

# Information Statistics Approach to Data Stream and Communication Complexity

IBM Almaden Research Center

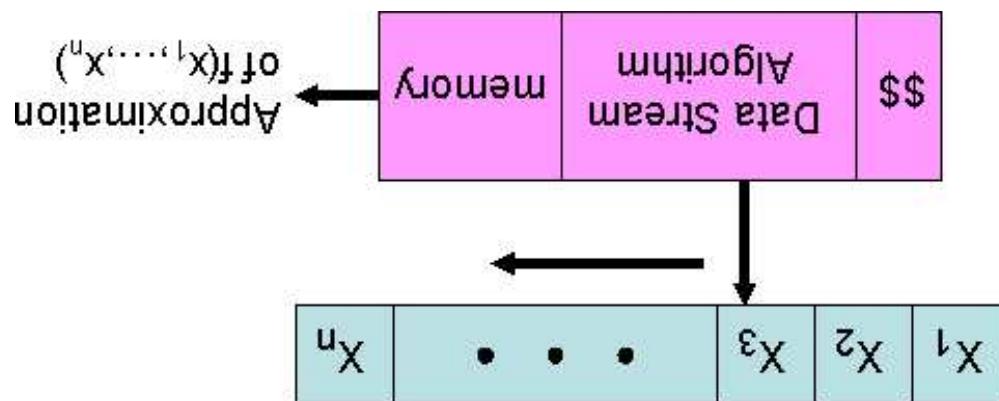
D. Sivakumar

Ravi Kumar

T.S. Jayram

Ziv Bar-Yossef

- Generalization:  $\ell$ -pass data stream algorithms
- Main measure of complexity: space
- Algorithms are allowed to be randomized and approximate
- Input arrives in a one-way stream in arbitrary order



The Data Stream Model  
[HR98, AMS96, FKS99]

- Motivation
  - Processing streams of IP packets
  - One-pass algorithms for large database relations
- Database
  - Processing search engine query logs
  - Web crawling
- Web Information Retrieval
  - Web crawling

Motivation

- Frequency statistics [FM83, AMS96, GT01, BKSO2, BJSTS02, CCF02, KPS02]
- Distances and norms [FKSV99, FS00, IOO]
- Histograms [GMP97, MVW00, GKS01, GGKMS02]
- Quantiles [ML98, MRL99, GK01]
- Clustering [GMMO00]
- Inversion counting [AJKS02]
- Triangle counting [BKSO2]

Algorithmic Results in Data Streams

$F_k$ , for any  $k > 2$ , requires polynomial space

Theorem 1 [This paper]

- $F_k$ , for  $k < 5$ , needs polynomial space
- $F_0, F_1, F_2$  can be approximated in logarithmic space

Theorem [Ailon, Matias, Szegedy '96]

- $F_k, k \geq 2$  = measure of  $k$ -wise collision probability
- $F_1 = n$
- $F_0 = \#$  of distinct data items

where,  $a_1, \dots, a_n \in [m]$ , and  $f_j = |\{i \in [n] \mid a_i = j\}|$

$$F_k(a_1, \dots, a_n) = \sum_{j=1}^m f_j^k$$

Example 1: Frequency Moments

$L^p$ , for  $p > 2$ , requires polynomial space, even for a constant number of passes over the input.

Theorem 2 [This paper]

- Any one pass algorithm for  $L^p$ , for  $p > 2$ , requires polynomial space [Feigenbaum et al. '99, Fong, Strauss '00, Indyk '00]
- $L^p$ , for any  $0 < p \leq 2$ , can be approximated in poly-logarithmic space [Saks, Sun '02]

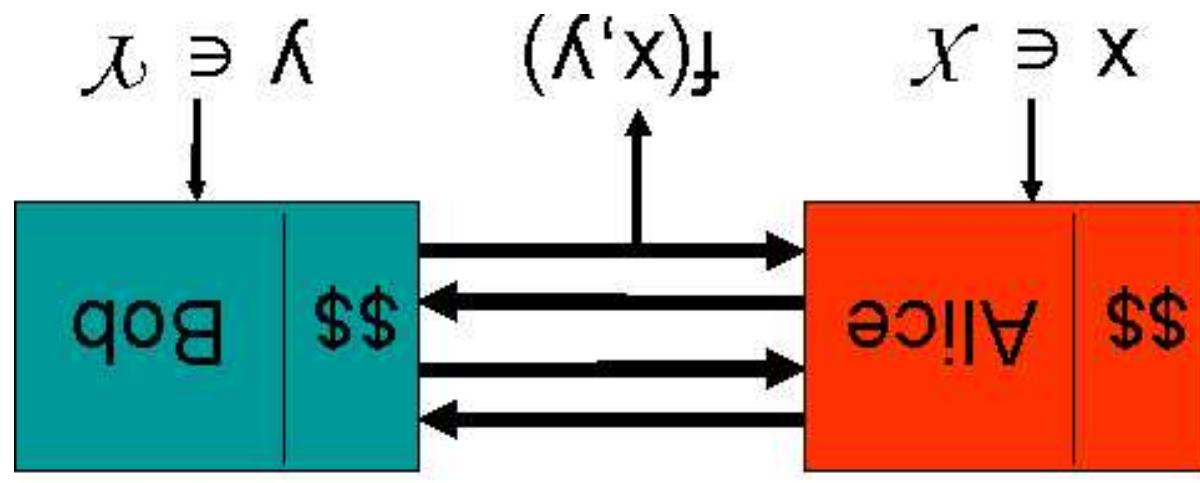
Theorem

where the entries of  $\underline{a}, \underline{b} \in [m]^n$  are given in arbitrary order

$$L^p(\underline{a}, \underline{b}) \stackrel{\text{def}}{=} \|\underline{a} - \underline{b}\|^p$$

Example 2:  $L^p$  Distance

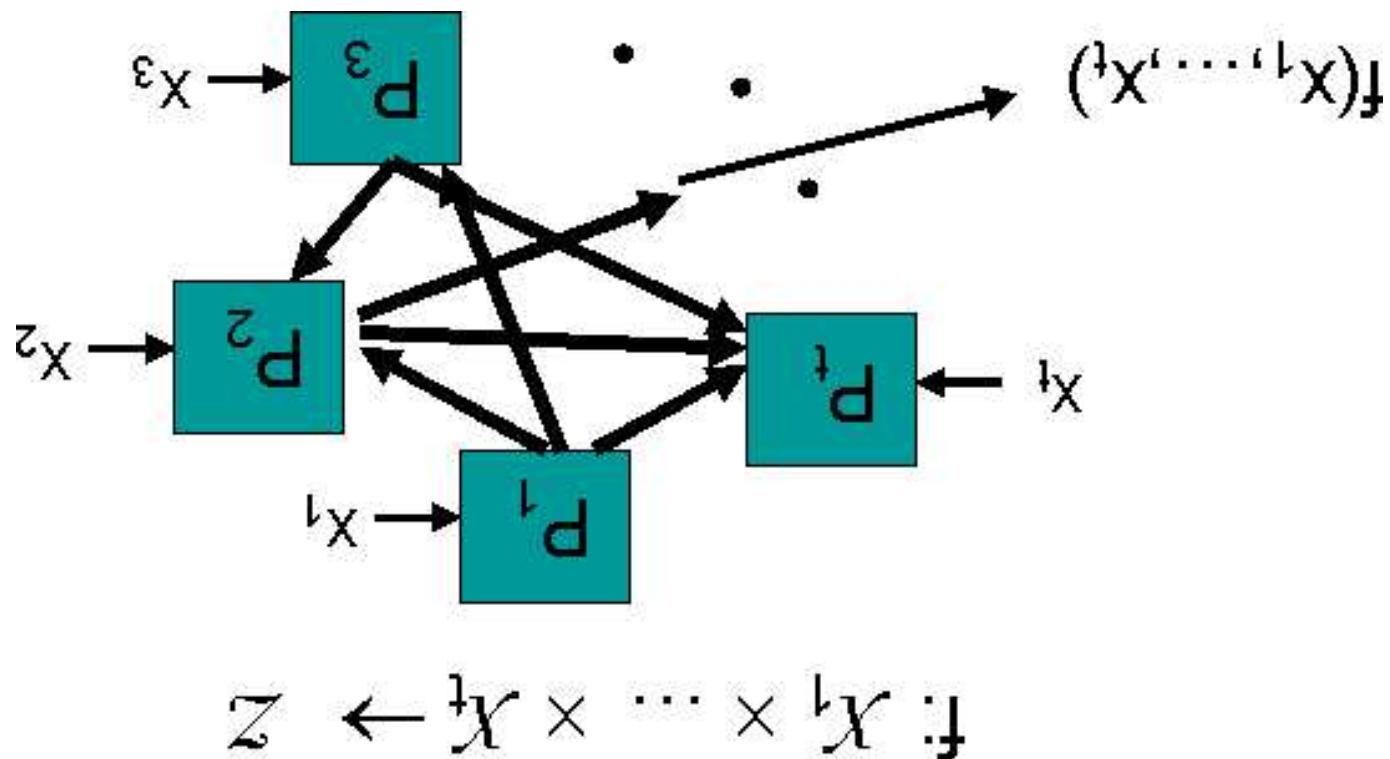
$R_\delta(f) =$  Randomized communication complexity of  $f$  with error  $\delta$



$$f: X \times \mathcal{Y} \rightarrow Z$$

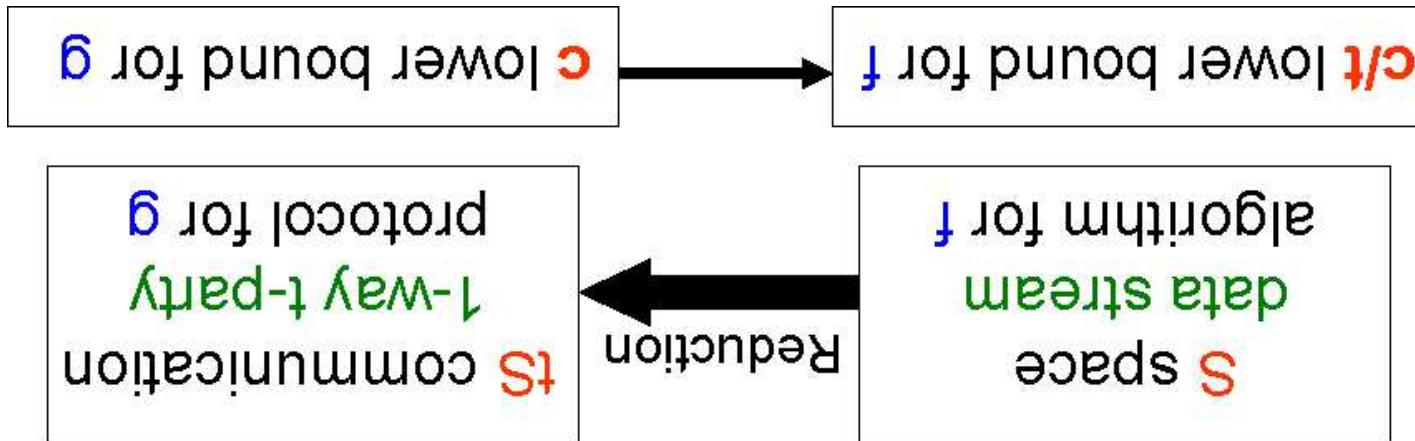
[Yao '79]

Communication Complexity



Multi-Party Communication Complexity  
("number in the hand")

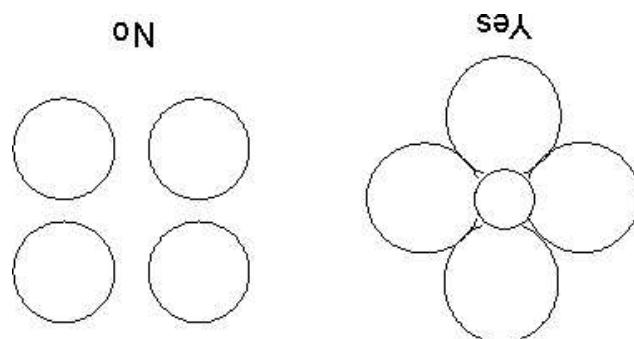
- Applications of Communication Complexity
- Circuit depth [Karchmer, Wigderson '88, Razborov '90, Raz, Wigderson '90]
  - Time-space tradeoffs [Chandra, Furst, Lipton '83, Alon, Maass '86, Babai, Nisan, Szegedy '89]
  - VLSI theory [Leengauer '90]
  - Decision tree complexity [Nisan '93]
  - Data structure complexity [Miltersen '95, Miltersen et al. '95]
  - Pseudorandom generators for logspace [Babai, Nisan, Szegedy '89, Impagliazzo, Nisan, Wigderson '94]
  - Computational economics [Nisan, Segal '01, Deng, Papadimitriou, Safra '02]
  - Data stream space complexity [Alon, Matias, Szegedy '96, Saks, Sun '02]



Easy reduction:



One-Way Communication Complexity



Input:  $S_1, \dots, S_t \subseteq [n]$  are either pairwise disjoint or mutually intersecting,  $\text{Disj}^{n,t}(S_1, \dots, S_t) = 1$  iff  $\bigcup_i S_i \neq \emptyset$ .

Multi-party set-disjointness



Input:  $S, T \subseteq [n]$ ,  $\text{Disj}^{n,2}(S, T) = 1$  iff  $S \cup T \neq \emptyset$ .

Set-disjointness

Example 1: Set-Disjointness

$F_k$  requires  $n^{1-2/k}$  space in the data stream model

Corollary

Using the reduction from  $\text{DISJ}^{n,n^{1/k}}$  to  $F_k$  [AMS96]:

- $R_{1\text{-way}}^g(\text{DISJ}^{n,t}) = \mathcal{O}(n/t^{1+\epsilon})$ , for any  $\epsilon < 0$
- $R_g^g(\text{DISJ}^{n,t}) = \mathcal{O}(n/t^2)$

Theorem 3 [This paper]

$R_{\text{Sim}}^g(\text{DISJ}^{n,t}) = \Theta(n/t)$   
 $R_g^g(\text{DISJ}^{n,t}) = \mathcal{O}(n/t_4)$   
 $R_g^g(\text{DISJ}^{n,2}) = \mathcal{O}(n)$

[Kalyanasundaram, Schnitger '87, Razborov '90]  
[Alon, Matias, Szegedy '96]  
[B, Jayram, Kumar, Sivakumar '02]

Example 1: Set-Disjointness (cont.)

Theorem

stream model, even for a constant number of passes over the input.  
 Estimating  $L_p$  to within  $n_\epsilon$  requires  $n^{1-4\epsilon-2/p}$  space in the data

**Corollary**

$$R_g(L^\infty) = \mathcal{O}(n/m^2)$$

**Theorem 4 [This Paper]**

$$R_{1\text{-way}}(L^\infty) = \Theta(n/m^2)$$

**Theorem [Saks, Sun '02]**

- $L^\infty(\underline{\hat{a}}, \underline{\hat{b}}) \geq m$  (NO instance)
- $L^\infty(\underline{\hat{a}}, \underline{\hat{b}}) \leq 1$  (YES instance)

$\underline{\hat{a}}, \underline{\hat{b}}$  are two vectors in  $[m]^n$  satisfying one of the following:

**Example 2:  $L^\infty$  Promise Problem**

## Outline of the Lower Bound Technique

1. Generalization of information complexity

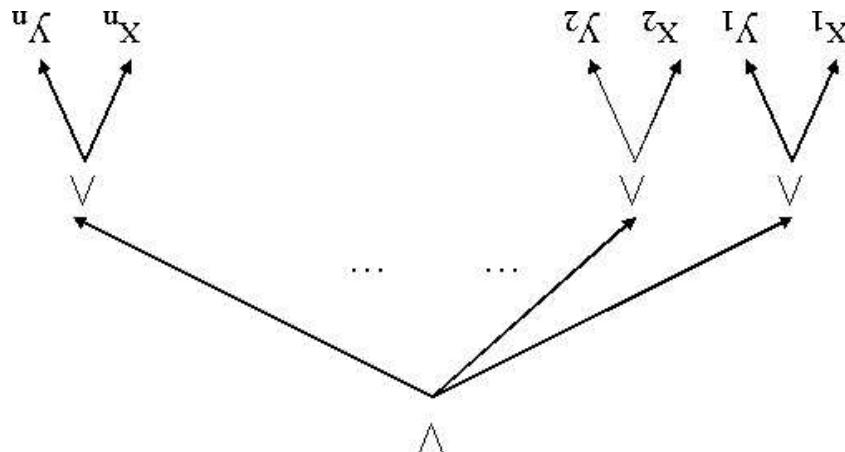
2. A **direct sum** theorem for information complexity

3. Statistical lower bounds on information complexity of "primitive"

functions

- Use information theory
- Cannot use communication length

Is a protocol for  $\text{DISJ}_{n,2}$  the "sum" of  $n$  protocols for AND?



where  $\underline{x}, \underline{y} \in \{0, 1\}^n$  are the characteristic vectors of the two sets

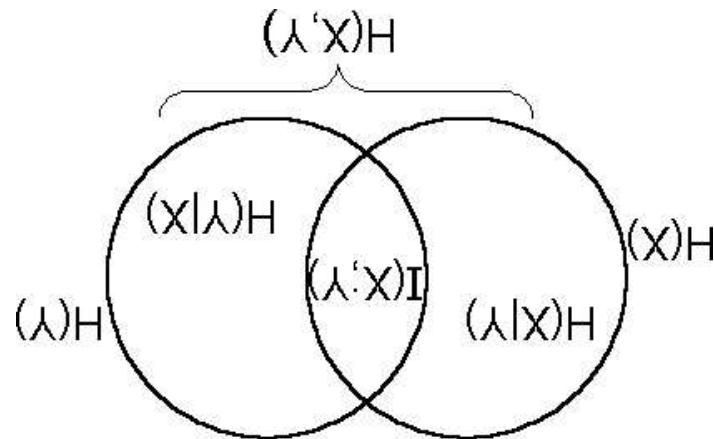
$$\text{DISJ}_{n,2}(\underline{x}, \underline{y}) \stackrel{\text{def}}{=} \bigwedge_{i=1}^n (x_i \vee y_i)$$

Example

The Direct Sum Question

$H(X, Y) \leq H(X) + H(Y)$ . Equality if  $X, Y$  are independent.

**Proposition [Subadditivity of entropy]**



$$(X | X) H - (X) H = (X | X) H - (X) H \stackrel{\text{def}}{=} (X; X) I$$

- $H(X | Y)$  – the amount of uncertainty left in  $X$  after knowing  $Y$
- $H(X)$  – the amount of “uncertainty” in  $X$

## Entropy and Mutual Information

**Proposition** For every  $\mu$ ,  $IC_{\mu,\delta}(f) \leq R_\delta(f)$

“Amount of information a protocol that computes  $f$  has to reveal about its inputs”

$$IC_{\mu,\delta}(f) \stackrel{\text{def}}{=} \min_{\Pi: \text{error}(\Pi) \leq \delta} I(X, Y; \Pi(X, Y))$$

**Definition [Information Complexity]**

$(X, Y) \sim \mu$  – a distribution over inputs of  $f$

[Chakrabarti, Shi, Wirth, Yao '01, Abolayev '93, Saks, Sun '02]

Information Complexity

**Proposition** For every  $u, D$ ,  $IC_{u,g}(f | D) \leq IC_{u,g}(f)$

$$IC_{u,g}(f | D) \stackrel{\text{def}}{=} \min_{H : \text{error}(H) \leq g} I(X, Y; H(X, Y) | D)$$

**Definition [Conditional Information Complexity]**

- Conditioned on  $\{D = d\}$ ,  $X, Y$  are independent
- Let  $D$  be a random variable such that  $\Pr[D = d] = \chi^d$

Express  $u$  as convex combination of product distributions  $\sum_d \chi^d u_d$

set-disjointness), we need  $u$  to be **non-product**.

- To create hard instances for some problems (e.g.,
- For direct sum, we need  $u$  to be **product**.

$u$  is **product** if  $X, Y$  are independent.

**Conditional Information Complexity**

- $\mu$  is concentrated on 0's of  $D_{Sj_n,2}$
  - Conditioned on  $\underline{D} = D^n$ ,  $\mu$  is product
  - $\mu$  is non-product
- $\mu \stackrel{\text{def}}{=} \nu_n$
- 
- If  $D = B$ , let  $X = 0$  and  $Y \in_R \{0, 1\}$
  - If  $D = A$ , let  $X \in_R \{0, 1\}$  and  $Y = 0$
  - $D \in_R \{A, B\}$
- $(X, Y) \sim \nu$ : a distribution on  $\{0, 1\}$  defined as follows:

Input Distribution for Set-Disjointness

$$(D \mid (\sum_{i=1}^j X_i, Y_i; \Pi(X_i, Y_i) \text{ AND } D)) \leq I(X_i, Y_i; \Pi(X_i, Y_i) \text{ AND } D)$$

2. Reduction step:

$$(D \mid (\sum_{i=1}^j X_i, Y_i; \Pi(X_i, Y_i) \text{ AND } D)) \leq I(X_i, Y_i; \Pi(X_i, Y_i) \text{ AND } D)$$

1. Decomposition step:

$$\text{Let } (X, Y) \sim \mu.$$

**Proof.** Let  $\Pi$  be a protocol for  $\text{DISJ}_{n,2}$ .

$$\text{Theorem } IC_{\mu, \delta}(\text{DISJ}_{n,2} \mid D) \geq n \cdot IC_{\mu, \delta}(\text{AND} \mid D)$$

Direct Sum for Information Complexity

$$\left( \underline{D} \mid (\underline{X}, \underline{X})_{\text{II}} : \underline{X}, \underline{X} \right) I \cdot \underline{\mathcal{Z}} =$$

and subadditivity of entropy)

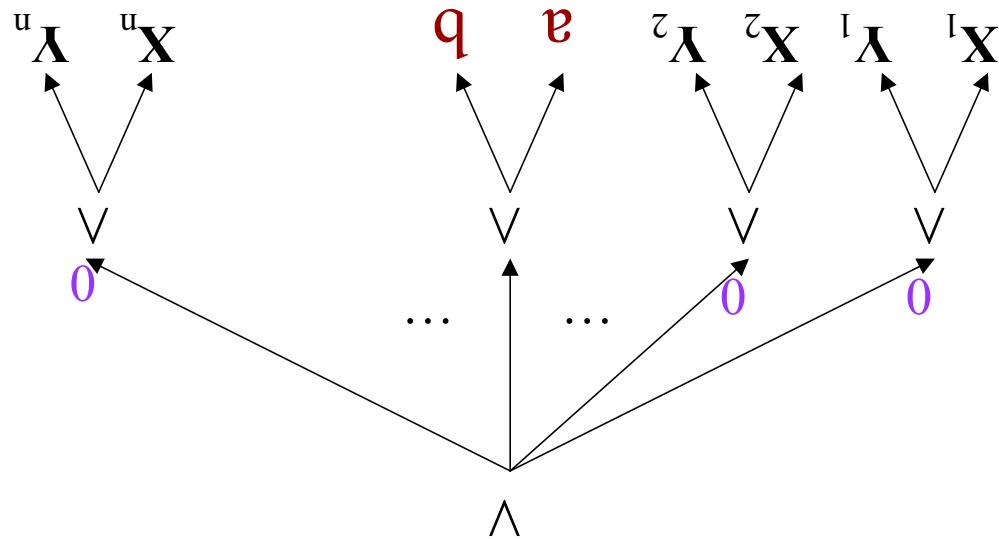
(by independence of  $\{\underline{X}, \underline{X}\}$ ,

$$\left( \underline{D}, (\underline{X}, \underline{X})_{\text{II}} \mid \underline{X}, \underline{X} \right) H \cdot \underline{\mathcal{Z}} - \left( \underline{D} \mid \underline{X}, \underline{X} \right) H \cdot \underline{\mathcal{Z}} \leq$$

$$\left( \underline{D}, (\underline{X}, \underline{X})_{\text{II}} \mid \underline{X}, \underline{X} \right) H - \left( \underline{D} \mid \underline{X}, \underline{X} \right) H =$$

$$\left( \underline{D} \mid (\underline{X}, \underline{X})_{\text{II}} : \underline{X}, \underline{X} \right) I$$

Proof of Decomposition Step



2. Create a protocol for computing  $\text{AND}(a, b)$

$$\sum_{\vec{D}} \Pr_{\vec{D}^{-j}} \left( \vec{D}^{-j} = \vec{d}^{-j} \mid I(X_j, Y_j; \Pi(\vec{X}, \vec{Y}) \mid \vec{D}) \cdot \prod_{i \neq j} I(X_i, Y_i; \Pi(\vec{X}, \vec{Y}) \mid \vec{D}_i) \right)$$

$$I(X_j, Y_j; \Pi(\vec{X}, \vec{Y}) \mid \vec{D}) \geq IC^{a, b}(\text{AND} \mid D)$$

Proof of Reduction Step

$$\left( \{U \in R \mid U \in \{0,1\} \right) = \frac{1}{2} \cdot [I(U; P(U, 0)) + I(U; P(U, 1))] =$$

$$[(B = D \mid (X, Y)P(X, Y)I + A = D \mid (X, Y)P(X, Y)I) =$$

$$I(X, Y \mid D)$$

Suppose  $P$  is a protocol for AND. Then,

- If  $D = B$ , let  $X = 0$  and  $Y \in R \setminus \{0,1\}$

- If  $D = A$ , let  $X \in R \setminus \{0,1\}$  and  $Y = 0$

- $D \in R \setminus \{A, B\}$

Recall the distribution  $\nu$ :

Lower Bound on  $IC_{\nu, g}(\text{AND} \mid D)$

$$\frac{1}{2} \cdot JS(P(1,0), P(0,1)) \leq$$

$$[((1,0)P,(0,0)P)SI + ((0,1)P,(0,0)P)(0,1)] \cdot \frac{2}{1} =$$

$$[((U,0)P;U)I + ((0,U)P;U)I] \cdot \frac{2}{1} =$$

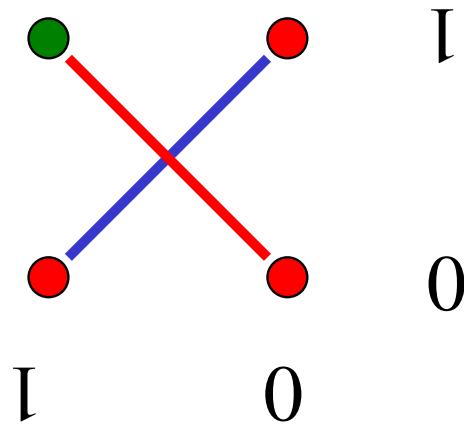
$$(D | (X,X)P; X,X)I$$

- Define  $JS(O^0, O^1) \stackrel{\text{def}}{=} I(U; O^U)$

- $U \in \{0, 1\}$

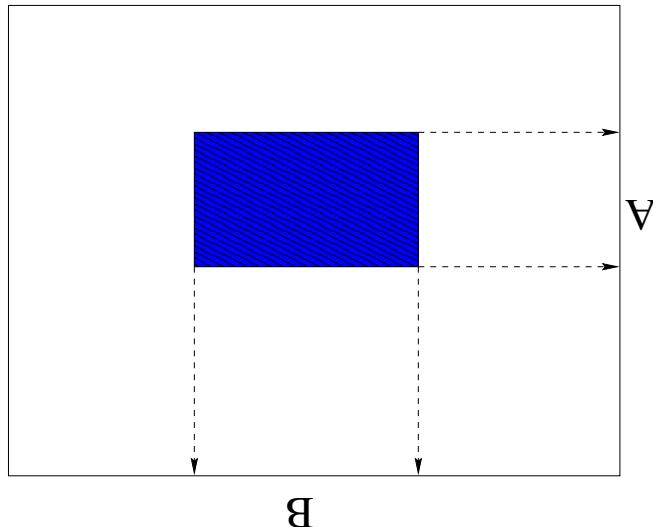
- $O^0, O^1$ : two distributions

Jensen-Shannon Divergence



- The **large** distance is on the other diagonal!
  - AND is 0 on both of these inputs
- If  $P$  computes AND, why should  $P(0, 1)$  be **far** from  $P(1, 0)$ ?

A Point to Ponder



Then,  $\Pi_{-1}(\tau)$  a combinatorial rectangle:

$$\Pi_{-1}(\tau) = \{(x, y) : \Pi(x, y) = \tau\}.$$

Fix a transcript  $\tau$

Let  $\Pi$  be a deterministic protocol.

Rectangle Property of Communication Complexity

$$\Pr[\Pi(x, y) = \tau] = d_\tau(x) \cdot b_\tau(y),$$

that

Then, there exist functions  $d_\tau : \mathcal{X} \rightarrow [0, 1]$  and  $b_\tau : \mathcal{Y} \rightarrow [0, 1]$  such

Fix a transcript  $\tau$

Let  $\Pi$  be a randomized protocol.

A Probabilistic Analog

## Hellinger Distance

$$\begin{aligned} \frac{1}{2} \cdot h_2(P(0,1), P(1,0)) &\leq \\ \frac{1}{2} \cdot [h_2(P(0,0), P(0,1)) + h_2(P(0,0), P(1,0))] &\leq \\ \frac{1}{2} \cdot [JS(P(0,0), P(0,1)) + JS(P(0,0), P(1,0))] &= \end{aligned}$$

$(D \mid (X, X)_D : X, X) I$

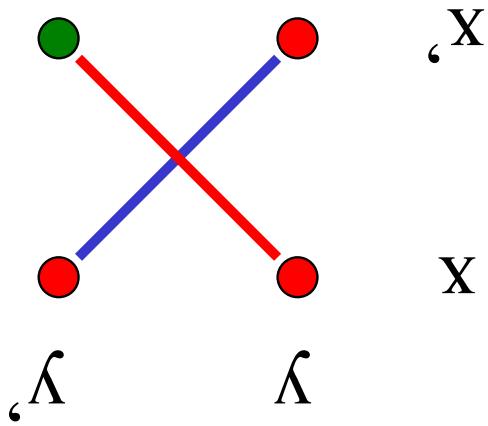
- $h(\cdot, \cdot)$  is a metric

- $JS(P, Q) \leq h_2(P, Q)$

$$h_2(P, Q) = \sqrt{\sum_{x \in \mathcal{X}} P(x)Q(x)}$$

Let  $P$  and  $Q$  be two probability distributions

$$\begin{aligned}
 & ((\mathcal{Y}, x)_{\mathcal{D}}, (\mathcal{Y}, x)_{\mathcal{D}}) \cup = \\
 & \underline{(\perp = (\mathcal{Y}, x)_{\mathcal{D}}) \sqcup} \cdot \underline{(\perp = (\mathcal{Y}, x)_{\mathcal{D}}) \sqcup} \wedge \perp \sqcup = \\
 & \underline{((\mathcal{Y})^{\perp} b \cdot (x)^{\perp} d \cdot (\mathcal{Y})^{\perp} b \cdot (x)^{\perp} d) \wedge \perp \sqcup} = \\
 & \underline{(\perp = (\mathcal{Y}, x)_{\mathcal{D}}) \sqcup} \cdot \underline{(\perp = (\mathcal{Y}, x)_{\mathcal{D}}) \sqcup} \wedge \perp \sqcup = \\
 & 1 - h_2((\mathcal{Y}, x)_{\mathcal{D}}, (\mathcal{Y}, x)_{\mathcal{D}})
 \end{aligned}$$



A Cut & Paste Lemma

$$R^g(\text{Disj}^{n,2}) \leq IC_{u,g}(\text{Disj}^{n,2} \mid D) \leq n \cdot IC_u(\text{AND} \mid D)$$

Therefore:

$$(Correctness of P) \leq \frac{1}{1 - 2\sqrt{\delta}}$$

$$(Cut \& Paste) = h_2(P(0,0), P(1,1))$$

$$\leq h_2(P(0,1), P(1,0))$$

$$(D \mid (X, X)P : X, X) I$$

Summary of Proof

divergences

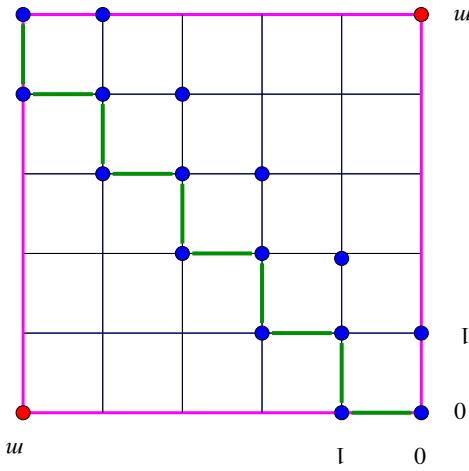
- Uses generalizations of the Hellinger distance – Rényi
  - Exploits the Markovian structure of one-way protocols.
- $\mathcal{O}(1/t_1^{\epsilon})$  bound for one-way communication:

- $t$  repeated applications of triangle inequality + cut & paste.
  - $\forall$ : pick  $D \in \mathbb{R}^{[t]}$  and then set  $X \in \mathbb{R}^{\{e_D, 0\}}$
- $\mathcal{O}(1/t_2)$  bound for general communication:
- Need a lower bound for  $t$ -bit AND.
  - Same direct sum argument

$t$ -Party Set-Disjointness

$$\frac{1}{2} [h_2(P(0,0), P(m,0)) + h_2(P(m,0), P(m,m))] \geq h_2(P(0,0), P(m,m))$$

A Pythagorean lemma:



$\Omega(1/m^2)$  lower bound:

decide whether  $|a - b| \leq 1$  or  $|a - b| \geq m$ .

- Need a lower bound for the difference problem: for  $a, b \in [m]$ ,
- Same direct sum argument

$L^\infty$  Promise Problem

- A powerful lower bound methodology in communication complexity
- Conclusions
- Method gives strong results even for **promise** problems
- particularly useful for data stream lower bounds
- Several novel ideas introduced:
- Conditional information complexity
  - Reduction to proving lower bounds for “simple” functions
  - Crisp connections between statistical distance measures and communication complexity

- $\Omega(n/t^{3/2})$  bound for  $t$ -party set-disjointness [Khot '02].  
[Jayram, Kumar, Sivakumar '02]
- Generalization to AND-OR trees of depth 3  
[Jayram '02]
- Distributional communication complexity lower bounds

## Subsequent Work

*Thank You!*