

Processing Data Streams



Andrew McGregor
University of Pennsylvania

Data Stream Model

[Morris '78] [Munro, Paterson '78] [Flajolet, Martin '85] [Alon, Matias, Szegedy '96]
[Henzinger, Raghavan, Rajagopalan '98] [Feigenbaum, Kannan, Strauss, Viswanathan '99]



Data Stream Model

[Morris '78] [Munro, Paterson '78] [Flajolet, Martin '85] [Alon, Matias, Szegedy '96]
[Henzinger, Raghavan, Rajagopalan '98] [Feigenbaum, Kannan, Strauss, Viswanathan '99]



Data Stream Model

[Morris '78] [Munro, Paterson '78] [Flajolet, Martin '85] [Alon, Matias, Szegedy '96]
[Henzinger, Raghavan, Rajagopalan '98] [Feigenbaum, Kannan, Strauss, Viswanathan '99]



- Constraints: **Sequential** access to data
Small working memory
Fast processing of elements

Data Stream Model

[Morris '78] [Munro, Paterson '78] [Flajolet, Martin '85] [Alon, Matias, Szegedy '96]
[Henzinger, Raghavan, Rajagopalan '98] [Feigenbaum, Kannan, Strauss, Viswanathan '99]



- Constraints: **Sequential** access to data
Small working memory
Fast processing of elements





Single Pass

e.g. network monitoring

Multi-Pass

e.g. huge log files

PASSES



Streams

Poly-log

e.g. statistics

Sub-linear

e.g. functions

SPACE

Single Pass

e.g. network monitoring

Multi-Pass

e.g. huge log files

PASSES

Time Series

e.g. stock prices

Adversarial

e.g. safety-critical

ORDER

Poly-log

e.g. statistics

Sub-linear

e.g. functions

SPACE



Streams

Single Pass

e.g. network monitoring

Multi-Pass

e.g. huge log files

PASSES

Time Series

e.g. stock prices

Adversarial

e.g. safety-critical

ORDER

Poly-log

e.g. statistics

Sub-linear

e.g. functions

SPACE



Streams

Single Pass

e.g. network monitoring

Multi-Pass

e.g. huge log files

PASSES

Numerical Data

e.g. quantiles, distances

Geometric Data

e.g. clusters, convex hulls

DATA

*Does it matter
how streams are
ordered?*

ORDER

Poly-log
e.g. statistics

Sub-linear
e.g. functions

SPACE



Streams

Single Pass
e.g. network monitoring

Multi-Pass
e.g. huge log files

PASSES

Numerical Data
e.g. quantiles, distances

Geometric Data
e.g. clusters, convex hulls

DATA

*Does it matter
how streams are
ordered?*

ORDER

Poly-log
e.g. statistics

Sub-linear
e.g. functions

SPACE



Streams

*How valuable is
an **extra pass**
over the data?*

PASSES

Numerical Data
e.g. quantiles, distances

Geometric Data
e.g. clusters, convex hulls

DATA

*Does it matter
how streams are
ordered?*

ORDER

*Can we create
small-space
sketches?*

SPACE



Streams

*How valuable is
an **extra pass**
over the data?*

PASSES

Numerical Data
e.g. quantiles, distances

Geometric Data
e.g. clusters, convex hulls

DATA

*Does it matter
how streams are
ordered?*

ORDER

*Can we create
small-space
sketches?*

SPACE



Streams

*How valuable is
an extra pass
over the data?*

PASSES

*What about a
stochastic-stream
or graph-stream?*

DATA

1. Thesis Overview
2. Some Algorithms
3. Some Lower-Bounds

I. Thesis Overview

Questions

Stochastic Streams

Graph Streams

Overview of Chapters

Questions

Questions

- Does it matter how streams are **ordered**?

Questions

- Does it matter how streams are **ordered**?



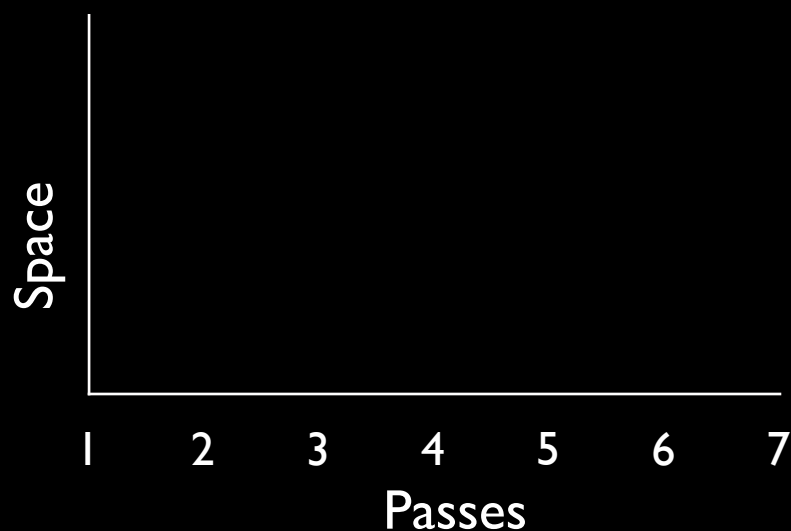
Yes!

Questions

- Does it matter how streams are **ordered**?
- How valuable is are extra **pass** over the data?

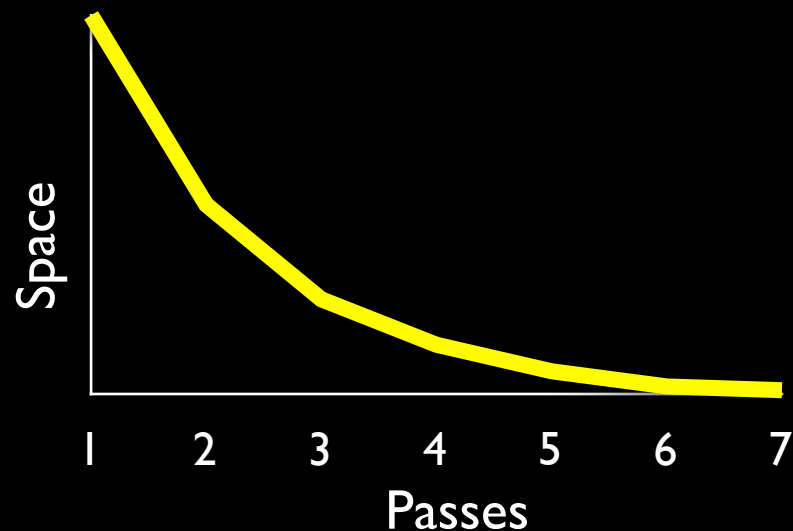
Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?



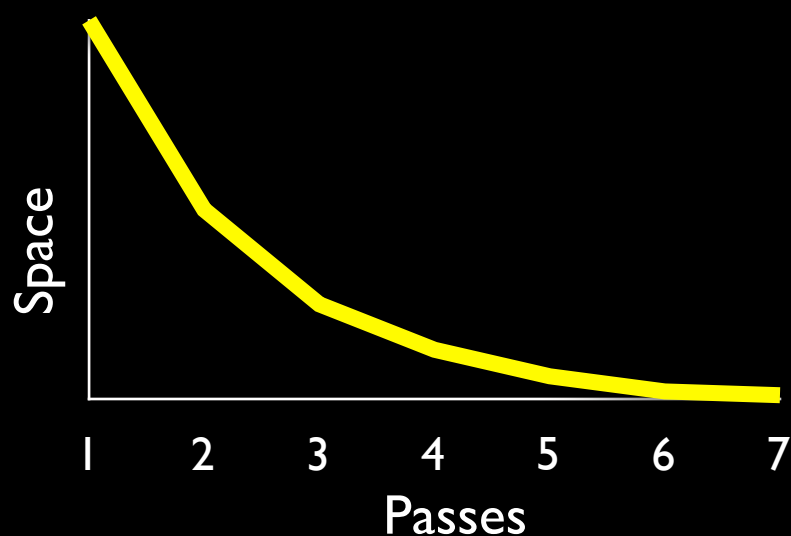
Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?

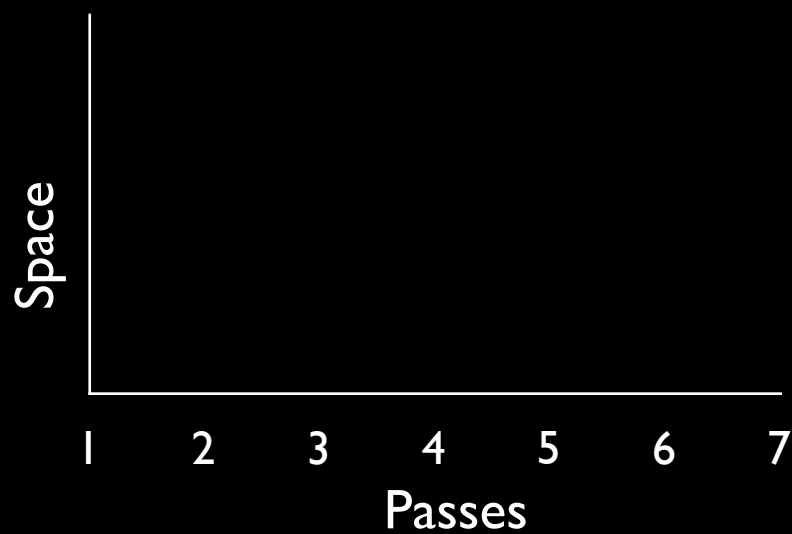


Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?

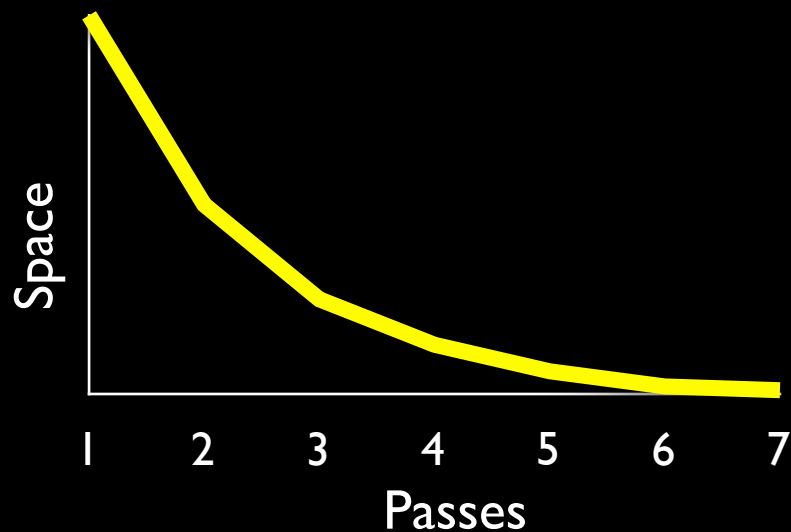


?

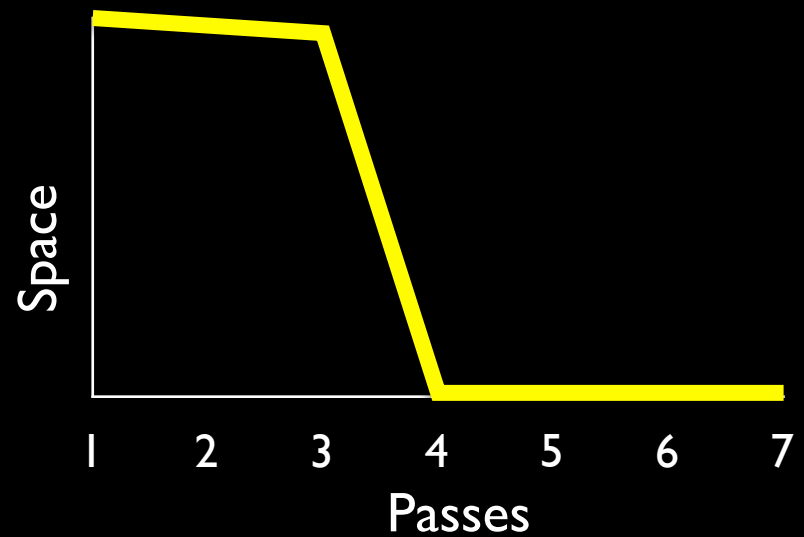


Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?



?



Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?
- Can we create **small-space** sketches?

Questions

- Does it matter how streams are **ordered**?
- How valuable is are extra **pass** over the data?
- Can we create **small-space** sketches?



New Stream

Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?
- Can we create **small-space** sketches?



Sketch of Old Stream

Can estimate L_1 distance.
[Indyk '00]



New Stream

Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?
- Can we create **small-space** sketches?



Sketch of Old Stream

Can estimate L_1 distance.
[Indyk '00]

What about the **relative distance** or **Hellinger distance**?



New Stream

Questions

- Does it matter how streams are **ordered**?
- How valuable is an extra **pass** over the data?
- Can we create **small-space** sketches?



Sketch of Old Stream

Can estimate L_1 distance.
[Indyk '00]

What about the **relative distance** or **Hellinger distance**?

No!



New Stream

Stochastic Streams

Stochastic Streams

- Stream: a_1, a_2, \dots, a_m where $a_j \in [n]$
 - Frequency:** $m_i =$ number of occurrences of i
 - Empirical Probability:** $p_i = m_i/m$

Stochastic Streams

- Stream: a_1, a_2, \dots, a_m where $a_j \in [n]$
 - Frequency:** $m_i =$ number of occurrences of i
 - Empirical Probability:** $p_i = m_i/m$
- Compute some function of (p_1, p_2, \dots, p_n)

Stochastic Streams

- Stream: a_1, a_2, \dots, a_m where $a_j \in [n]$
Frequency: $m_i =$ number of occurrences of i
Empirical Probability: $p_i = m_i/m$
- Compute some function of (p_1, p_2, \dots, p_n)

e.g. $\sum p_i^k$ [Alon, Matias, Szegedy '96]

$$\sum |p_i - q_i|$$

[Indyk '00]

$$\sum_{i \leq j} p_i \approx 1/2$$

[Greenwald, Khanna '01]

Stochastic Streams

- Stream: a_1, a_2, \dots, a_m where $a_j \in [n]$
Frequency: $m_i =$ number of occurrences of i
Empirical Probability: $p_i = m_i/m$
- Compute some function of (p_1, p_2, \dots, p_n)

e.g. $\sum -p_i \log p_i$ **Entropy?**

$\sum p_i \log(p_i/q_i)$ **Kullback-Leibler?**

$\sum (\sqrt{p_i} - \sqrt{q_i})^2$ **Hellinger?**

REAL

Stochastic Streams

- Who cares about the **empirical distribution?!**

REAL

Stochastic Streams

- Who cares about the **empirical distribution?!?**
- What can we infer about the **stream's source?**

REAL

Stochastic Streams

- Who cares about the **empirical distribution?!?**
- What can we infer about the **stream's source?**
- Assume the stream consists of **independent samples** or an evolving **Markov chain**.



REAL

Stochastic Streams

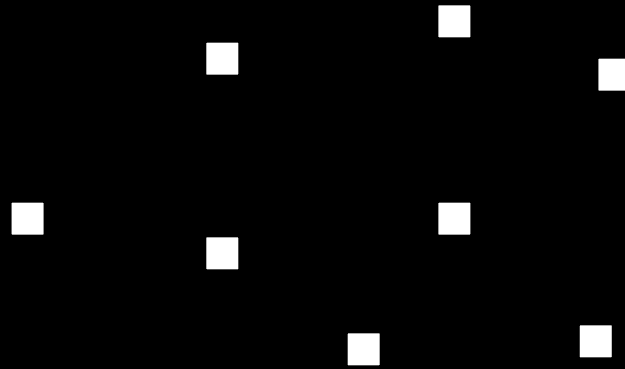
- Who cares about the **empirical distribution?!?**
- What can we infer about the **stream's source?**
- Assume the stream consists of **independent samples** or an evolving **Markov chain**.
- Can we **learn in small-space?**



Graph Streams

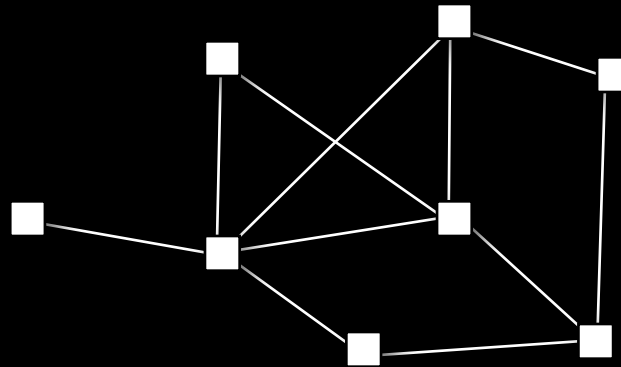
Graph Streams

- **Stream:** a_1, a_2, \dots, a_m where $a_j \in [n]^2$



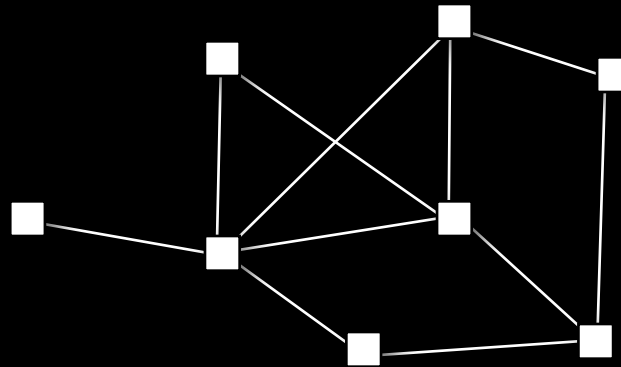
Graph Streams

- **Stream:** a_1, a_2, \dots, a_m where $a_j \in [n]^2$



Graph Streams

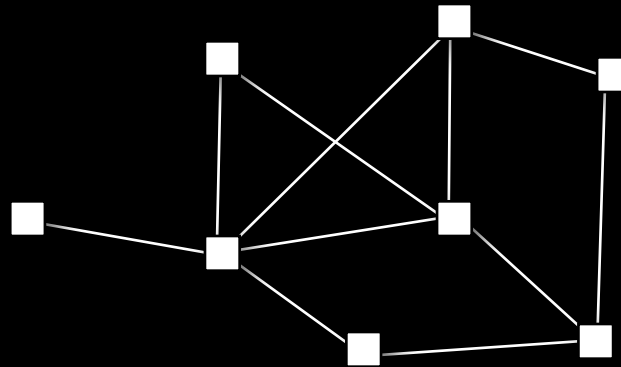
- **Stream:** a_1, a_2, \dots, a_m where $a_j \in [n]^2$



- E.g. **Web-graph** or **Call-graph**

Graph Streams

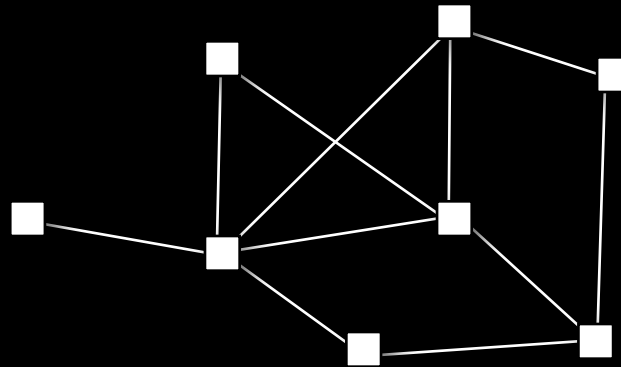
- **Stream:** a_1, a_2, \dots, a_m where $a_j \in [n]^2$



- E.g. **Web-graph** or **Call-graph**
- Previously: **Triangles and Common Neighborhoods**
[Bar-Yossef, Kumar, Sivakumar '02] [Buchsbaum, Giancarlo, Westbrook '03]

Graph Streams

- **Stream:** a_1, a_2, \dots, a_m where $a_j \in [n]^2$



- E.g. **Web-graph** or **Call-graph**
- Previously: Triangles and Common Neighborhoods
[Bar-Yossef, Kumar, Sivakumar '02] [Buchsbaum, Giancarlo, Westbrook '03]
- Connectivity, diameter, girth, matchings?

I. Oracles & Ordering

Does it matter how streams are ordered?



Quantiles
Histograms
Frequency Moments

Guha, M. (06)

Guha, M. (PODS 06)

Guha, M., Venkatasubramanian (SODA 06)

I. Oracles & Ordering

Does it matter how streams are ordered?



Quantiles
Histograms
Frequency Moments

Guha, M. (06)
Guha, M. (PODS 06)
Guha, M., Venkatasubramanian (SODA 06)

II. Entropy & Distances

What distances are “sketchable”?



Entropy
 f -Divergences
Bregman-Divergences

Guha, M., Indyk (07)
Chakrabarti, Cormode, M. (SODA 07)

I. Oracles & Ordering

Does it matter how streams are ordered?



Quantiles
Histograms
Frequency Moments

Guha, M. (06)
Guha, M. (PODS 06)
Guha, M., Venkatasubramanian (SODA 06)

II. Entropy & Distances

What distances are “sketchable”?

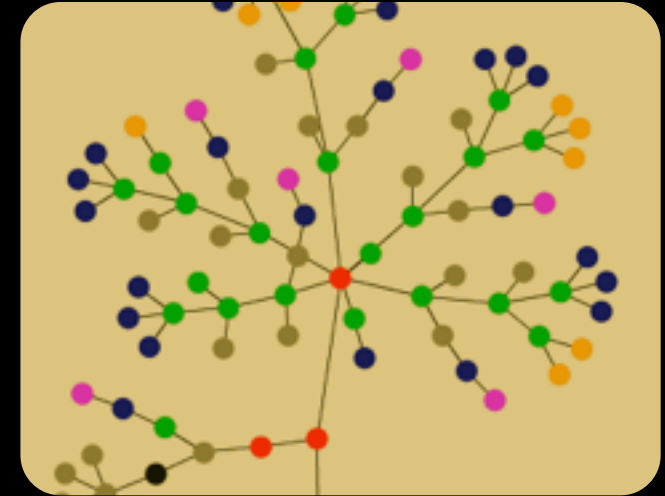


Entropy
 f -Divergences
Bregman-Divergences

Guha, M., Indyk (07)
Chakrabarti, Cormode, M. (SODA 07)

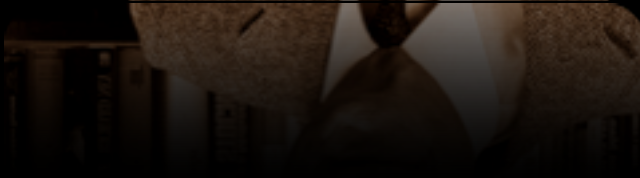
III. Graph Streams

How valuable is an extra pass?



Girth
Diameter
Large Matchings

M. (APPROX 05)
Feigenbaum et al. (SODA 05)
Feigenbaum et al. (ICALP 04)



II. Algorithms for Estimating Entropy



II. Algorithms for Estimating Entropy



III. Use Entropy for Lower Bounds

1. Thesis Overview
- 2. Some Algorithms**
3. Some-Lower Bounds

1. Thesis Overview

2. Some Algorithms

Estimating the entropy of a data stream

Higher-order entropy and random-walks

Entropy

- **Entropy:** $H(p) = \sum -p_i \lg p_i$

Upper-bound: $O(\varepsilon^{-2} \log m \log \delta^{-1})$

Lower-bound: $\Omega(\varepsilon^{-2} / \log^2 \varepsilon^{-1})$

[Guha et al. '06], [Chakrabarti et al. '06],
[Lall et al. 06], [Bhuvanagiri, Ganguly '06].

- **Generalizations:**

Higher-order entropy

Random-walk entropy

A:

● c a a b a a b a c a a a c b

A: ● c ● a ● a ● b ● a ● a ● b ● a ● c ● a ● a ● a ● a ● c ● b

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$

A: ● ● ● ● ● ● ● ● ● ● ● ● ● ●
 c a a b a a b a c a a a c b

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A: ● ● ● ● ● ● ● ● ● ● ● ● ● ●
 c a a b a a b a c a a a c b
Tags: 0.408

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A: ● ● ● ● ● ● ● ● ● ● ● ● ● ●
 c a a b a a b a c a a a c b
 Tags: 0.408

(min-tag, item, repeats) = (0.408, c, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815												

(min-tag, item, repeats) = (0.408, c, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217											

(min-tag, item, repeats) = (0.408, c, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217											

(min-tag, item, repeats) = (0.217, a, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191										

(min-tag, item, repeats) = (0.217, a, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191										

(min-tag, item, repeats) = (0.191, b, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770									

(min-tag, item, repeats) = (0.191, b, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082								

(min-tag, item, repeats) = (0.191, b, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082								

(min-tag, item, repeats) = (0.082, a, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366							

(min-tag, item, repeats) = (0.082, a, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228						

(min-tag, item, repeats) = (0.082, a, 1)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228						

(min-tag, item, repeats) = (0.082, a, 2)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549					

(min-tag, item, repeats) = (0.082, a, 2)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173				

(min-tag, item, repeats) = (0.082, a, 2)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173				

(min-tag, item, repeats) = (0.082, a, 3)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627			

(min-tag, item, repeats) = (0.082, a, 3)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627			

(min-tag, item, repeats) = (0.082, a, 4)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202		

(min-tag, item, repeats) = (0.082, a, 4)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202		

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	c	a	a	b	a	a	b	a	c	a	a	a	c	b
	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$
- **Basic-Estimator:** $X = f(j) - f(j - 1)$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

- **Basic-Estimator:** $X = f(j) - f(j - 1)$
- Correct in expectation so repeat and average

$$E[X] = \sum_{i \in [m]} \frac{m_i}{m} \sum_{j \in [m_i]} \frac{1}{m_i} (f(j) - f(j - 1))$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

- **Basic-Estimator:** $X = f(j) - f(j - 1)$
- Correct in expectation so repeat and average

$$3\left(1 + \frac{a}{E[X]}\right)^2 \epsilon^{-2} \ln(2\delta^{-1}) \frac{a + b}{a + E[X]} \text{ where } -a \leq X \leq b$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

(min-tag, item, repeats) = (0.082, a, 5)

- Estimate: $\frac{1}{m} \sum f(m_i)$ where $f(x) = x \lg \frac{m}{x}$
- Define a random variable j :

$$i \in_R A \text{ and } j \in_R [m_i]$$

- **Basic-Estimator:** $X = f(j) - f(j - 1)$
- Correct in expectation so repeat and average

$$O((1 + H^{-1})^2 \epsilon^{-2} \log m \log \delta^{-1})$$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

- If H is small, there exists a **majority element** x

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

- If H is small, there exists a **majority element** x
- Let A' be the sub-stream ignoring x

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

- If H is small, there exists a **majority element** x
- Let A' be the sub-stream ignoring x
- Define a random variable $j: i \in_R A'$ and $j \in_R [m_i]$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

- If H is small, there exists a **majority element** x
- Let A' be the sub-stream ignoring x
- Define a random variable $j: i \in_R A'$ and $j \in_R [m_i]$
- **Basic-Estimator:** $X = f(j) - f(j - 1)$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

- If H is small, there exists a **majority element** x
- Let A' be the sub-stream ignoring x
- Define a random variable $j: i \in_R A'$ and $j \in_R [m_i]$
- **Basic-Estimator:** $X = f(j) - f(j - 1)$
- Then, $H(p) = p_x \lg p_x^{-1} + (1 - p_x)E[X]$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

- If H is small, there exists a **majority element** x
- Let A' be the sub-stream ignoring x
- Define a random variable $j: i \in_R A'$ and $j \in_R [m_i]$
- **Basic-Estimator:** $X = f(j) - f(j - 1)$
- Then, $H(p) = p_x \lg p_x^{-1} + (1 - p_x) E[X]$

[Misra, Gries '82]

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
	c	a	a	b	a	a	b	a	c	a	a	a	c	b
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274

- If H is small, there exists a **majority element** x
- Let A' be the sub-stream ignoring x
- Define a random variable $j: i \in_R A'$ and $j \in_R [m_i]$
- **Basic-Estimator:** $X = f(j) - f(j - 1)$
- Then, $H(p) = p_x \lg p_x^{-1} + (1 - p_x) E[X]$

[Misra, Gries '82]
Constant!

Tags:

A:

0.408 **c** ●

0.815 **a** ●

0.217 **a** ●

0.191 **b** ●

0.770 **a** ●

0.082 **a** ●

0.366 **b** ●

0.228 **a** ●

0.549 **c** ●

0.173 **a** ●

0.627 **a** ●

0.202 **a** ●

0.549 **c** ●

0.274 **b** ●

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

(min-tag₁, item₁, repeats₁)

(min-tag₂, item₂, repeats₂)

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

(min-tag₁, item₁, repeats₁)

(min-tag₂, item₂, repeats₂)

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

(min-tag₁, item₁, repeats₁)

(min-tag₂, item₂, repeats₂)

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

$(\text{min-tag}_1, \text{item}_1, \text{repeats}_1) = (0.082, a, 5)$

$(\text{min-tag}_2, \text{item}_2, \text{repeats}_2) = (0.191, b, 2)$

A:	●	●	●	●	●	●	●	●	●	●	●	●	●	●
Tags:	0.408	0.815	0.217	0.191	0.770	0.082	0.366	0.228	0.549	0.173	0.627	0.202	0.549	0.274
	c	a	a	b	a	a	b	a	c	a	a	a	c	b

$(\text{min-tag}_1, \text{item}_1, \text{repeats}_1) = (0.082, a, 5)$

$(\text{min-tag}_2, \text{item}_2, \text{repeats}_2) = (0.191, b, 2)$

Update triples on seeing item y and tag t :

if $y = \text{item}_1$:

if $t < \text{min-tag}_1$ then $(\text{min-tag}_1, \text{item}_1, \text{repeats}_1) = (t, y, l)$ else repeats_1++

else:

if $y = \text{item}_2$ then repeats_2++

if $t < \text{min-tag}_1$:

$(\text{min-tag}_2, \text{item}_2, \text{repeats}_2) = (\text{min-tag}_1, \text{item}_1, \text{repeats}_1)$

$(\text{min-tag}_1, \text{item}_1, \text{repeats}_1) = (t, y, l)$

else:

if $t < \text{min-tag}_2$ then $(\text{min-tag}_2, \text{item}_2, \text{repeats}_2) = (t, y, l)$

Summary

Summary

- **Thm:** Single pass (ϵ, δ) -approx for entropy using $O(\epsilon^{-2} \log m \log \delta^{-1})$ space.

Summary

- **Thm:** Single pass (ϵ, δ) -approx for entropy using $O(\epsilon^{-2} \log m \log \delta^{-1})$ space.
- **Thm:** Higher-order entropy:

$$H_1 = \sum_i \frac{m_i}{m} \sum_j \frac{m_{ij}}{m_i} \lg \frac{m_i}{m_{ij}}$$

Needs $\Omega(n/\log n)$ space but can additive approx.

Summary

- **Thm:** Single pass (ϵ, δ) -approx for entropy using $O(\epsilon^{-2} \log m \log \delta^{-1})$ space.

- **Thm:** Higher-order entropy:

$$H_1 = \sum_i \frac{m_i}{m} \sum_j \frac{m_{ij}}{m_i} \lg \frac{m_i}{m_{ij}}$$

Needs $\Omega(n/\log n)$ space but can additive approx.

- **Thm:** Drunkards-walk entropy:

$$H_G = \frac{1}{2m} \sum_v d_v \lg d_v$$

Single pass (ϵ, δ) -approx using $O(\epsilon^{-2} \log^2 n \log^2 \delta^{-1})$ space.

1. Thesis Overview
2. Some Algorithms
- 3. Some Lower-Bounds**

1. Thesis Overview
2. Some Algorithms
- 3. Some Lower-Bounds**

Communication Complexity

Pointer Chasing: Multi-Value and Multi-Party

Constructing BFS Trees and Finding the Median

Communication Complexity

[Yao '78] [Nisan, Wigderson '93] [Jain, Radhakrishnan, Sen '03]

Communication Complexity

[Yao '78] [Nisan, Wigderson '93] [Jain, Radhakrishnan, Sen '03]



Alice

x



Bob

y

Communication Complexity

[Yao '78] [Nisan, Wigderson '93] [Jain, Radhakrishnan, Sen '03]



Alice

x

How many bits need to be sent for the final recipient to learn $f(x,y)$?



Bob

y

Communication Complexity

[Yao '78] [Nisan, Wigderson '93] [Jain, Radhakrishnan, Sen '03]



Alice
x

How many bits need to be sent for the final recipient to learn $f(x,y)$?



Bob
y

- **One-round** or **Multi-round**?
- **Two-party** or **Multi-party**?
- **Randomized** or **Deterministic**?



Alice

binary string x



Bob

binary string y



Alice

binary string x

Is Hamming distance $\Delta(x,y)$ less than $n/2$ or at least $n/2 + \sqrt{n}$?



Bob

binary string y



Alice

binary string x

Is Hamming distance $\Delta(x,y)$ less than $n/2$ or at least $n/2 + \sqrt{n}$?

Requires $\Omega(n)$ bits. [Indyk, Woodruff 03]



Bob

binary string y



Alice

binary string x

Is Hamming distance $\Delta(x,y)$ less than $n/2$ or at least $n/2 + \sqrt{n}$?

Requires $\Omega(n)$ bits. [Indyk, Woodruff 03]



Bob

binary string y

- Assume \exists an (ϵ, δ) -approx for entropy using S space

● ● ● ● ● ● ●
(a,x₁) (b,x₂) (c,x₃) (d,x₄) (e,x₅) (f,x₆) (g,x₇)



Alice

binary string x

Is Hamming distance $\Delta(x,y)$ less than $n/2$ or at least $n/2 + \sqrt{n}$?

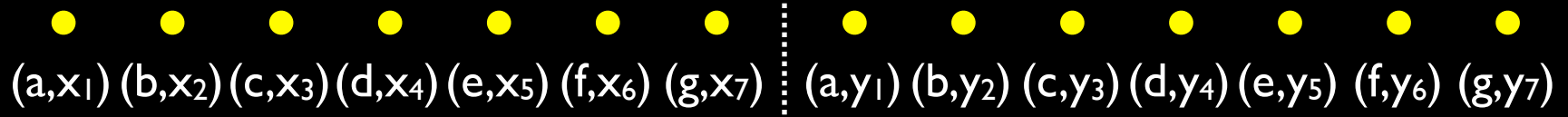
Requires $\Omega(n)$ bits. [Indyk, Woodruff 03]



Bob

binary string y

- Assume \exists an (ϵ, δ) -approx for entropy using S space



Alice

binary string x

Is Hamming distance $\Delta(x,y)$ less than $n/2$ or at least $n/2 + \sqrt{n}$?

Requires $\Omega(n)$ bits. [Indyk, Woodruff 03]



Bob

binary string y

- Assume \exists an (ϵ, δ) -approx for entropy using S space

(a,x₁) (b,x₂) (c,x₃) (d,x₄) (e,x₅) (f,x₆) (g,x₇) | (a,y₁) (b,y₂) (c,y₃) (d,y₄) (e,y₅) (f,y₆) (g,y₇)



Alice

binary string x

$$H(p) = \lg n + \frac{\Delta(x, y)}{n}$$



Bob

binary string y

- Assume \exists an (ϵ, δ) -approx for entropy using S space

(a,x₁) (b,x₂) (c,x₃) (d,x₄) (e,x₅) (f,x₆) (g,x₇) | (a,y₁) (b,y₂) (c,y₃) (d,y₄) (e,y₅) (f,y₆) (g,y₇)



Alice

binary string x

$$\begin{aligned} \Delta(x, y) \leq n/2 &\Leftrightarrow H(p) \leq \lg n + 1/2 \\ \Delta(x, y) \geq n/2 + \sqrt{n} &\Leftrightarrow H(p) \geq \lg n + 1/2 + 1/\sqrt{n} \end{aligned}$$



Bob

binary string y

- Assume \exists an (ϵ, δ) -approx for entropy using S space

(a,x₁) (b,x₂) (c,x₃) (d,x₄) (e,x₅) (f,x₆) (g,x₇) | (a,y₁) (b,y₂) (c,y₃) (d,y₄) (e,y₅) (f,y₆) (g,y₇)



Alice

binary string x



Bob

binary string y

$$\begin{aligned} \Delta(x, y) \leq n/2 &\Leftrightarrow H(p) \leq \lg n + 1/2 \\ \Delta(x, y) \geq n/2 + \sqrt{n} &\Leftrightarrow H(p) \geq \lg n + 1/2 + 1/\sqrt{n} \end{aligned}$$

- Assume \exists an (ϵ, δ) -approx for entropy using S space
- Setting $n = 3\epsilon^{-2} / \lg^2 \epsilon^{-1}$ ratio is more than $\frac{1 + \epsilon}{1 - \epsilon}$

(a,x₁) (b,x₂) (c,x₃) (d,x₄) (e,x₅) (f,x₆) (g,x₇) | (a,y₁) (b,y₂) (c,y₃) (d,y₄) (e,y₅) (f,y₆) (g,y₇)



Alice

binary string x



Bob

binary string y

$$\Delta(x, y) \leq n/2 \iff H(p) \leq \lg n + 1/2$$

$$\Delta(x, y) \geq n/2 + \sqrt{n} \iff H(p) \geq \lg n + 1/2 + 1/\sqrt{n}$$

MEMORY STATE OF ALGORITHM

- Assume \exists an (ϵ, δ) -approx for entropy using S space
- Setting $n = 3\epsilon^{-2} / \lg^2 \epsilon^{-1}$ ratio is more than $\frac{1 + \epsilon}{1 - \epsilon}$

(a,x₁) (b,x₂) (c,x₃) (d,x₄) (e,x₅) (f,x₆) (g,x₇) | (a,y₁) (b,y₂) (c,y₃) (d,y₄) (e,y₅) (f,y₆) (g,y₇)



Alice

binary string x



Bob

binary string y

$$\Delta(x, y) \leq n/2 \iff H(p) \leq \lg n + 1/2$$

$$\Delta(x, y) \geq n/2 + \sqrt{n} \iff H(p) \geq \lg n + 1/2 + 1/\sqrt{n}$$

MEMORY STATE OF ALGORITHM

- Assume \exists an (ϵ, δ) -approx for entropy using S space
- Setting $n = 3\epsilon^{-2} / \lg^2 \epsilon^{-1}$ ratio is more than $\frac{1 + \epsilon}{1 - \epsilon}$
- **Thm:** $S = \Omega(\epsilon^{-2} / \log^2 \epsilon^{-1})$

Pointer-Chasing



Alice

function

$f_A: [m] \rightarrow [m]$



Bob

function

$f_B: [m] \rightarrow [m]$

Pointer-Chasing



Alice

function

$f_A: [m] \rightarrow [m]$

Compute

$f_A(f_B(\dots (f_A(f_B(f_A(l)))) \dots))$

← k →



Bob

function

$f_B: [m] \rightarrow [m]$

Pointer-Chasing



Alice
function
 $f_A: [m] \rightarrow [m]$

Compute

$$f_A(f_B(\dots(f_A(f_B(f_A(l))))\dots))$$



If Bob speaks first:

k messages: $O(k \log m)$ bits.

$k-1$ messages: $\Omega(m)$ bits.

[Nisan, Wigderson '93]



Bob
function
 $f_B: [m] \rightarrow [m]$

Multi-Value Pointer-Chasing



Alice

d-valued function
 $f_A: [m] \rightarrow [m]$

Compute

$$f_A(f_B(\dots (f_A(f_B(f_A(I)))) \dots))$$



If Bob speaks first:

k messages: $O(k d^k \log m)$ bits.

k-1 messages: $\Omega(d m)$ bits.



Bob

d-valued function
 $f_B: [m] \rightarrow [m]$

Multi-Value Pointer-Chasing



Alice

d-valued function

$f_A: [m] \rightarrow [m]$

Compute

$f_A(f_B(\dots(f_A(f_B(f_A(I))))))\dots))$

\longleftrightarrow
k

If Bob speaks first:

k messages: $O(k d^k \log m)$ bits.

k-1 messages: $\Omega(d m)$ bits.



Bob

d-valued function

$f_B: [m] \rightarrow [m]$

- Proof: Relate to solving d copies of single-valued problem

Multi-Value Pointer-Chasing



Alice

d-valued function
 $f_A: [m] \rightarrow [m]$

Compute

$$f_A(f_B(\dots(f_A(f_B(f_A(l))))\dots))$$



k

If Bob speaks first:

k messages: $O(k d^k \log m)$ bits.

k-1 messages: $\Omega(d m)$ bits.



Bob

d-valued function
 $f_B: [m] \rightarrow [m]$

- Proof: Relate to solving d copies of single-valued problem

Appeal to “Direct-sum” Theorem

[Jain, Radhakrishnan, Sen 03]

Multi-Value Pointer-Chasing



Alice

d-valued function
 $f_A: [m] \rightarrow [m]$

Compute

$$f_A(f_B(\dots (f_A(f_B(f_A(I)))) \dots))$$

← k →

If Bob speaks first:

k messages: $O(k d^k \log m)$ bits.

k-1 messages: $\Omega(d m)$ bits.



Bob

d-valued function
 $f_B: [m] \rightarrow [m]$

- Proof: Relate to solving d copies of single-valued problem
Appeal to “Direct-sum” Theorem [Jain, Radhakrishnan, Sen 03]
- Can be reduced to the problem of finding BFS trees

Multi-Party Pointer-Chasing



Alice

function

$$f_A: [m] \rightarrow [m]$$



Bob

function

$$f_B: [m] \rightarrow [m]$$



Claude

function

$$f_C: [m] \rightarrow [m]$$

...



Zebedee

function

$$f_Z: [m] \rightarrow [m]$$

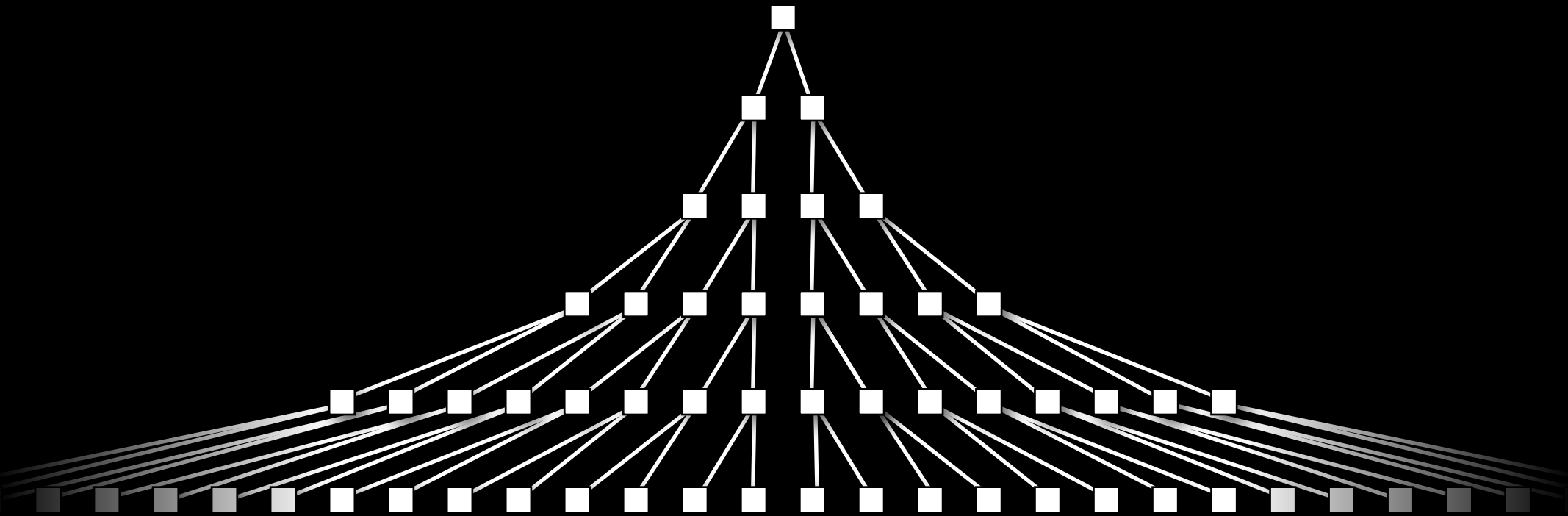
Compute $f_Z(\dots f_C(f_B(f_A(I))) \dots)$

Order of speaking “Z, ..., C, B, A”

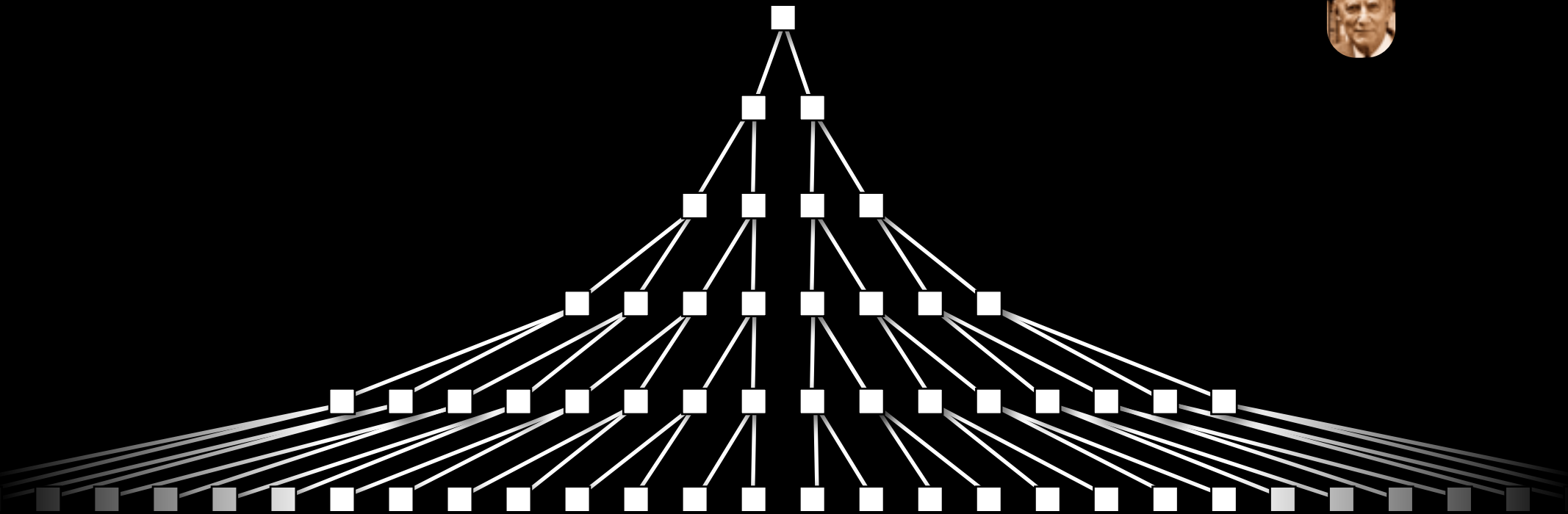
k rounds: $O(k \log m)$ bits.

$k-1$ rounds: $\Omega(m)$ bits.

Proof Sketch



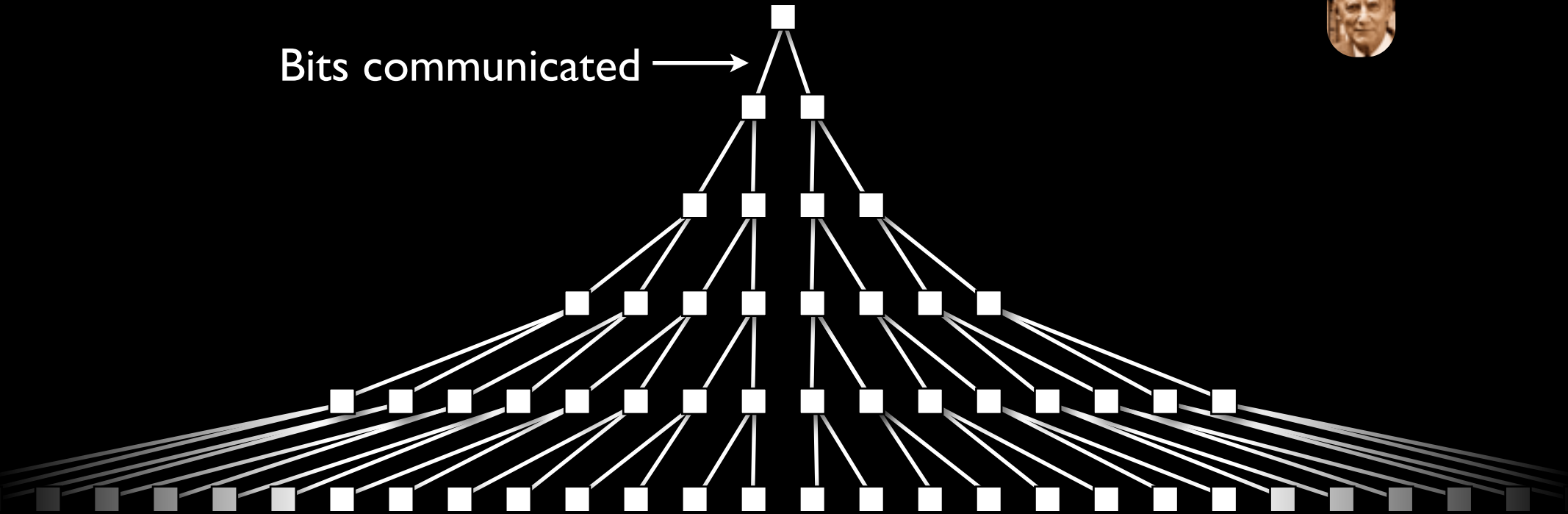
Proof Sketch



Proof Sketch

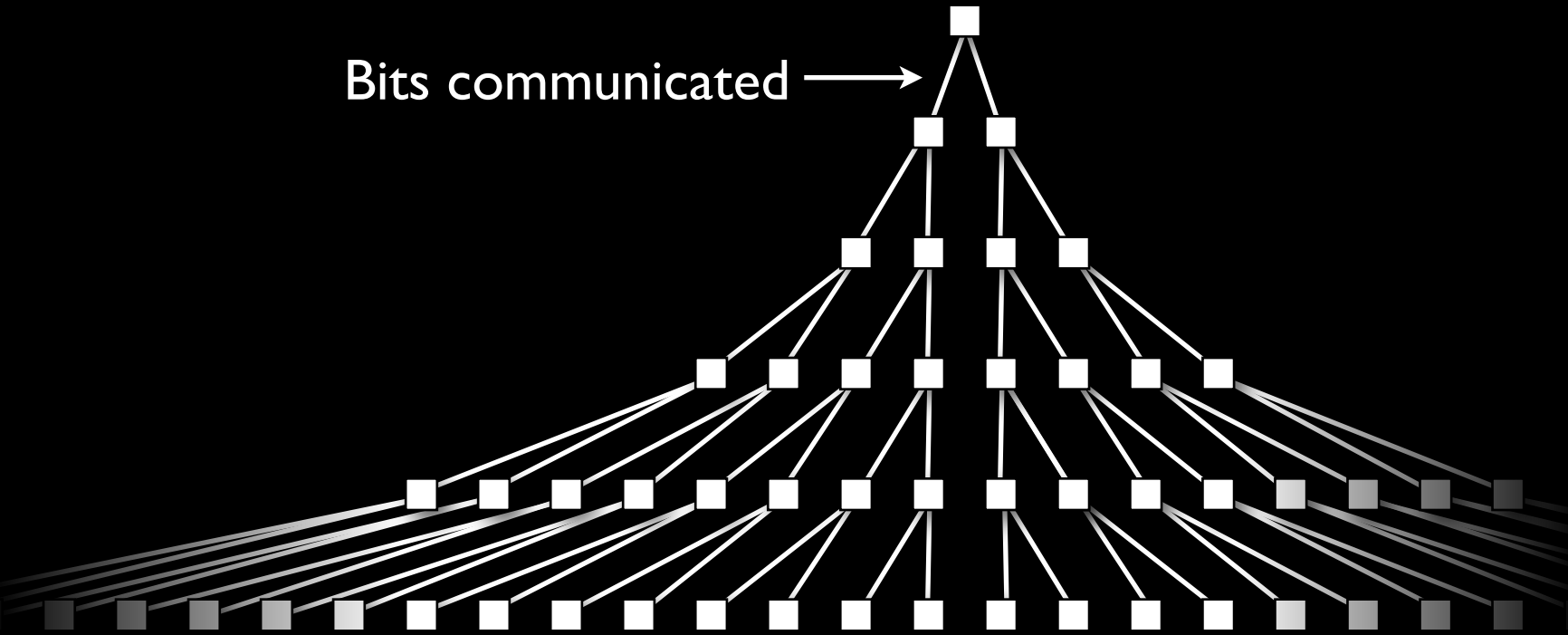


Bits communicated →

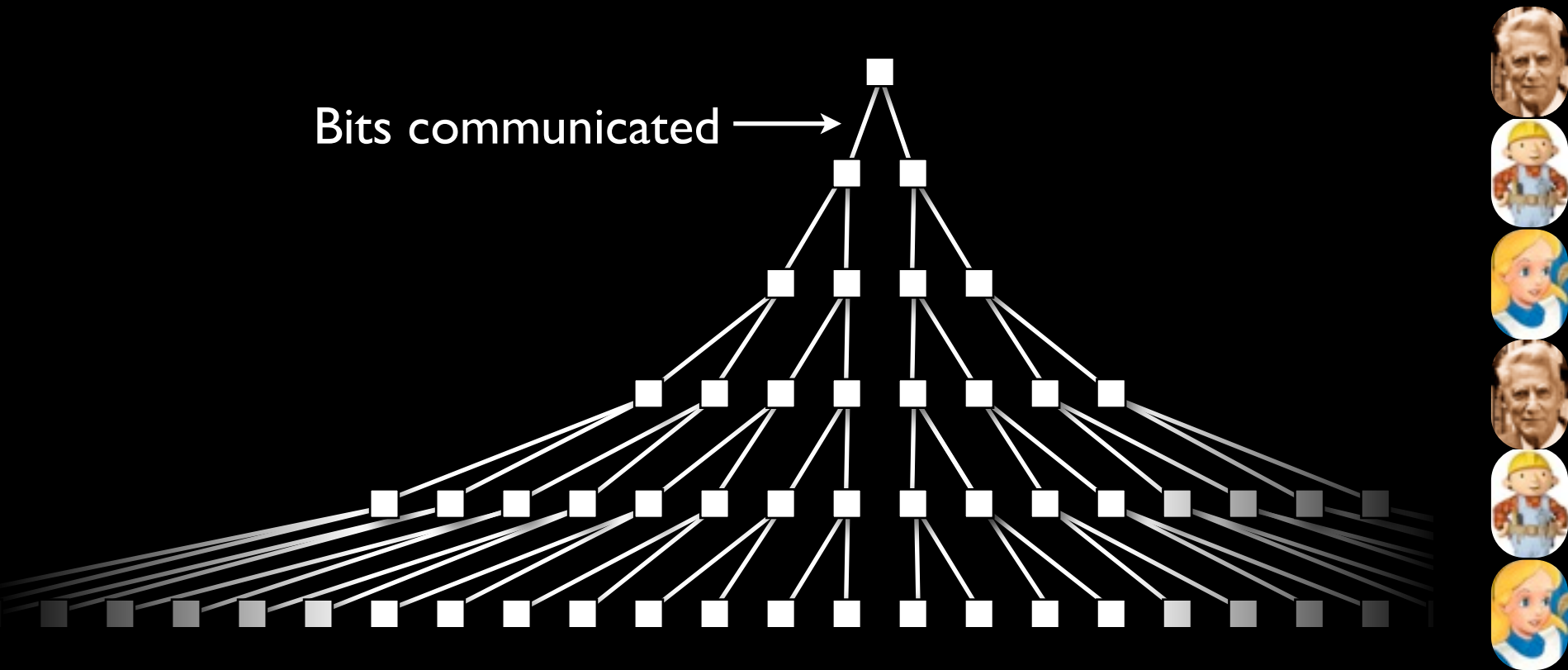


Proof Sketch

Bits communicated →

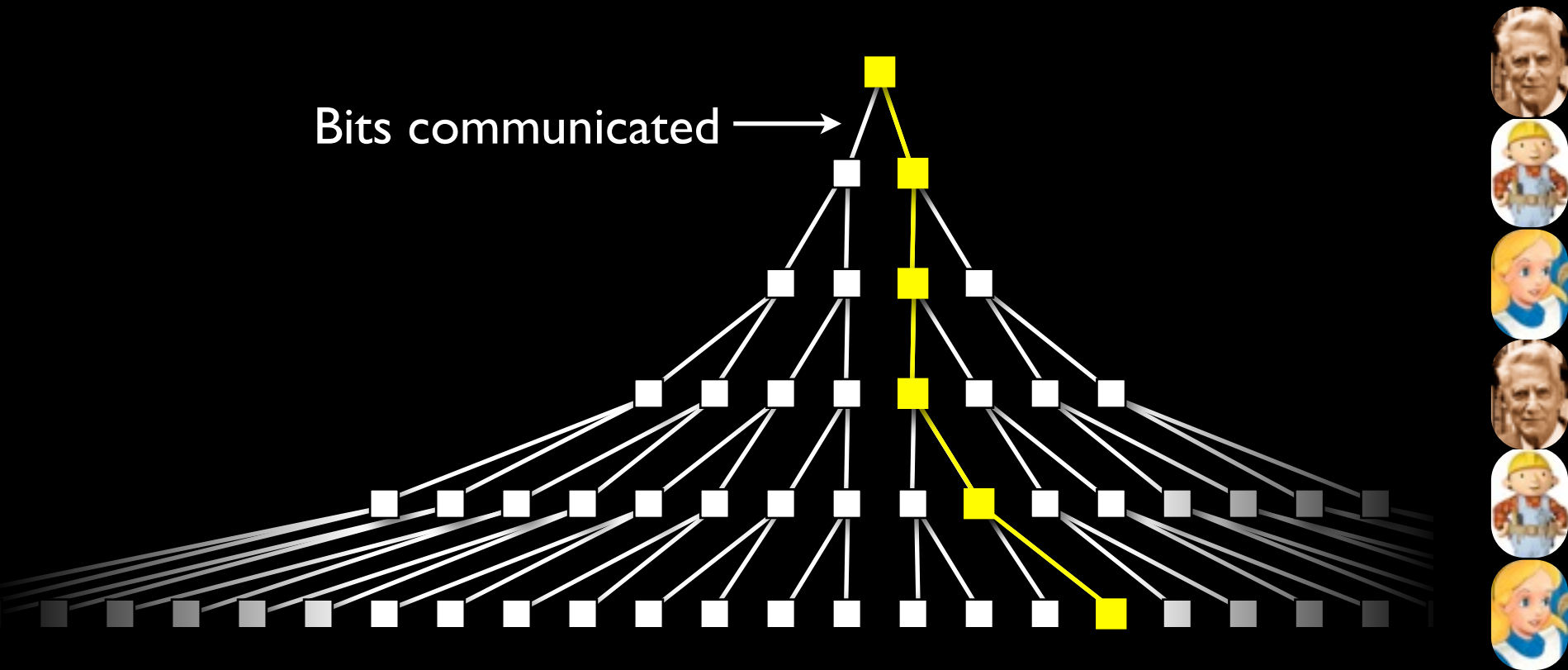


Proof Sketch



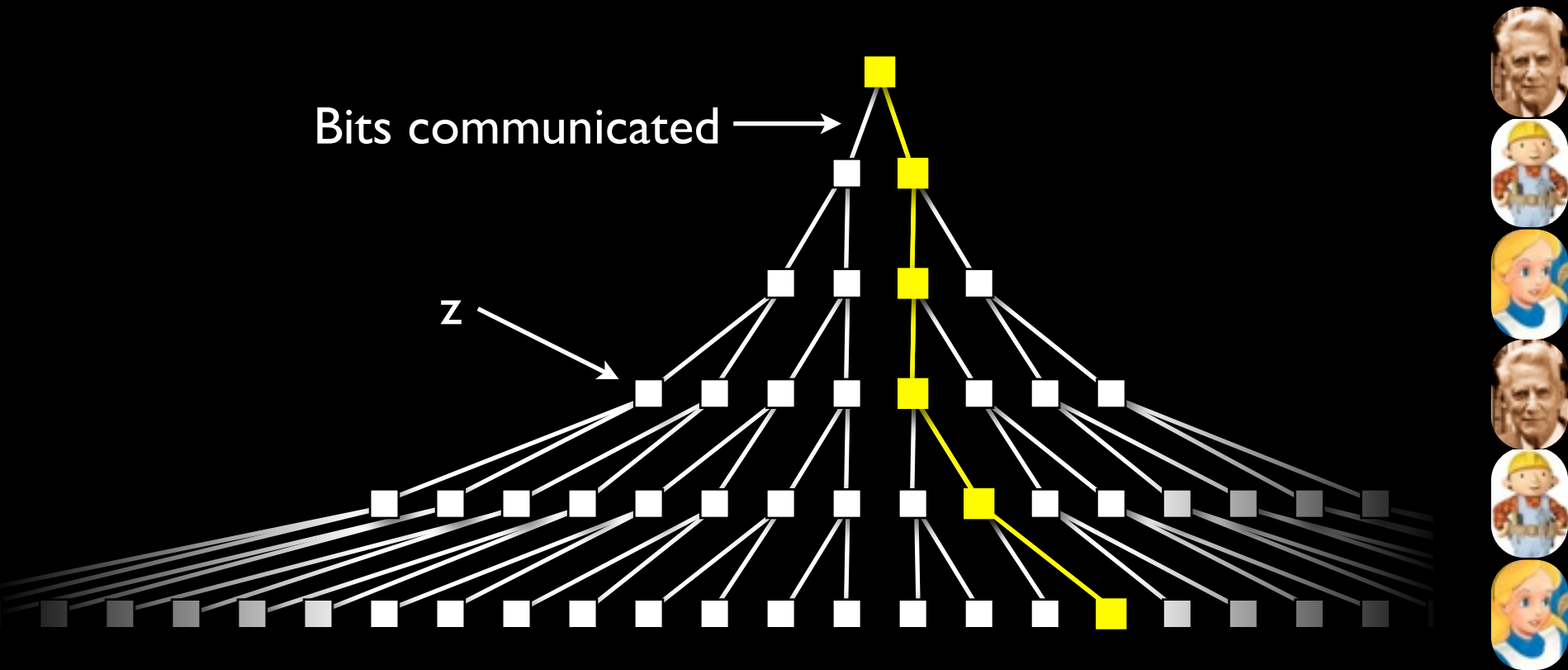
- Consider **deterministic** protocols and **random** f_A, f_B, f_C

Proof Sketch



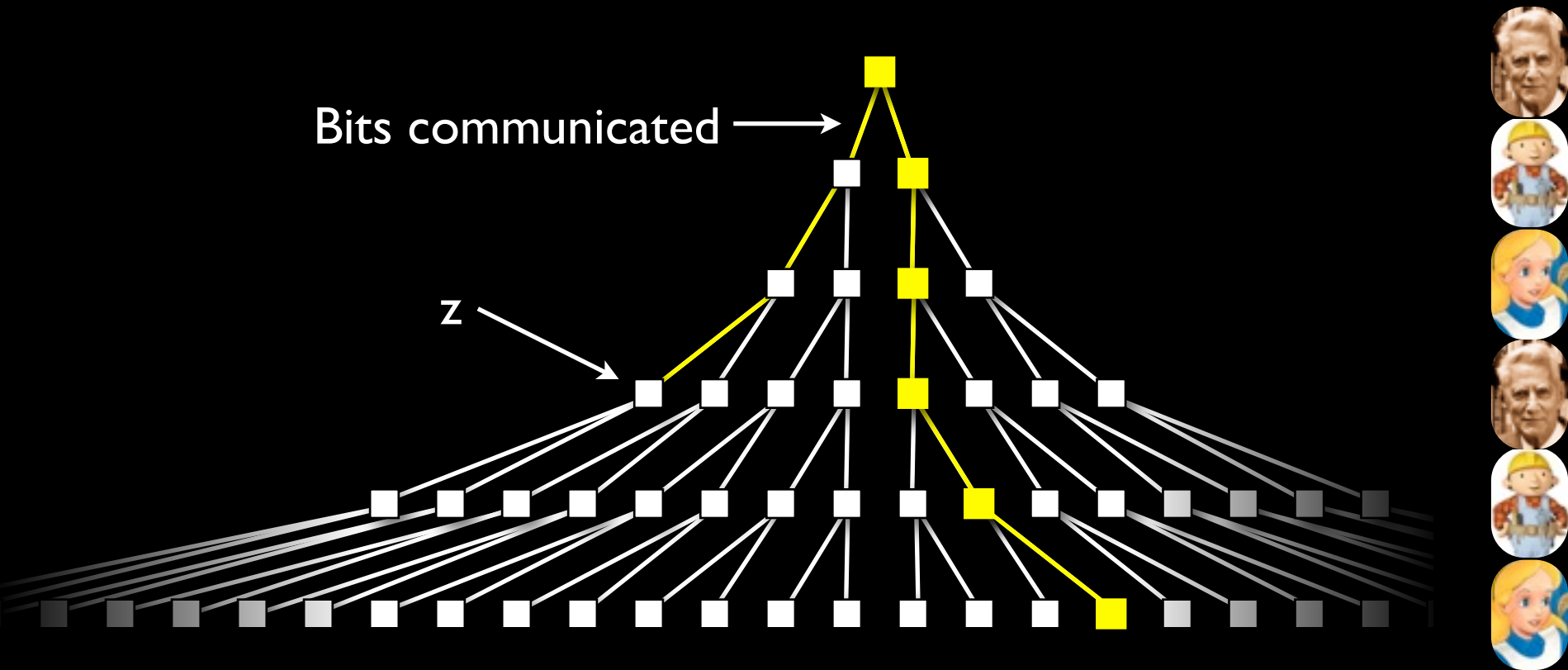
- Consider **deterministic** protocols and **random** f_A, f_B, f_C

Proof Sketch



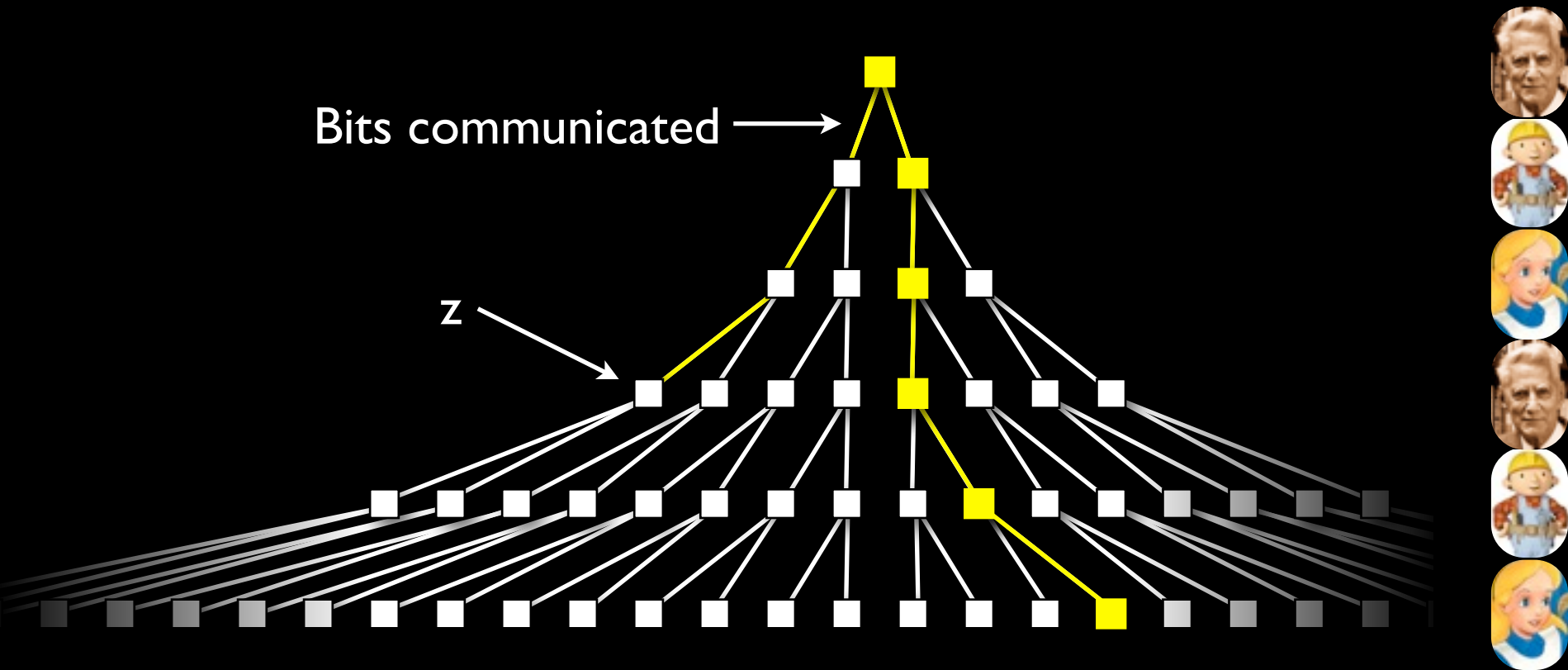
- Consider **deterministic** protocols and **random** f_A, f_B, f_C
- At each node z defines random variables $z(f_A), z(f_B), z(f_C)$

Proof Sketch



- Consider **deterministic** protocols and **random** f_A, f_B, f_C
- At each node z defines random variables $z(f_A), z(f_B), z(f_C)$

Proof Sketch



- Consider **deterministic** protocols and **random** f_A, f_B, f_C
- At each node z defines random variables $z(f_A), z(f_B), z(f_C)$
- Induction: For each z , entropy of variables is high.

Reduction to Selection



$f_A: [m] \rightarrow [m]$



$f_B: [m] \rightarrow [m]$



$f_C: [m] \rightarrow [m]$

Reduction to Selection



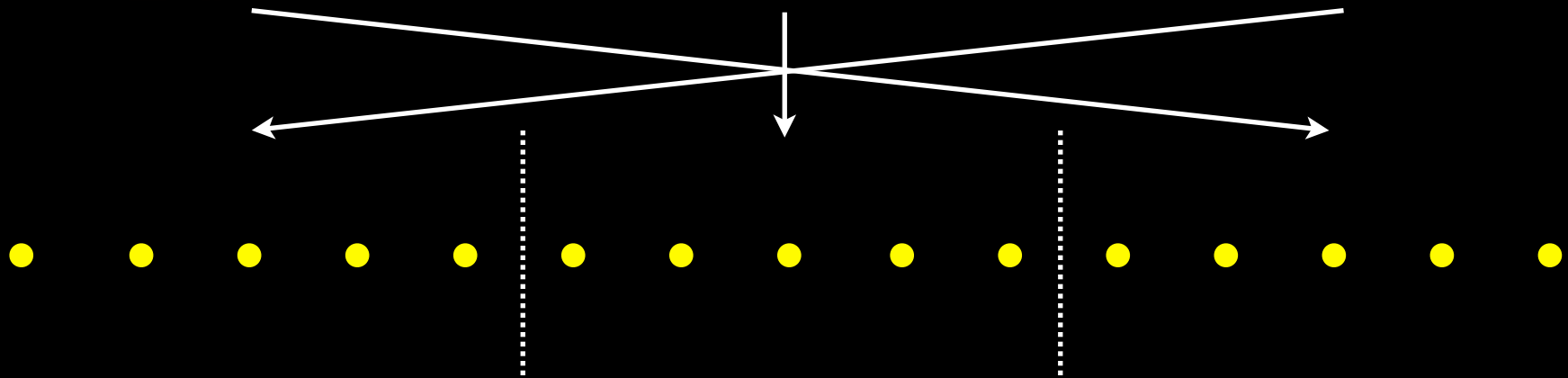
$f_A: [m] \rightarrow [m]$



$f_B: [m] \rightarrow [m]$



$f_C: [m] \rightarrow [m]$



Reduction to Selection



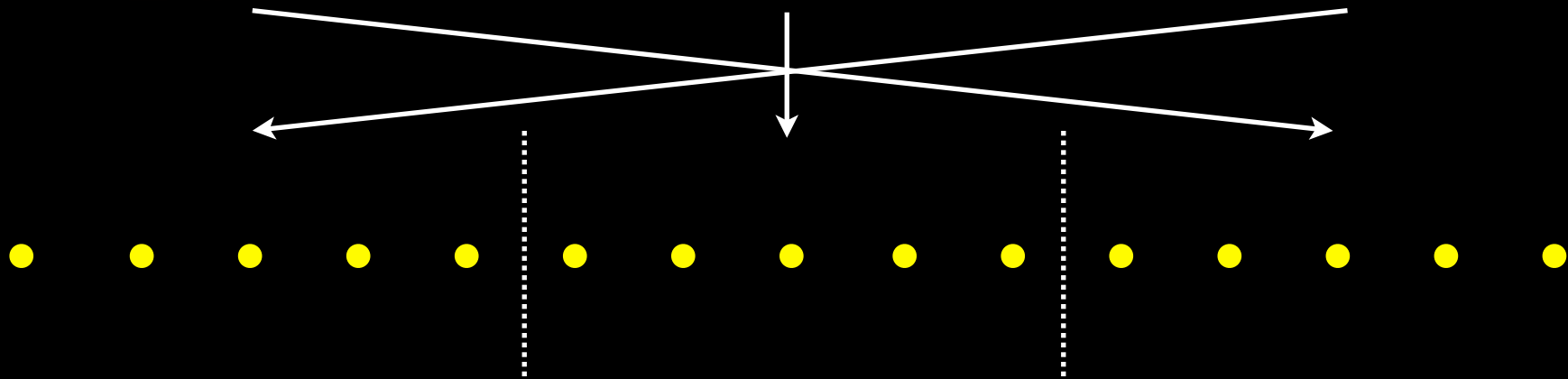
$f_A: [m] \rightarrow [m]$



$f_B: [m] \rightarrow [m]$



$f_C: [m] \rightarrow [m]$



Median = $f_A(I) f_B(f_A(I) f_C(f_B(f_A(I))))$

Reduction to Selection



$$000 \times (3-f_A(1)) \times 5$$

$$100 \times (3-f_B(1))$$

$$11 f_C(1)$$

$$12 f_C(2)$$

$$13 f_C(3)$$

$$140 \times (f_B(1)-1)$$

$$200 \times (3-f_B(2))$$

$$21 f_C(1)$$

$$22 f_C(2)$$

$$23 f_C(3)$$

$$240 \times (f_B(2)-1)$$

$$300 \times (3-f_B(3))$$

$$31 f_C(1)$$

$$32 f_C(2)$$

$$33 f_C(3)$$

$$340 \times (f_B(3)-1)$$

$$400 \times (f_A(1)-1) \times 5$$

VALUE



Summary

- **Thm:** Any single pass $(\epsilon, 1/4)$ -approx for entropy requires $\Omega(\epsilon^{-2} / \log^2 \epsilon^{-1})$ space.
- **Thm:** Computing first t layers of a BFS-tree requires either $t/2$ passes or $\Omega(n^{1+1/t})$ space.
- **Thm:** Finding the median requires either p passes or $\Omega(n^{1/p})$ space.

Future Directions

Future Directions

- **Little things:** Factors in space and passes

Future Directions

- **Little things:** Factors in space and passes
- **Medium things:** Random-order BFS trees, Lower-bounds for random-order streams, Fully characterize the “sketch-able” distances.

Future Directions

- **Little things:** Factors in space and passes
- **Medium things:** Random-order BFS trees,
Lower-bounds for random-order streams,
Fully characterize the “sketch-able” distances.
- **Big thing:** Space-Efficient Sampling!

Thanks

Thanks

Advisor and Mentors:

Sampath, Sudipto,
Sasha, & Bruce

Thanks

Advisor and Mentors:

Sampath, Sudipto,
Sasha, & Bruce

Committee:

Sudipto, Piotr,
Michael, & Sanjeev

Thanks

Advisor and Mentors:

Sampath, Sudipto,
Sasha, & Bruce

Committee:

Sudipto, Piotr,
Michael, & Sanjeev

Theory Group (Past & Present):

Stan, Milan, Yael, Boulos, Niel,
Kuku, Sid, & Mirko



Thanks

Advisor and Mentors:

Sampath, Sudipto,
Sasha, & Bruce

Committee:

Sudipto, Piotr,
Michael, & Sanjeev

Theory Group (Past & Present):

Stan, Milan, Yael, Boulos, Niel,
Kuku, Sid, & Mirko

Co-Authors:

Deepak, Stan, Sasha, Tugkan,
Amit, Graham, Joan, Peter,
Boulos, Piotr, Sampath, Sanjeev,
Keshav, Eduardo, Jeff, Bruce, Sid,
Suresh, Jian, & Zhengyuan



Thanks

Advisor and Mentors:

Sampath, Sudipto,
Sasha, & Bruce

Committee:

Sudipto, Piotr,
Michael, & Sanjeev

Theory Group (Past & Present):

Stan, Milan, Yael, Boulos, Niel,
Kuku, Sid, & Mirko

Co-Authors:

Deepak, Stan, Sasha, Tugkan,
Amit, Graham, Joan, Peter,
Boulos, Piotr, Sampath, Sanjeev,
Keshav, Eduardo, Jeff, Bruce, Sid,
Suresh, Jian, & Zhengyuan

“Normal” People:

Friends, family, visitors,
& especially Christie



Questions?