

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Streaming Lower Bounds for the Needle Problems

Jiapeng Zhang

University of Southern California

October 16, 2023

Joint works with Shachar Lovett; Qian Li and Shuo Wang

Overview

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- 1 The Needle Problem
- 2 Lower Bounds
- 3 Asymmetric Disjointness
- 4 Information Complexity Approaches

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

The Needle Problem

The Needle Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

Let $n > 1$ be a large integer and let $p > 0$ be a parameter.

- **Uniform distribution D_0** : each sample is uniformly sampled from $[n]$.
- **Needle distribution D_1** : First sample a needle $x \in [n]$. Each sample equals x with probability p , and uniformly otherwise.

The Needle Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

Let $n > 1$ be a large integer and let $p > 0$ be a parameter.

- **Uniform distribution D_0** : each sample is uniformly sampled from $[n]$.
- **Needle distribution D_1** : First sample a needle $x \in [n]$. Each sample equals x with probability p , and uniformly otherwise.

Question:

given a bounded memory of s bits, how many samples t are needed to distinguish these two distributions?

Simple Algorithms

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Algorithm 1

Put all samples in the memory, and find the most frequent element.

Simple Algorithms

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Algorithm 1

Put all samples in the memory, and find the most frequent element. It needs $\Theta(1/p)$ samples and $\Theta((\log n)/p)$ space.

Simple Algorithms

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Algorithm 1

Put all samples in the memory, and find the most frequent element. It needs $\Theta(1/p)$ samples and $\Theta((\log n)/p)$ space.

Algorithm 2

Keep the most recent two samples in the memory, and check the consecutive identical elements.

Simple Algorithms

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Algorithm 1

Put all samples in the memory, and find the most frequent element. It needs $\Theta(1/p)$ samples and $\Theta((\log n)/p)$ space.

Algorithm 2

Keep the most recent two samples in the memory, and check the consecutive identical elements. It needs $\Theta(1/p^2)$ samples and $\Theta(\log n)$ space.

Simple Algorithms

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Algorithm 1

Put all samples in the memory, and find the most frequent element. It needs $\Theta(1/p)$ samples and $\Theta((\log n)/p)$ space.

Algorithm 2

Keep the most recent two samples in the memory, and check the consecutive identical elements. It needs $\Theta(1/p^2)$ samples and $\Theta(\log n)$ space.

For general space s , we need $t \approx \Theta((\log n)/(s \cdot p^2))$ samples.

Our Results

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Theorem (Lovett-Z, Li-Wang-Z)

Any streaming algorithm that distinguishes the needle distribution needs $t = \Omega(1/(s \cdot p^2))$ samples.

Our Results

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Theorem (Lovett-Z, Li-Wang-Z)

Any streaming algorithm that distinguishes the needle distribution needs $t = \Omega(1/(s \cdot p^2))$ samples. For ℓ -pass streaming algorithm, it needs $t = \Omega(1/(\ell \cdot s \cdot p^2))$ samples.

Frequency Estimation

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Corollary

In the random order setting. It requires $\Omega(n^{1-2/k})$ space for a streaming algorithm to approximate k -th frequency moment of a data stream.

Frequency Estimation

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Corollary

In the random order setting. It requires $\Omega(n^{1-2/k})$ space for a streaming algorithm to approximate k -the frequency moment of a data stream.

Theorem (Andoni-McGregor-Onak-Panigrahy)

In the random order setting. It requires $\Omega(n^{1-2.5/k})$ space for a streaming algorithm to approximate k -the frequency moment of a data stream.

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Lower Bounds of the Needle Problems

Disjointness Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

There are k players with each of them holds a (random) set $S_i \subseteq [n]$. It is promised that,

- **Disjoint:** the sets S_1, \dots, S_k are pairwise disjoint.
- **Unique intersection:** there is an $x \in S_1 \cap \dots \cap S_k$, and the sets $S_1 \setminus \{x\}, \dots, S_k \setminus \{x\}$ are pairwise disjoint.

Disjointness Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

There are k players with each of them holds a (random) set $S_i \subseteq [n]$. It is promised that,

- **Disjoint:** the sets S_1, \dots, S_k are pairwise disjoint.
- **Unique intersection:** there is an $x \in S_1 \cap \dots \cap S_k$, and the sets $S_1 \setminus \{x\}, \dots, S_k \setminus \{x\}$ are pairwise disjoint.

Theorem

The randomized communication complexity of the disjointness problem is $\Omega((|S_1| + \dots + |S_k|)/k)$

A Simple Case of the Needle Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

In expectation, there are $\Theta(p \cdot t)$ needles in the stream.

A Simple Case of the Needle Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

In expectation, there are $\Theta(p \cdot t)$ needles in the stream.

$$X_1, \dots, X_{i_1}, \dots, X_{(1/p)}$$

$$X_{(1/p)+1}, \dots, X_{i_2}, \dots, X_{(2/p)}$$

...

$$X_{t+1-(1/p)}, \dots, X_{i_\ell}, \dots, X_t$$

The symmetric case.

Needle Algorithm to Communication Protocol

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Communication protocols

Let \mathcal{A} be an algorithm that distinguishes needles. Recall that each communication player i has a set S_i .

Needle Algorithm to Communication Protocol

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Communication protocols

Let \mathcal{A} be an algorithm that distinguishes needles. Recall that each communication player i has a set S_i .

- The first player randomly order S_1 , and sends $M_1 := \mathcal{A}(S_1)$ to the second player.

Needle Algorithm to Communication Protocol

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Communication protocols

Let \mathcal{A} be an algorithm that distinguishes needles. Recall that each communication player i has a set S_i .

- The first player randomly order S_1 , and sends $M_1 := \mathcal{A}(S_1)$ to the second player.
- The second player randomly order S_2 , and sends $M_2 := \mathcal{A}(M_1, S_2)$ to the third player.

Needle Algorithm to Communication Protocol

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Communication protocols

Let \mathcal{A} be an algorithm that distinguishes needles. Recall that each communication player i has a set S_i .

- The first player randomly order S_1 , and sends $M_1 := \mathcal{A}(S_1)$ to the second player.
- The second player randomly order S_2 , and sends $M_2 := \mathcal{A}(M_1, S_2)$ to the third player.
- ...

Needle Algorithm to Communication Protocol

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Communication protocols

Let \mathcal{A} be an algorithm that distinguishes needles. Recall that each communication player i has a set S_i .

- The first player randomly order S_1 , and sends $M_1 := \mathcal{A}(S_1)$ to the second player.
- The second player randomly order S_2 , and sends $M_2 := \mathcal{A}(M_1, S_2)$ to the third player.
- ...
- The last player receives M_{k-1} and outputs $\mathcal{A}(M_{k-1}, S_k)$.

Needle Algorithm to Communication Protocol

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Communication protocols

Let \mathcal{A} be an algorithm that distinguishes needles. Recall that each communication player i has a set S_i .

- The first player randomly order S_1 , and sends $M_1 := \mathcal{A}(S_1)$ to the second player.
- The second player randomly order S_2 , and sends $M_2 := \mathcal{A}(M_1, S_2)$ to the third player.
- ...
- The last player receives M_{k-1} and outputs $\mathcal{A}(M_{k-1}, S_k)$.

The total communication cost is $(k \cdot s)$ bits.

From Communication to Needle Lower Bounds

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- The total communication cost is $(k \cdot s)$ bits.

From Communication to Needle Lower Bounds

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- The total communication cost is $(k \cdot s)$ bits.
- From the communication lower bounds of disjointness, we have that $(k \cdot s) = \Omega((|S_1| + \dots + |S_k|)/k) = \Omega(t/k)$.

From Communication to Needle Lower Bounds

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- The total communication cost is $(k \cdot s)$ bits.
- From the communication lower bounds of disjointness, we have that $(k \cdot s) = \Omega((|S_1| + \dots + |S_k|)/k) = \Omega(t/k)$.
- Recall that $k = t \cdot p$, hence $s \cdot t = \Omega(1/p^2)$

Jiapeng Zhang

The Needle
Problem

Lower Bounds

**Asymmetric
Disjointness**

Information
Complexity
Approaches

Asymmetric Disjointness

Asymmetric Case

Jiapeng Zhang

The Needle
Problem

Lower Bounds

**Asymmetric
Disjointness**

Information
Complexity
Approaches

We still assume there are $p \cdot t$ needles.

Asymmetric Case

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

We still assume there are $p \cdot t$ needles.

$$X_1, \dots, X_{i_1}, \dots, X_{i_2}, \dots, X_{(1/p)}$$

$$X_{(1/p)+1}, \dots, X_{(2/p)}$$

...

$$X_{t+1-(1/p)}, \dots, X_{i_\ell}, \dots, X_t$$

Then the first player would know the answer. In expectation, there are two needles with distance $O(1/(p^2 \cdot t))$.

Asymmetric Disjointness

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

There are k players with each of them holds a (random) set $S_i \subseteq [n]$ of size at most s_i . It is promised that either these sets are pairwise disjoint or have a unique intersection

Asymmetric Disjointness

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

There are k players with each of them holds a (random) set $S_i \subseteq [n]$ of size at most s_i . It is promised that either these sets are pairwise disjoint or have a unique intersection

Theorem (Lovett-Z)

Let Π be a randomized protocol that solves the asymmetric disjointness. Let c_i be the communication bits by the i -th player. Then we have that,

$$\sum_{i \in [k]} \frac{c_i}{s_i} = \Omega(1).$$

Needle Bounds from Asymmetric Disjointness

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Theorem (Lovett-Z)

Any algorithm that distinguishes the needle distribution needs $t = \Omega(1/(s \cdot p^2 \cdot \log n))$ samples.

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Information Complexity Approaches

The Needle Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

Let $n > 1$ be a large integer and let $p > 0$ be a parameter.

- **Uniform distribution D_0 :** each sample is uniformly sampled from $[n]$.
- **Needle distribution D_1 :** Sample a needle x . Each sample equals x with probability p , and uniformly otherwise.

The Needle Problem

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

Let $n > 1$ be a large integer and let $p > 0$ be a parameter.

- **Uniform distribution D_0** : each sample is uniformly sampled from $[n]$.
- **Needle distribution D_1** : Sample a needle x . Each sample equals x with probability p , and uniformly otherwise.
- **Local needle distribution D^S** : Sample a needle x . Each sample **in S** equals $x \in S$ with probability p , and uniformly otherwise.

Proof Strategy

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- We consider $|S| \approx 2 \cdot p \cdot t$

Proof Strategy

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- We consider $|S| \approx 2 \cdot p \cdot t$
- A half of elements from S are the needle

Proof Strategy

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- We consider $|S| \approx 2 \cdot p \cdot t$
- A half of elements from S are the needle
- $\mathbf{D}_1 = \sum_S \alpha_S \cdot \mathbf{D}^S$, where $\sum_S \alpha_S = 1$.

Proof Strategy

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- We consider $|S| \approx 2 \cdot p \cdot t$
- A half of elements from S are the needle
- $\mathbf{D}_1 = \sum_S \alpha_S \cdot \mathbf{D}^S$, where $\sum_S \alpha_S = 1$.
- If \mathcal{A} distinguishes \mathbf{D}_0 and \mathbf{D}_1 , then it distinguishes \mathbf{D}_0 and \mathbf{D}^S for many S .

Proof Strategy

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- We consider $|S| \approx 2 \cdot p \cdot t$
- A half of elements from S are the needle
- $\mathbf{D}_1 = \sum_S \alpha_S \cdot \mathbf{D}^S$, where $\sum_S \alpha_S = 1$.
- If \mathcal{A} distinguishes \mathbf{D}_0 and \mathbf{D}_1 , then it distinguishes \mathbf{D}_0 and \mathbf{D}^S for many S .
- The information cost of distinguishing \mathbf{D}^S and \mathbf{D}^0 is $\Omega(1)$

Proof Strategy

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

- We consider $|S| \approx 2 \cdot p \cdot t$
- A half of elements from S are the needle
- $\mathbf{D}_1 = \sum_S \alpha_S \cdot \mathbf{D}^S$, where $\sum_S \alpha_S = 1$.
- If \mathcal{A} distinguishes \mathbf{D}_0 and \mathbf{D}_1 , then it distinguishes \mathbf{D}_0 and \mathbf{D}^S for many S .
- The information cost of distinguishing \mathbf{D}^S and \mathbf{D}^0 is $\Omega(1)$
- The information cost of distinguishing \mathbf{D}_1 and \mathbf{D}^0 is $\Omega(1/p^2)$

Information Complexity

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition (Braverman-Garg-Woodruff)

Let \mathcal{A} be a streaming algorithm. We define its information complexity by,

$$\text{IC}(\mathcal{A}, \mathbf{D}_0) := \sum_{i=1}^t \sum_{k=1}^i I(\mathbf{M}_i; \mathbf{X}_k \mid \mathbf{M}_{k-1})$$

Information Complexity

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition (Braverman-Garg-Woodruff)

Let \mathcal{A} be a streaming algorithm. We define its information complexity by,

$$IC(\mathcal{A}, \mathbf{D}_0) := \sum_{i=1}^t \sum_{k=1}^i I(\mathbf{M}_i; \mathbf{X}_k \mid \mathbf{M}_{k-1})$$

Lemma

$$IC(\mathcal{A}, \mathbf{D}_0) \leq t \cdot s$$

Local Information Complexity

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Definition

Let \mathcal{A} be a streaming algorithm and let $S = \{p_1, \dots, p_m\}$ be a set. We define the local information complexity by,

$$IC^S(\mathcal{A}, \mathbf{D}_0) := \sum_{i=1}^m \sum_{k=1}^i I(M_{p_{i+1}-1}; \mathbf{X}_{p_k} \mid M_{p_k-1}).$$

Lemma

- If \mathcal{A} distinguishes \mathbf{D}_0 and \mathbf{D}^S , then $IC^S(\mathcal{A}, \mathbf{D}_0) = \Omega(1)$
- $IC(\mathcal{A}, \mathbf{D}_0) \approx \mathbb{E}_S[IC^S(\mathcal{A}, \mathbf{D}_0)]/p^2$

Tight Needle Lower Bounds

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Theorem (Li-Wang-Z)

Any algorithm that distinguishes the needle distribution needs $t = \Omega(1/(s \cdot p^2))$ samples.

Tight Needle Lower Bounds

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Theorem (Li-Wang-Z)

Any algorithm that distinguishes the needle distribution needs $t = \Omega(1/(s \cdot p^2))$ samples.

Lower bounds can be extended to the multi-pass setting by a multi-pass information complexity notion.

Jiapeng Zhang

The Needle
Problem

Lower Bounds

Asymmetric
Disjointness

Information
Complexity
Approaches

Thank you!