

Beyond Set Disjointness: The Communication Complexity of Finding the Intersection

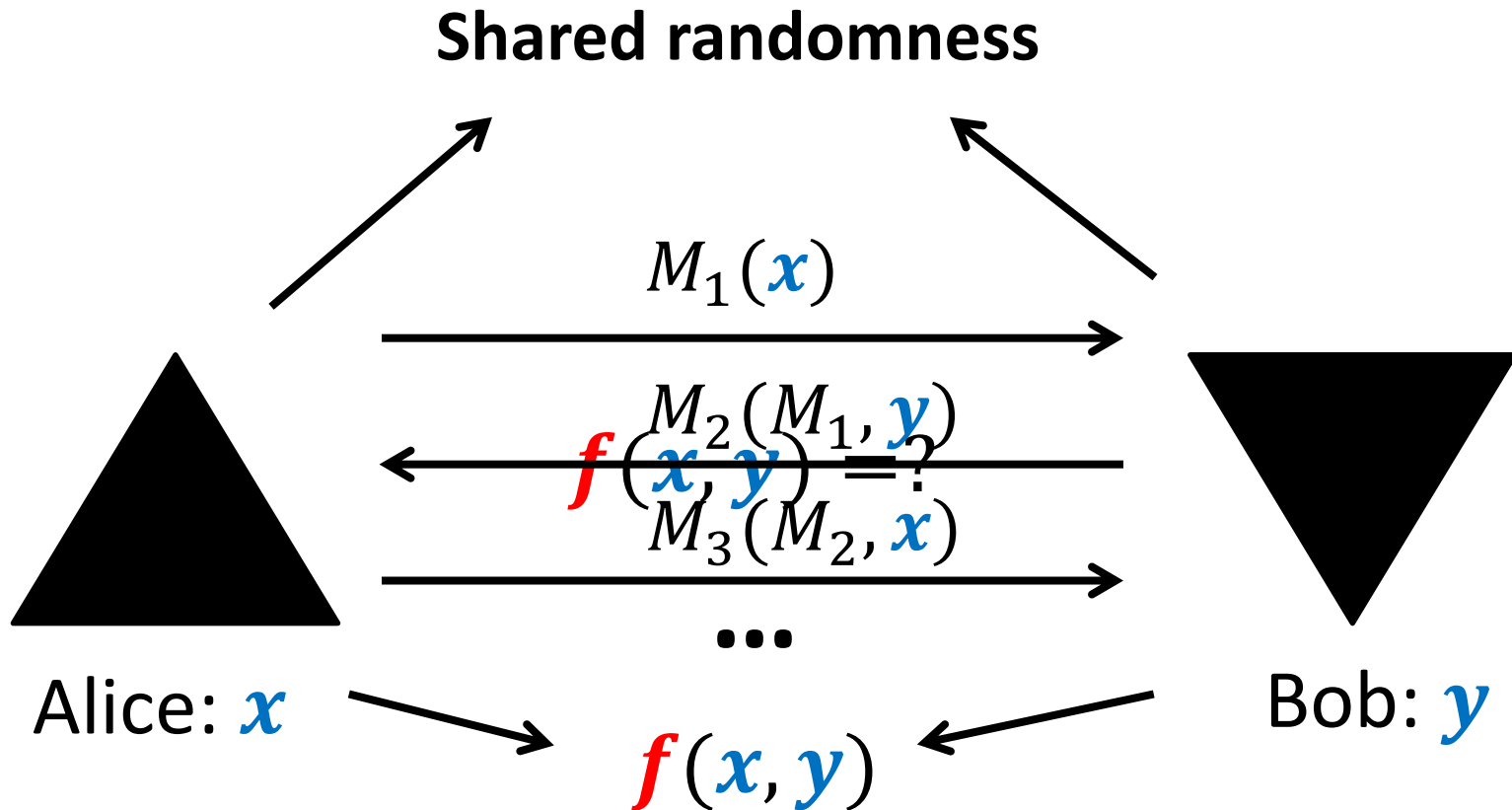
Grigory Yaroslavtsev

<http://grigory.us>



Joint with Brody, Chakrabarti, Kondapally and Woodruff

Communication Complexity [Yao'79]



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

Set Intersection

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

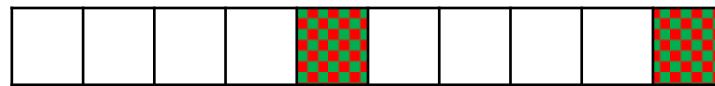


$$S \subseteq [n], |S| \leq k$$



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

$$S \cap T = ?$$



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

k is big, n is **huge**, where **huge** \gg big

Our results

Let $i\log^r k = \log \underbrace{\log \dots \log}_r k$
 r times

- $R^r(\mathbf{k}\text{-Intersection}) = O(\mathbf{k} i\log^{\beta^r} \mathbf{k})$

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

- $R^r(\mathbf{k}\text{-Intersection}) = \Omega(\mathbf{k} i\log^r \mathbf{k})$

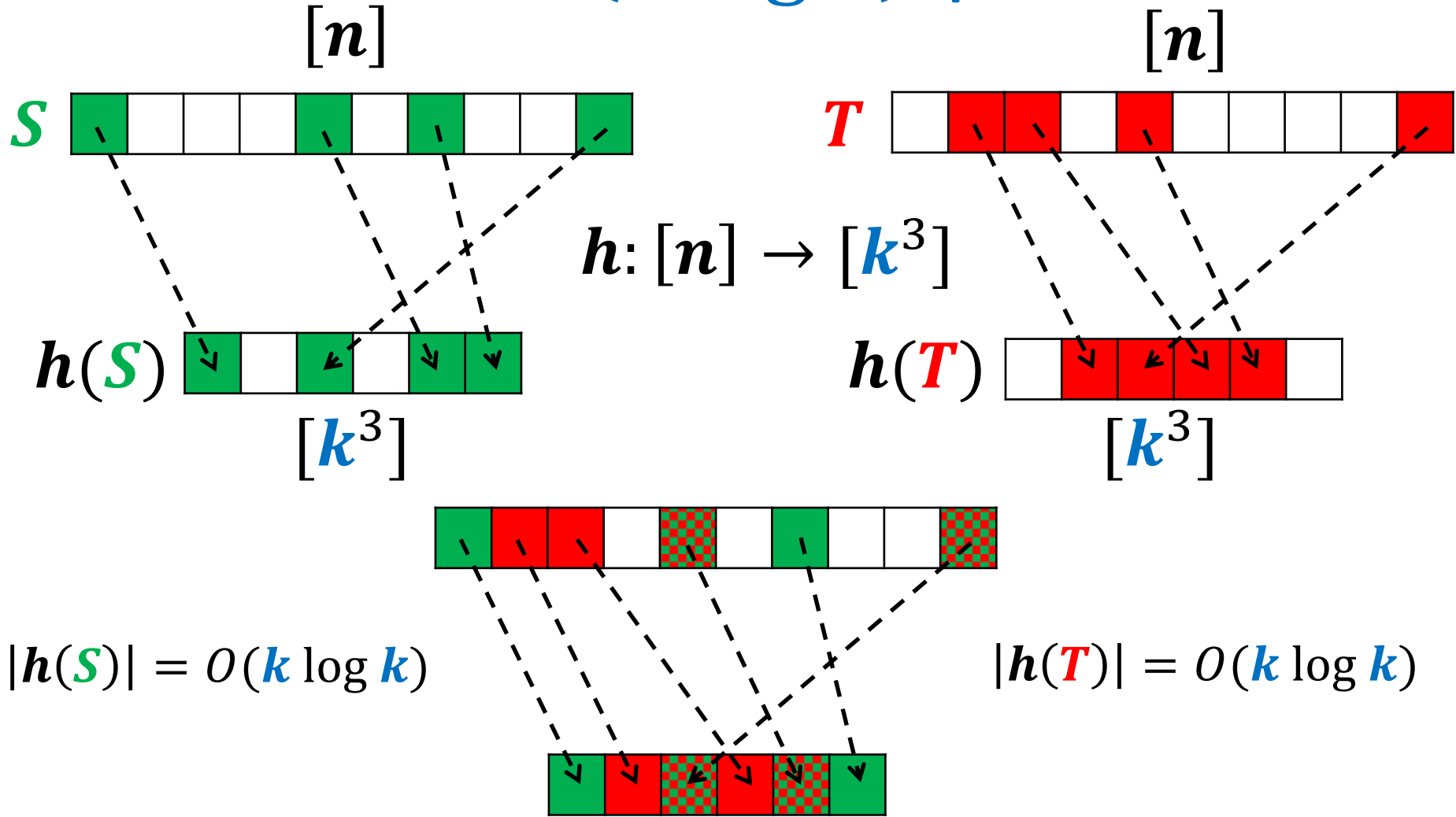
[Saglam-Tardos FOCS'13; Brody, Chakrabarti, Kondapally, Woodruff, Y.; RANDOM'14]

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

Applications

- **Exact** Jaccard index $J(\mathbf{S}, \mathbf{T}) = \frac{|\mathbf{S} \cap \mathbf{T}|}{|\mathbf{S} \cup \mathbf{T}|}$
(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: $R(\mathbf{S} \cup \mathbf{T}) = R(\mathbf{S} \Delta \mathbf{T}) = \Theta\left(k \log \frac{n}{k}\right)$

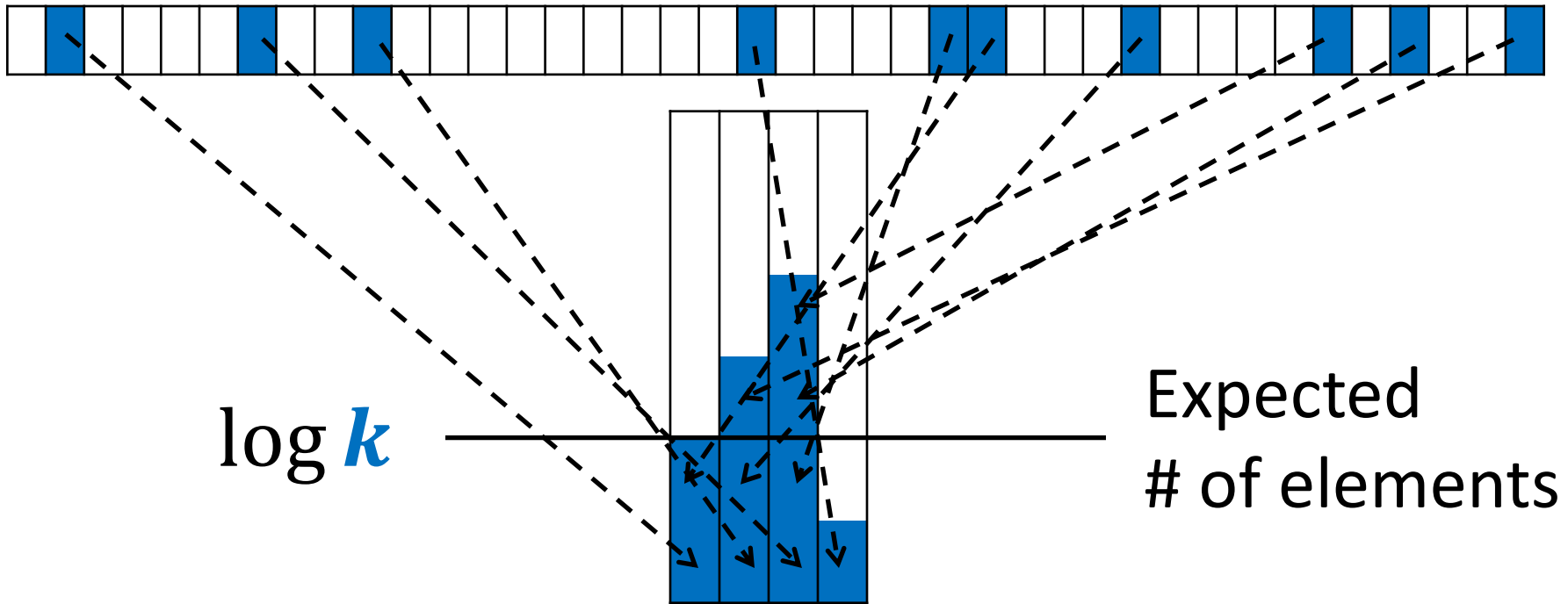
1-round $O(k \log k)$ -protocol



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

Hashing

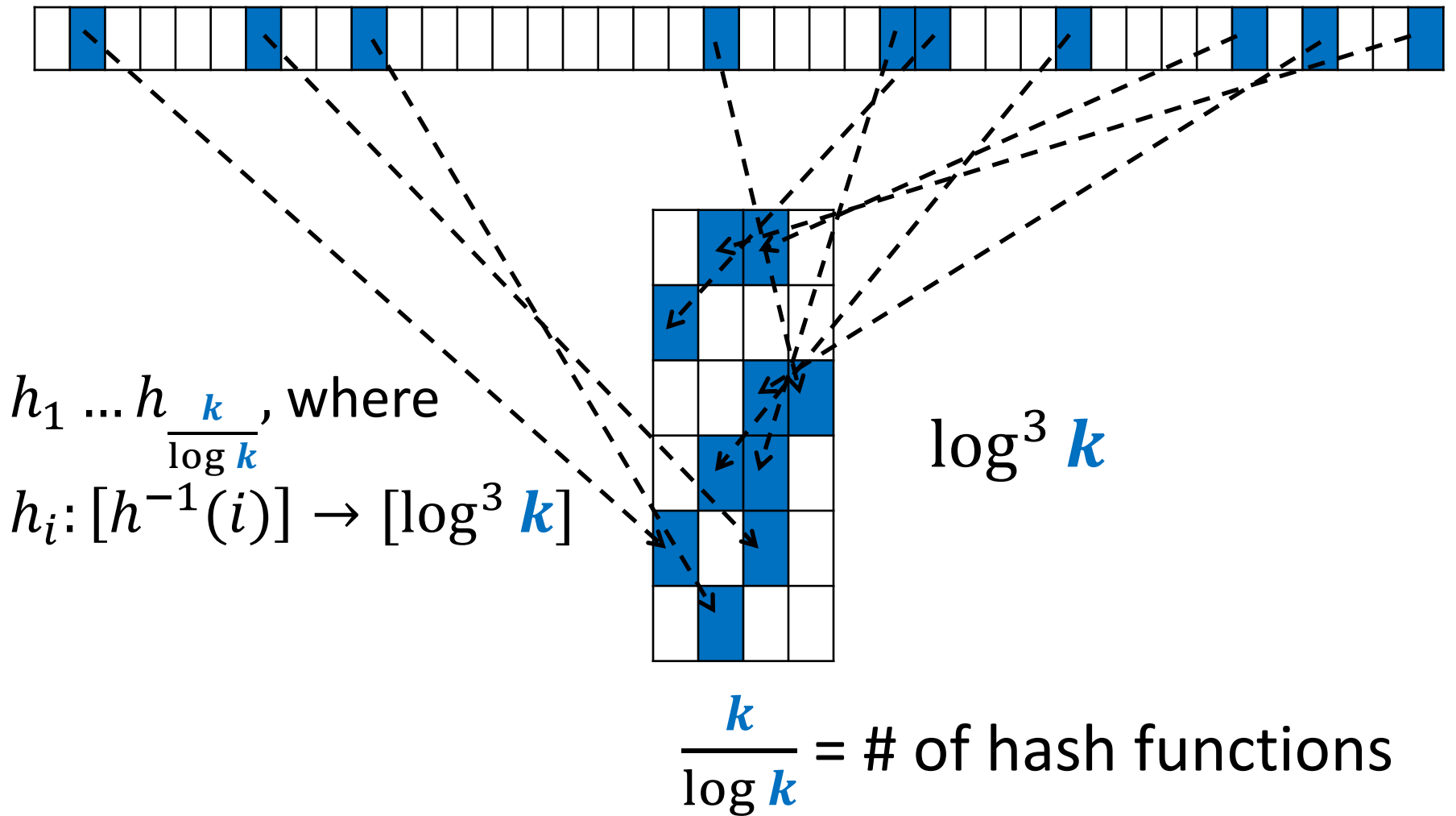
$$h: [n] \rightarrow [k / \log k]$$



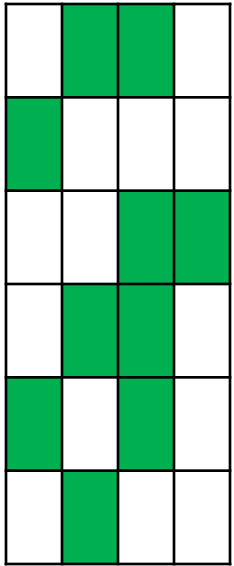
Expected
of elements

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

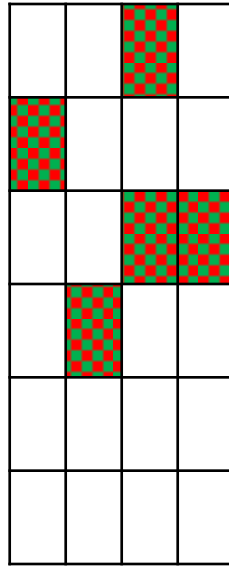
Secondary Hashing



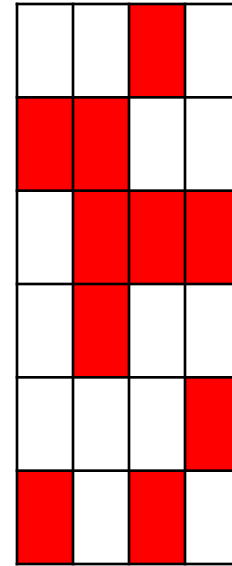
2-Round $O(k \log \log k)$ -protocol



$\log^3 k$



$\log^3 k$



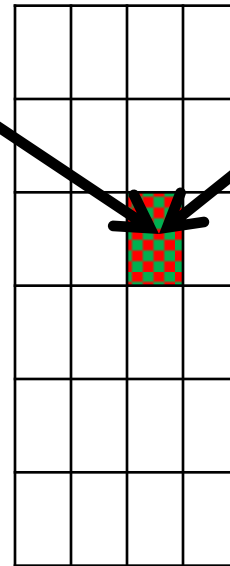
$$\frac{k}{\log k}$$

$$|h_i(\mathbf{S})|, |h_i(\mathbf{T})| = O(\log k \log \log k)$$

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

$$\text{Total communication} = \frac{k}{\log k} O(\log k \log \log k) = O(k \log \log k)$$

Collisions

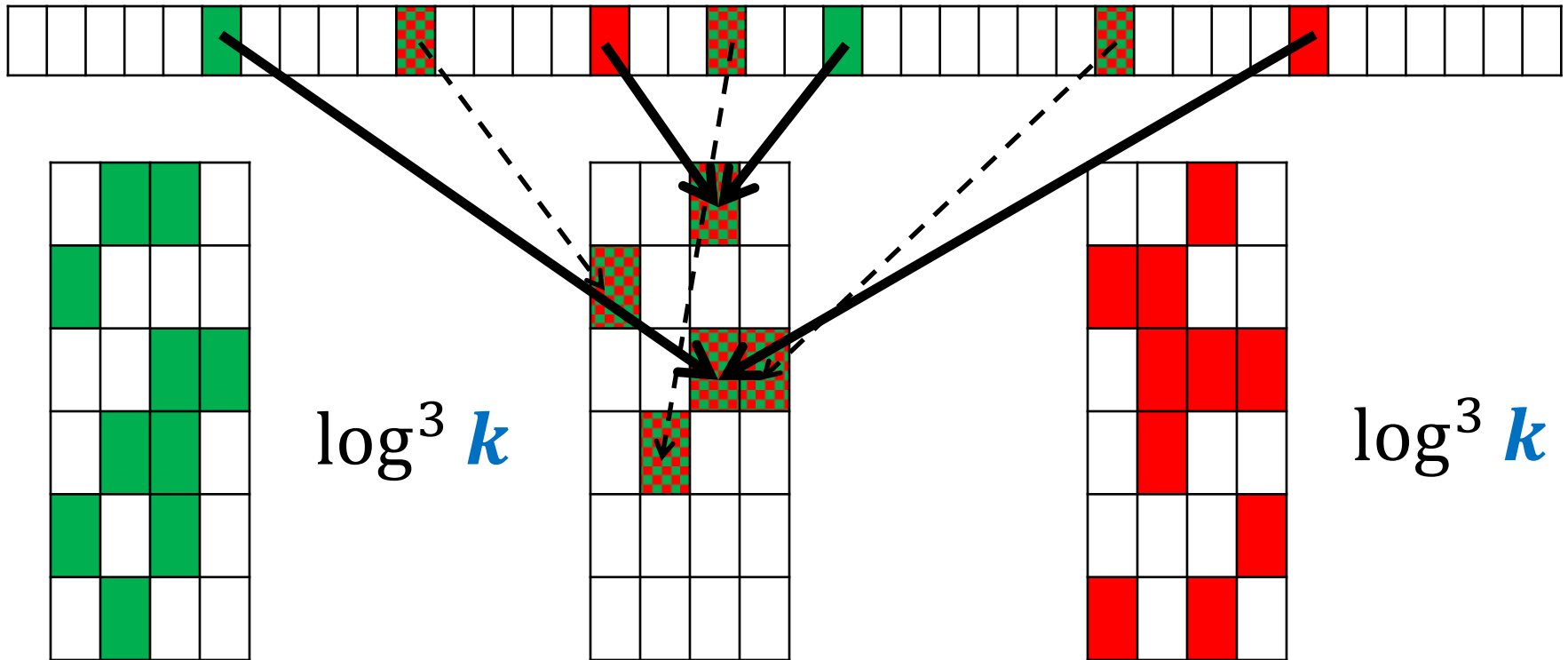


$$\Pr[\textit{collision}] = o\left(\frac{1}{\log k}\right)$$

$\log^3 k$

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

Collisions



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

$\log^3 k$

$$S \cap T \subseteq S \cap H^{-1}(\text{checkered})$$

$$S \cap T \subseteq T \cap H^{-1}(\text{checkered})$$

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

$\log^3 k$

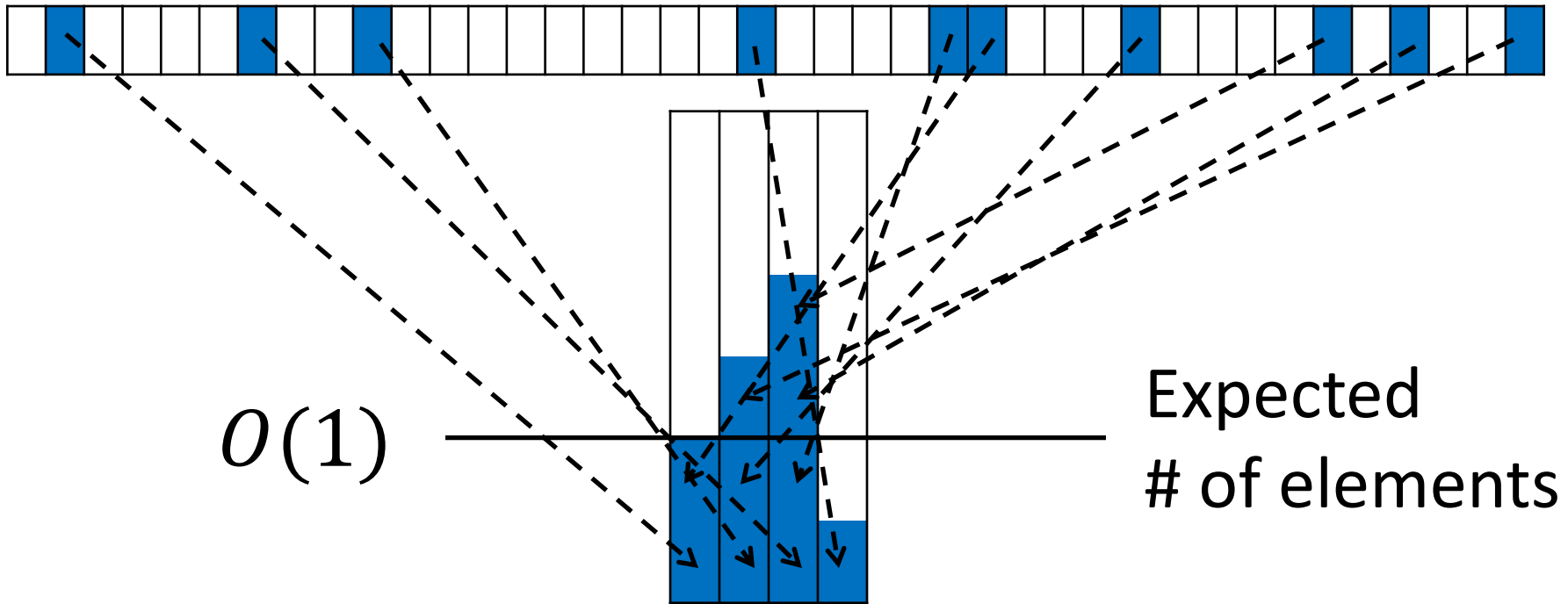
Key fact: If $S \cap H^{-1}(\text{checkered}) = T \cap H^{-1}(\text{checkered})$ then also = $S \cap T$

Collisions

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

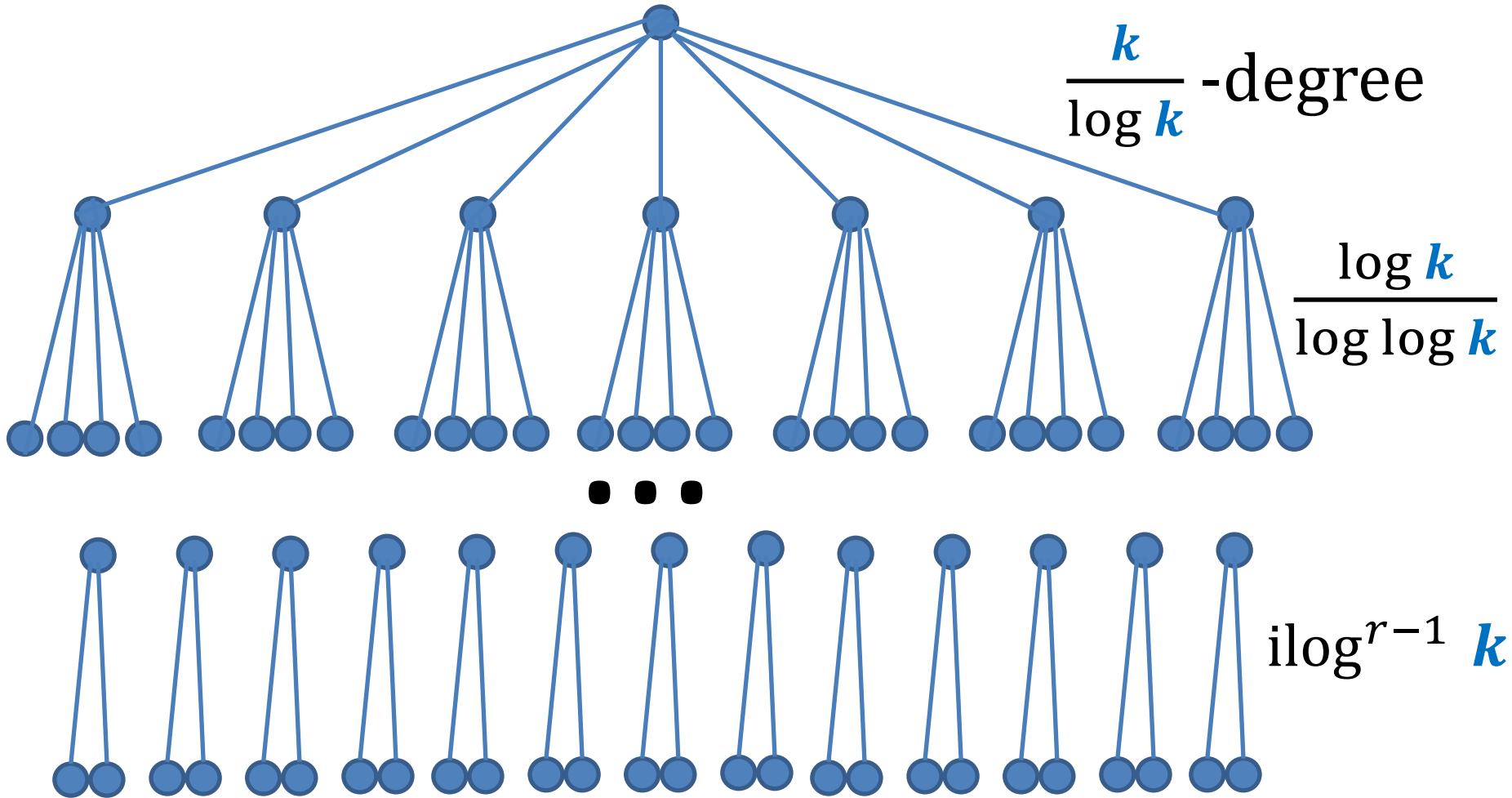
Main protocol

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



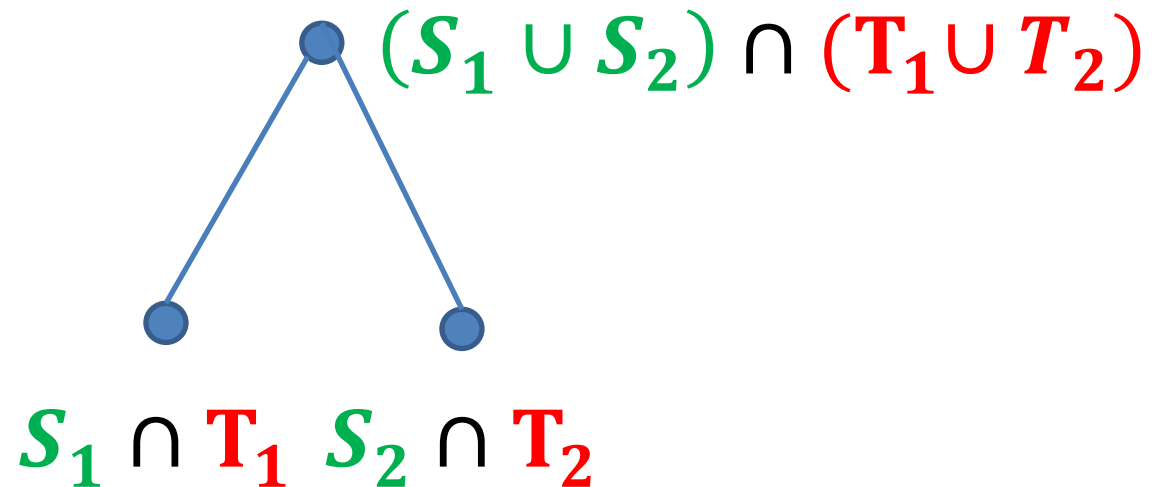
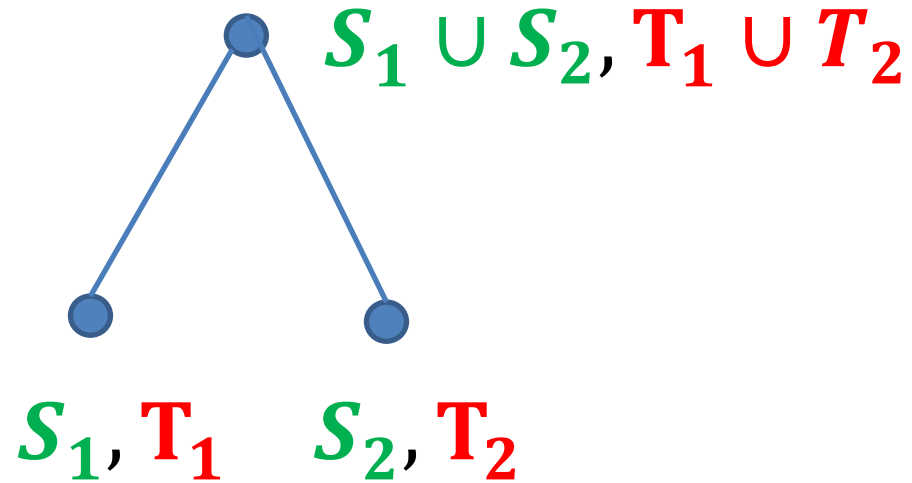
$k = \#$ of buckets

Verification tree

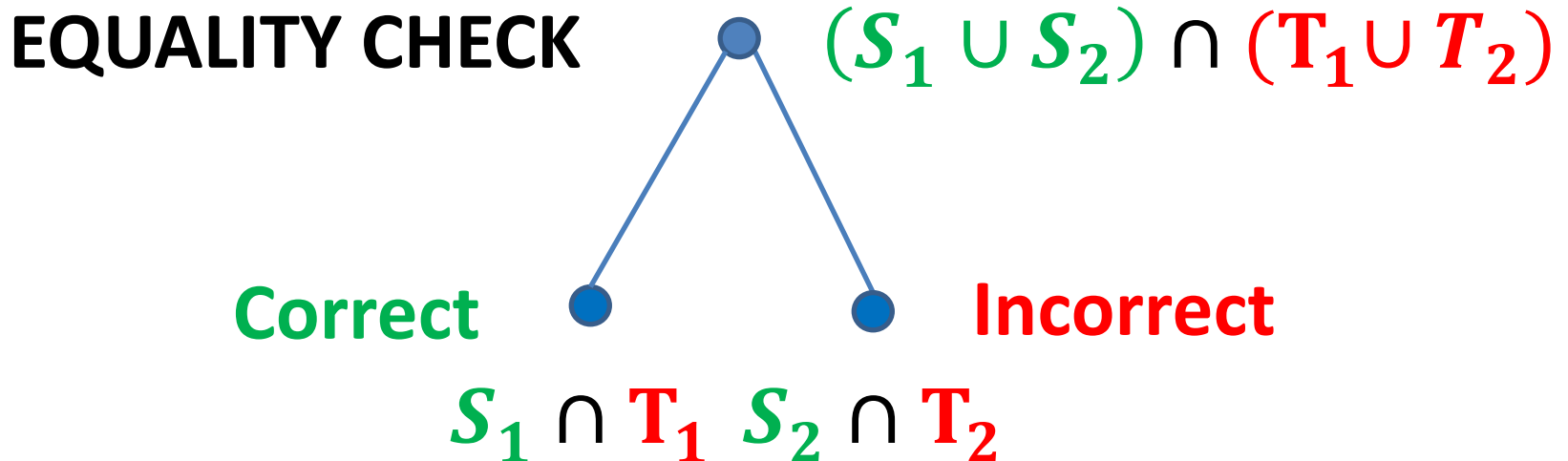
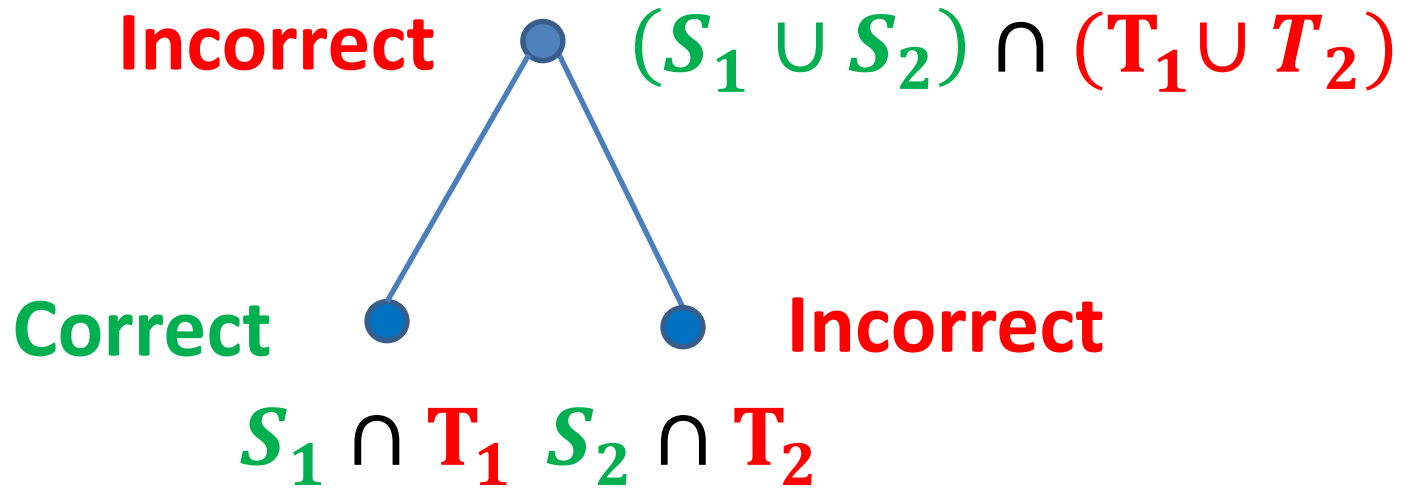


k buckets = leaves of the verification tree

Verification bottom-up

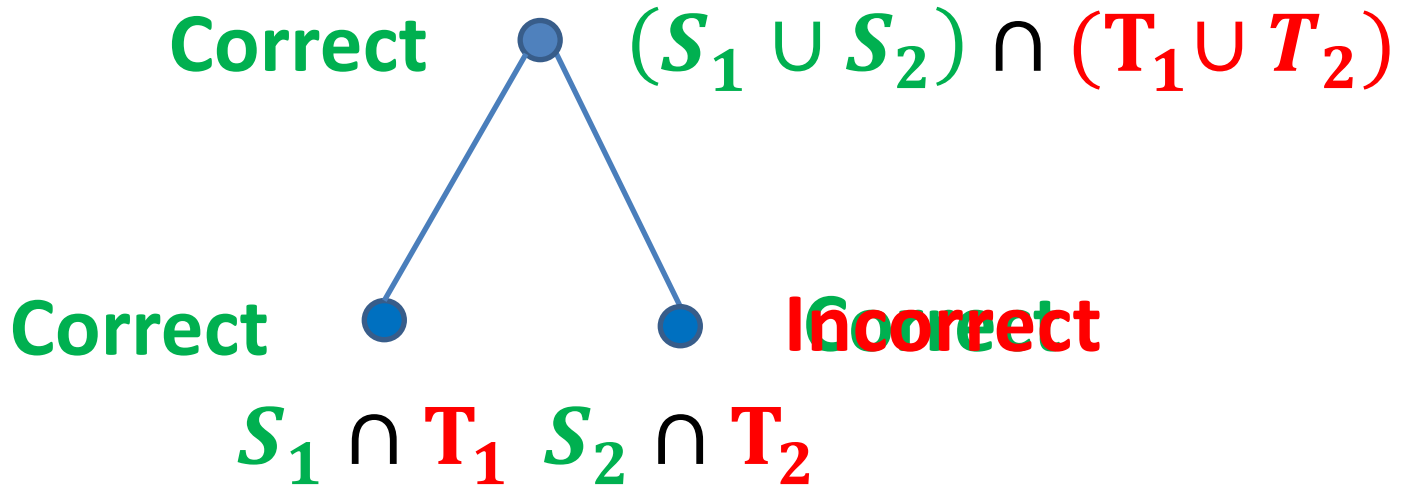
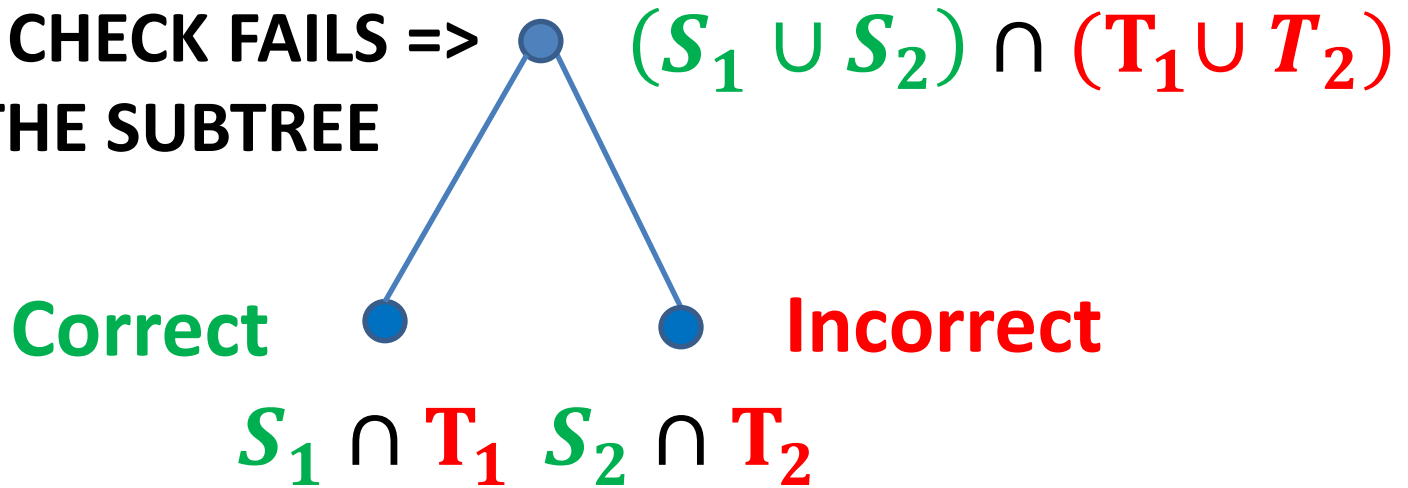


Verification bottom-up

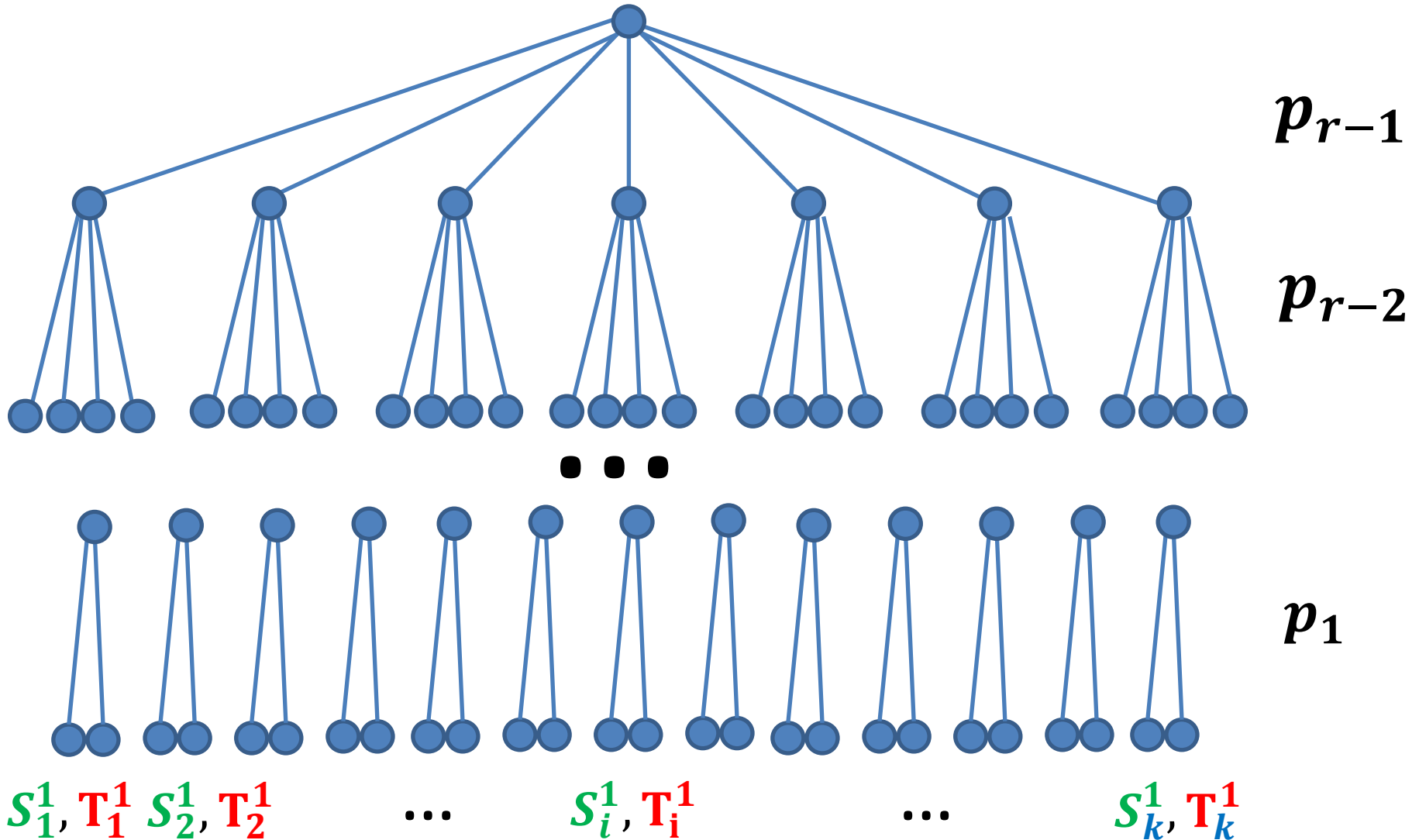


Verification bottom-up

EQUALITY CHECK FAILS =>
RESTART THE SUBTREE



Verification bottom-up



Analysis of Stage i

- $p_i = Pr[\text{node at stage } i \text{ computed correctly}]$
- Set $p_i = 1 - \frac{1}{(i \log^{r-i} k)^4}$
 - Run equality checks and basic intersection protocols with success probability p_i
 - **Key lemma:** $\mathbb{E}[\# \text{ of restarts per leaf}] = O(1) \Rightarrow$
Cost of Intersection in leafs = $O(k)$
 - Cost of Equality = $O(k i \log^r k)$
- $p_{r-1} = Pr[\text{protocol succeeds}] = 1 - 1/k^4$

Multi-party extensions

m players: S_1, \dots, S_m , where $|S_i| \leq k$

- $S = S_1 \cap \dots \cap S_m = ?$
- Boost error probability of 2-player protocol to $1 - \frac{1}{2^k}$
- Average per player (using coordinator):
 $O(k \text{ilog}^r k)$ in $O\left(r \max\left(1, \frac{\log m}{k}\right)\right)$ rounds
- Worst-case per player (using a tournament)
 $O\left(k^2 \text{ilog}^r k \max\left(1, \frac{\log m}{k}\right)\right)$ in $O\left(rk \max\left(1, \frac{\log m}{k}\right)\right)$ rounds

Open Problems

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

k -Disjointness

- $f(\mathbf{S}, \mathbf{T}) = 1$, iff $|\mathbf{S} \cap \mathbf{T}| = 0$
- $R(k\text{-Disjointness}) = \Theta(k)$ [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]
- $R^r(k\text{-Disjointness}) = \Theta(k \operatorname{ilog}^r k)$ [Saglam, Tardos'13]
- $R(k\text{-Disjointness}) = \alpha k + o(k)$ [Braverman, Garg, Pankratov, Weinstein'13]