

# Nearly Optimal Sparse Group Testing

Venkata Gandikota, *Purdue University*

Elena Grigorescu, *Purdue University*

Sidhart Jaggi, *The Chinese University of Hong Kong*

Samson Zhou, *Purdue University*

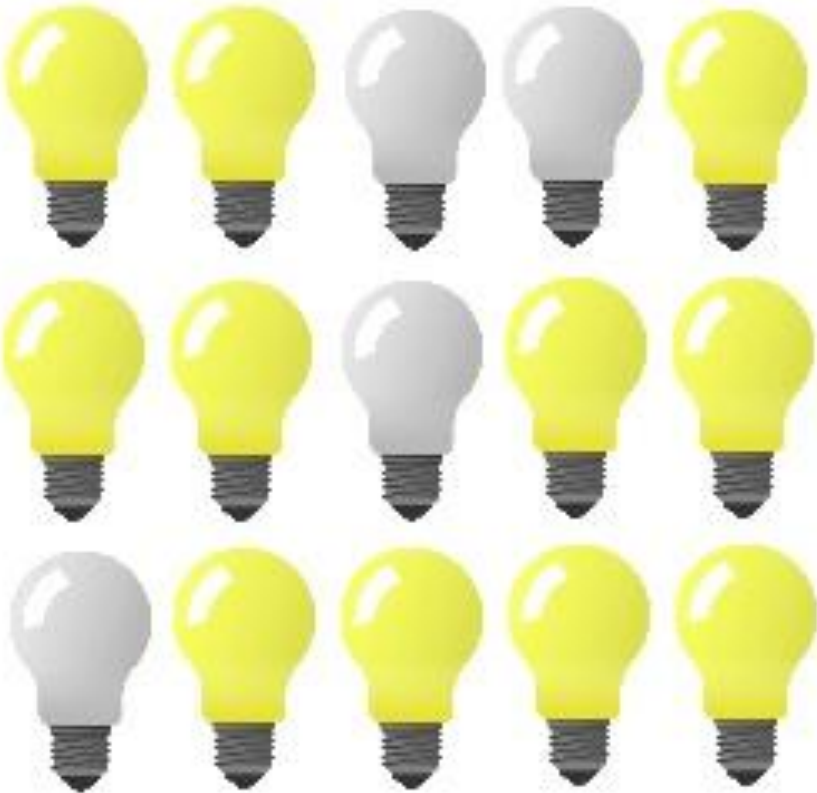
# Group Testing

**Task:** Identify a small set of defects among a larger population using tests

Test = a subset of the items

A test is **positive** if a defective item is included; **negative** otherwise

**Goal:** Minimize the number of tests



# Background



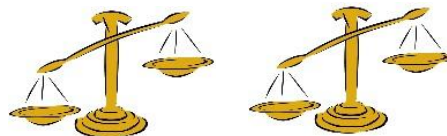
- **Motivation:** Identify WW2 draftees with syphilis
- Each test is expensive, but blood may be pooled
- Blood from individual can be used in multiple tests
- How to identify infected individuals?

# Adaptive Group Testing



$$T = O(d \log n)$$

$$T = \Omega\left(d \log \frac{n}{d}\right)$$



# Non-adaptive Group Testing (NAGT)



# Non-adaptive Group Testing (NAGT)



# Non-adaptive Group Testing (NAGT)

													
	0	 1	0	0	0	0	 1	0	0	 1	0	0	 1
	 1	 1	0	 1	0	0	0	 1	 1	0	0	0	0
	0	0	0	0	 1	 1	 1	0	0	 1	 1	 1	0
	 1	 1	 1	 1	0	0	0	0	0	0	0	 1	 1
	0	0	 1	 1	0	0	0	 1	0	0	 1	0	0
	0	0	0	0	0	 1	 1	0	0	0	 1	0	 1
	0	0	 1	0	 1	0	0	 1	0	 1	0	0	0
	 1	0	0	0	0	0	 1	 1	0	 1	 1	0	0

# Non-adaptive Group Testing

0	1	0	0	0	0	1	0	0	1	0	0	1
1	1	0	1	0	0	0	1	1	0	0	0	0
0	0	0	0	1	1	1	0	0	1	1	1	0
1	1	1	1	0	0	0	0	0	0	0	1	1
0	0	1	1	0	0	0	1	0	0	1	0	0
0	0	0	0	0	1	1	0	0	0	1	0	1
0	0	1	0	1	0	0	1	0	1	0	0	0
1	0	0	0	0	0	1	1	0	1	1	0	0



=

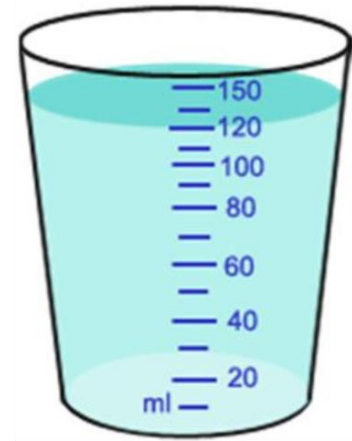




# “Classical” NAGT

- Upper Bound:  $O(d^2 \log n)$  (Du, Hwang '00)
- Explicit:  $O(d^2 \log n)$  (Porat, Rothschild '08)
- Lower Bound:  $\Omega(d^2 \log_d n)$  (D'yachkov, Rykov '82)
- Noisy tests (CheraghchiHKV '11)
- Efficient Decoding (IndyKRN'10)
- Graph-Constrained (CheraghchiKMS'11)
- Phase-Transitions (Scarlett, Cehver'16)

# Real World Limitations



Each item can be included  
in at most  $\gamma$  tests  
( $\gamma$  -divisible items)

Each test can include  
at most  $\rho$  items  
( $\rho$  -sized tests)

# Our Results: $\gamma$ -divisible items

**Theorem.** Given  $n$  items, with  $d$  defects:  
 $\Omega(\gamma d (n/d)^{1/\gamma})$  tests are needed in the NAGT,  
 $\gamma$  -divisible items model.

**Theorem.**  $\exists$  a randomized algorithm:  $T = O\left(\gamma d \left(\frac{n-d}{\varepsilon}\right)^{1/\gamma}\right)$

**Theorem.**  $\exists$  a deterministic algorithm:  $T = O\left(\frac{d^2 \gamma}{\varepsilon} \left(\frac{n\varepsilon}{d^2}\right)^{1/\gamma}\right)$

# Our Results: $\rho$ -sized tests

**Theorem.** Given  $n$  items, with  $d$  defects:

$\Omega\left(\frac{n \log(n/d)}{\rho \log(n/\rho d)}\right)$  tests are needed in the

NAGT,  $\rho$ -sized tests model.

**Theorem.**  $\exists$  a randomized algorithm:  $T = O\left(\frac{n}{\rho} \log\left(\frac{n}{\varepsilon}\right)\right)$

**Theorem.**  $\exists$  a deterministic algorithm:  $T = O\left(\frac{n}{\rho} \left(\frac{d^2 \log n}{\varepsilon \log(n/\rho)}\right)\right)$

# Structure of Talk

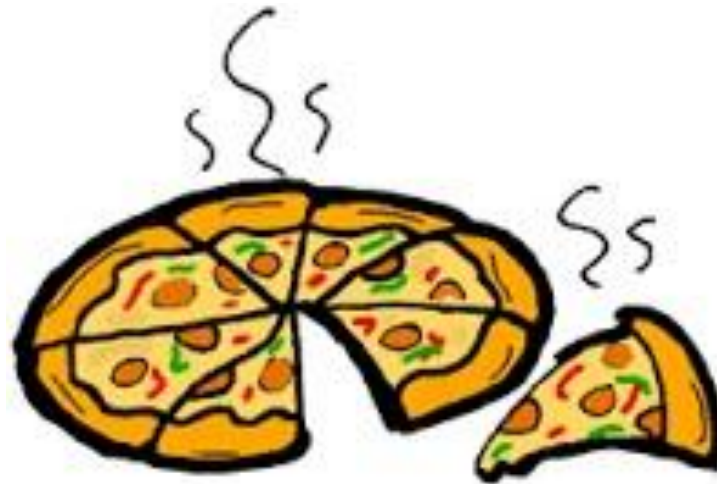
1. Background
2. Achievability: Randomized Construction
3. Achievability: Deterministic Construction
4. Lower Bounds

# Structure of Talk

1. Background
2. **Achievability: Randomized Construction**
3. Achievability: Deterministic Construction
4. Lower Bounds

# Upper Bounds: Randomized

- Consider the  $\gamma$ -divisible items model



# Upper Bounds: Randomized

- **Recall:** Total items tested is at most  $\gamma n$
- **Idea:** Have each test include roughly  $\frac{\gamma n}{T}$  items.
- Columns of test matrix **M** are uniformly sampled from  $\{0,1\}^T$  with weight  $\gamma$ .
- How to decode?





# Decoding Algorithm Philosophy



Innocent until  
proven guilty



Defective until  
proven innocent

# Decoding Algorithm

- Algorithm: Marks item  $i$  non-defective if some test which includes  $i$  is negative
- Observations:
  - **CANNOT** incorrectly identify defective items
  - Incorrectly marks non-defective item  $i$  if all tests which include  $i$  are positive



# Upper Bounds: Randomized

- **Recall:** Total number of positive tests is at most  $d\gamma$ .
- Probability an item is included only in positive tests:  $\binom{d\gamma}{\gamma} / \binom{T}{\gamma}$
- Union bound over all  $(n-d)$  non-defective items, so we require:

$$(n - d) \binom{d\gamma}{\gamma} / \binom{T}{\gamma} < \epsilon$$

# Upper Bounds: Randomized

- **Recall:** Total number of positive tests is at most  $d\gamma$ .
- Probability an item is included only in positive tests:  $\binom{d\gamma}{\gamma} / \binom{T}{\gamma}$
- Union bound over all  $(n-d)$  non-defective items, so we require:

$$(n - d) \binom{d\gamma}{\gamma} / \binom{T}{\gamma} < \epsilon \Rightarrow T > (e\gamma d) \left(\frac{n-d}{\epsilon}\right)^{1/\gamma}$$

# Structure of Talk

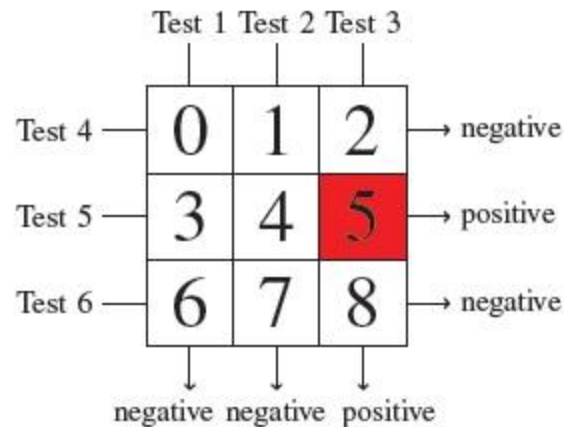
1. Background
2. **Achievability: Randomized Construction**
3. Achievability: Deterministic Construction
4. Lower Bounds

# Structure of Talk

1. Background
2. Achievability: Randomized Construction
3. Achievability: Deterministic Construction
4. Lower Bounds

# Upper Bounds: Deterministic

- **Intuition:** Should be able to “encode” each item so that test outcomes uniquely identify defects
- **Idea:** Use  $\gamma$ -dimensional hypergrid, represent each item by base- $b$  representation, where  $b = n^{1/\gamma}$



If  $n = 9$ ,  $\gamma = 2$ ,  $d = 1$ , the above test uniquely determines that item 5 is defective.

# Upper Bounds: Deterministic

- What if we have multiple defects?

	Test 1	Test 2	Test 3	
Test 4	0	1	2	→ positive
Test 5	3	4	5	→ positive
Test 6	6	7	8	→ negative
	↓	↓	↓	
	negative	positive	positive	

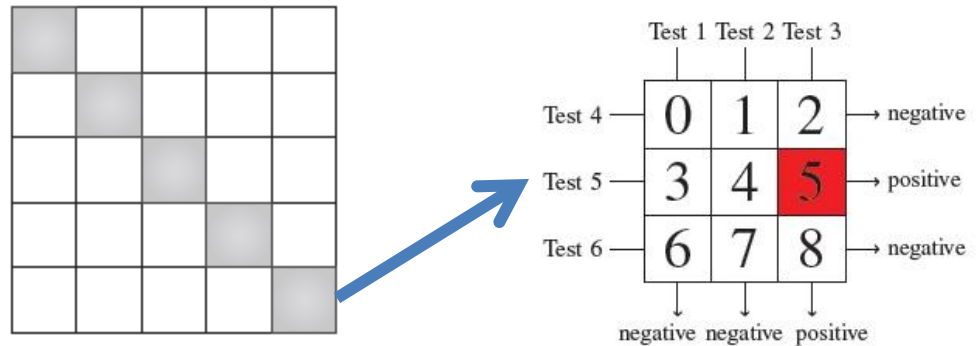
- These tests cannot distinguish whether the red items or the blue items are defective.





# Upper Bounds: Deterministic

- What if we have multiple defects?
- **Idea:** Divide and Conquer!



- Each gray box catches one defective item with high probability

# Upper Bounds: Deterministic

- How many blocks are necessary?  $\frac{d^2}{\epsilon}$
- Probability that no two defects fall into the same block:

$$\begin{aligned} & 1 \left(1 - \frac{1}{cd^2}\right) \left(1 - \frac{2}{cd^2}\right) \cdots \left(1 - \frac{d-1}{cd^2}\right) \\ & \geq \left(1 - \frac{d}{cd^2}\right)^d = \left(1 - \frac{1}{cd}\right)^d \\ & \geq 1 - \frac{1}{c} = 1 - \epsilon \quad (\text{by Bernoulli's Inequality}) \end{aligned}$$

- **In summary:**  $\frac{d^2}{\epsilon}$  blocks, each requiring  $\gamma \left(\frac{n\epsilon}{d^2}\right)^{1/\gamma}$  tests, for a total of

$$T = \frac{d^2}{\epsilon} \gamma \left(\frac{n\epsilon}{d^2}\right)^{1/\gamma}$$

# Structure of Talk

1. Background
2. Achievability: Randomized Construction
3. Achievability: Deterministic Construction
4. Lower Bounds

# Structure of Talk

1. Background
2. Achievability: Randomized Construction
3. Achievability: Deterministic Construction
4. Lower Bounds

# Classical NAGT Design Philosophy



- Each test gives ~1 bit of information

# Lower Bounds: $\gamma$ -divisible items

- Total number of items tested is at most  $\gamma n$
- “Light” tests: each includes less than  $\frac{n}{\varepsilon d \log\left(\frac{T}{\gamma d}\right)}$  items.
- “Heavy” tests: each includes at least  $\frac{n}{\varepsilon d \log\left(\frac{T}{\gamma d}\right)}$  items.

# Lower Bounds: $\gamma$ -divisible items

- “Light” tests: each includes less than  $\frac{n}{\epsilon d \log\left(\frac{T}{\gamma d}\right)}$  items.
- “Heavy” tests: each includes at least  $\frac{n}{\epsilon d \log\left(\frac{T}{\gamma d}\right)}$  items.

$$H(Y) = \sum_{i \in S_1} H(Y_i) + \sum_{i \in S_2} H(Y_i)$$

$$X \rightarrow Y \rightarrow \hat{X}$$

$$\leq T \left( \frac{\gamma d}{T} + 3\delta \right) \log \left( \frac{T}{\gamma d} \right) + \epsilon \gamma d \log \left( \frac{T}{\gamma d} \right)$$

$$\leq (1 + 2\epsilon) \gamma d \log \left( \frac{T}{\gamma d} \right)$$

(for appropriate choice of  $\delta$ ).

# Lower Bounds: $\gamma$ -divisible items

$$\begin{aligned}
 H(Y) &= \sum_{i \in \mathcal{S}_1} H(Y_i) + \sum_{i \in \mathcal{S}_2} H(Y_i) \\
 &\leq T \left( \frac{\gamma d}{T} + 3\delta \right) \log \left( \frac{T}{\gamma d} \right) + \epsilon \gamma d \log \left( \frac{T}{\gamma d} \right) \\
 &\leq (1 + 2\epsilon) \gamma d \log \left( \frac{T}{\gamma d} \right)
 \end{aligned}$$

$$X \rightarrow Y \rightarrow \hat{X}$$

$$\begin{aligned}
 H(X) &= H(X|\hat{X}) + I(X; \hat{X}) \\
 &\leq H(\epsilon) + \epsilon \log(|\mathcal{X}| - 1) - H(Y)
 \end{aligned}$$

$$\log \binom{\binom{n}{d}}{d} \leq -2\epsilon \log \epsilon + \epsilon \log \binom{\binom{n}{d}}{d} \quad \text{Data processing inequality}$$

$$+ (1 + 2\epsilon) \gamma d \log \left( \frac{T}{\gamma d} \right). \quad \text{Fano's inequality}$$

Entropy over all possible combinations of  $d$  defects



# Lower Bounds: $\gamma$ -divisible items

$$H(Y) = \sum_{i \in S_1} H(Y_i) + \sum_{i \in S_2} H(Y_i)$$

$$\leq T \left( \gamma d \log \left( \frac{T}{\gamma d} \right) + \epsilon \log \binom{n}{d} \right)$$

Implies  $T = \Omega\left(\gamma d (n/d)^{1/\gamma}\right)$

$$\log \binom{\binom{n}{d}}{d} \leq -2\epsilon \log \epsilon + \epsilon \log \binom{\binom{n}{d}}{d} \quad \text{Data processing inequality}$$

$$+ (1 + 2\epsilon) \gamma d \log \left( \frac{T}{\gamma d} \right) \quad \text{Fano's inequality}$$



Entropy over all possible combinations of  $d$  defects

# Structure of Talk

1. Background
2. Achievability: Randomized Construction
3. Achievability: Deterministic Construction
4. Lower Bounds

# Structure of Talk

1. Background
  2. Achievability: Randomized Construction
  3. Achievability: Deterministic Construction
  4. Lower Bounds
- ?? Results for zero-error tests and noisy tests

# Further Results

- Noisy tests: Tests can give incorrect result with probability  $\sigma$ .

**Theorem:** CANNOT recover defectives with probability at least  $1 - \epsilon$  for arbitrary  $\epsilon < 1/2$  and  $\gamma = o(\log n)$ .



# Questions?

