

# CS250/EE387 - LECTURE 14 - LOCALITY!

## AGENDA

- ① LOCALLY CORRECTABLE CODES
- ② RM CODES as LCCs
- ③ HIGH-RATE LCCs [sketch]

## TODAY'S OCTOPUS FACT

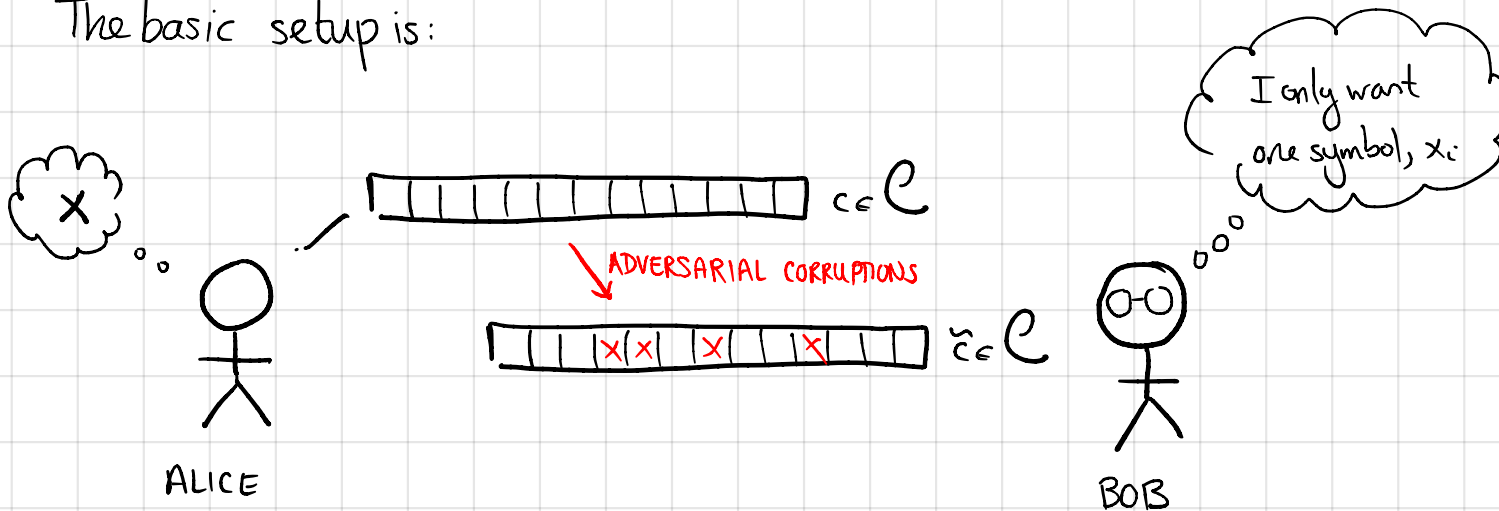
The mimic octopus not only camouflages itself to blend in with its surroundings, but it also can mimic other species, like a lion fish, sea snake, jelly fish or zebra sole.



No octopus here...

Today we will talk about **LOCALLY DECODABLE CODES**.

The basic setup is:



If Bob only wants one symbol of Alice's message (or her codeword), then he **COULD** decode the whole thing and figure out  $x_i$ .

But that seems wasteful...

The idea of **LOCAL DECODING** is to allow Bob to figure out  $x_i$  in **SUBLINEAR TIME**. In particular, he won't even have enough time to look at all of  $\tilde{c}$ !

# ① LOCALLY CORRECTABLE CODES.

Let us try to formalize this goal:

NOT THE CORRECT DEF.

$\mathcal{C} \subseteq \mathbb{F}_q^n$  is a  $(\delta, Q)$ -LOCALLY CORRECTABLE CODE (LCC) if there is an algorithm  $A$  so that the following holds:

for all  $w \in \mathbb{F}_q^n$  so that  $\exists c \in \mathcal{C}$  s.t.  $\Delta(c, w) \leq \delta n$ , and  $\forall i \in [n]$ ,

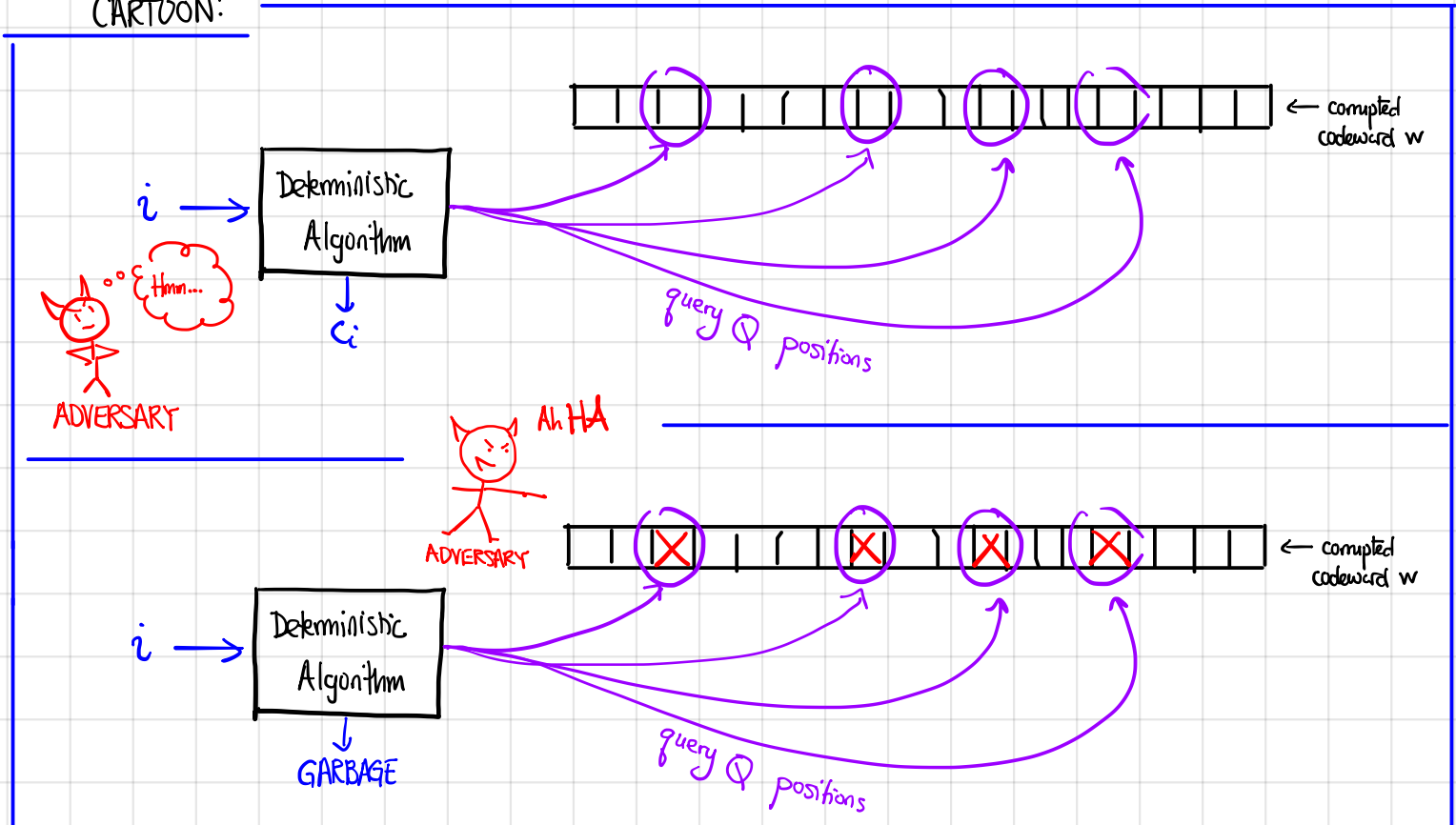
$A^{(w)}(i)$  makes at most  $Q$  queries to  $w$  and returns  $c_i$ .

input is  $i$   
 $A$  has oracle access to  $w$

Does this make sense? NO.

Sure, it parses, but if  $Q = o(n)$  then this is a vacuous definition.

CARTOON:



That is, if the queries are deterministic, and  $Q < \delta n$ , then the adversary can COMPLETELY mess up the algorithm's view.

Instead, we will need to RANDOMIZE the queries if we want to deal with an adversary.

DEF.  $\mathcal{C} \subseteq \mathbb{F}_q^n$  is a  $(\delta, Q, \gamma)$ -LOCALLY CORRECTABLE CODE (LCC) if there is a randomized algorithm  $A$  so that the following holds:

for all  $w \in \mathbb{F}_q^n$  so that  $\exists c \in \mathcal{C}$  s.t.  $\Delta(c, w) \leq \delta n$ , and  $\forall i \in [n]$ ,

- $A^{(w)}(i)$  makes at most  $Q$  queries to  $w$   
input is  $i$   
 $A$  has oracle access to  $w$
- $A^{(w)}(i) = c_i$  with probability at least  $1 - \gamma$ .

OTHER NOTIONS of LOCALITY:

- If you only want to recover a MESSAGE SYMBOL  $x_i$  instead of a CODEWORD SYMBOL  $c_i$ , it's called a LOCALLY DECODABLE CODE.
- If there's no adversary and you just want to be able to recover any symbol in SOME local way (not including that symbol) it's a LOCALLY REPAIRABLE CODE. (or LOCALLY RECOVERABLE CODE)
- Also: REGENERATING CODES, LOCALLY TESTABLE CODES, RELAXED LCCs, MAXIMALLY RECOVERABLE CODES, ...

## BRIEF LIT. REVIEW on LCC's.

$Q$	$n$ (as a function of $k$ )	Comments
2	$n = \Theta(2^k)$	Matching upper + lower bounds here.
3	$k^2 \leq n \leq \exp(\exp(\lg \lg(k)^{0.99}))$	The upper bd is actually an LDC
$O(\log(n))$	$k \leq n \leq \text{poly}(k)$	
$O(n^\epsilon)$	$k \leq n \leq (1+\alpha)k$ for any $\alpha > 0$	

Today we'll see how RM codes fit in, starting at  $Q=2$  and ending at  $Q=n^\epsilon$ .

② RM codes as LCCs. First let's recall the def. of REED-MULLER CODES:

Recall that  $\mathbb{F}_q[X_1, \dots, X_m]$  is the space of  $m$ -variate polynomials over  $\mathbb{F}_q$ .

The (total) DEGREE of a monomial  $X_1^{i_1} X_2^{i_2} \dots X_m^{i_m}$  is  $\sum_{j=1}^m i_j$ .

The DEGREE of  $f \in \mathbb{F}_q[X_1, \dots, X_m]$  is the largest degree of any monomial in  $f$ .

DEF. The  $m$ -VARIATE REED-MULLER CODE of DEGREE  $r$  over  $\mathbb{F}_q$  is

$$\text{RM}_q(m, r) = \left\{ (f(\alpha_1), \dots, f(\alpha_m)) : f \in \mathbb{F}_q[X_1, \dots, X_m], \deg(f) \leq r \right\}$$

REMARK. Note that we may assume that each  $X_i$  has degree  $< q$ , since  $\alpha = \alpha^q$  for all  $\alpha \in \mathbb{F}_q$ .

We saw BINARY RM CODES back in Lecture 6 when we were trying to figure out how to get good binary codes.

Let's start with  $RM_2(m, 1)$ : that is, codewords are just the evaluations of LINEAR polynomials

$$f(x_1, x_2, \dots, x_m) = \sum_i a_i x_i$$

This is also called the HADAMARD CODE.  
You saw it on your HW.

NOTE: Technically for  $RM_2(m, 1)$  we should also have a constant term here, but it will be convenient for us to ignore it...

Consider the following algorithm for locally decoding the Hadamard code:

ALG. Input: Query access to  $g: \mathbb{F}_2^m \rightarrow \mathbb{F}_2$  s.t.  $\Delta(g, f) < 2^{m-2}$  for some  $f \in RM_2(m, 1)$ ,  
and an index  $\alpha \in \mathbb{F}_2^m$   
Output: A guess for  $f(\alpha)$

Choose  $\beta \in \mathbb{F}_2^m$  at random.  
RETURN  $g(\beta) + g(\beta + \alpha)$

CLAIM:  $RM_2(m, 1)$  is a  $(\delta, 2, 1-2\delta)$ -LCC for any  $\delta < 1/4$ .

proof.

If  $g(\beta) = f(\beta)$  and  $g(\beta + \alpha) = f(\beta + \alpha)$  (\*)

then  $g(\beta) + g(\beta + \alpha) = f(\beta) + f(\beta + \alpha) = f(\alpha)$  since  $\deg(f) = 1$ .

← Here I'm using the assm that  $f$  has no constant term.

(\*) happens with probability  $\geq 1 - 2\delta$ , since

$\mathbb{P}\{g(\beta) \neq f(\beta)\} = \mathbb{P}\{g(\beta + \alpha) \neq f(\beta + \alpha)\} \leq \delta$ , since  $\beta$  and  $\alpha + \beta$  are both uniformly random.

(Notice that they are NOT jointly uniform, but each marginal is uniform).

GREAT! Now we have a 2-query LCC. But the rate is not great:  $m/2^m$ .

QUESTIONS:

① Can we do better for  $Q=2$ ?

NO. See [Kerenidis+Wolf]

② What if  $Q = \omega(1)$ ?

YES! Coming up next.

②B  $Q = \log(n)$

We'd like to use the same idea, but there's a problem.

If our strategy is "hope that our  $\log(n)$  queries completely avoid the errors," we'll be in trouble. Indeed, w.h.p. there will be about  $\delta \cdot \log(n)$  errors in our  $\log(n)$  queries.

The idea will be to make our queries themselves somewhat robust to error.

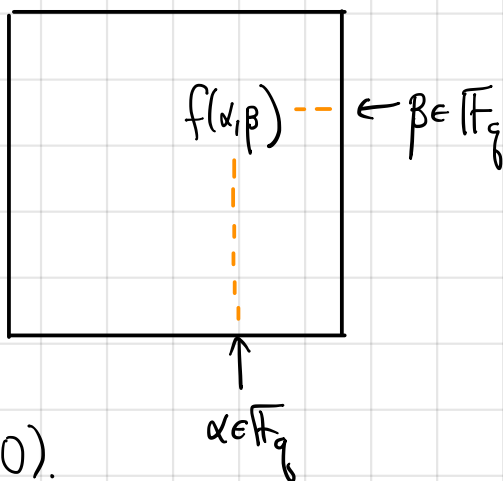
For motivation, consider  $RM_q(z, r)$ .

That is, the codewords of  $RM_q(z, r)$  are evaluations of bivariate polynomials

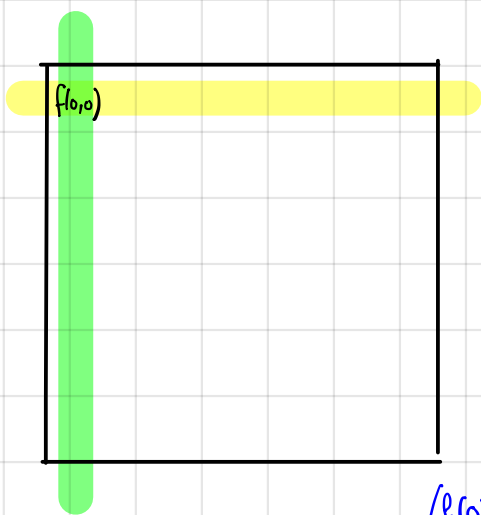
$$f(X, Y) = \sum_{i+j \leq r} c_{ij} X^i Y^j.$$

GOAL: Recover a single symbol (say,  $f(\alpha, \beta)$ ) given query access to  $g: \mathbb{F}_q^2 \rightarrow \mathbb{F}_q$  with  $\Delta(g, f) \leq \delta$ .

We can think of codewords as  $q \times q$  grids of evaluation points.



Suppose I want to recover  $f(0,0)$ .  
As before, we want to find a bunch of LOCAL, LINEAR relationships involving  $f(0,0)$ .



← This row is  $(f(0,0), f(0,\gamma), \dots, f(0,\gamma^{r-1}))$   
 $=: (g(0), g(\gamma), \dots, g(\gamma^{r-1}))$

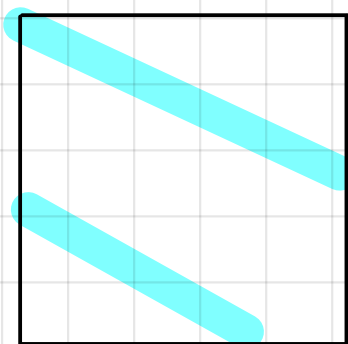
$$\text{where } g(\gamma) := f(0,\gamma) = \sum_{i+j \leq r} c_{ij} 0^i \cdot \gamma^j \\ = \sum_{j \leq r} c_{0j} \gamma^j$$

↗ Similarly, this column is  $\begin{pmatrix} h(0) \\ h(\gamma) \\ \vdots \\ h(\gamma^{r-1}) \end{pmatrix}$  where  $h(x) = \sum_{j \leq r} c_{j0} x^j$

Hey, those are univariate polynomials! (aka, RS codewords).

Moreover, the restriction of  $f$  to ANY line is an RS codeword!

Consider the line  $L(z) = (a_1 z + b_1, a_2 z + b_2)$ ,  
 $a_i, b_i \in \mathbb{F}_q$ .

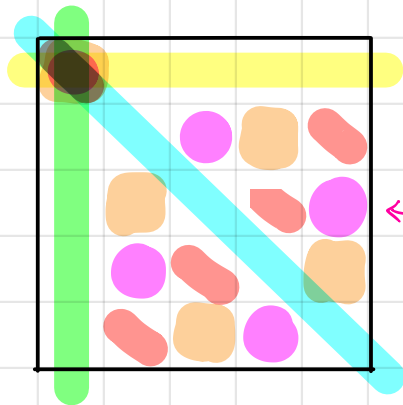


Then  $f(L(z)) = \sum_{i+j \leq r} (a_1 z + b_1)^i (a_2 z + b_2)^j$   
 $= \text{some degree } \leq r \text{ polynomial in } z$ .

The lines through  $f(0,0)$  have the properties we want:

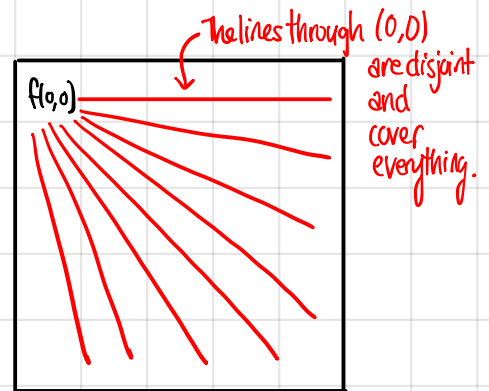
- There are not too many (only  $q$ ) points per line.
- Any two lines through  $f(0,0)$  don't intersect anywhere else.

PICTURE(S):



These pink dots form a line in  $\mathbb{F}_5$

Confusing but more accurate



Less accurate but hopefully more clear.

This inspires an algorithm:

ALG. (Let  $r < q$  and suppose  $\delta < \frac{1}{2}(1 - r/q)$ .)

Input: Query access to  $g: \mathbb{F}_q^2 \rightarrow \mathbb{F}_q$  s.t.  $\Delta(g, f) < \delta \cdot q^2$  for some  $f \in \text{RM}_q(2, r)$ .

and an index  $(\alpha, \beta) \in \mathbb{F}_q^2$

Output: A guess for  $f(\alpha, \beta)$ .

Choose  $(\sigma, \tau) \in \mathbb{F}_q^2 \setminus \{(0,0)\}$  at random, let  $L(Z) = (\sigma \cdot Z + \alpha, \tau \cdot Z + \beta)$ .

Query  $g(L(\lambda))$  for all  $\lambda \in \mathbb{F}_q$  and let  $\tilde{h}(Z) := g(L(Z))$ .

Use RS decoding to find an  $h \in \mathbb{F}_q[Z]$ ,  $\deg(h) \leq r$ , so that  $\Delta(h, \tilde{h}) < \frac{q-r}{2}$

half the distance of the RS code.

RETURN  $h(0)$ .



CLAIM. For any  $\delta > 0$ ,  $RM_q$  is correct with prob.  $\geq 1 - \left(\frac{2\delta q}{q-r-1}\right)$

Proof. The RS decoder will successfully find  $h(Z) = f(L(Z))$  as long as the number of errors on  $\{L(x) : x \in \mathbb{F}_q\}$  is  $< \lfloor \frac{q-r}{2} \rfloor$ , since  $f(L(Z)) \in RS_q(q, r+1)$ .

$\mathbb{E}\{\text{\#errors on a line}\} = \delta q$ , so by Markov's inequality,

$$\mathbb{P}\{\text{\#errors on a line} \geq \lfloor \frac{q-r}{2} \rfloor\} \leq \frac{\delta q}{\lfloor \frac{q-r}{2} \rfloor} < \frac{2\delta q}{q-r-1}$$

NOTE: If  $\delta < \frac{1}{2}(1 - r/q) = \frac{1}{2} \text{dist}(RM_q(2, r))$ , then the failure probability above is interesting, otherwise it reads "with prob.  $\geq 0$ ."

For example:

COR.  $RM_q(2, r = q/2) \subseteq \mathbb{F}_q^N$ ,  <sup>$N = q^2$  here</sup> is a  $(Q = \sqrt{N}, \delta, 4\delta)$ -LCC for any  $\delta < \frac{1}{4}$ . The rate is  $\approx 1/8$  and the distance is  $1/2$ .

We can do EXACTLY the same thing with  $m > 2$ .

LARGE-but-CONSTANT m:

Then we get  $Q = q = N^{1/m}$ , since  $N = q^m$ .

However, as  $m \uparrow$  then the rate  $\downarrow$ .  $\leftarrow$  Recall  $R = \binom{q+m}{m} / q^m \leq \left(\frac{c}{m}\right)^m \rightarrow 0$  as  $m \rightarrow \infty$ .

But this does give us a constant-rate code w/  $Q = N^{1/100}$  (say).

EVEN LARGER m:

Choose  $m = q / \log(q)$ .

Then  $N = q^{q/\log(q)} = 2^q$  so  $Q = q = \log(N)$

But the rate is even worse, and in fact goes to 0 like  $1/\text{poly}(n)$ .

$\leftarrow$  This simple construction is the state-of-the-art for  $\log(n)$  queries. Can you do better???

So far, we have seen how to use RM codes to get:

$Q$	$n$ (as a function of $k$ )	Code
2	$n = \Theta(2^k)$	$RM_2(m, 1)$
$\log(n)$	$n = \text{poly}(k)$	$RM_q(m, r)$ for $m \approx \frac{8}{\log(q)}$ , $q > r$
$n^\epsilon$	$n = \Theta_\epsilon(k)$	$RM_q(m, r)$ for $m = \frac{1}{\epsilon}$ , $q > r$
$\sqrt{n}$	$n = 8k$	$RM_q(2, \frac{q}{2})$

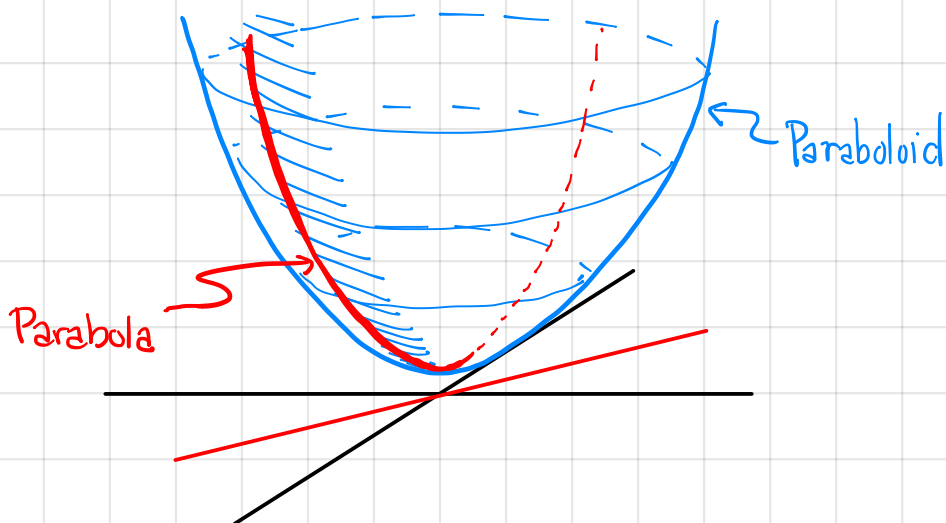
All of these have pretty low rate. Could we get an LCC with rate  $\rightarrow 1$ ?

For  $Q = n^\epsilon$  (and even a bit smaller), the answer is YES.

There are several constructions. Here's a sketch of one based on RM codes.

## (2) HIGH-RATE LCCs.

The thing we needed from RM codes are that restrictions to lines are low-deg polys.



To make the rate better, we might try

$$\mathcal{C} = \left\{ (f(\alpha_1), \dots, f(\alpha_{q^m})) : f \in \mathbb{F}_q[X], \text{ AND } \deg(f(L(z))) \leq r \forall \text{ lines } L \right\}$$

This would be a win as long as  $|C| \geq |RM_q(m,r)|$ , aka, as long as there are high-degree polynomials whose restrictions to lines are low-degree.

QUESTION. Does there exist a polynomial  $f: \mathbb{F}_q^m \rightarrow \mathbb{F}_q$  of degree  $> r$  so that,  
 $\forall$  lines  $L: \mathbb{F}_q \rightarrow \mathbb{F}_q^m$ ,  $\deg(f(L(z))) \leq r$ ? (for  $r < q-1$ )

This means  $\exists g(z)$  w/  $\deg(z) \leq r$  s.t.  $g(\lambda) = f(L(\lambda))$   
 $\forall \lambda \in \mathbb{F}_q$

ANSWER.

Over  $\mathbb{R}$  or  $\mathbb{C}$ : **NO.** (fun exercise!)

Over  $\mathbb{F}_p$ , for prime  $p$ : **NO.** (See [Rabinfeld-Sudan '96])

Over  $\mathbb{F}_q$ , and  $q > 2r$ : **NO.** (" " ")

Over  $\mathbb{F}_q$ , and  $q \approx r(1+\epsilon)$ : **YES**, and there are **LOTS** of them.

[Guo, Kopparty, Sudan '13]

EXAMPLE. Consider  $f(X,Y) = X^2 Y^2$  over  $\mathbb{F}_4$ .

The degree of  $f$  is 4.

Any restriction of  $f$  to a line is equivalent to a polynomial of  $\deg \leq 3 = q-1$ . (Not too helpful).

CLAIM  $\forall$  lines  $L: \mathbb{F}_4 \rightarrow \mathbb{F}_4^2$ ,  $\deg(f(L(z))) \leq 2$ .

pf. Say  $L(z) = (\sigma z + \alpha, \tau z + \beta)$ .

$$\begin{aligned} f(L(z)) &= (\sigma z + \alpha)^2 (\tau z + \beta)^2 \\ &= (\sigma^2 z^2 + \alpha^2) (\tau^2 z^2 + \beta^2) \quad [(a+b)^2 = a^2 + b^2 \text{ in } \mathbb{F}_2] \\ &= \sigma^2 \tau^2 z^4 + (\alpha^2 \tau^2 + \sigma^2 \beta^2) z^2 + \alpha^2 \beta^2 \quad [\text{algebra}] \\ &\equiv (\alpha^2 \tau^2 + \sigma^2 \beta^2) z^2 + \sigma^2 \tau^2 z + \alpha^2 \beta^2. \quad [\alpha^4 = \alpha \text{ in } \mathbb{F}_4] \end{aligned}$$

That's just one example, but it turns out there are actually LOTS, enough so that

$$\mathcal{C} = \left\{ (f(\alpha_1), \dots, f(\alpha_{q^m})) : f \in \mathbb{F}_q[X], \text{ AND } \deg(f(L(z))) \leq r \ \forall \text{ lines } L \right\}$$

has  $|\mathcal{C}| \geq q^{(1-\epsilon) \cdot (q^m)}$ , aka  $\text{RATE}(\mathcal{C}) \geq 1-\epsilon$ .

$\mathcal{C}$  is called a "LIFTED CODE."

Thm (Guo, Kopparty, Sudan)

$\forall m > 0, q = 2^t, \forall \epsilon > 0, \exists \epsilon' > 0$  s.t. the set

$$S = \left\{ f: \mathbb{F}_q^m \rightarrow \mathbb{F}_q \mid f \text{ has degree } \leq (1-\epsilon') \cdot q \text{ restrictions to ALL lines} \right\}$$

has  $\dim(S) \geq (1-\epsilon) \cdot q^m$

COR.  $\forall \epsilon, \alpha > 0, \exists \delta > 0$  and  $\gamma > 0$  s.t. there exists a family of codes  $\mathcal{C} \subseteq \mathbb{F}_q^n$  so that  $\mathcal{C}$  is a  $(n^\alpha, \delta, \gamma)$ -LCC of rate  $1-\epsilon$ .

(One can do a bit better than this: see [Kopparty, Meir, Ron-Zewi, Saraf, 2015].)

- RECAP:
- RM codes have nice local structure
  - They are LCCs with  $Q = 2$ ,  $\log(n)$ ,  $n^{1/100}$  although the rate gets bad.
  - To get rate  $1 - \epsilon$  with  $Q = n^{1/100}$ , we can "lift" RM codes.
- In general, there are TONS of open questions about LCCs!

## QUESTIONS TO PONDER

- ① Can you show that  $k$  must be at least  $n^{2+\epsilon}$  for 3-query LCC's?
- ② Can you beat RM codes for  $Q = \log(n)$ ?
- ③ Can you do anything with the Hadamard code when  $\delta > 1/4$ ?