

# Lower Bounds via the Cell-Sampling Method

Omri Weinstein  
Columbia

---



# Locality in TCS

- Locality/Sparsity is central to TCS and Math:
  - PCP Theorems
  - Locally-Decodable Codes (LDCs)
  - Data Structures
  - Derandomization (expanders, k-wise independence)
  - Matrix Rigidity
  - Compressed sensing
  - Graph decompositions (LLL)
  - ...

## Limits of Local Computation ?

- Typically, locality comes at price (e.g. blowup in size of input)
- How can we prove lower bounds on this tradeoff?
- This Tutorial:
  - “Cell Sampling”: A simple & unified technique. Proves highest known **unconditional** lower bounds in various computational models.

# Plan

- Cell-Sampling technique
- Applications:
  - I) Time-Space Tradeoffs in Data Structures (near-neighbor search)
  - II) LB for Locally Decodable Codes (rate vs. locality)
  - III) Matrix Rigidity (sparsity vs. rank)
- Limits of cell-sampling method

# Information Theory I01

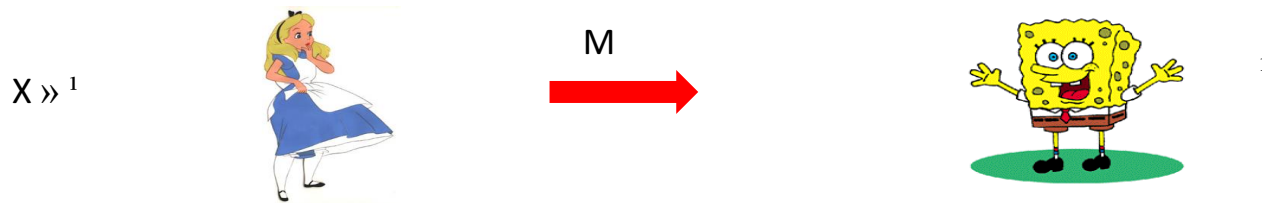
- **Entropy** : For random variable  $X \gg 1$ ,

$$H_1(X) := \sum_{x \in \mathcal{X}} p(x) \lg(1/p(x)) = \mathbf{E}_X[\lg 1/p(X)]$$

Captures how “unpredictable”  $X$  is – E.g.,  $H(\text{Ber}(1/2)) = 1$  bit,  $H(\text{Unif}_n) = \lg(n)$  bits .

- **Conditional Entropy** :  $H_1(X|Y) := \mathbf{E}_y[H_1(X | Y=y)]$

**Thm** (Shannon '48):  $\mathbf{E}_1[\text{cost of sending } X]$ ,  $H_1(X)$  bits . (tight by Huffman code)



# Cell Sampling

- LBs on “locality” via **compression** argument.
- **High-level idea:**  
Too-good-to-be-true “**local**” Algorithm → impossible compression of **input**.

(Typically not enough by itself – need to combine argument with extra features/  
structure of problem, e.g., geometric/ combinatorial etc – more on this soon)

- Let’s exemplify this method by proving **time-space** tradeoffs for **data structures**.

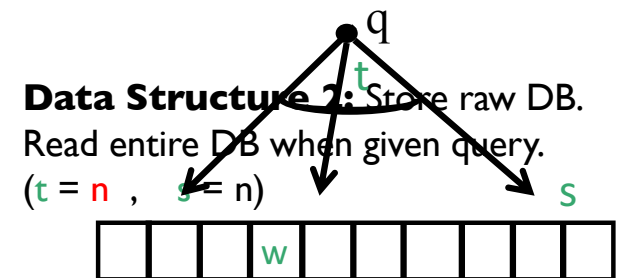
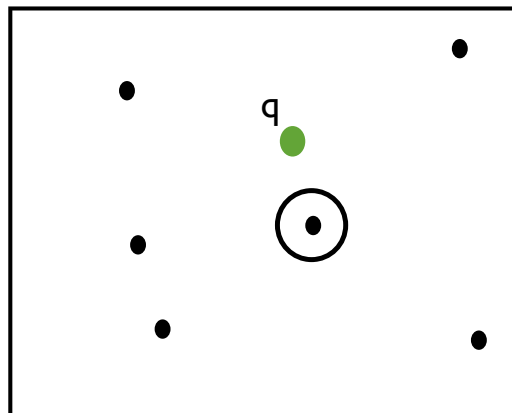


## Data Structures LBs (“Cell-Probe” model)

- DS = “compact” representation of info in **database**, so that **queries** about data can be answered quickly.
- **Static Data Structures**: Given **data**  $X$  of  $n$  elements in advance (e.g., graph, string, set of pts, etc.), *preprocess* it into small memory  $s$  so that  $\exists$  **query**  $q \in Q$  can be computed fast with  $t$  memory accesses (computations free of charge!).

- **NNS**: Data =  $n$  pts in  $\{0,1\}^d$

**Data Structure 1**: Precompute and store all answers in lookup table.  
 ( $t = 1$ ,  $s = 2^d$ )



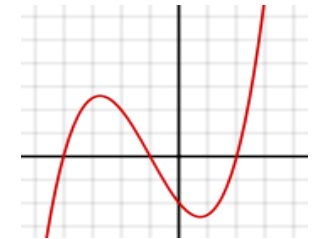
**Data Structure 2**: Store raw DB.  
 Read entire DB when given query.  
 ( $t = n$ ,  $s = n$ )

- **Data Structure LBs**: Is there anything in between? Study **time-space** tradeoffs ( $s$  vs.  $t$ ).

## Ex: Polynomial Evaluation

- PolyEval:

- Input: Random degree- $n$  polynomial  $P \in \mathbb{F}_m$  ( $m = n^2$ ).
- Query: Element  $x \in \mathbb{F}_m \rightarrow$  Return  $P(x)$ .
- $H(P) = (n+1)\lg(m)$  ( $n+1$  random coeff  $\in \mathbb{F}_m$ , word size  $w = \lg m$ )



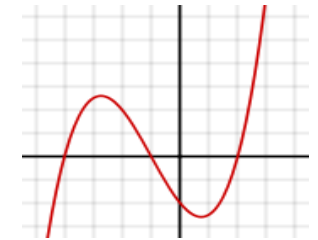
- Trivial:  $s = n+1$ ,  $t = n+1$  (read all coefficients)

**Thm** : Any  $\mathbf{D}$  with space  $s = O(n)$  must have query time  $t, \Omega(\lg n)$ .



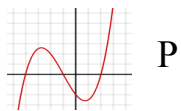
# Proof (via cell-sampling)

- Assume toward contradiction  $\exists \mathbf{D}$  with space  $s = 10n$  and query time  $t = o(\lg n)$ . (recall word-size  $w = \lg m$ ).
- Use data structure to **encode**  $P$  using less than  $(n+1)\lg m$  bits (!)
- Alice **Encodes**  $P$  :
  - Build DS  $\mathbf{D}$  on  $P$ .
  - Alice picks a *random sample*  $C$  of mem cells:
    - Include each cell w.p  $p := 1/100$ .
  - Sends Bob  $C$  (addresses + contents).
    - $\mathbf{E}[\text{message length}] = (s / 100) 2w < 10n 3\lg n / 100 < (n+1)\lg(m)$  bits =  $H(P)$  !



$|\mathbf{D}(P)| = 10n$  (words)

	222	j41	If...	Else	389	j4#	\$y	j13	
--	-----	-----	-------	------	-----	-----	-----	-----	--



- Expected cost  $< H(P)$  bits.

- Decoding  $P$  :

- Bob iterates through all  $x \in \mathbf{F}_m$ :
  - Run query algorithm of DS on  $x$ :
    - If read **outside**  $C$ , discard  $x$ .

- Probability recover answer to fixed  $x = p^t = (1/100)^t$

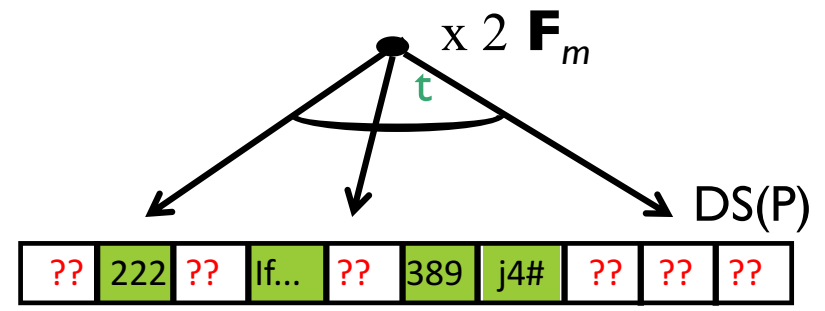
- $E_C[\# \text{ surviving queries } x] = m \cdot (1/100)^t = \Omega(n^2 \cdot 2^{-t}) = n^{2-o(1)}$

But **every**  $n+1$  queries **determine**  $P$ , hence Bob learns  $H(P) \sim n \log m$  bits of info from  $< n \log m$  bits of CC (!)

)  $t = \Omega(\lg n)$ .

[ More generally :

$t, \Omega(\lg(n)/\lg(s/n)) ]$



- Special property of polynomials: **Any large enough** set of answers recovers entire input (“ $n$ -wise independence”). Most natural problems don’t have this feature..

# Time-Space Lower Bounds for Near-Neighbor Search



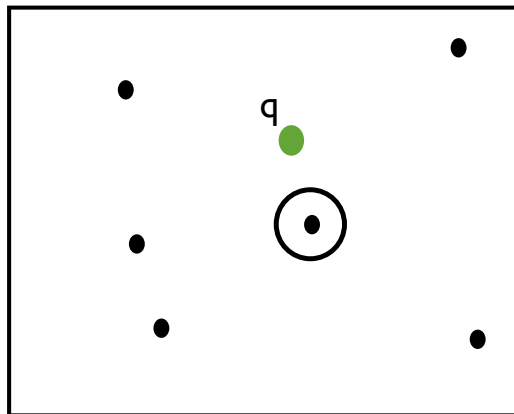
# Nearest-Neighbor Search

- **NNS:** Preprocess dataset  $X = x_1, \dots, x_n$  in metric space (say  $\mathbb{R}^d$  with  $l_1$  norm), s.t given a query  $q \in \mathbb{R}^d$ , closest point in  $X$  to  $q$  can be retrieved as fast as possible.

## Data Structure 1:

Precompute and store all answers in lookup table.

$$(t = 1, s = 2^d = n^{100})$$



$$\{0,1\}^d \quad (d = 100 \log(n))$$

## Data Structure 2:

Store DB (graph) as is. Read entire DB when given query (linear scan).

$$(t = n, s = n)$$

- Better time-space tradeoffs ( $s$  vs.  $t$ ) ?
- Probably not... (“Curse of dimensionality”)

## Approximate NNS

- **(c,r)-ANN**: Relaxed requirement: Given *radius*  $r$  and apx parameter  $c > 1$ , if  $\exists x_i$  s.t  $|q - x_i| < r$ , return  $x_j$  s.t  $|q - x_j| \leq cr$ .

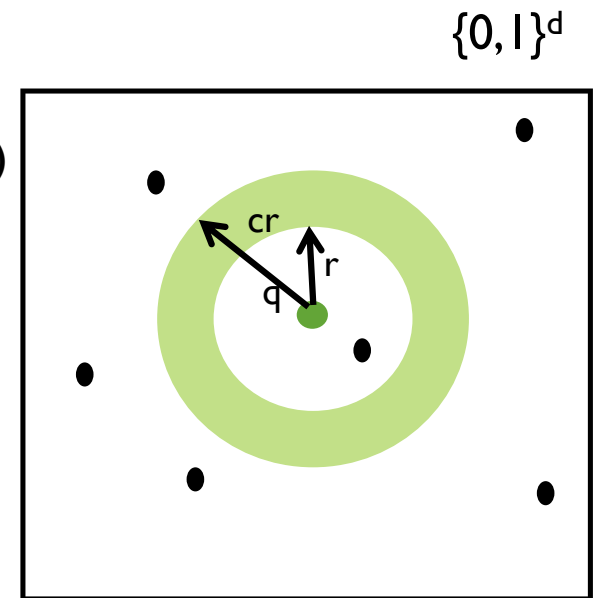
- “Robust” version has dramatic consequences (LSH) :  
 $s = n^{1+2}$ ,  $t = O(n^2)$  for  $c=(1/2)$ -apx. ( $l_1, l_2$  [IM98, Pan06, AR15..])

- Is this optimal? Can we get near-linear space and  $t = n^{o(1)}$ ?

- **Thm** [PTW'10, LMWY'19] :  
 $\exists$  DS  $\mathbf{D}$  for  $(1/2)$ -ANN over  $d$ -dim Hamming space,

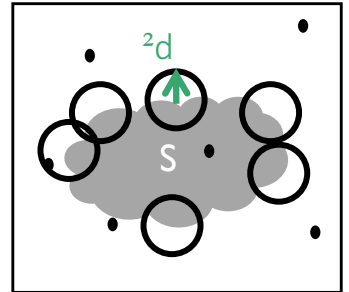
$$t, \Omega(d / \lg(dw)) \text{ for } s = O(n) \text{ space.}$$

- For  $d = \Omega(\lg n) \rightarrow \Omega(\lg n / \lg \lg n)$ .



# Proof

$$2^d = n^{10}$$



- Consider  $X = x_1, \dots, x_n \gg U(2^d)$ ,  $d = 10 \log(n)$  (whp,  $B_{2^{-d}}(x_i)$  are all unique)

- Isoperimetric Fact :  $|S| = 2^{(1-2^{-d})d} \cdot j_2(S)$ ,  $2^{d-1}$   
 (Harper's Inequality: least-expanding subset of hypercube = ball)

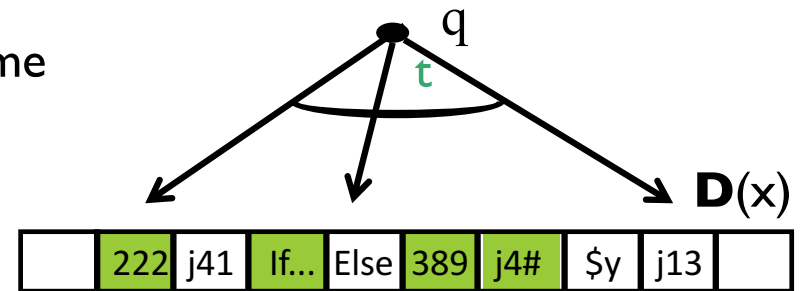
(\*) Cor:  $\delta$  fixed subset of  $|S|$  &  $2^{(1-2^{-d})d}$  r-ANN queries ( $r=2^{-d}$ ),  $\Pr_{x_i \gg U} [x_i \in S] = 1/2$

)  $n/4$  data points  $x_i$  fall into **any** such  $S$  whp.

- Consider (0-err)  $\mathbf{D}$  solving ANN with  $s=10n$  space (say),  $t = o(2^{2d} / \lg w)$  query time.

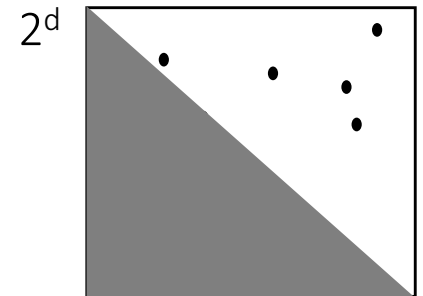
- $\mathbf{D} \rightarrow$  too good to be true (rand) compression scheme for encoding  $n/8$   $x_i$ 's using  $o(n \log d) = o(n \lg n)$  bits !

- Alice samples  $\delta$  cell  $c \in \mathbf{D}(x)$  iid w/  $p := 1/100w$ .



- Alice sends Bob contents + addresses of **sampled** cells

$$\mathbf{E}[|C| \mid \phi, w] = 2^d p w < n/10 \text{ bits (recall } p = 1/100w)$$



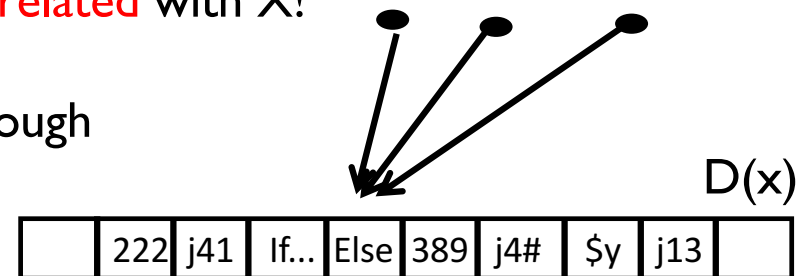
- $\mathbf{E}_{C,X}[\# \text{ surviving queries } Q] = 2^d \phi \Pr_{C,X}[\mathbf{D}(q) \leq 1/2 C] = 2^d \phi p^t$

$$\& 2^d \phi (1/100w)^{o(2^2 d/\lg w)} > 2^d - o(2^2 d) > 2^{(1-2^2)d} \quad (p = 1/100w)$$

- **If** it were the case that  $Q \perp X$  (i.e.,  $(X|Q) \gg U(2^d)$ ) By (\*)  $(i_2(Q), 2^{d-1})$ , Bob would have been able to recover  $n/4$  (say)  $x_i$ 's just from  $C \rightarrow$  contradiction!

- But  $\mathbf{D}$  is **adaptive**  $\rightarrow$  surviving queries heavily depend on *content* of cells (function of  $X$ )  $\rightarrow$  Surviving set  $Q = Q(X)$  **correlated** with  $X$ !

- In principle, all  $x_i$ 's could all fall into  $i_2(Q)$  (even though it has covers  $1/2$  the space).

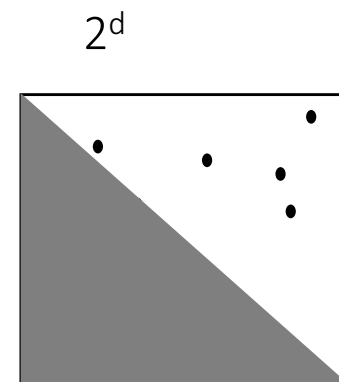


- Obs:  $Q(X)$  determined by only  $|C|w < o(n)$  bits (actually  $n/10$  but good enough)

(DPI)  $H(X|_{i_2(Q(X))}) > nd - o(n)$  bits :  $X$  is still “close” to  $U(2^d)$ ...

- Formalize this using simple “geometric packing” argument: Suppose  $f_{soc}$  that  $> n/4$   $x_i$ 's fall **outside**  $i_2(Q)$   $\rightarrow$  these pts are (essentially) **(d-1)**-dim.

$\rightarrow$  Can save 1 bit for their encoding, already gives impossible compression...



- So may assume  $> n/4$   $x_i$ 's indeed fall into  $i_2(Q)$  as desired, in which case prev “naive” analysis goes through.

- Cell-Sampling also used in highest  $(\sim \lg^2 n)$  **dynamic** data structure lower bounds... [Lar12, LWY18]



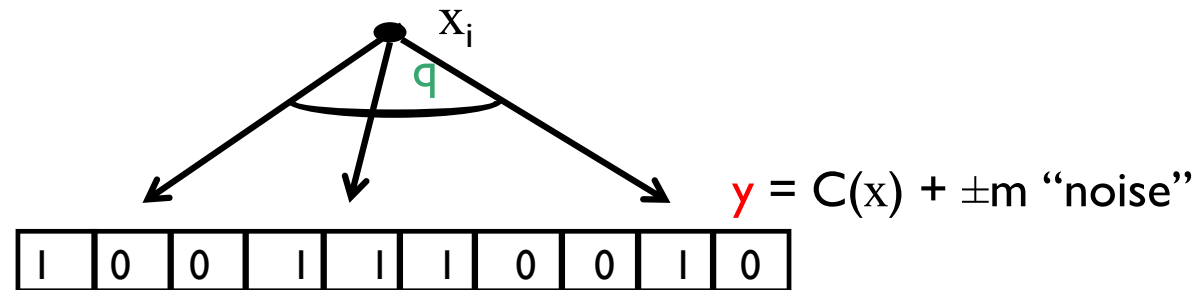
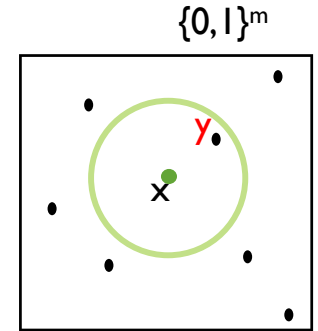
# Lower Bounds on Locally Decodable Codes



# Error Correcting Codes

$\pm = (\text{frac})$  distance of  $C$

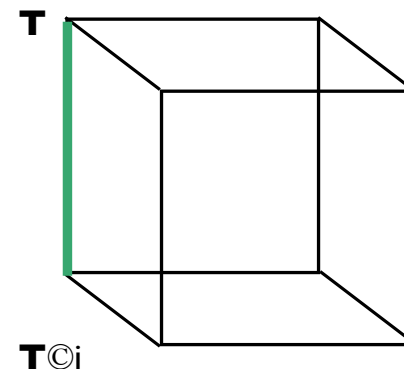
- ECC  $C : \mathbf{F}^n \mapsto \mathbf{F}^m$  ( $m > n$ ) s.t.  $\exists x, y |C(x) - y| < \pm \Rightarrow x$  recovered from  $y$ .
- For  $\pm = 1/4$  (say), 9 ECCs with constant **rate**  $m = O(n)$ .
- But decoding requires reading **entire** codeword  $C(x)$ , even if just want  $x_i$ .
- If interested in decoding only  $x_i$ , can hope to read **few** (ideally  $O(1)$ ) bits of  $C(x)$ ?
- **q-LDC**  $C : \{0,1\}^n \mapsto \{0,1\}^m$  s.t.  $\exists x, y |C(x) - y| < \pm \Rightarrow$  recover  $x_i$  by reading only  $q$  bits of  $y$ .



# Locally-Decodable Codes

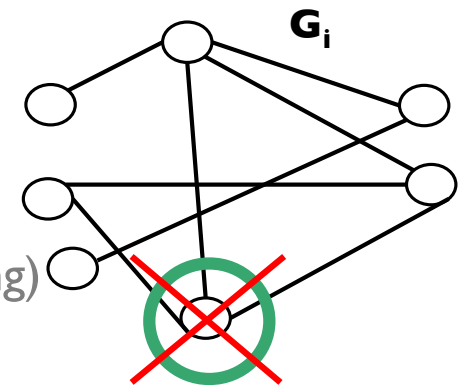
LDCs must randomize!

- **q-LDC**  $C : \{0,1\}^n \mapsto \{0,1\}^m$  s.t.  $\forall x, d(C(x),y) < 1/4$  ) recover  $x_i$  by reading  $q$  bits (whp).
- Tradeoff b/w  $q$  and  $m$  ? Is  $q=O(1)$  possible with  $m=O(n)$  ??
- Claim: for  $q=1$ , impossible (intuition: some bit  $j \in [m]$  must convey info on  $\Omega_{\pm}(n)$   $x_i$ 's)
- $q=2$  ? Possible with  $m = 2^n$  :
- To Encode  $x \in \{0,1\}^n$ , store  $\mathcal{T} \subseteq [n]$   $C(x)_{\mathcal{T}} := \bigoplus_{i \in \mathcal{T}} x_i$  ( $m = 2^n$ )
- To Decode  $x_i$  from  $y$ , pick  $\mathcal{T} \subseteq [n]$  & query  $y_{\mathcal{T}} \oplus y_{\mathcal{T} \ominus i}$
- $\Pr_{\text{random } \mathcal{T}}[\text{both } y_{\mathcal{T}} \text{ \& } y_{\mathcal{T} \ominus i} \text{ uncorrupted}] \approx 1/4$  ☺
- LB on tradeoff b/w  $q$  and  $m$  ? Is  $q=O(1)$  possible with  $m=O(n)$  ??



**Thm** [KatzTrevisan'00] :  $\exists$   $q$ -LDC ,  $m$  ,  $\sim \Omega_{\pm}(n^{1+1/q})$

- **Proof:** For  $q$ -LDC  $C$ , the **query graph**  $\mathbf{G}_i$  of  $C$  is the  $q$ -hypergraph containing possible  $q$ -tuples from which  $x_i$  can be recovered ( $|\mathbf{G}_i|=m$ ).
- **“Smoothness”:** Intuitively,  $q$ -edges of  $\mathbf{G}_i$  are  $1/q$  uniformly distributed: No vertex  $j \in \mathbf{G}_i$  has (weighted) degree  $\phi > q/\pm m$  (o.w adversary can **corrupt** it. Avg deg =  $q/m$  (Markov)).
- Corollary:  $\exists$   $q$ -LDC,  $\mathbf{G}_i$  contains **Matching**  $|M_i|$  ,  $\pm m/q^2$ .
- **Proof:** Max  $|Matching(\mathbf{G}_i)|$  in  $q$ -hypergraph  
 , Min  $|VC(\mathbf{G}_i)| / q$  (any VC must pick  $1$  v from max matching)  
 ,  $1/\phi \leq 1/q$  (each vertex covers  $\cdot \phi$  “mass”)  
 ,  $1 / (q\phi q/\pm m) > \pm m/q^2$  (max-deg  $\phi < q/\pm m$ )



- Use LDC to **compress** input (via Cell-Sampling) :
- Alice builds  $C(x)$ , samples  $\delta \approx 2 \lceil \log m \rceil$  w.p  $p := n/10m \rightarrow \mathbf{E}[S] = m \cdot p = \mathbf{n/10}$  bits.
- $\delta \approx 2 \lceil \log n \rceil \Pr_S[\text{Bob can recover } x_i] \leq \Pr_S[\bigcup_{e \in M_i} \text{"e survives"}]$   
 $= \sum_{e \in M_i} \Pr_S[\text{"e survives"}]$  (disjoint events since  $M_i = \text{matching!}$ )  
 $= |M_i| \cdot p^q = \boxed{\pm m/q^2 \cdot (n/10m)^q < 3/4}$  for  $>n/2$  i's (o.w. recover  $n/2$   $x_i$ 's from  $< n/10$  bits!)  
 $\leq m^{q-1} \cdot n^q / q^2 \leq \Omega_{\pm}(n^{1+1/q})$  ■



# Matrix Rigidity



## Matrix Rigidity

$$M \neq A + B$$

$\swarrow$                        $\searrow$   
 $t$ -sparse                   $\text{rk} \leq \epsilon n$

- **Def:** A matrix  $M \in \mathbf{F}^{m \times n}$  is  **$t$ -Rigid** if decreasing its rank  $n \rightarrow n/2$  requires modifying  $t$  entries in some row.

- $M$  “ $t$ -far” from any low-rank matrix. (assume  $m = \text{poly}(n)$ )

- **Thm:**  $n^2 \times n$  **Vandermonde** matrix is  $\Omega(\lg n)$ -Rigid.

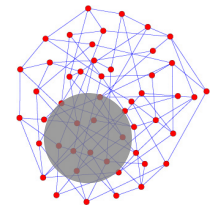
$$\begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & & x_n^n \end{bmatrix}$$

- Cell-Sampling + Subadditivity of Rank :

- Sketch: Sample  $\delta$  column of  $A$  w.p  $1/10$  (call it  $S$ )  $\rightarrow$  If every row is  $< \lg(n)/100$  sparse  $\rightarrow$   $\exists$   $n$  rows  $R$  s.t  $|S|=n/10$  “covers” **all** 1’s in these rows  $\rightarrow \text{rk}(A_R + B_R) \leq \text{rk}(A_R) + \text{rk}(B_R) \leq n/10 + n/2 < n$ . But every submatrix of  $V$  is also full rank ( $n$ ) !

## Limits of Cell Sampling

- Cell Sampling relies on simple fact : In every graph of size  $n$  with  $m$  edges, there is a small set ( $\sim pn$ ) containing “nontrivial” ( $\sim m/2^p$ ) edges.
- Tight for *expanders*... (any  $o(\lg n)$  subset contains  $o(n)$  edges)
- $\log(n)$  is a fundamental limit of cell sampling ☹️
- Still very useful technique, that unifies/explains current barrier in complexity holy-grails (LDC/Rigidity/DS...)





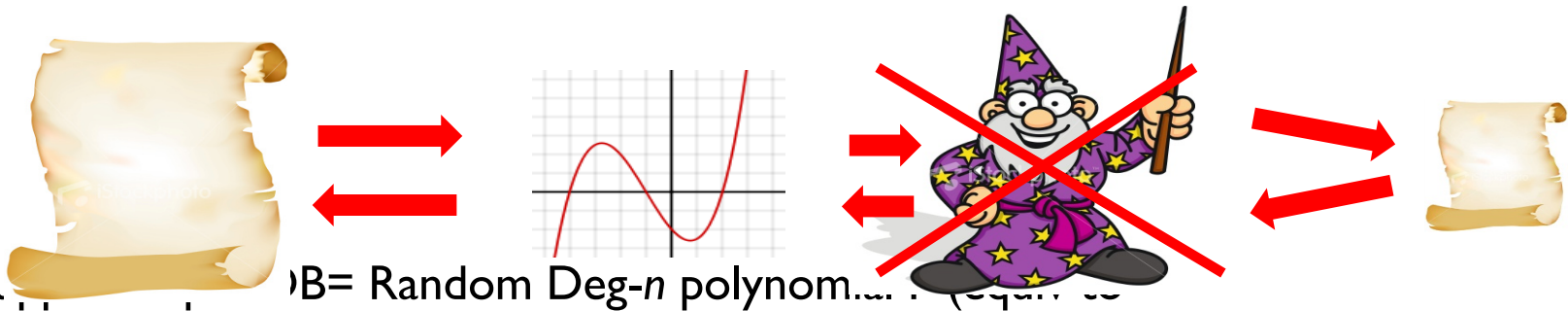
Thanks!

---



# Cell Sampling: Time-Space LBs via Compression

- Assume for contradiction that a “too-good-to-be-true” DS exists for PolyEval with  $t < \lg n$  and linear space ( $s = O(n)$ ).



- Step 1:  $S = \{T_i\}$  where  $T_i$  are  $n$ -letter text  $T$  w. random symbols :  $P(i) = T_i$ .
- Step 2:  $B = \text{Random Deg-}n \text{ polynomial}$ ...

- Use magic data structure to **encode** and **decode** the input set using less than  $y$  symbols. **A contradiction!**