# Beyond Set Disjointness:
# The Communication Complexity of Finding the Intersection

**Grigory Yaroslavtsev**
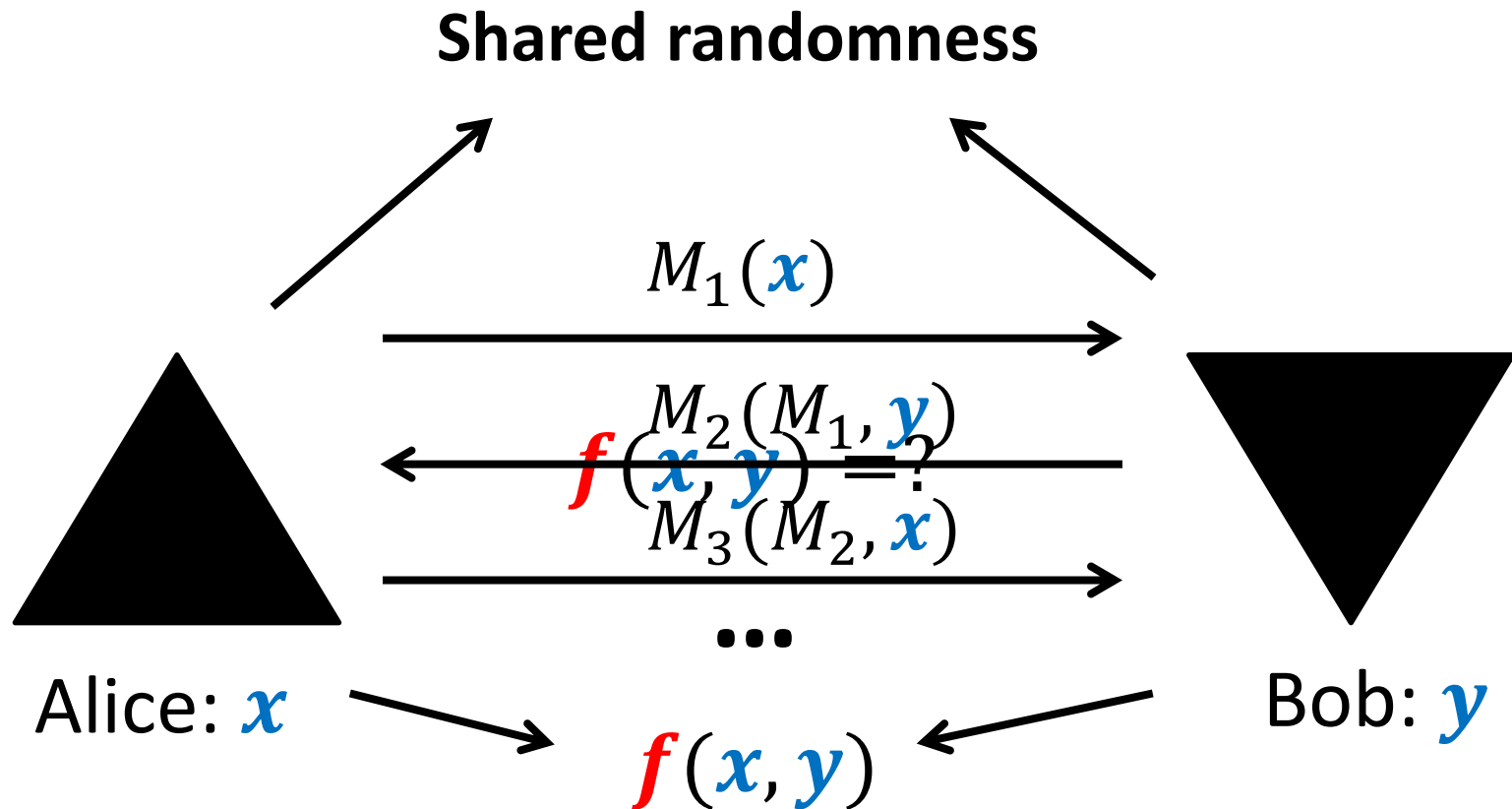
**http://grigory.us**

Joint with Brody, Chakrabarti, Kondapally and Woodruff

# Communication Complexity [Yao'79]

**Shared randomness**



$M_1(x)$

$M_2(M_1, y)$

$f(x, y) = ?$

$M_3(M_2, x)$

$\cdots$

Alice: $x$

Bob: $y$

$f(x, y)$

- $R(f)$ = min. communication (error 1/3)
- $R^k(f)$ = min. $k$-round communication (error 1/3)

# Set Intersection

- $x = S, y = T, f(x, y) = S \cap T$

$S \subseteq [n], |S| \leq k$

$T \subseteq [n], |T| \leq k$

$S \cap T = ?$

$R^r(k\text{-Intersection}) = ?$

# This talk

Let $ilog^r k = \underbrace{\log \log \ldots \log}_{r \text{ times}} k$

- $R^r(k\text{-Intersection}) = O\left(k \; ilog^{\boldsymbol{\beta} r} k\right)$

[Brody, Chakrabarti, Kondapally, Woodruff, Y.; PODC'14]

- $R^r(k\text{-Intersection}) = \Omega(k \; ilog^r k)$

[Saglam-Tardos FOCS'13; Brody, Chakrabarti, Kondapally, Woodruff, Y.'13]

$R^r(k\text{-Intersection}) = \Theta(k)$ for $r = O(\log^* k)$

# $k$-Disjointness

- $f(S, T) = 1$, iff $|S \cap T| = 0$

- $R(k\text{–Disjointness}) = \Theta(k)$ [Razborov'92; Hastad-Wigderson'96]

- $R^1(k\text{–Disjointness}) = \Theta(k \log k)$

[Folklore + Dasgupta, Kumar, Sivakumar; Buhrman'12, Garcia-Soriano, Matsliah, De Wolf'12]

- $R^r(k\text{–Disjointness}) = \Theta(k \, \text{ilog}^r k)$ [Saglam, Tardos'13]

- $R(k\text{–Disjointness}) = \alpha k + \text{o}(k)$ [Braverman, Garg, Pankratov, Weinstein'13]

# Applications

- $J(\textcolor{green}{S}, \textcolor{red}{T}) = \frac{|\textcolor{green}{S} \cap \textcolor{red}{T}|}{|\textcolor{green}{S} \cup \textcolor{red}{T}|}$ : exact Jaccard index

( for $(1 \pm \epsilon)$-approximate use MinHash [Broder'98; Li-Konig'11; Path-Strokel-Woodruff'14])

- Rarity, distinct elements, joins,…

- Multi-party set intersection (later)

- Contrast: $R(\textcolor{green}{S} \cup \textcolor{red}{T}) = R(\textcolor{green}{S} \, \Delta \, \textcolor{red}{T}) = \Theta(k \log \frac{n}{k})$

# 1-round $O(k \log k)$-protocol

$[n]$

$S$

$[n]$

$T$

$h: [n] \to [k^3]$
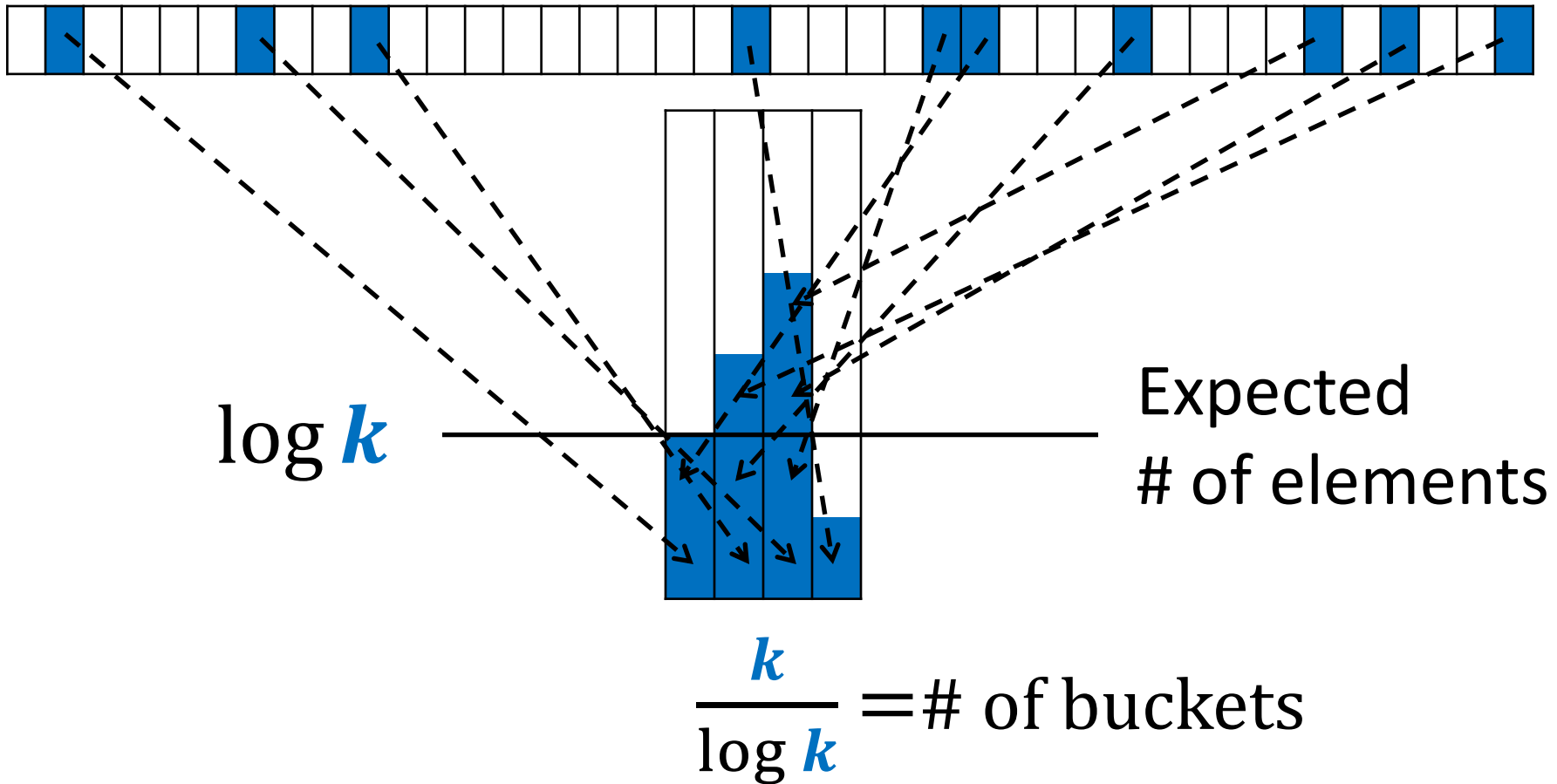
$h(S)$

$[k^3]$

$h(T)$
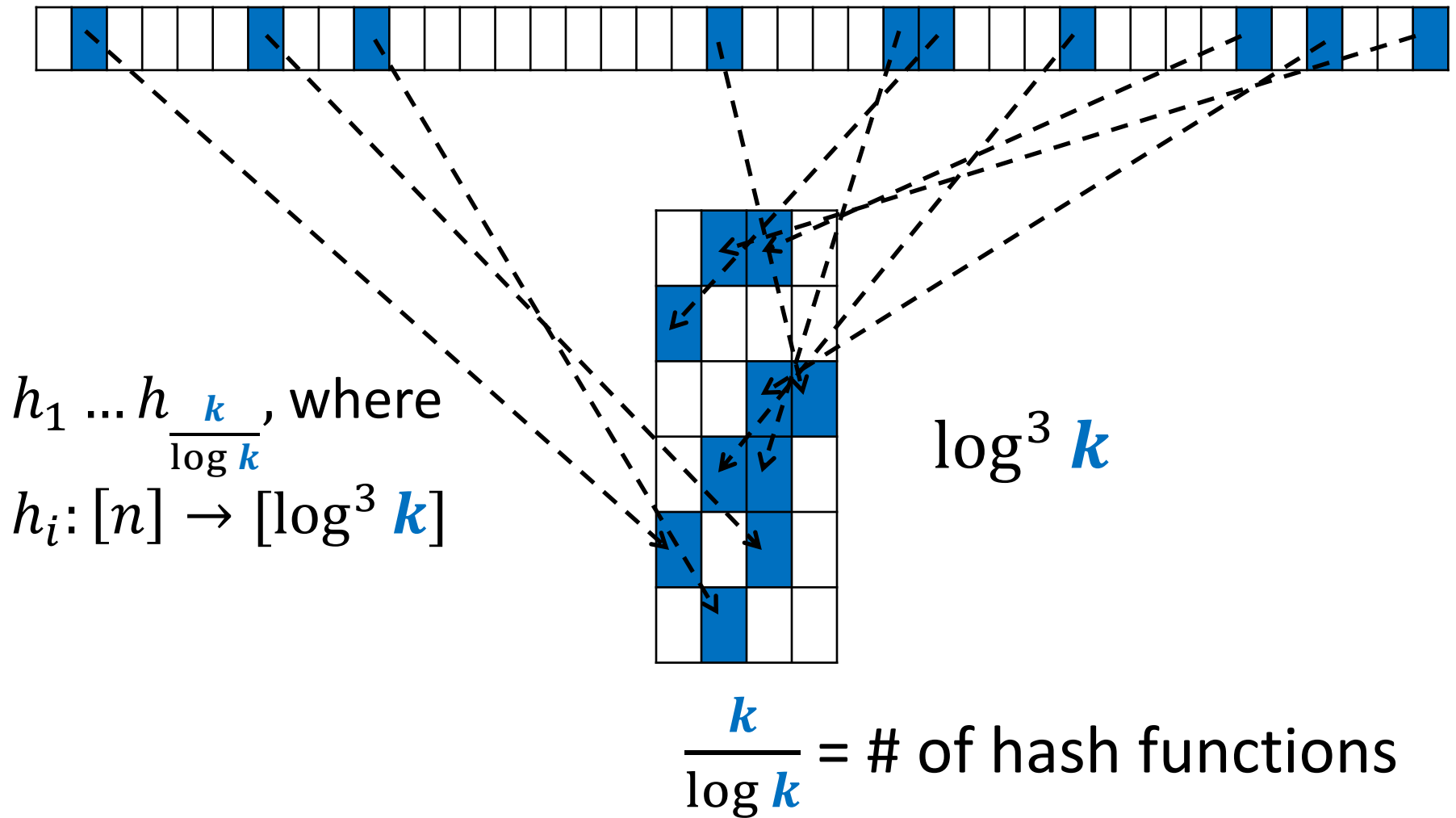
$[k^3]$

$|h(S)| = O(k \log k)$

$|h(T)| = O(k \log k)$

$$S \cap T = S \cap h^{-1}(h(T)) = h^{-1}(h(S)) \cap T$$

# Hashing

$$h: [n] \to [k / \log k]$$



$\log k$

Expected
# of elements

$$\frac{k}{\log k} = \text{\# of buckets}$$

# Secondary Hashing



$h_1 \ldots h_{\frac{k}{\log k}}$, where

$h_i : [n] \rightarrow [\log^3 k]$

$\log^3 k$

$\dfrac{k}{\log k}$ = # of hash functions

# 2-Round $O(k \log \log k)$-protocol



$\log^3 k$

$\log^3 k$

$\dfrac{k}{\log k}$  $\quad |h_i(S)| = |h_i(T)| = O(\log k \log \log k) \quad$  $\dfrac{k}{\log k}$
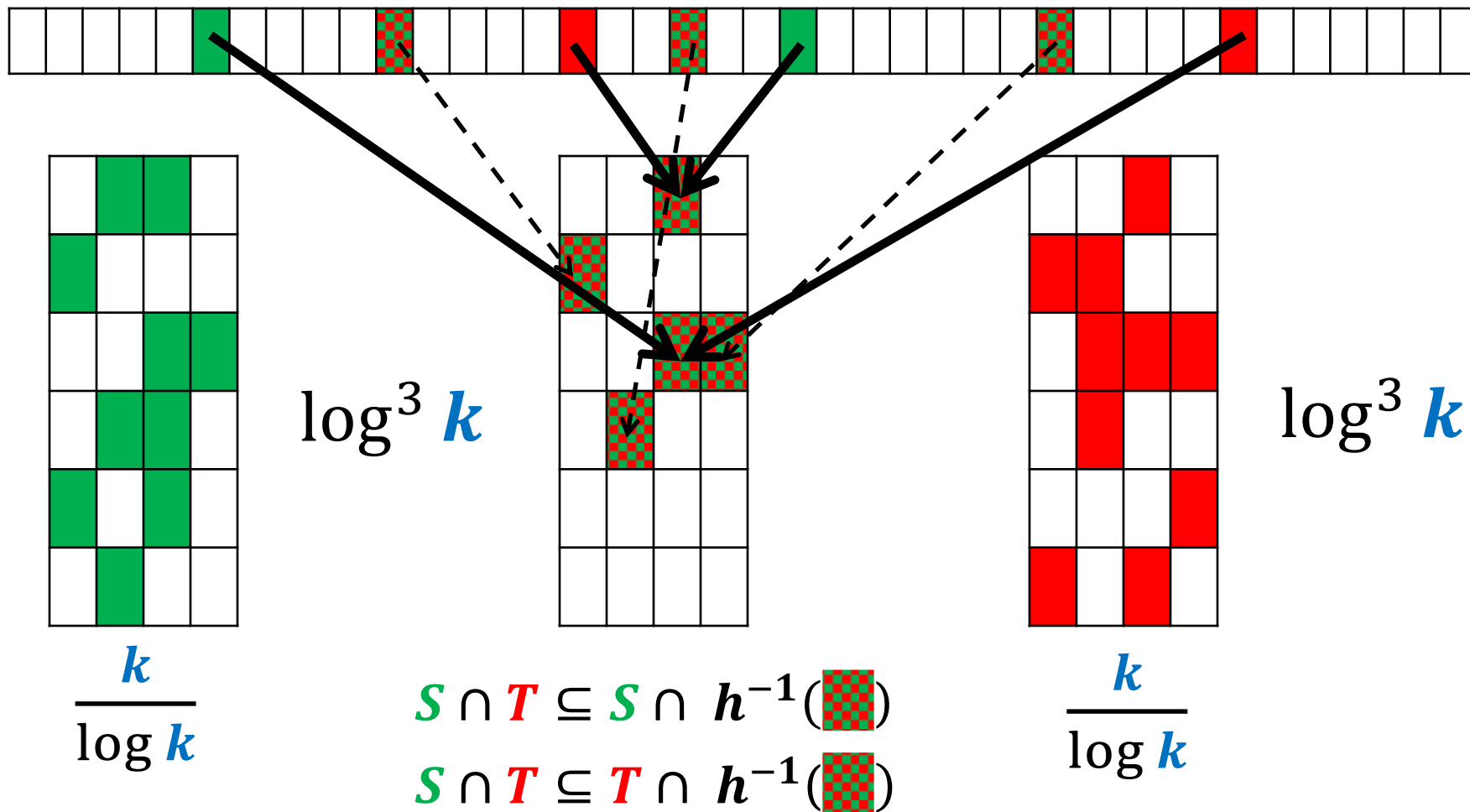
Total communication $= \dfrac{k}{\log k} O(\log k \log \log k) = O(k \log \log k)$

# Collisions



$$\Pr[collision] = O(\frac{1}{\log k})$$

$$\log^3 k$$

$$\frac{k}{\log k}$$

# Collisions



$$\log^3 k$$

$$\frac{k}{\log k}$$

$$S \cap T \subseteq S \cap h^{-1}(\ \blacksquare\ )$$
$$S \cap T \subseteq T \cap h^{-1}(\ \blacksquare\ )$$
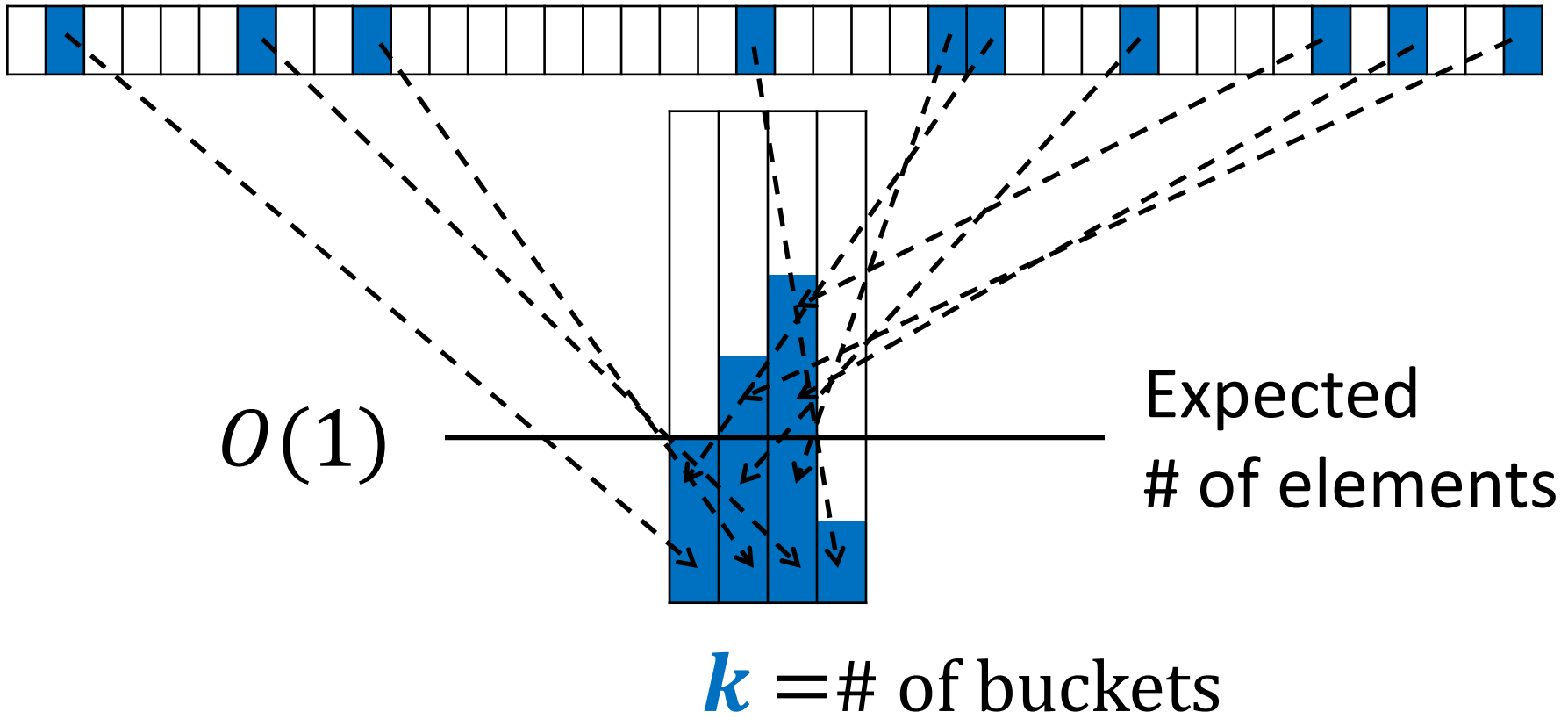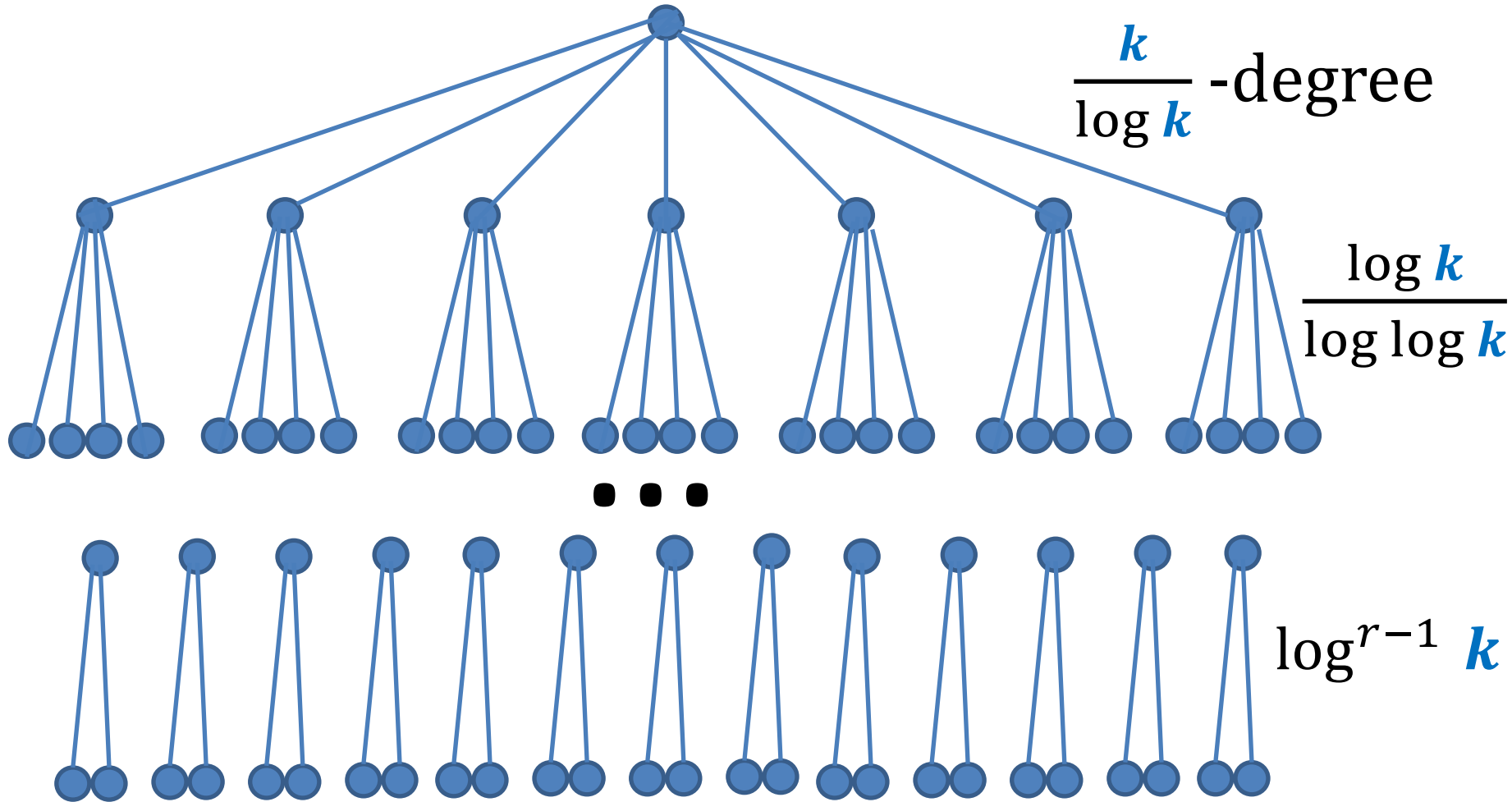
$$\log^3 k$$

$$\frac{k}{\log k}$$

# Collisions

- Second round:
  - For each bucket send $O(\log k)$-bit equality check (total $O(k)$-communication)
  - Correct intersection computed in buckets $i$ where
  $$\boldsymbol{S} \cap \boldsymbol{h_i^{-1}}(\blacksquare) = \boldsymbol{T} \cap \boldsymbol{h_i^{-1}}(\blacksquare)$$
  - Expected # of items in incorrect buckets $O(k / \log k)$
  - Use 1-round protocol for incorrect buckets
  - Total communication $O(k \log \log k)$

# Main protocol

$$h : [n] \to [k]$$
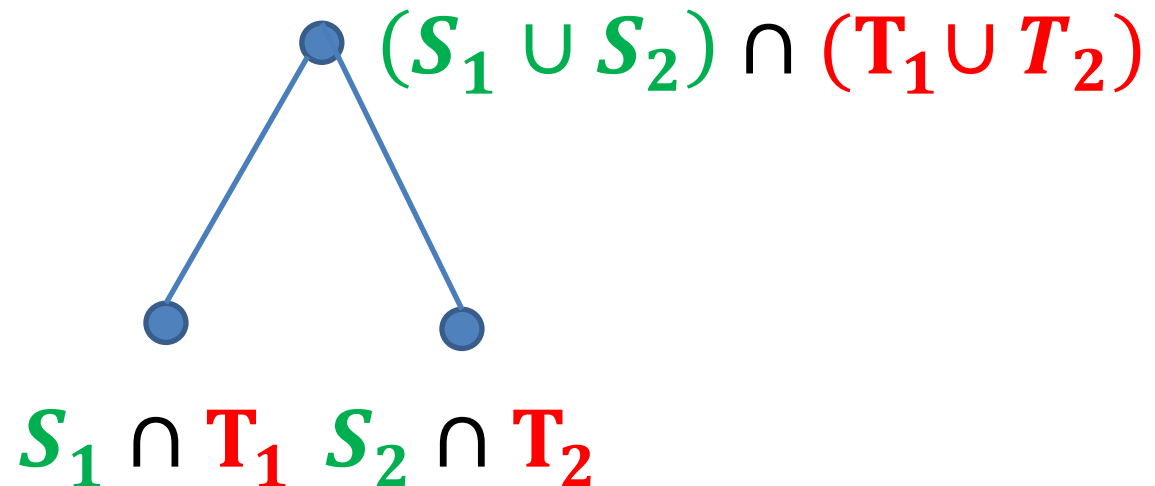


$O(1)$ — Expected # of elements

$k =$ # of buckets

# Verification tree



$\frac{k}{\log k}$ -degree

$\frac{\log k}{\log \log k}$

$\log^{r-1} k$

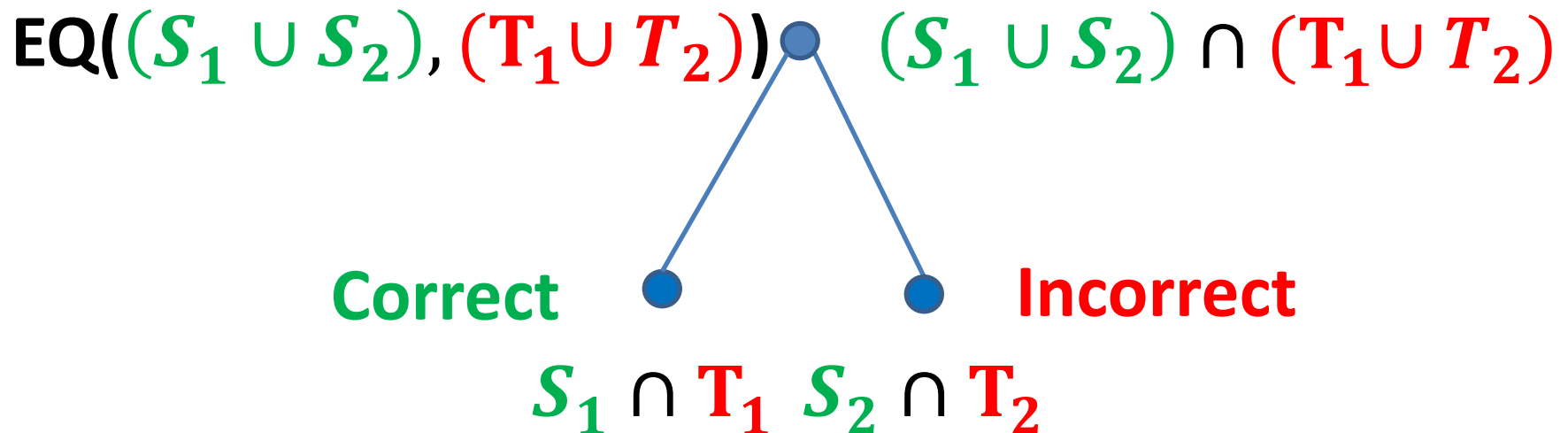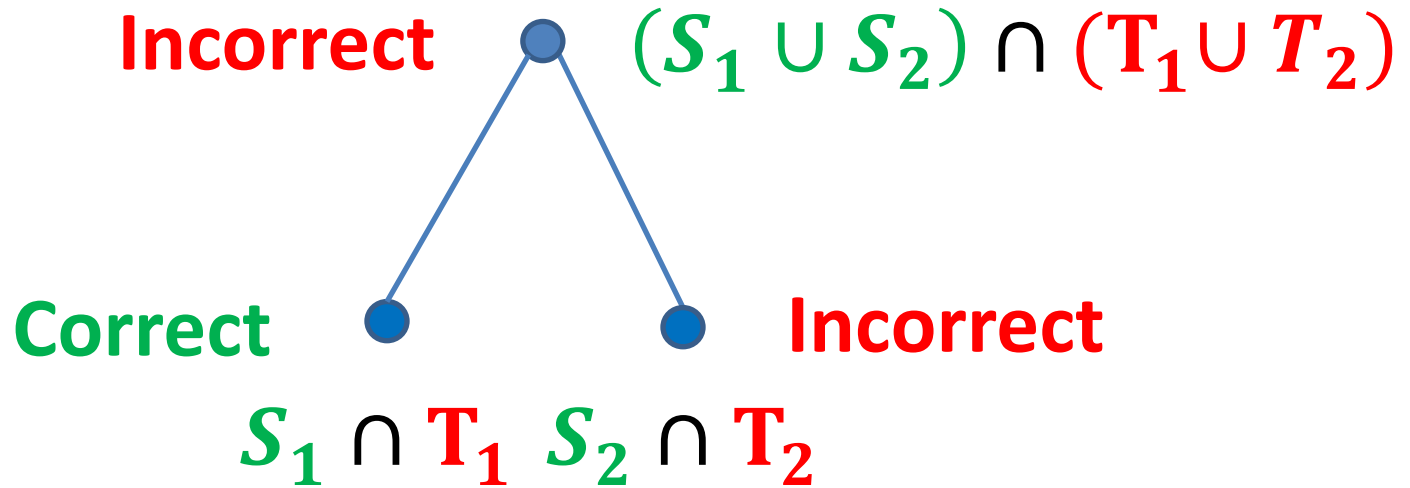$k$ buckets = leaves of the verification tree

# Verification bottom-up

# Verification bottom-up

**Incorrect** $(S_1 \cup S_2) \cap (T_1 \cup T_2)$

**Correct** **Incorrect**

$S_1 \cap T_1 \quad S_2 \cap T_2$

$\textbf{EQ}((S_1 \cup S_2), (T_1 \cup T_2))$ $(S_1 \cup S_2) \cap (T_1 \cup T_2)$

**Correct** **Incorrect**

$S_1 \cap T_1 \quad S_2 \cap T_2$

# Verification bottom-up

**Incorrect** $(S_1 \cup S_2) \cap (\mathbf{T_1} \cup \mathbf{T_2})$

**Correct** **Incorrect**

$$S_1 \cap \mathbf{T_1} \quad S_2 \cap \mathbf{T_2}$$

$\mathbf{EQ}((S_1 \cup S_2), (\mathbf{T_1} \cup \mathbf{T_2}))$ **Correct** $(S_1 \cup S_2) \cap (\mathbf{T_1} \cup \mathbf{T_2})$

**Correct** **Incorrect**

$$S_1 \cap \mathbf{T_1} \quad S_2 \cap \mathbf{T_2}$$

# Verification bottom-up

$p_{r-1}$

$p_{r-2}$

$p_1$

$S_1^1, T_1^1$ $S_2^1, T_2^1$ $\ldots$ $S_i^1, T_i^1$ $\ldots$ $S_k^1, T_k^1$

# Analysis of Stage $i$

- $p_i = Pr[\text{node at stage } i \text{ computed correctly}]$

- Set $p_i = 1 - \dfrac{1}{\left(ilog^{r-i-1}k\right)^4}$

  - Run equality checks and basic intersection protocols with success probability $p_i$

  - **Key lemma**: $\mathbb{E}[\text{\# of restarts per leaf}] = O(1)$

  - Cost of Equality = $O(k\ ilog^r k)$

  - Cost of Intersection in leafs = $O(k)$

- $p_{r-1} = Pr[\text{protocol succeeds}] = 1 - 1/k^4$

# Lower Bound

- $R^r(k\text{-Intersection}) = \Omega(k\ ilog^r k)$

[Brody, Chakrabarti, Kondapally, Woodruff, Y.'13]

- $EQ_m(x, y) = 1$ iff $x = y$, where $x, y \in \{0,1\}^m$

- $EQ_m^k =$ solving $k$ independent instances of $EQ_m$

- $EQ_m^k$ reduces to $k$-Intersection:
  - Given $(x_1, \ldots, x_k)$ and $(y_1, \ldots, y_k)$
  - Construct sets with elements $(1, x_1), \ldots, (k, x_k)$ and $(1, y_1), \ldots, (k, y_k)$

# Communication Direct Sums

"Solving **m** copies of a communication problem requires **m** times more communication":

$$R^r(\boldsymbol{f^m}) = \Omega(\boldsymbol{m})R^r(\boldsymbol{f})$$

- For arbitrary $\boldsymbol{f}$ [… Braverman, Rao 10; Barak Braverman, Chen, Rao 11, ….]

- In general, can't go beyond

$$\boldsymbol{R}(EQ_{\boldsymbol{m}}) = O(1)$$
$$\boldsymbol{R}\big(EQ^{\boldsymbol{m}}_{,\boldsymbol{m}}\big) = O(\boldsymbol{m})$$

# Specialized Communication Direct Sums

Information cost $\leq$ Communication complexity

- $R(\text{Disjointness}) = \Omega(n)$ [Bar Yossef, Jayram, Kumar, Sivakumar'01]

$$\text{Disjointness}(x, y) = \wedge_i (\neg x_i \vee \neg y_i)$$

- Stronger direct sum for bounded-round complexity of Equality-type problems (a.k.a. "union bound is optimal") [Molinaro, Woodruff, Y.'13]

$$\boldsymbol{R}^1 (EQ^{\boldsymbol{k}}) = \Omega(\boldsymbol{k} \log \boldsymbol{k}) \boldsymbol{R}(EQ)$$
$$\boldsymbol{R}^r (EQ^{\boldsymbol{k}}) = \Omega(\boldsymbol{k}\, i \log^r \boldsymbol{k}) \boldsymbol{R}(EQ)$$

# Extensions

- Multi-party: $m$ players, $S_1, \ldots, S_m$, where $|S_i| \leq k$

  $-\ S = S_1 \cap \cdots \cap S_m = ?$

  $-$ Boost error probability to $1 - 1/2^k$

  $-$ Average per player (using coordinator):

  $$O(k\ ilog^r k) \text{ in } O\left(r \max\left(1, \frac{\log m}{k}\right)\right) \text{ rounds}$$

  $-$ Worst-case per player (using a tournament)

  $$O\left(k^2\ ilog^r k \max\left(1, \frac{\log m}{k}\right)\right) \text{ in } O\left(rk \max\left(1, \frac{\log m}{k}\right)\right) \text{ rounds}$$

# Open Problems

- $R^r(k\text{-Intersection}) = O(k \; ilog^r k)$ ?
- Better protocols for the multi-party setting