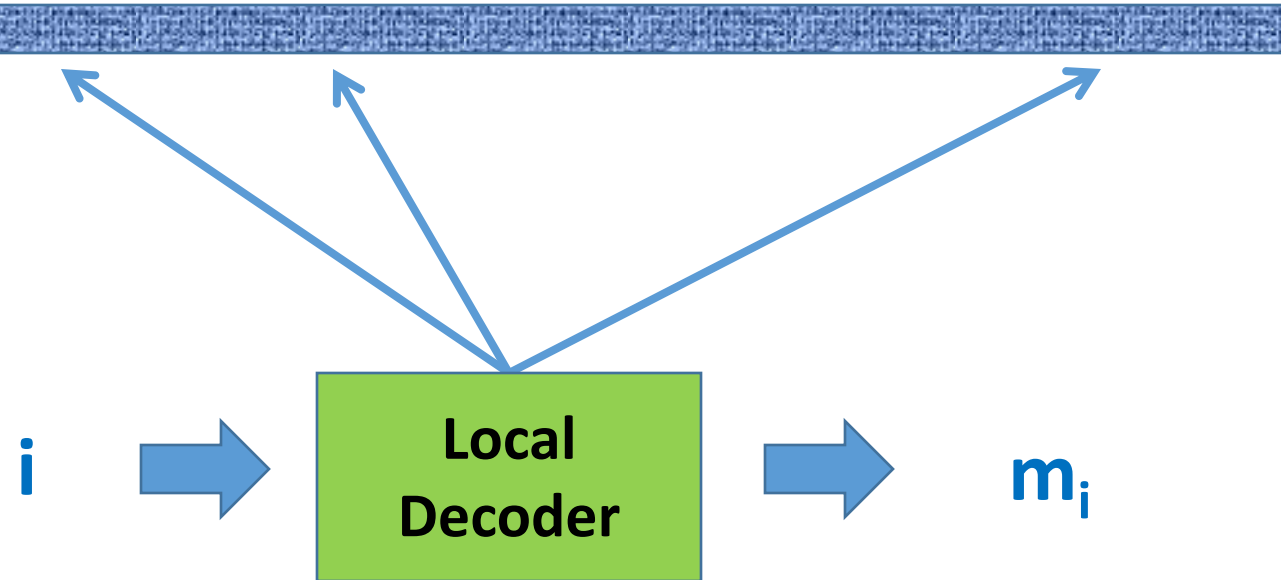
Given: $r \in \Sigma^n$



Locally Decodable Codes

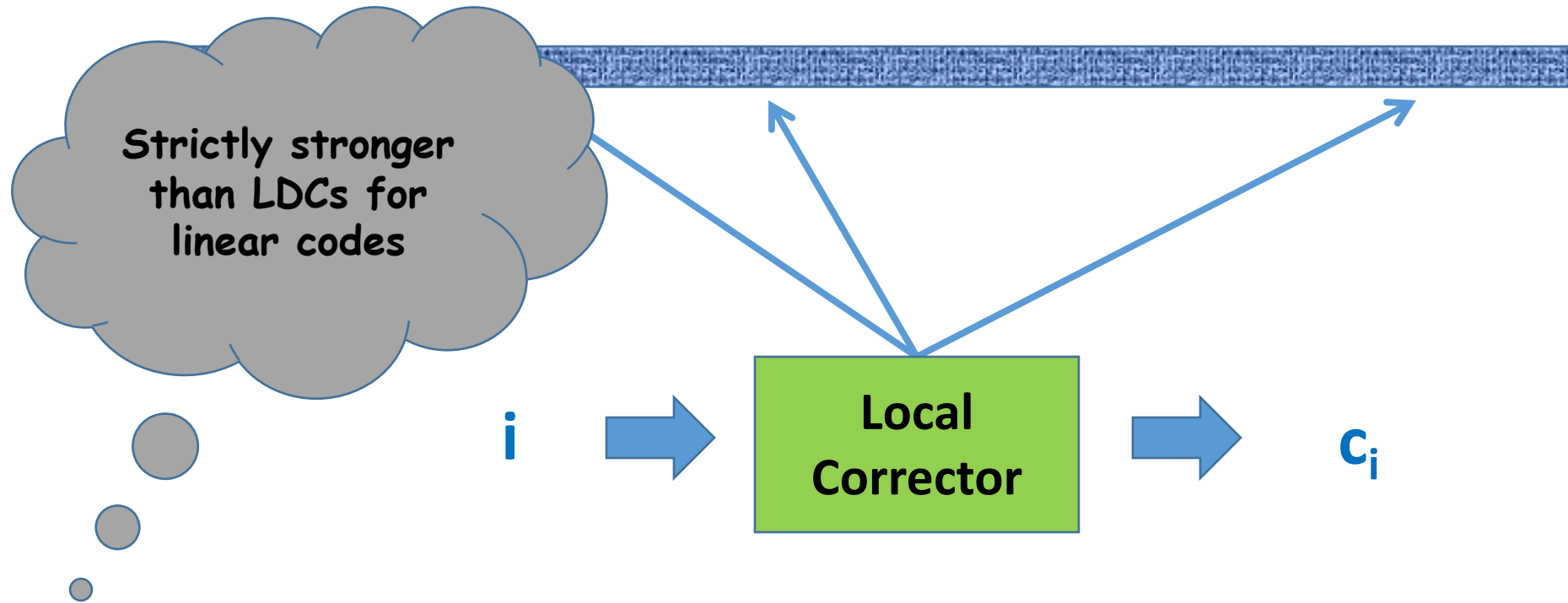
Given: $r \in \Sigma^n$ such that $\Delta(r, C) < \epsilon n$

Given: $i \in [k]$



Locally Correctable Codes

Given: $r \in \Sigma^n$ such that $\Delta(r, C) < \epsilon n$



Motivation for Local Decoding/Local Correcting

Many applications to cryptography and complexity theory

- Worst case to Average Case reductions
- Constructions of PRGs from One-Way functions
- Connections to Polynomial Identity Testing, Matrix Rigidity, Circuit Lower bounds
- Private information retrieval
- Learning theory

- Mathematically very interesting

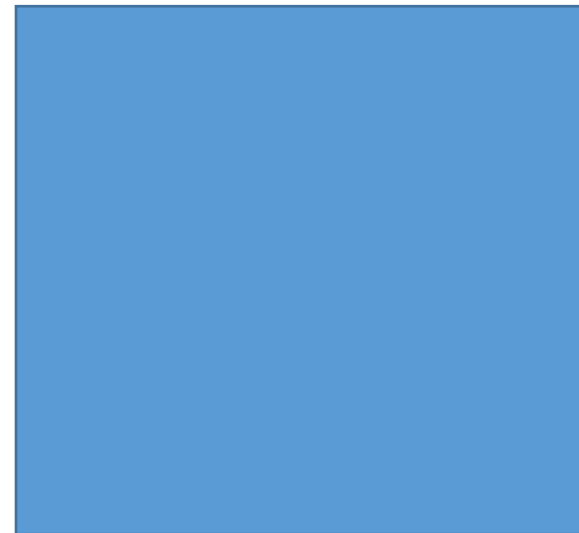
- Interesting for coding theory in practice?

Motivation for Local Testing

- Implicit connections to the PCP theorem
 - Advances have led to improved PCPs
 - Limitations should lead to an understanding of limitations of PCPs
- Applications to Unique Games conjecture and hardness of approximation
- Many relations to testing of functions
 - Original [Blum-Luby-Rubinfeld] linearity tester \approx testability of the Hadamard Code which led to the proof checking revolution

A nice local code

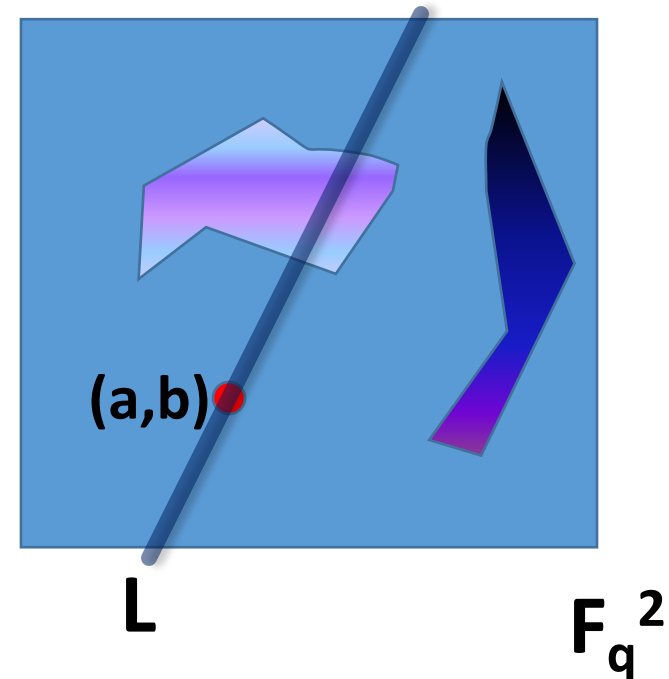
- Reed-Muller codes (multivariate polynomial evaluation codes)
 - constant rate, constant distance
 - $O(n^\epsilon)$ query locally testable
 - $O(n^\epsilon)$ query locally decodable
- Large finite field \mathbf{F}_q of size q
- Interpret original data as a polynomial $P(X,Y)$
 - $\text{degree}(P) = d = 0.1 q$
- Encoding:
 - Evaluate P at each point of \mathbf{F}_q^2
- Rate = $\Omega(1)$
- Distance = 0.9
 - Two low degree polynomials cannot agree on many points of \mathbf{F}_q^2



\mathbf{F}_q^2

Local testing/correcting RM codes

- Main idea:
 - Restricting a low-degree multivariate polynomial to a line gives a low-degree univariate polynomial
- Local testing:
 - Check that restriction to a random line is a low-degree univariate polynomial
 - Analysis highly nontrivial [Rubinfeld-Sudan + others]
- Local correcting:
 - To recover $P(a,b)$:
 - Pick random line L through (a,b)
 - Fit univariate polynomial through $r|_L$
 - Use it to recover value at (a,b)
- Query complexity
 - # points on a line = $q = O(\sqrt{n})$



Local codes of constant rate

- Reed-Muller codes (multivariate polynomial evaluation codes)
 - constant rate, constant distance
 - $O(n^\epsilon)$ query locally testable
 - $O(n^\epsilon)$ query locally decodable
- Since the 2010s, several improved codes:
 - Local testing:
 - tensor codes [BS, V], lifted codes [GKS]
 - Local decoding:
 - multiplicity codes [KSY], lifted codes [GKS], expander codes [HOW]
- rate $\rightarrow 1$, better rate vs. distance vs. queries

Plan of talk

- Survey of some known results
- [Kopparty-Meir-RonZewi-S `16]
 - High rate LTCs/LCCs with improved query complexity
- [Gopi-Kopparty-Oliveira-RonZewi-S `17]
 - LTCs and LCCs approaching* Gilbert-Varshamov bound
- [Kopparty-RonZewi-S-Wootters `18]
 - Capacity achieving locally list decodable codes
- Some proofs

Locally decodable/correctable codes

Two regimes

- **Low query regime:**

- Number of queries is small (2, 3, constant)
- What is the best rate?

- Theoretically very interesting
 - applications to *Cryptography, average-case complexity*
- Too inefficient for codes in practice

Extensively studied
Many deep and amazing
results (upper and lower
bounds)
Many basic problems
unanswered

- **High rate regime**

- Let the rate be high (constant rate or rate ≈ 1)
- What is the best query complexity that can be achieved?

- Focus of more recent work.
- Relevant regime for data storage and retrieval.
- Even mild lower bounds would have very interesting consequences to rigidity, lower bounds [Dvir]

Low Query Regime (LCCs, LDCs)

- $\ell = 2$: Hadamard Code is best possible $n = 2^{\Omega(k)}$ [Goldreich-Karloff]
- $\ell = 3$: $n = 2^{\sqrt{k}}$ (till not very long ago ...)
- For any constant ℓ : Reed Muller code best known construction: $n =$ (till not very long ago)

Matching Vector Codes:
LDCs with $n = \exp(\exp(o(\log k)))$
[Yekhanin, Efremenko, Dvir-Gopalan-Yekhanin]

- Lower bounds:

- $\ell = 3$: $n = \Omega(k^2)$ [Woodruff]
 - [Dvir-S-Wigderson] Over Real numbers, if code is linear then for LCCs $n = \Omega(k^{2+\epsilon})$
- General ℓ : $n \geq k^{1+\frac{1}{\ell}}$ (too inefficient for codes in practice)

Open question:

Can one get LDCs/LCCs with $O(1)$ queries and polynomial rate?

ong

High rate regime (LCCs, LDCs)

- Till about 8 years ago:
 - Reed-Muller codes were the only example
 - To get query complexity $\ell = k^\epsilon$, Rate $R = \exp\left(\frac{1}{\epsilon}\right)$
- More recently:
 - [Kopparty-S-Yekhanin '11] Multiplicity Codes
 - [Guo-Kopparty-Sudan '13] Lifted Codes
 - [Hemenway-Ostrovsky-Wooters '13] Expander based codes
 - Query complexity $\ell = k^\epsilon$, Rate $R = 1 - \epsilon$
(locally decodable and correctable from a constant fraction)
- [Katz-Trevisan]:
 - Constant rate \Rightarrow must have query complexity $\Omega(\log n)$

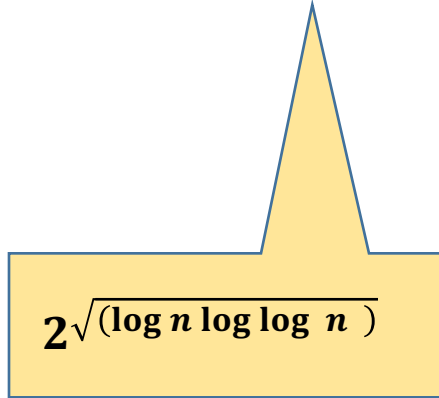
Interesting question:

What is the best rate/query complexity tradeoff?

Can one get LDCs/LCCs with rate $\Omega(1)$ or $1 - \epsilon$ and with query complexity $k^{o(1)}$

Somewhat recent result:

[Kopparty-Meir-RonZewi-S `16]: There exists a family of codes of rate $1 - \epsilon$ that is locally decodable and locally correctable with $n^{o(1)}$ queries from a constant fraction of errors.


$$2^{\sqrt{(\log n \log \log n)}}$$

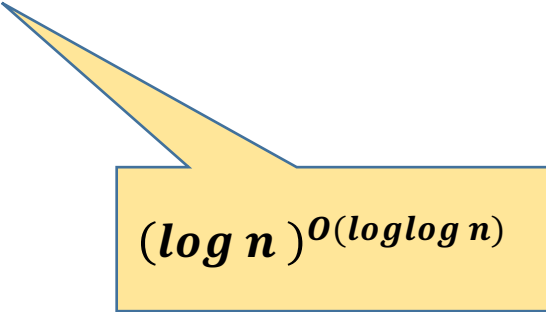
What we know about constant rate LTCs

- As far as we know,
 - there could be 3-query LTCs of constant rate
- RM codes achieve:
 - For all $R < 1/\exp(\frac{1}{\beta})$
 - Query complexity = $O(n^\beta)$
- Recent progress beyond Reed-Muller codes:
 - For all $R < 1$
 - For all $\beta > 0$
 - Query complexity = $O(n^\beta)$
 - Two families of codes achieving this!
 - Tensor codes [BenSasson-Sudan], [Viderman]
 - Lifted Reed-Solomon codes [Guo-Kopparty-Sudan]

Constructions known with 3-queries and Rate = $\frac{1}{\text{poly}(\log n)}$
[BenSasson-Sudan`05, Dinur`06]

More recently:

[Kopparty-Meir-RonZewi-S `16]: There exists a family of codes of rate $1 - \epsilon$ that are locally testable with $n^{o(1)}$ query complexity.


$$(\log n)^{o(\log \log n)}$$

KMRS Theorem for LCCs: There exists a family of codes of rate $1 - \epsilon$ that is locally decodable and locally correctable with $2^{\sqrt{(\log n \log \log n)}}$ queries from a constant fraction of errors

KMRS Theorem for LTCs: There exists a family of codes of rate $1 - \epsilon$ that is locally testable with $(\log n)^{O(\log \log n)}$ queries from a constant fraction of errors.

LTCs and LCCs approaching the GV bound

- Theorem [[Gopi-Kopparty-Oliveira-RonZewi-S `17](#)]

(informal) We can construct LTCs and LCCs which achieve the best possible rate-distance tradeoff that we know how to achieve with general (nonlocal) codes.

Main Result: LTCs

[Gopi-Kopparty-Oliveira-RonZewi-S `17]

Theorem:

For all R, δ with:

$$R < 1 - H(\delta)$$

there exists an infinite family of codes C_n

such that:

- $\text{length}(C_n) = n$
- $\text{Rate} \geq R$
- $\text{Distance} \geq \delta$
- C_n is locally testable with $(\log n)^{O(\log \log n)}$ queries

Local codes can be list decoded up to capacity

[Hemenway-RonZewi-Wootters`17, Kopparty-RonZewi-S-Wootters`18]

There exist codes that can be *locally list decoded* up to capacity

with query complexity $2^{(\log n)^{\frac{3}{4}}}$

[KMRS] result (and proof ideas) – an important ingredient in all these results.

Rest of talk – sketch of proof of KMRS result for LCCs

KMRS Theorem for LCCs: There exists a family of codes of rate $1 - \epsilon$ that is locally decodable and locally correctable with $2^{\sqrt{(\log n \log \log n)}}$ queries from a constant fraction of errors

KMRS Theorem for LTCs: There exists a family of codes of rate $1 - \epsilon$ that is locally testable with $(\log n)^{O(\log \log n)}$ queries from a constant fraction of errors.

Proof of KMRS result: *2 components*

- **Component 1:** High rate codes with *sub-polynomial query complexity* but only tolerating a tiny *sub-constant fraction of errors*
- **Component 2:** “Distance Amplification”
 - Takes code as above and transforms it to a code that can tolerate many more errors

Component 1

- High rate codes with *sub-polynomial query complexity* but only tolerating a tiny *sub-constant fraction of errors*

Can be achieved by Multiplicity Codes!

(In a regime of parameters not studied before)

Multiplicity Codes

[Kopparty-S-Yekhanin`11]

Theorem (original)

For every $\epsilon > 0$,
for inf. many k , there are codes encoding

k bits $\rightarrow (1 + \epsilon) k$ bits (symbols)
decodable in $O(k^\epsilon)$ time (+queries)
from $\delta(\epsilon) > 0$ fraction errors.

Theorem (sub-constant distance)

For every $\epsilon > 0$
for inf. many k , there are codes encoding

k bits $\rightarrow (1 + \epsilon) k$ bits (symbols)
decodable in $O(2^{\sqrt{\log k \log \log k}})$ time (+queries)
from $\approx \sqrt{(\log \log k) / \log k}$ fraction errors.

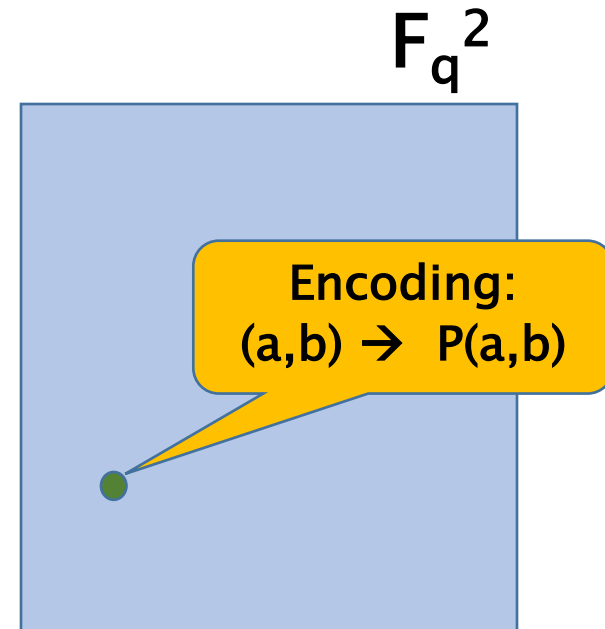
Construction of Mult. Codes

- Reed Muller Codes
- Augment it with “derivatives”

Reed-Muller Codes

Bivariate Reed-Muller

- Large finite field of size q
- Interpret original data as a polynomial $P(X,Y)$
 - $\text{degree}(P) \cdot d = (1 - \delta) q$
- *Encoding:* $\text{Enc}(P)$
 - At each point $(a,b) \in \mathbb{F}_q^2$, Evaluate $P(a,b)$



Key observations

- Schwartz-Zippel Lemma

- 2 polynomials of degree $< (1 - \delta)q$ differ on at least δ fraction of points

- So:

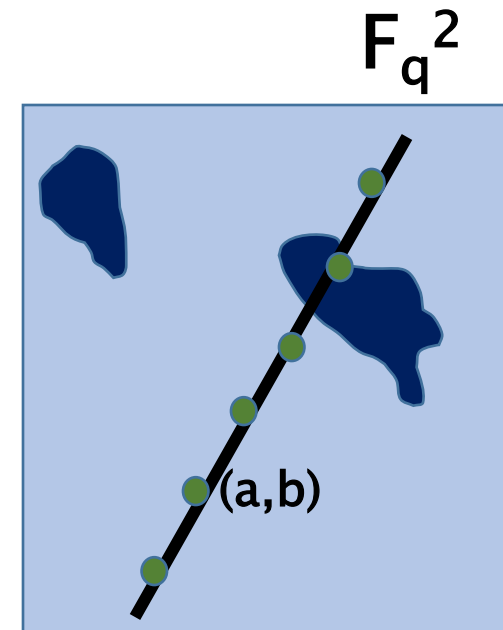
- Any two codewords are at least δn apart

Decoding Reed-Muller Codes

- **Given:**
 - noisy encoding of $P(X,Y)$
 - $\text{Deg}(P) = q(1 - \delta)$
 - point (a,b) in \mathbb{F}_q^2
- **Goal:** recover $P(a,b)$

Algorithm

- Take random line L through (a,b)
- Query points on L
 - Should have small error
 - Noisy encoding of $P|_L$ (univariate polynomial)
- Recover $P|_L$
 - “Reed Solomon” decoding
- Compute $P|_L(a,b)$
= $P(a,b)$



Parameters of Reed-Muller Codes

- Bivariate Reed Muller:

- $k = (d+2) \text{ choose } 2 \approx \frac{(1-\delta)^2 q^2}{2}$

- $n = q^2$

- **Rate** $\approx \frac{1}{2} - \delta$

- **# Queries:** $\ell \approx O(k^{1/2})$

- Improve query complexity \rightarrow increase # of variables

More variables

- Polynomials of deg $\cdot (1-\delta) q$ in ***m variables***
- $k = (d+m)$ choose $m \approx \frac{(1-\delta)^m q^m}{m!}$
- $n = q^m$
- **Rate** $\approx \frac{(1-\delta)^m}{m!}$
- **Queries** $= q \approx n^{1/m} \approx \mathbf{O(k^{1/m})}$
- Decodable from $\Omega(\delta)$ errors
- Bottleneck for rate: Degree needs to be small

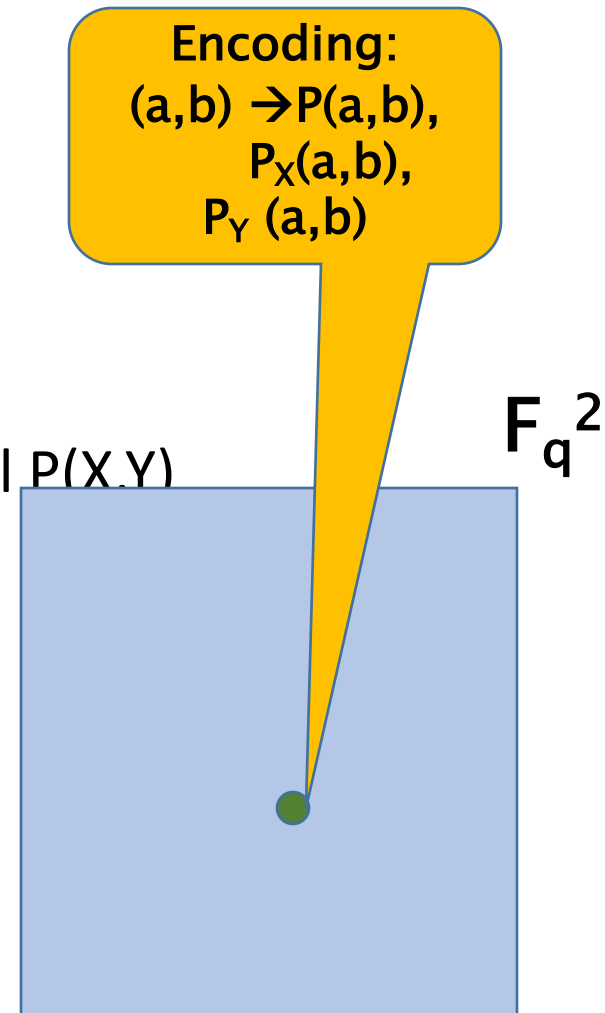
Multiplicity Codes

- Key idea: Derivatives
- Higher degree polynomials
 - (too high for Reed-Muller)

Multiplicity Codes

Bivariate Multiplicity codes

- Large finite field of size q
- Interpret original data as a (high) degree polynomial $P(X,Y)$
 - degree(P): $d = 2 \times (1 - \delta) q$
- *Encoding:* $\text{Enc}(P)$
 - At each point $(a,b) \in \mathbb{F}_q^2$, evaluate:
 - $\langle P(a,b), P_X(a,b), P_Y(a,b) \rangle$



Schwartz–Zippel with Multiplicities [Dvir–Kopparty–S–Sudan’10]

- 2 polynomials of degree $< 2q(1-\delta)$ cannot agree on their evaluations and evaluations of derivatives in more than $(1-\delta)$ fraction points
- # roots of P *counted with multiplicity* $\cdot \deg(P) \leq |F|^{n-1}$
- Multiplicity Codes have good distance

Decoding Multiplicity Codes

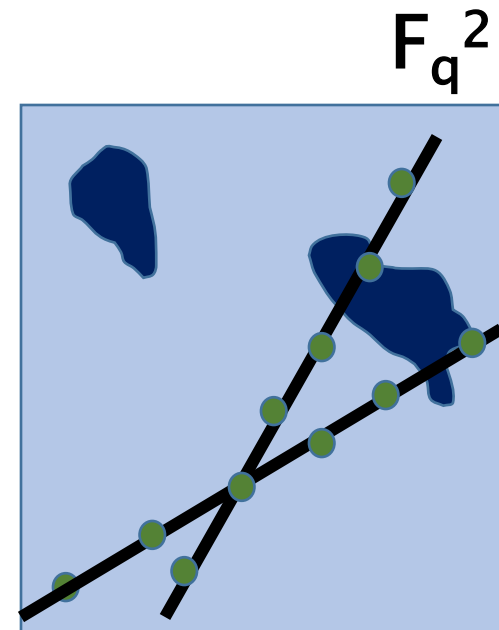
Given:

- noisy encoding of $\langle P, P_x, P_y \rangle$
 - $\text{Deg}(P) = 2 \times q(1-\delta)$
- point (a,b) in F_q^2

Goal: recover $\langle P(a,b), P_x(a,b), P_y(a,b) \rangle$

Algorithm

- Take random line L through (a,b)
 - Should have small error
- Query points on L
 - P_x, P_y give directional derivative of P along L
 - Noisy encoding of $P|_L$ (univariate polynomial), and of $\text{der}(P|_L)$
- Recover $P|_L$
- Repeat above steps
- We thus know $P(a,b), \text{der}(P|_{L_1})(a,b), \text{der}(P|_{L_2})(a,b)$
- This gives us $P(a,b), P_x(a,b), P_y(a,b)$



Parameters of Multiplicity Codes

- Bivariate Multiplicity Codes **of order 2**:
 - $k = (d+2) \text{ choose } 2 / 3 \approx (2(1-\delta)q)^2 / 6$
 - $n = q^2$
 - **Rate $\approx 2/3 - \delta$**
 - **# Queries: $\approx O(k^{1/2})$**
- Improve Rate \rightarrow increase order of derivatives
- Improve query complexity \rightarrow increase # variables


More variables, many derivatives

- **m – variate, derivatives up to order s**
- Polynomials of **degree $(1-\delta)sq$**
 - ▶ **Query Complexity: $\approx k^{1/m}$**
- **Rate $\approx (s/m+s)^m \times (1-\delta)^m$**
 - so if $s \gg m$, rate $\rightarrow 1$
- Decoding as before ...
 - (+ some “robustification”)

Reed-Muller Codes

- Messages: Low degree polynomials
- Encoding: Evaluation of polynomial on full domain
- #queries: Decreases with increase in # variables
- Rate: Decreases exponentially with increase in #vars

Multiplicity Codes

- Messages: High degree polynomials
- Encoding: Evaluation of polynomial *and its derivatives* on full domain
- #queries: Decreases with increase in # variables
- Rate: **1**


Multiplicity codes in low distance regime

To make *queries sub-polynomial*, choose m to be super-constant. For *constant rate* this *forces distance to be sub-constant*.

Theorem (sub-constant distance)

For every $\epsilon > 0$

for inf. many k , there are codes encoding

k bits $\rightarrow (1 + \epsilon) k$ bits (symbols)

decodable in $O(2^{\sqrt{\log k \log \log k}})$ time (+queries)

from $\approx \sqrt{(\log \log k) / \log k}$ fraction errors.

Component 2

- **Distance amplification**

- Similar technique used by [Alon-Luby'96] and then by others [GI'05, GR'08]

Theorem (sub-constant distance)

For every $\epsilon > 0$

for inf. many k , there are codes encoding

k bits $\rightarrow (1 + \epsilon) k$ bits (symbols)

decodable in $O(2^{\sqrt{\log k \log \log k}})$ time (+queries)

from $\approx \sqrt{(\log \log k) / \log k}$ fraction errors.

Component 2

- **Distance amplification**

- Similar technique used by [Alon-Luby'96] and then by others [GI'05, GR'08]

Theorem (~~sub-constant distance~~)

For every $\epsilon > 0$

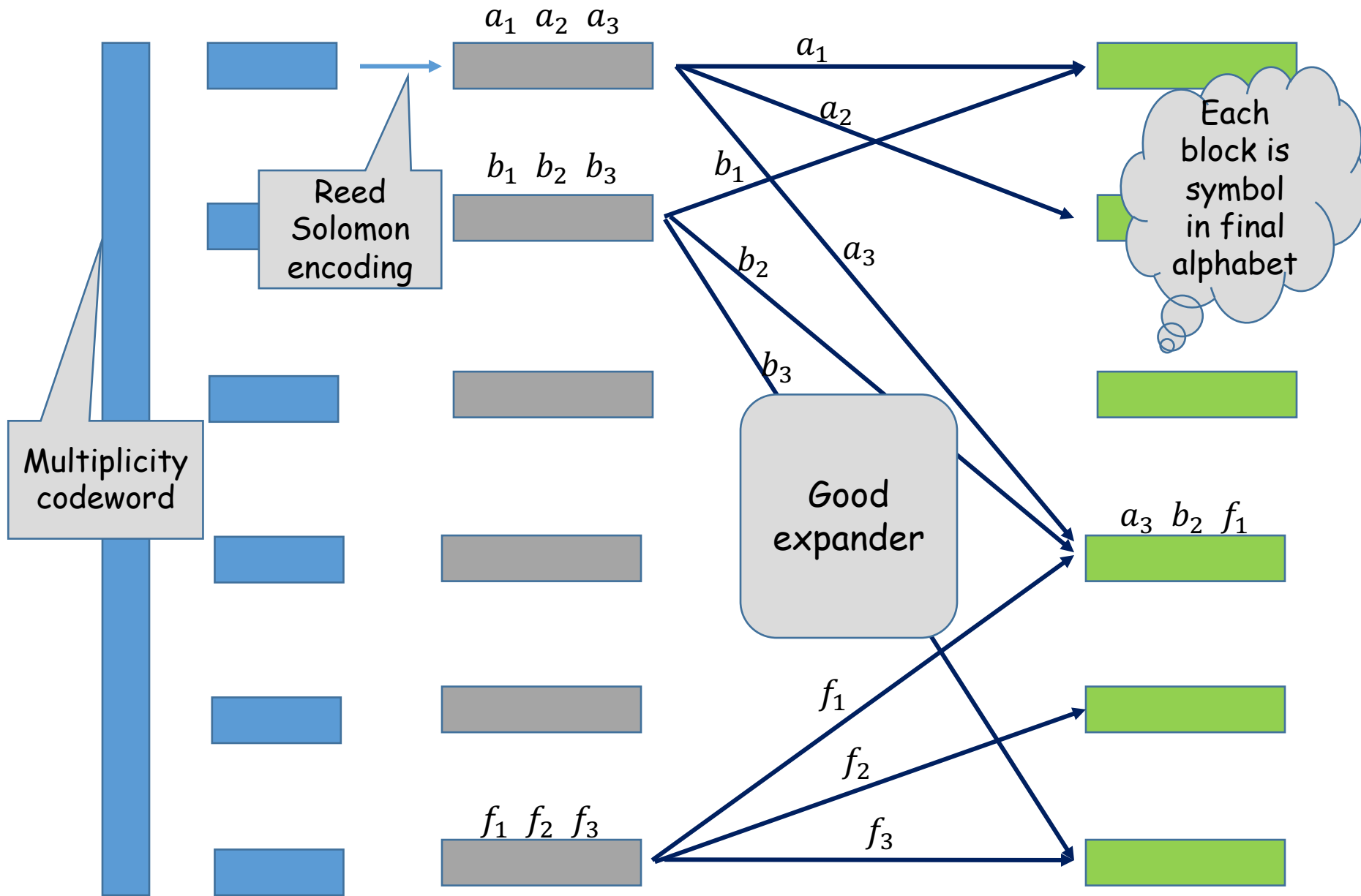
for inf. many k , there are codes encoding

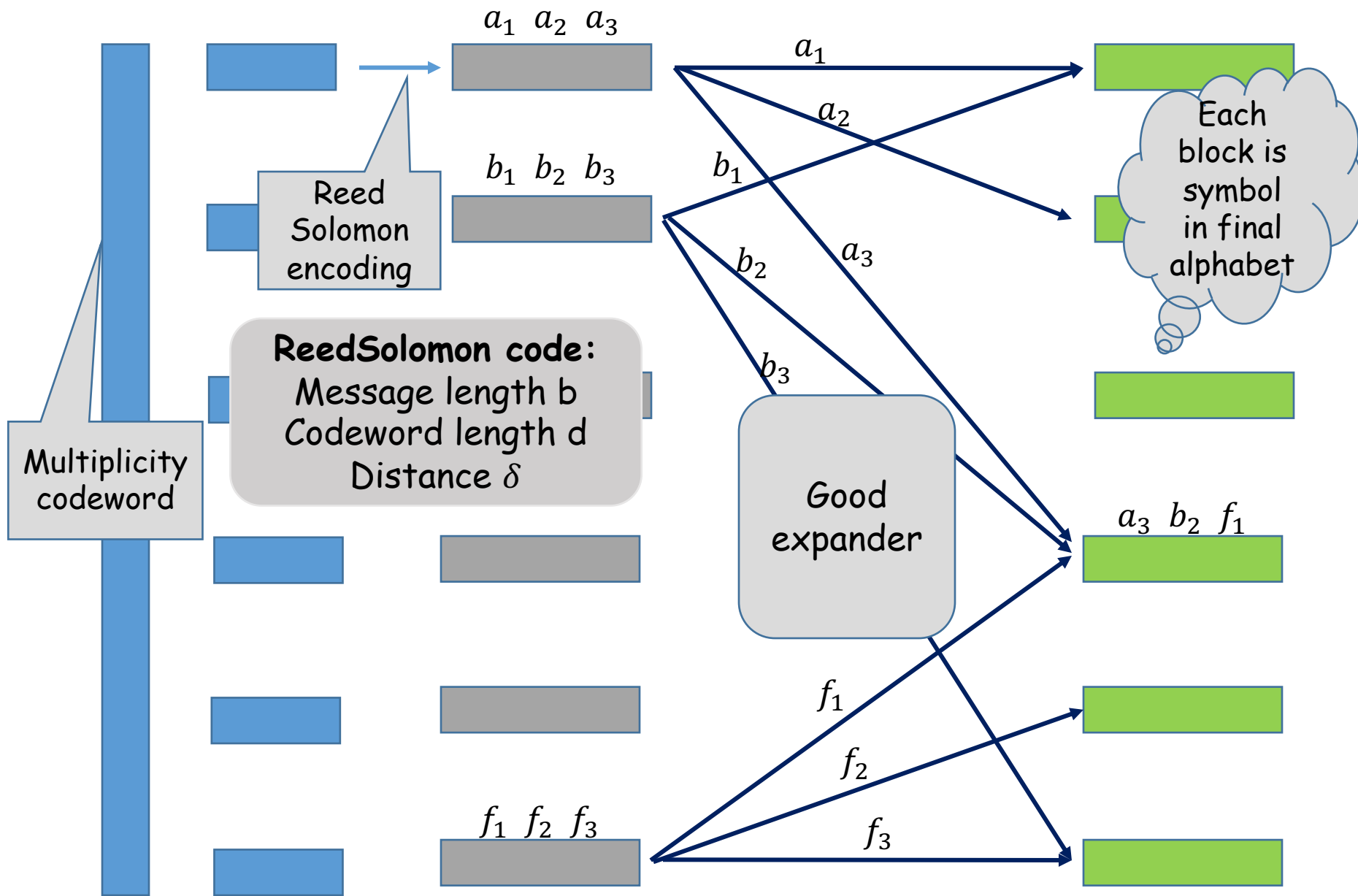
k bits $\rightarrow (1 + 2\epsilon) k$ bits (symbols)

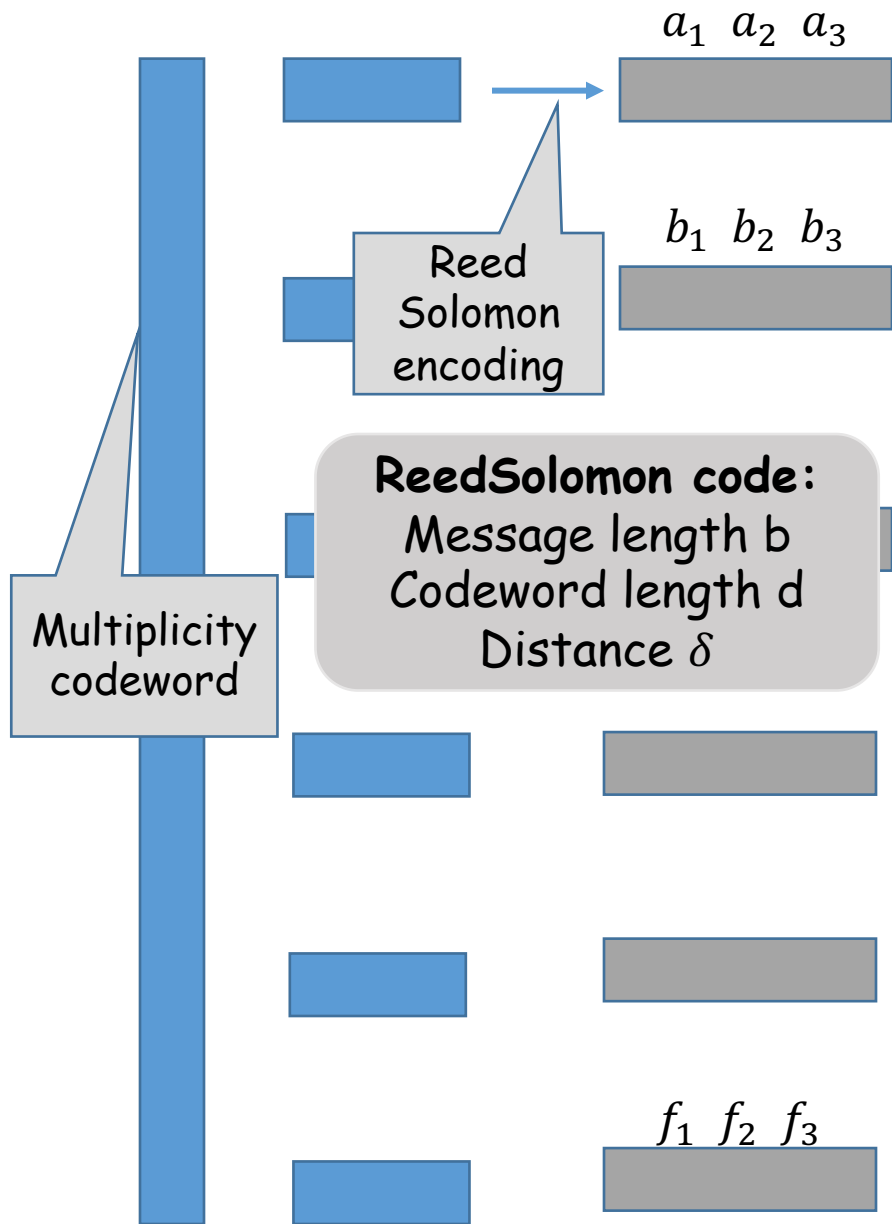
decodable in $O(2^{2\sqrt{\log k \log \log k}})$ time (+queries)

from $\approx \sqrt{\frac{\log \log k}{\log k}}$ fraction errors.

$\Omega(1)$







Decoding from random errors:

Suppose $\frac{\delta}{2} - \epsilon$ fraction of random errors

Most $(1-o(1))$ grey blocks have at most $\frac{\delta}{2}$ corruptions

Those Reed-Solomon codewords can be correctly decoded

Thus $1-o(1)$ fraction of the blue blocks can be correctly recovered.
 This is low enough error for multiplicity codes to handle

Everything can be done locally

Decoding from adversarial errors:

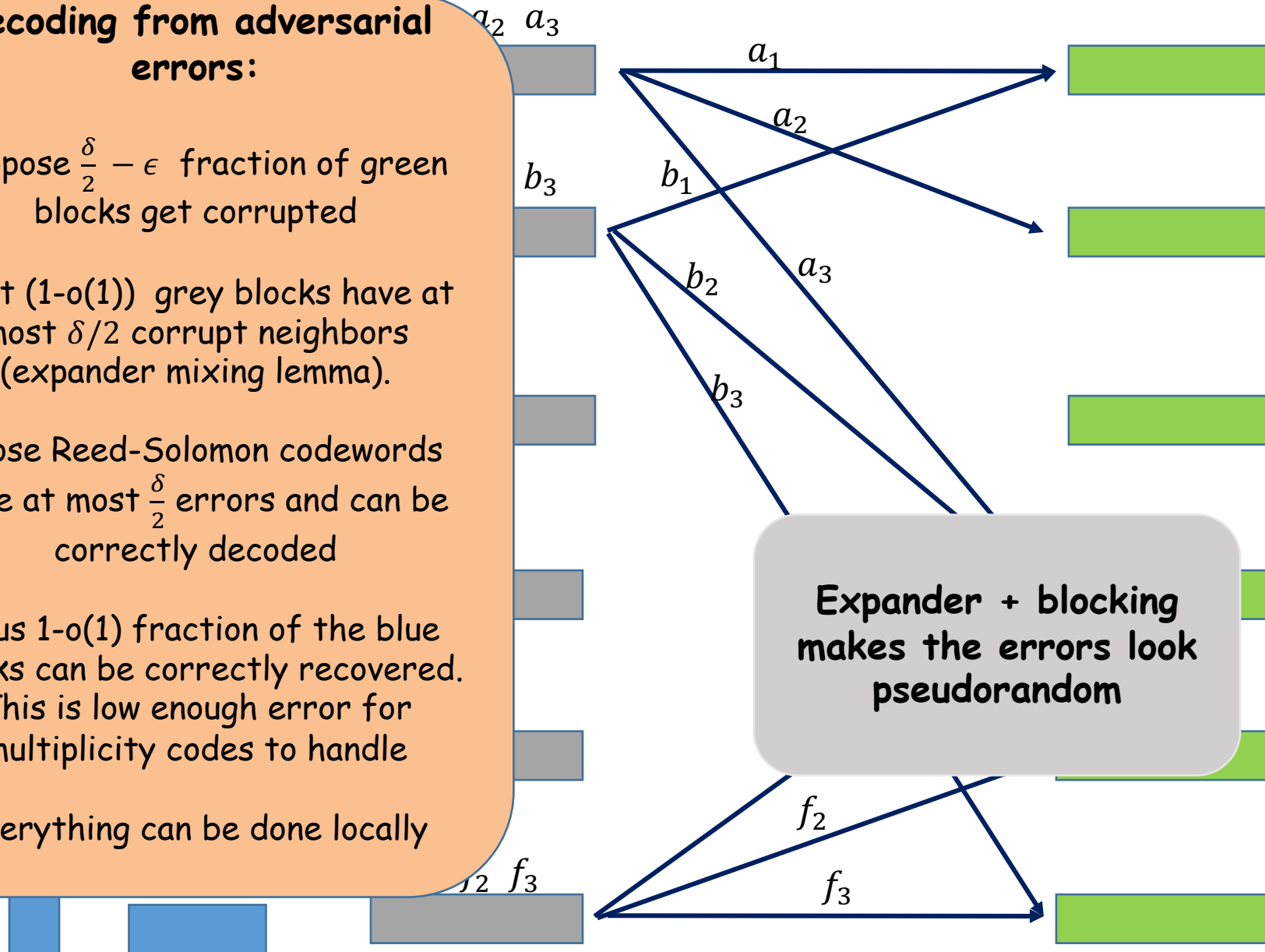
Suppose $\frac{\delta}{2} - \epsilon$ fraction of green blocks get corrupted

Most $(1-o(1))$ grey blocks have at most $\delta/2$ corrupt neighbors (expander mixing lemma).

Those Reed-Solomon codewords have at most $\frac{\delta}{2}$ errors and can be correctly decoded

Thus $1-o(1)$ fraction of the blue blocks can be correctly recovered. This is low enough error for multiplicity codes to handle

Everything can be done locally



Open questions

- Best possible query complexity for high rate LDCs and LTCs?
 - LTCS – potentially high rate 3 query LTCs!
 - LDCs/LCCs – potentially high rate $\log n$ query LCCs
- Explicit codes meeting the GV bound?
 - Almost solved by Ta-Shma!
- Is the GV bound tight?

Thanks!