

Course Summary

CS 598 DH

What is this course about?

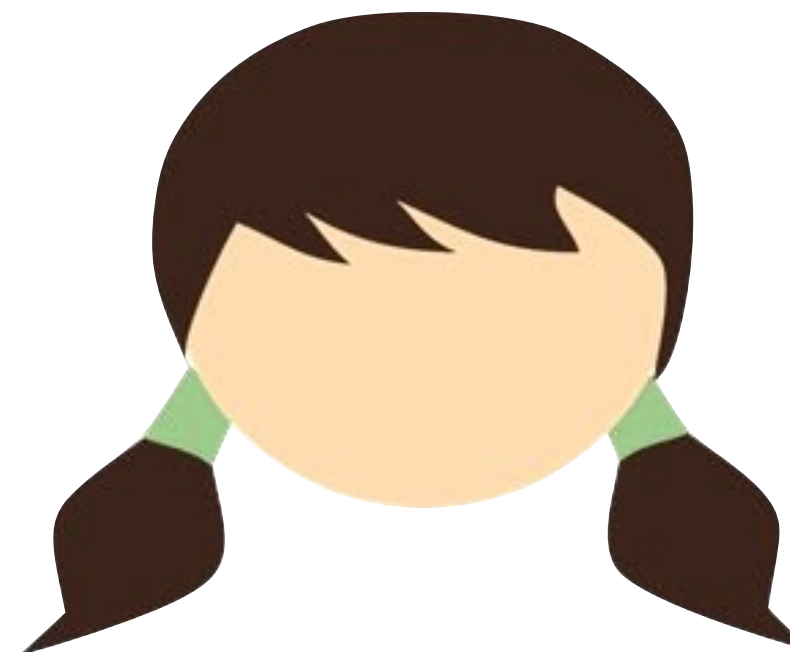
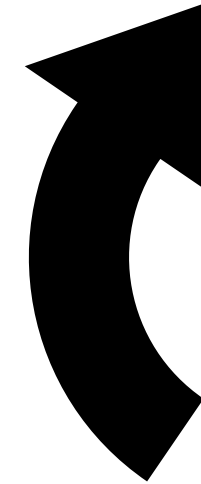
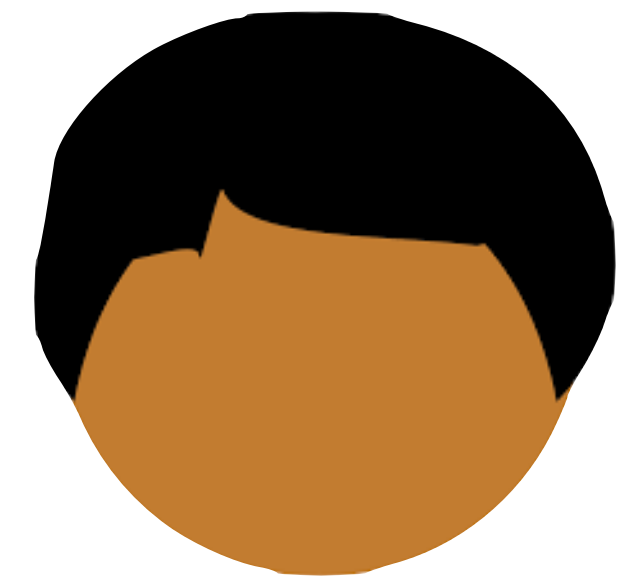
Running programs on secret-shared data.

Emulate the existence of a
trusted third party.

By doing so, we can solve
problems while keeping the
underlying data private.

What is this course *not* about?

classic cryptography setting



Privacy
Authenticity

Differentially Private Secure Multi-Party Computation for Federated Learning in Financial Applications

David Byrd
db@gatech.edu
School of Interactive Computing
Georgia Institute of Technology
Atlanta, Georgia

Antigoni Polychroniadou
antigoni.poly@jpmorgan.com
J.P. Morgan AI Research
New York, New York

ABSTRACT

Federated Learning enables a population of clients, working with a trusted server, to collaboratively learn a shared machine learning model while keeping each client's data within its own local systems. This reduces the risk of exposing sensitive data, but it is still possible to reverse engineer information about a client's private data set from communicated model parameters. Most federated learning systems therefore use differential privacy to introduce noise to the parameters. This adds uncertainty to any attempt to reveal private client data, but also reduces the accuracy of the shared model, limiting the useful scale of privacy-preserving noise. A system can further reduce the coordinating server's ability to recover private client information, without additional accuracy loss, by also including secure multiparty computation. An approach combining both techniques is especially relevant to financial firms as it allows new possibilities for collaborative learning without exposing sensitive client data. This could produce more accurate models for important tasks like optimal trade execution, credit origination, or fraud detection. The key contributions of this paper are: We present a privacy-preserving federated learning protocol to a non-specialist audience, demonstrate it using logistic regression on a real-world credit card fraud data set, and evaluate it using an open-source simulation platform which we have adapted for the development of federated learning systems.

KEYWORDS

federated learning, simulation, multiagent, finance, privacy

ACM Reference Format:

David Byrd and Antigoni Polychroniadou. 2020. Differentially Private Secure Multi-Party Computation for Federated Learning in Financial Applications. In *ACM International Conference on AI in Finance (ICAIF '20)*, October 15–16, 2020, New York, NY, USA. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3383455.3422562>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ICAIF '20, October 15–16, 2020, New York, NY, USA
© 2020 Association for Computing Machinery.
ACM ISBN 978-1-4503-7584-9/20/10...\$15.00
<https://doi.org/10.1145/3383455.3422562>

1 INTRODUCTION

Modern financial firms routinely need to conduct analysis of large data sets stored across multiple servers or devices. A typical response is to combine those data sets into a single central database, but this approach introduces a number of privacy challenges: The institution may not have appropriate authority or permission to transfer locally stored information, the owner of the data may not *want* it shared, and centralization of the data may worsen the potential consequences of a data breach.

For example, the mobile app ai.type collected personal data from its users' phones and uploaded this information to a central database. Security researchers gained access to the database and obtained the names, email addresses, passwords, and other sensitive information of 31 million users of the Android version of the app. Such incidents highlight the risks and challenges associated with centralized data solutions. [5]

In this section, we motivate our approach while providing an extensive non-technical overview of the underlying techniques.

1.1 Federated Learning

One approach to mitigate the mentioned privacy concerns is to analyze the multiple data sets separately and share only the resulting insights from each analysis. This approach is realized in a recently-introduced technique called federated analysis. [2] Federated *learning*, already adopted by large companies like Google, allows users to share insights (perhaps the parameters of a trained model) from the data on their laptops or mobile devices *without* ever sharing the data itself, typically as follows:

1. Users train a local model on their individual data.
2. Each user sends their model weights to a trusted server.
3. The server computes an average-weight shared model.
4. The shared model is returned to all of the users.
5. Users retrain a local model starting from the shared model.

For instance, email providers could use federated learning to reduce the amount of spam their customers receive. Instead of each provider using its own spam filter trained from its customers' reported spam email, the providers could combine their models to create a shared spam-detection mechanism, without sharing their individual customers' reported spam emails. For a survey of recent advances in federated learning, see Kairouz et al. [13]

It is still possible, however, for a malicious party to potentially compromise the privacy of the individual users by inferring details of a training data set from the trained model's weights or parameters [16, 19]. It is important to protect sensitive user information while still providing highly accurate inferences.

Secure Auctions

Privacy-preserving studies

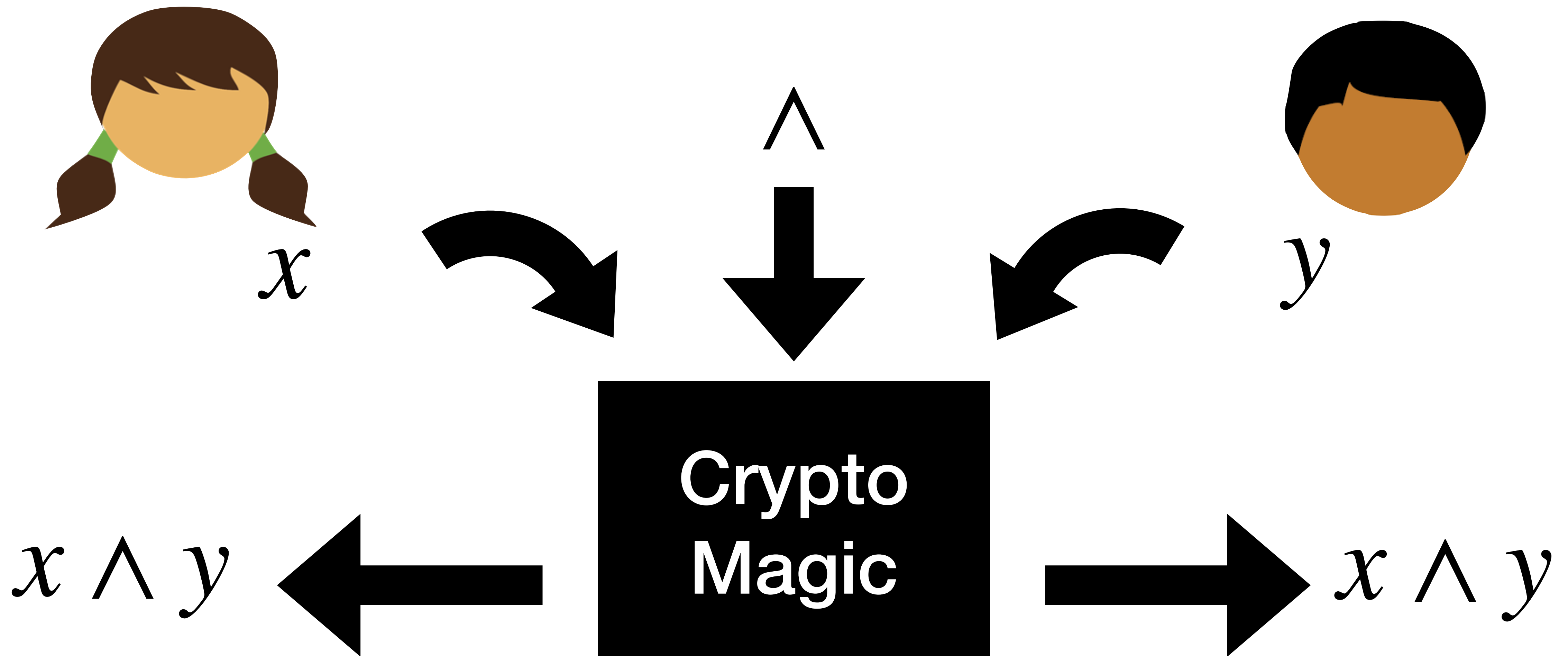
Privacy-preserving advertising

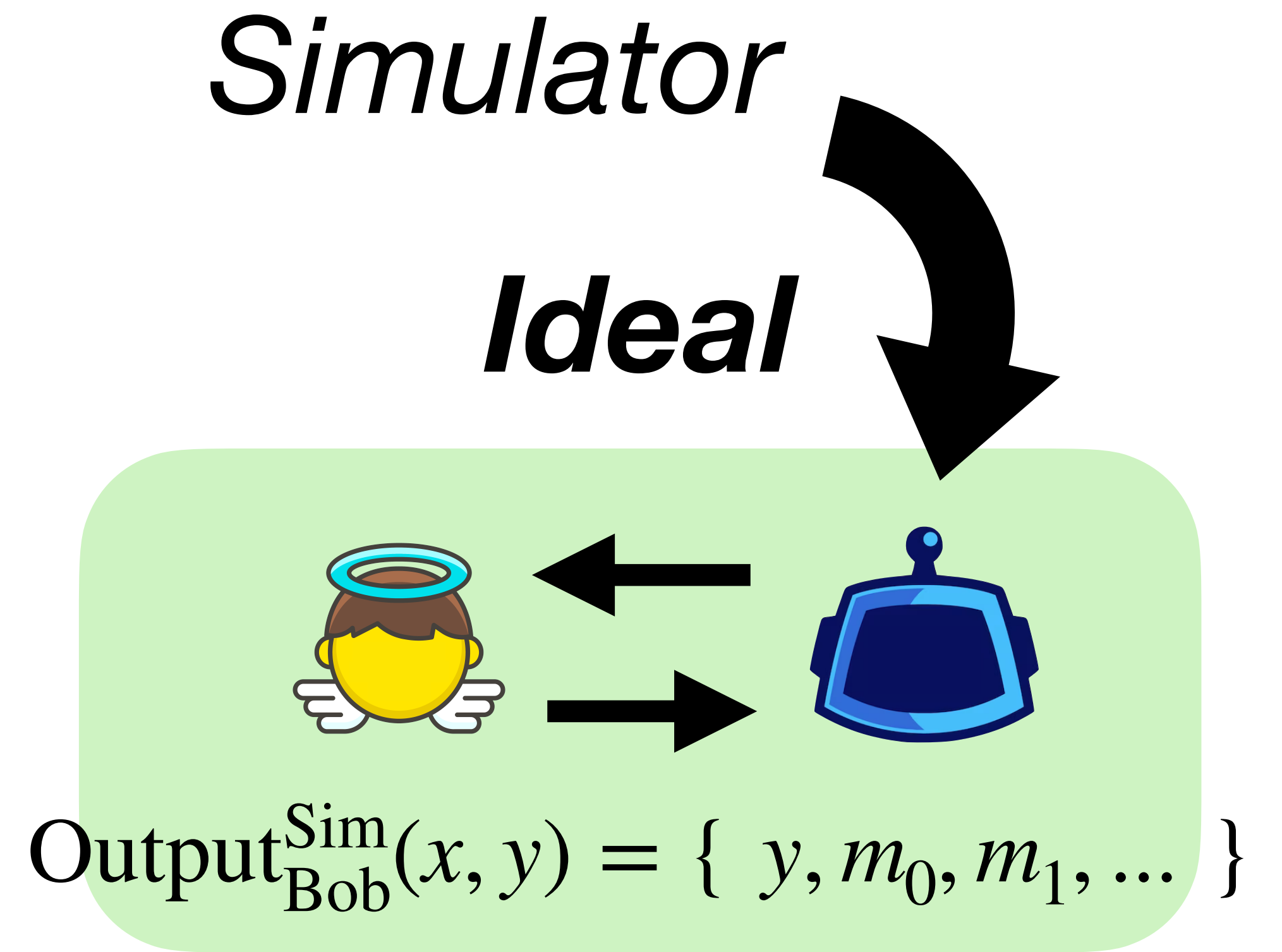
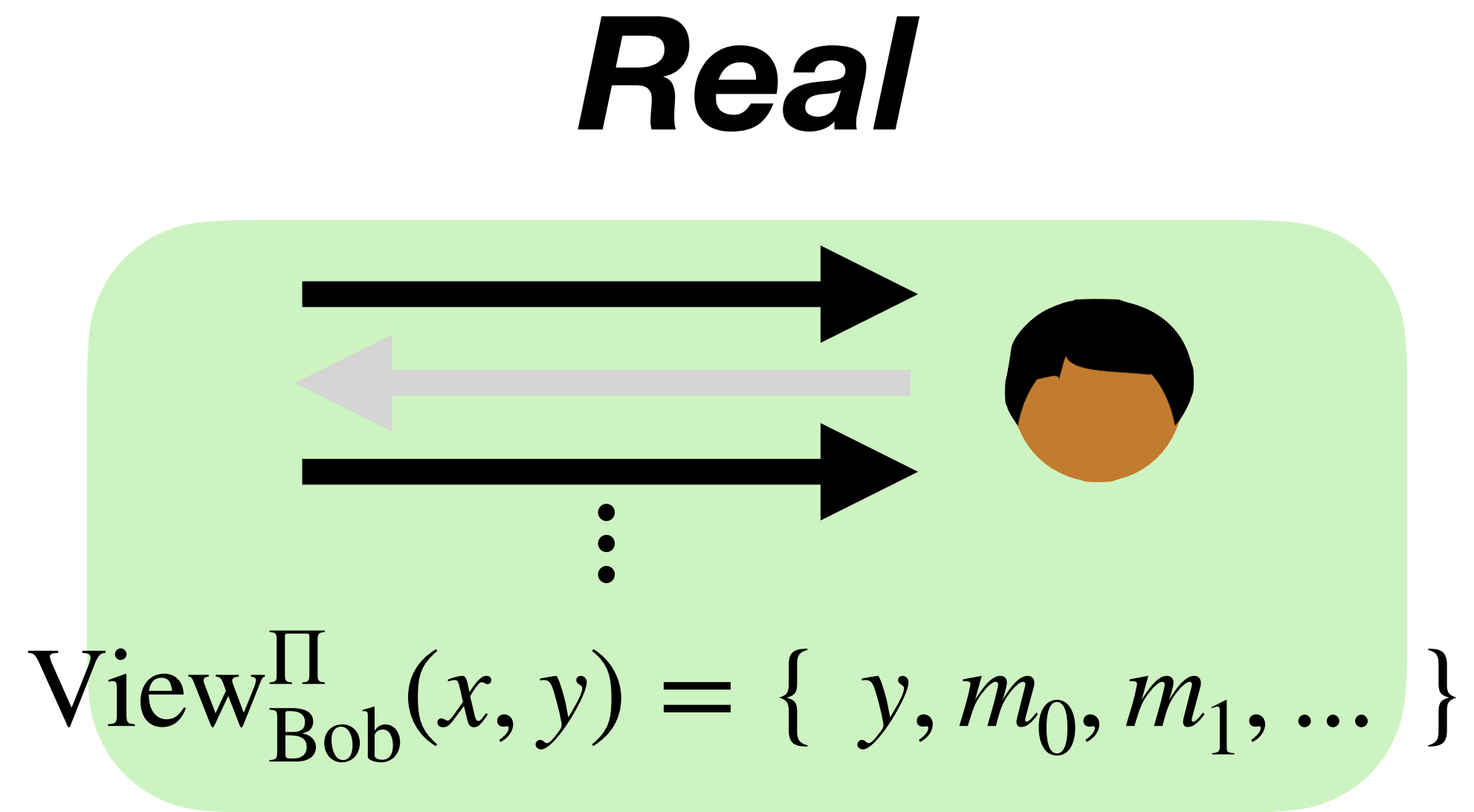
Privacy-preserving analytics

(Secure Machine Learning)

Financial Fraud Detection

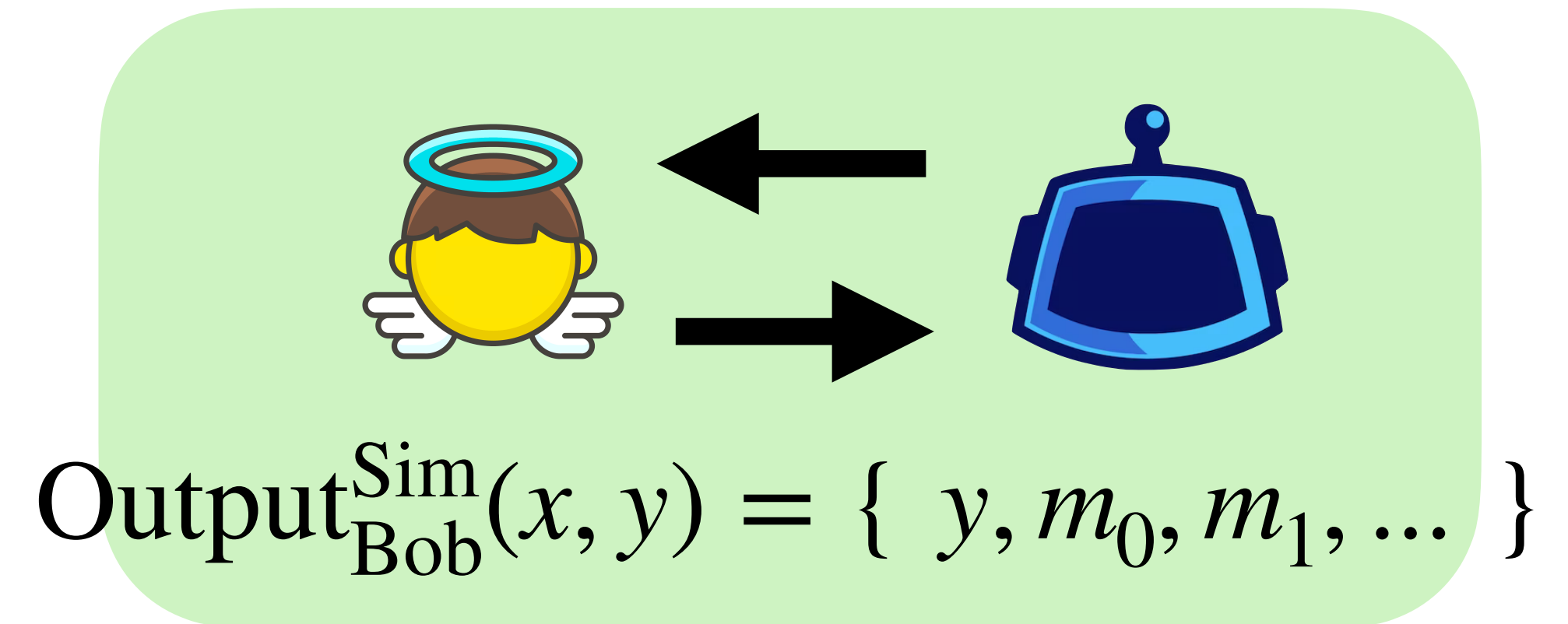
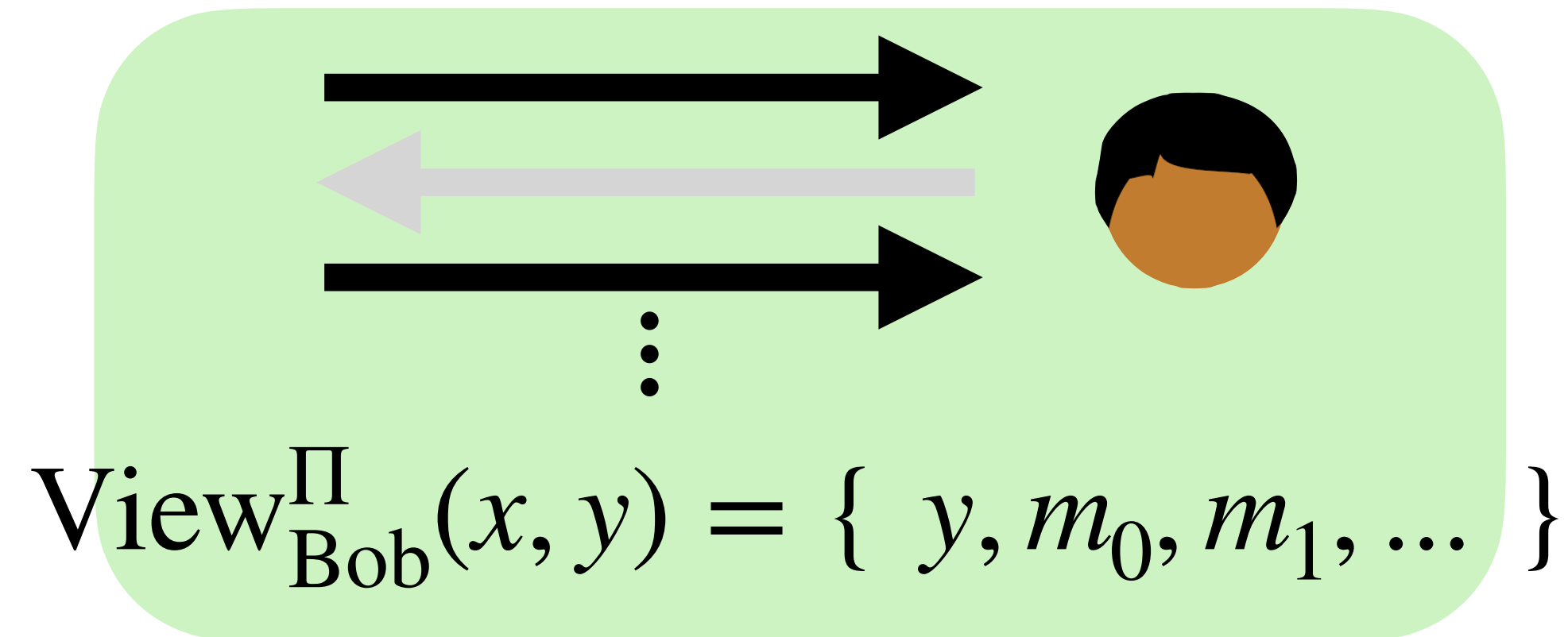
...and much more





These should “look the same”

“No efficient algorithm can tell these two things apart”



Three notions of “hard to tell apart”

$X \equiv Y$ Identically distributed

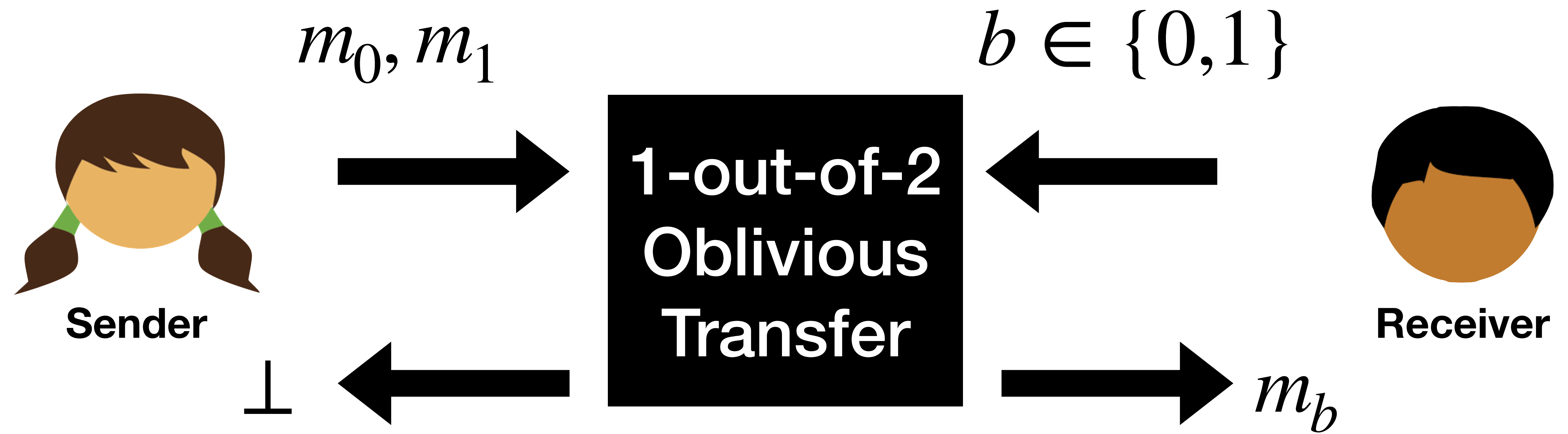
$X \approx Y$ Statistically close As we increase a parameter, the distributions **quickly** become close together.

$X \stackrel{c}{=} Y$ Indistinguishable As we increase a parameter, it **quickly** becomes difficult for programs to tell the distributions apart.

Two-Party Semi-Honest Security

Let f be a functionality. We say that a protocol Π securely computes f in the presence of a semi-honest adversary if for each party $i \in \{0,1\}$ there exists a polynomial time simulator \mathcal{S}_i such that for all inputs x_0, x_1 :

$$\begin{aligned} & \{ \text{View}_i^\Pi(x_0, x_1), \text{Output}^\Pi(x_0, x_1) \} \\ & \quad \stackrel{\mathcal{C}}{=} \\ & \{ \mathcal{S}_i(x_i, y_i), (y_0, y_1) \mid (y_0, y_1) \leftarrow f(x_0, x_1) \} \end{aligned}$$



OT Extension

In MPC (e.g., GMW), we need lots of short OTs

Can we turn a few OTs into a lot of OTs?

λ *base* OTs



n *extended* OTs

Public key

Symmetric key

S

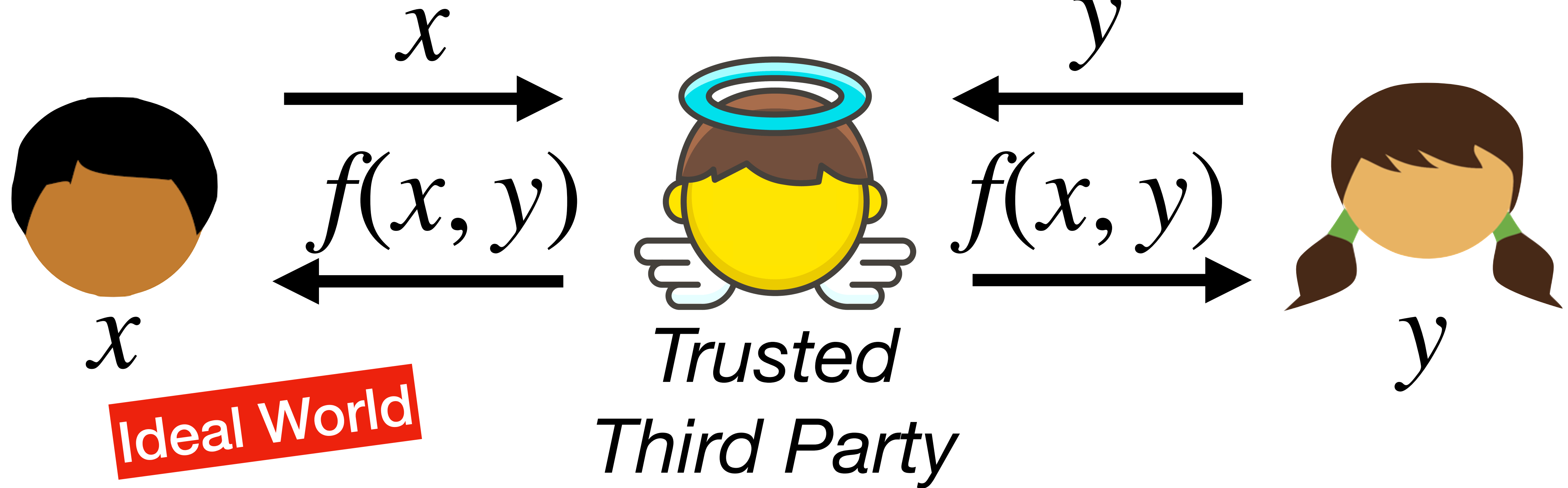
Δ	0	1	0	0	0	1	1	0
	<hr/>							
	1	0	0	1	1	1	1	1
	1	0	1	1	0	1	0	0
	0	1	1	0	0	0	1	1
	1	1	0	1	1	0	1	1
q	0	1	1	0	1	1	1	1
	1	0	1	0	1	0	0	0
	1	0	1	1	1	1	0	0
	1	0	0	0	0	0	1	1

R

r	1	1	0	1	1	0	0	1
	0	1	1	1	0	1	0	0
	0	1	1	0	0	0	1	1
	0	1	1	0	1	1	0	1
	1	0	0	1	0	1	0	0
	0	1	0	1	1	0	0	0
	1	1	1	1	1	0	1	0
	1	1	0	0	0	1	0	1

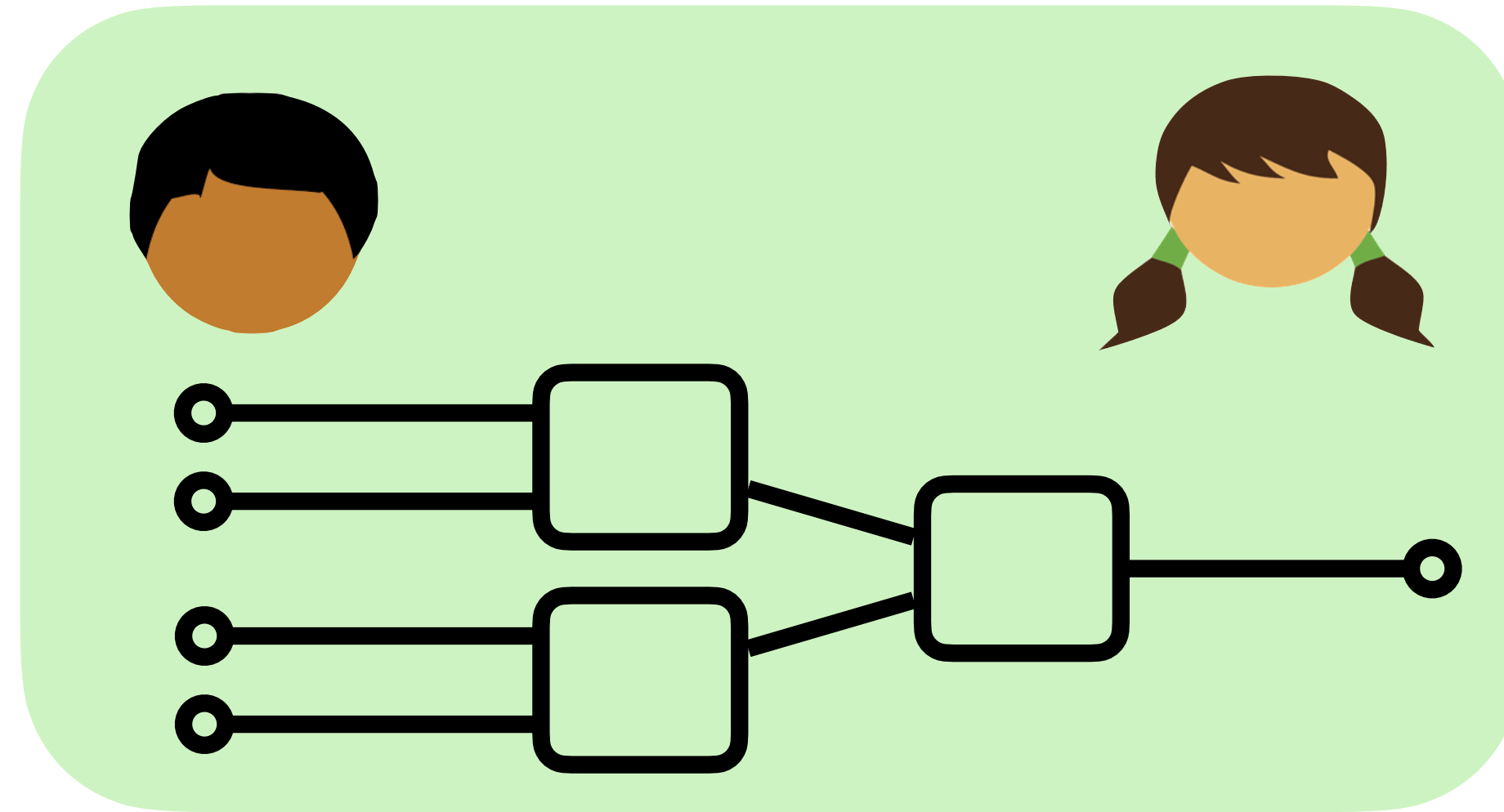
								t





GMW Protocol
Hint: Lots of OT

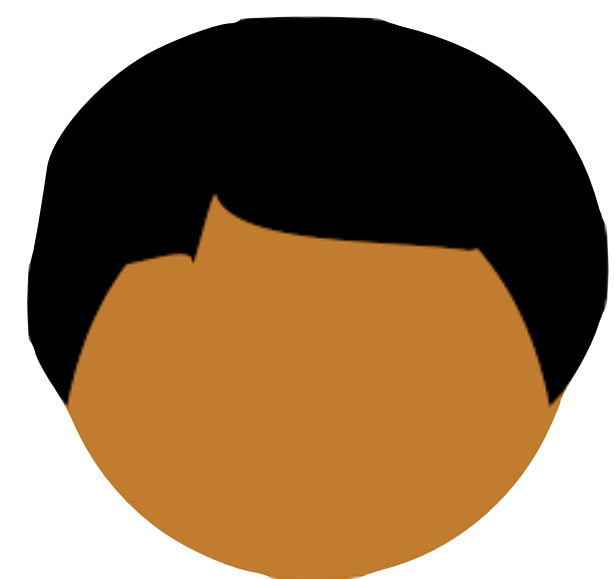
Real World



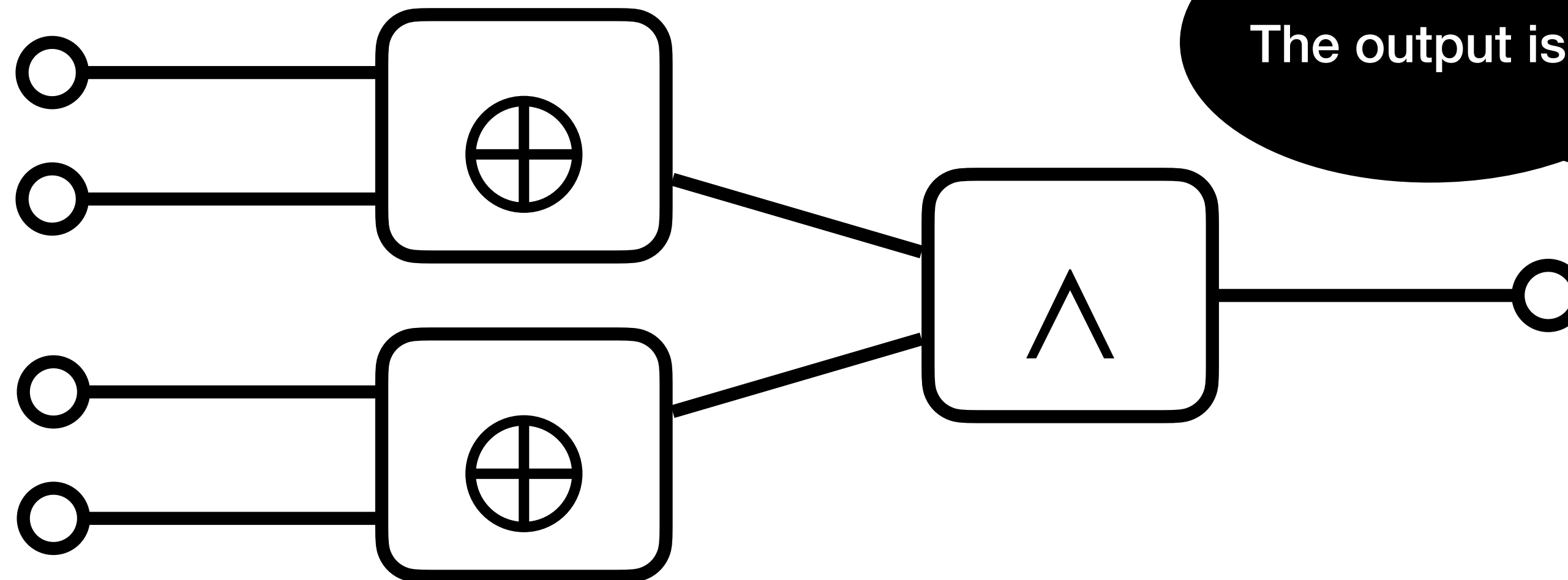
In GMW, Number of protocol rounds scales with multiplicative depth of C



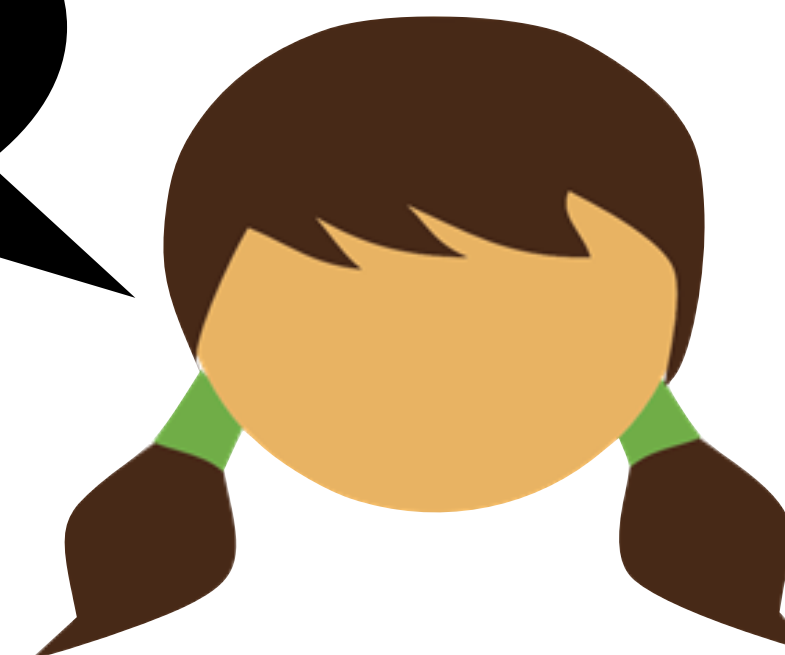
Our protocol's efficiency is fundamentally bounded by the speed of light



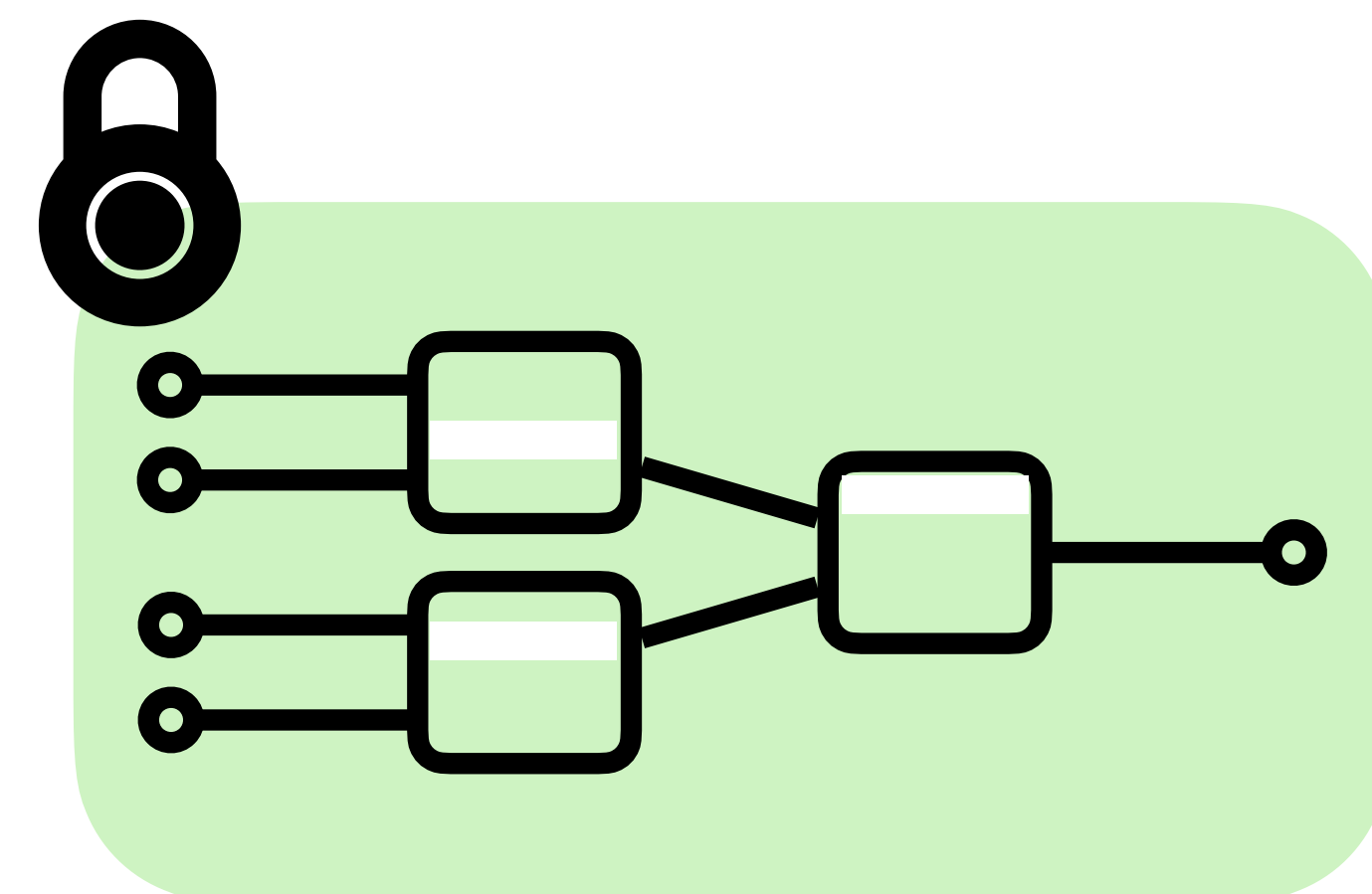
Garbler



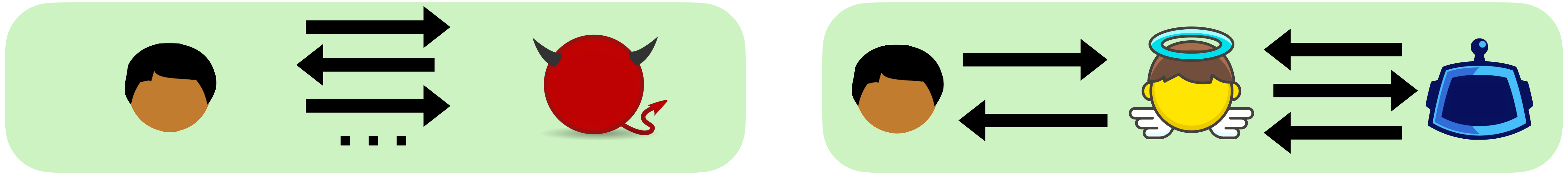
The output is 1



Evaluator



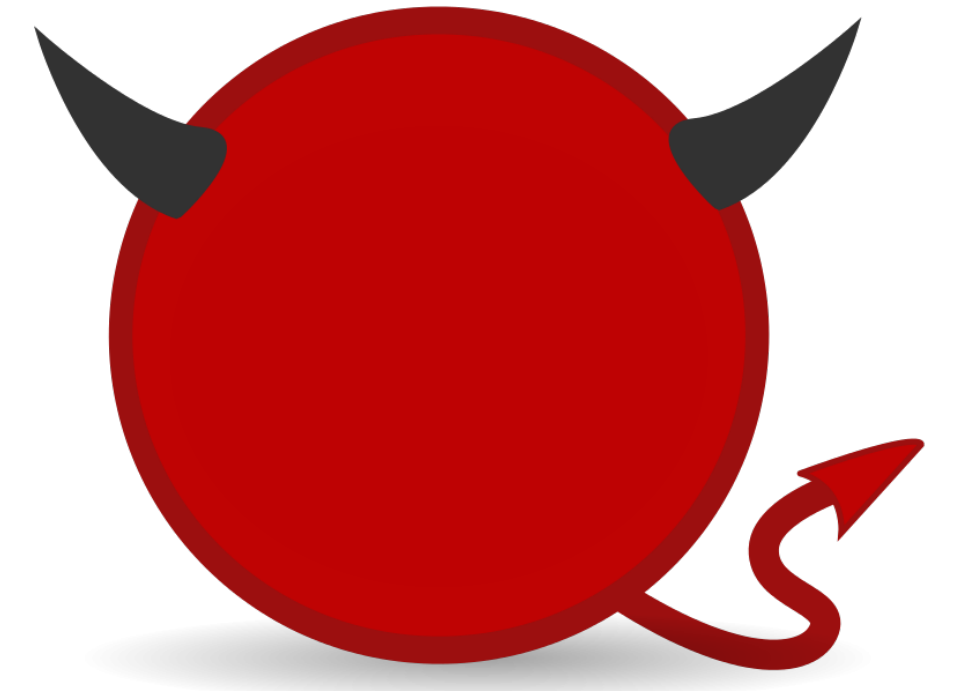
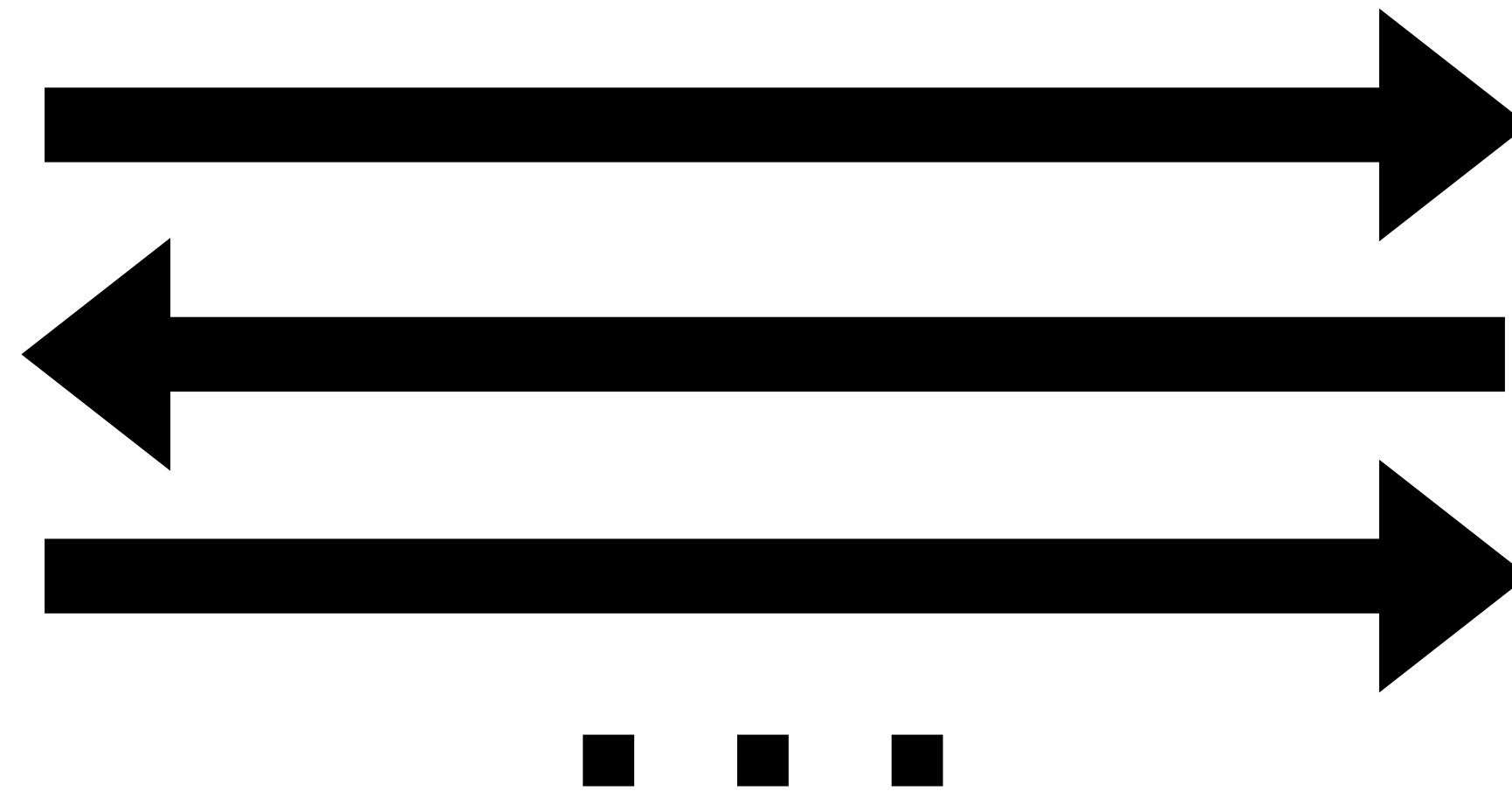
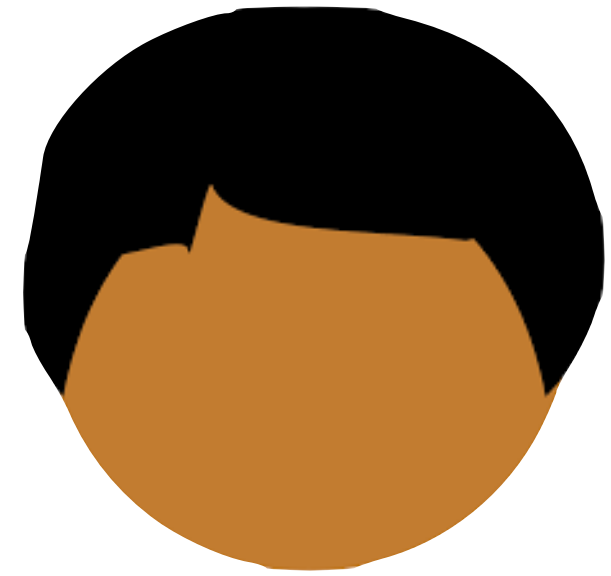
Malicious Security



*A protocol Π securely realizes a functionality f in the presence of a malicious adversary if for **every** real-world adversary \mathcal{A} corrupting party i , **there exists** an ideal-world adversary \mathcal{S}_i (a simulator) such that for all inputs x, y the following holds:*

$$\text{Real}_{\mathcal{A}}^{\Pi}(x, y) \approx \text{Ideal}_{\mathcal{S}_i}^f(x, y)$$

Ensemble of outputs of **each** party



What can go *in terms of outcomes*?

Cause honest party to output wrong answer



Learn too much information about other party's input



Prevent honest party from learning output



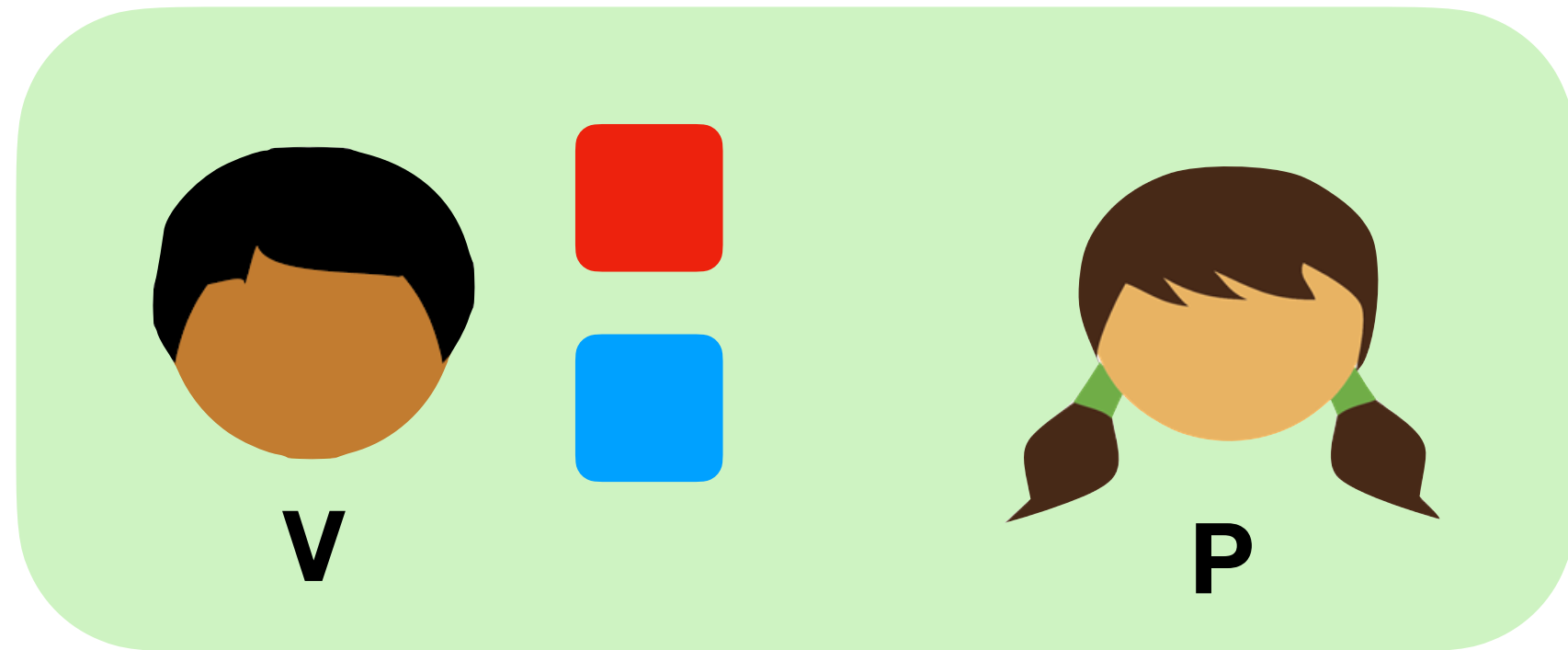
Malicious security ideal-world execution



honest party outputs
 $f(x, y')$

adversary outputs... ?
whatever it wants

What is a *zero-knowledge* proof?



Completeness: If $x \in \mathcal{L}$ and if P and V are honest, then V accepts the proof (except with negligible probability)

“P can prove true things”

Soundness: If $x \notin \mathcal{L}$, even malicious P cannot cause honest V to accept the proof

“P cannot prove false things”

Zero Knowledge: “V learns nothing except that the thing is true”

ZK from MPC in the Head

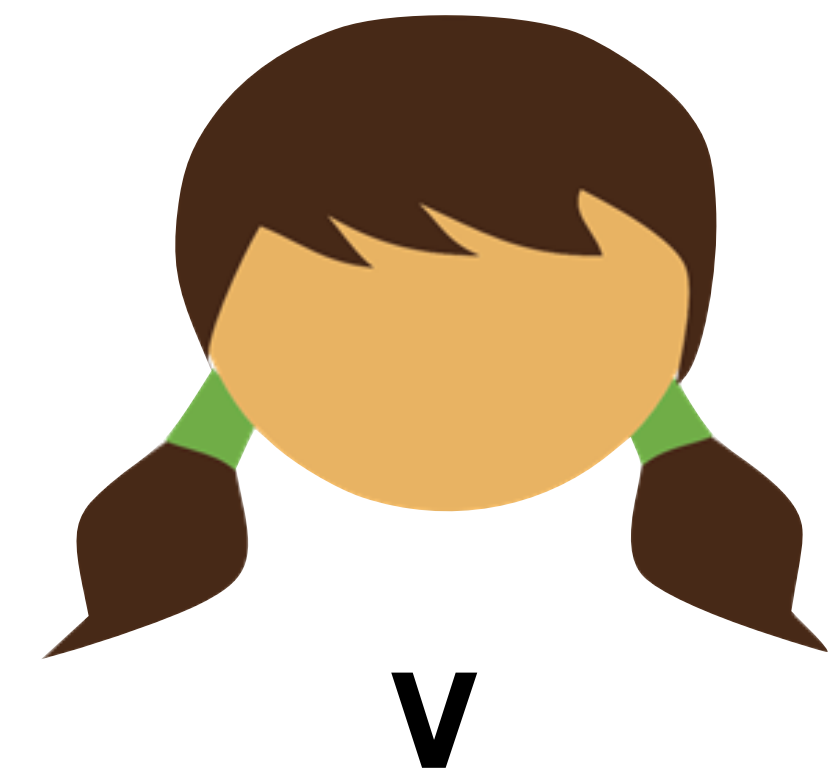
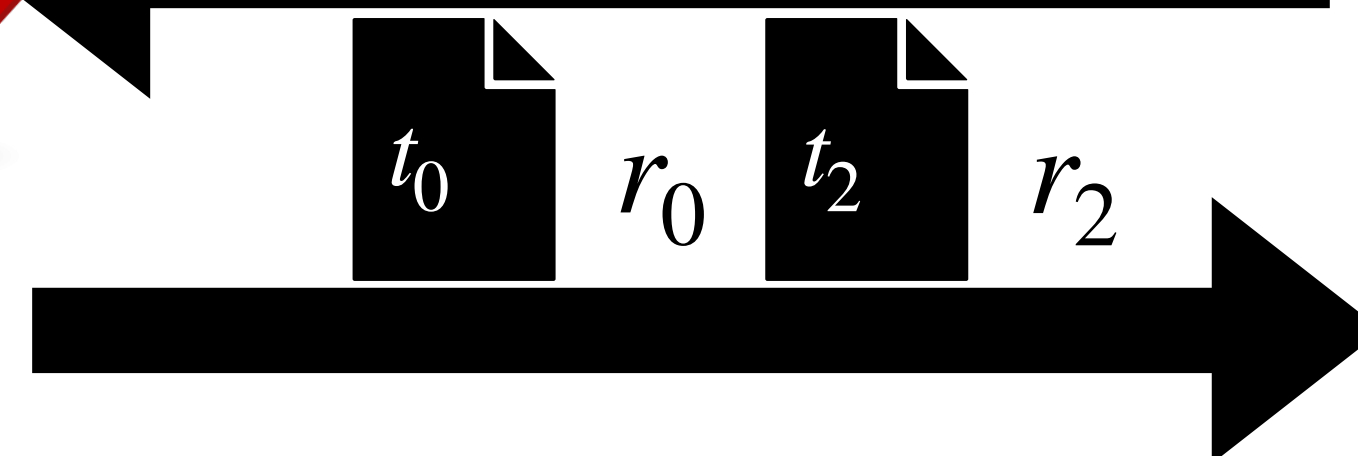
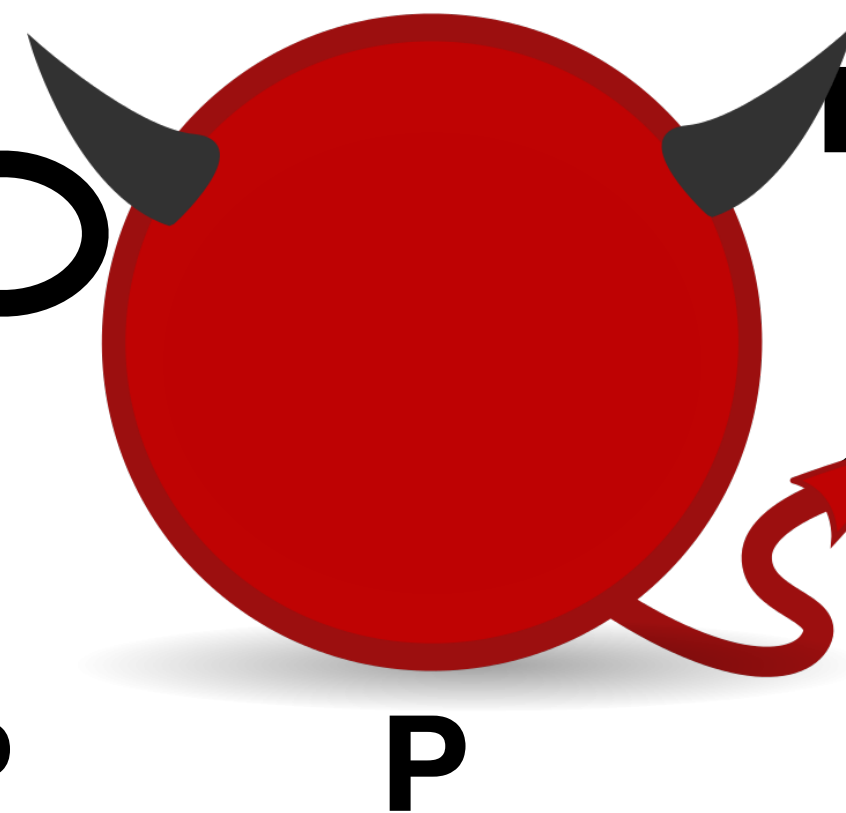
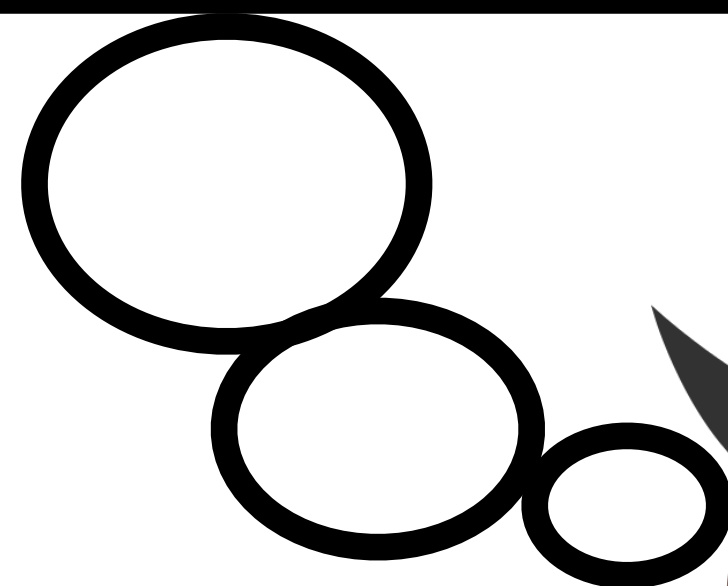
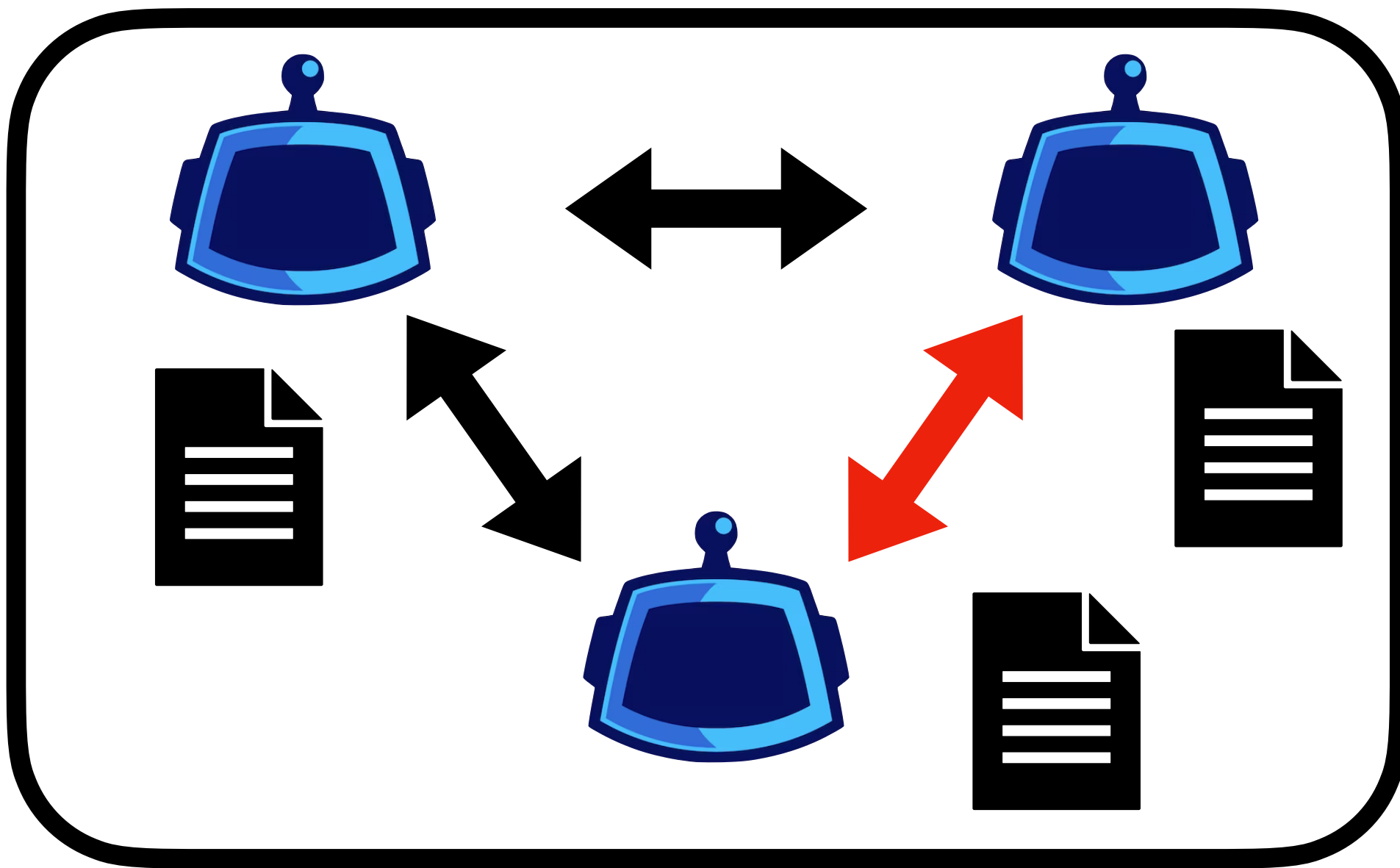
Soundness?

To cheat, P must corrupt at least one edge (i.e., one party receives a message that was not sent by the other)

$\text{com}(t_0, r_0)$
 $\text{com}(t_1, r_1)$
 $\text{com}(t_2, r_2)$

By opening an edge, V has probability at least 1/3 to catch cheating P

Repeat to obtain desired soundness

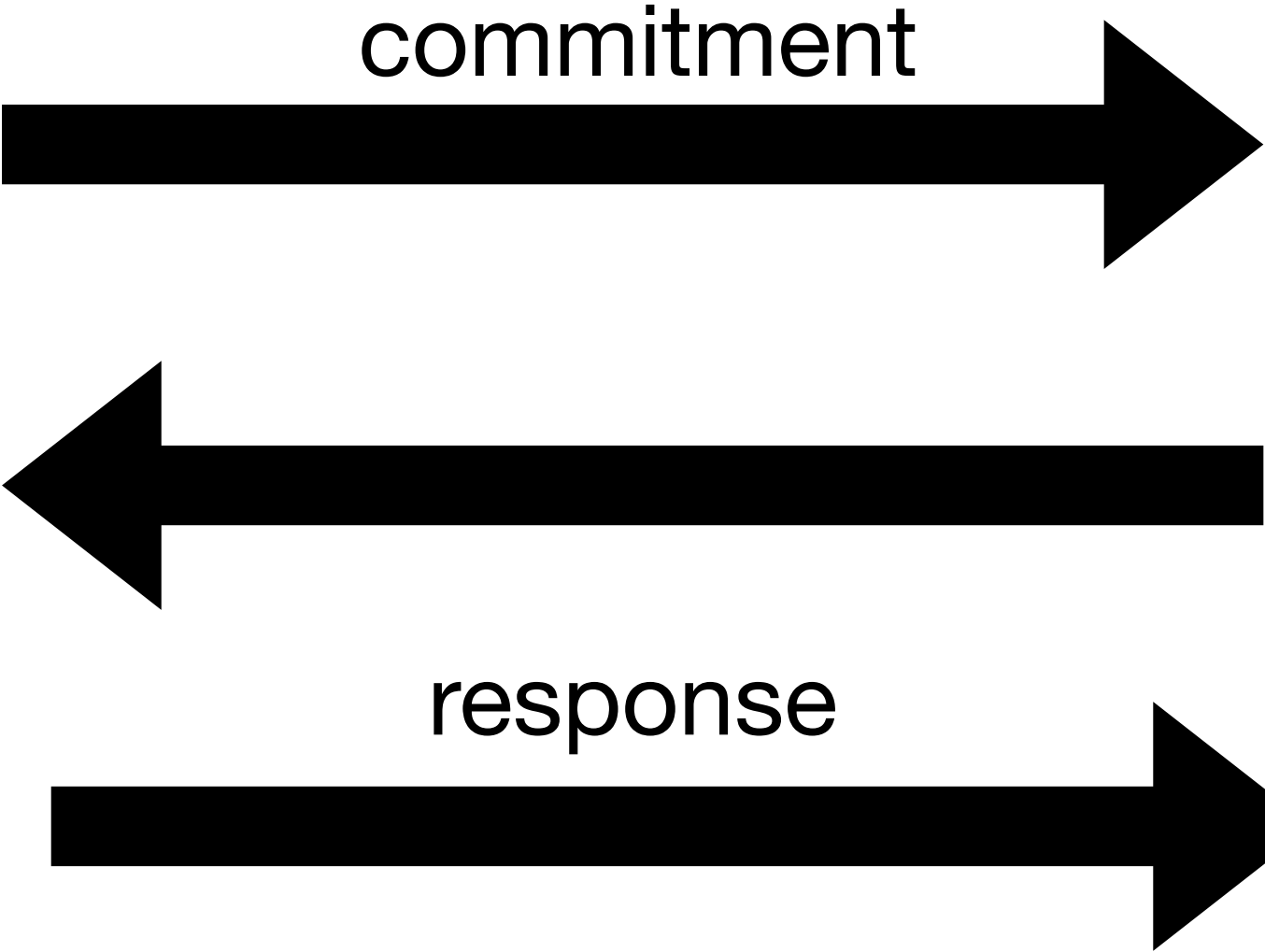
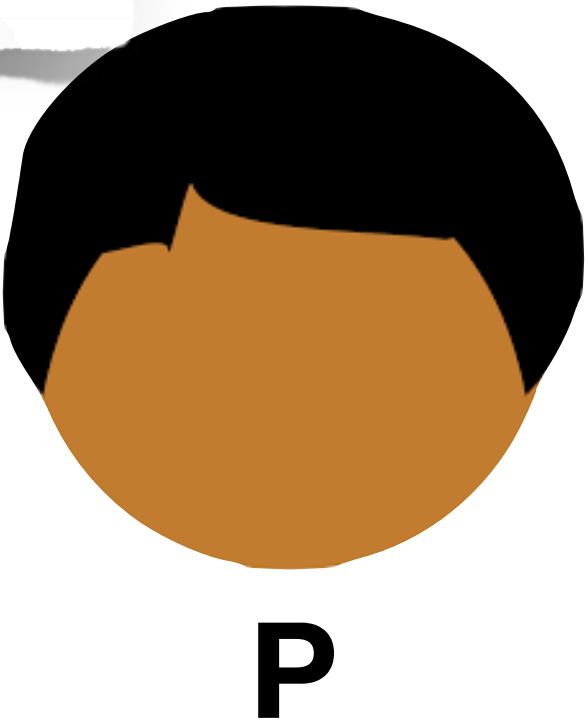




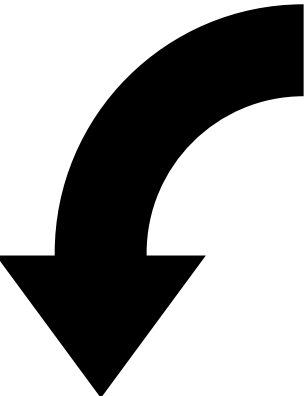
Fiat Shamir Heuristic

Public coin ZK can be made **non-interactive**

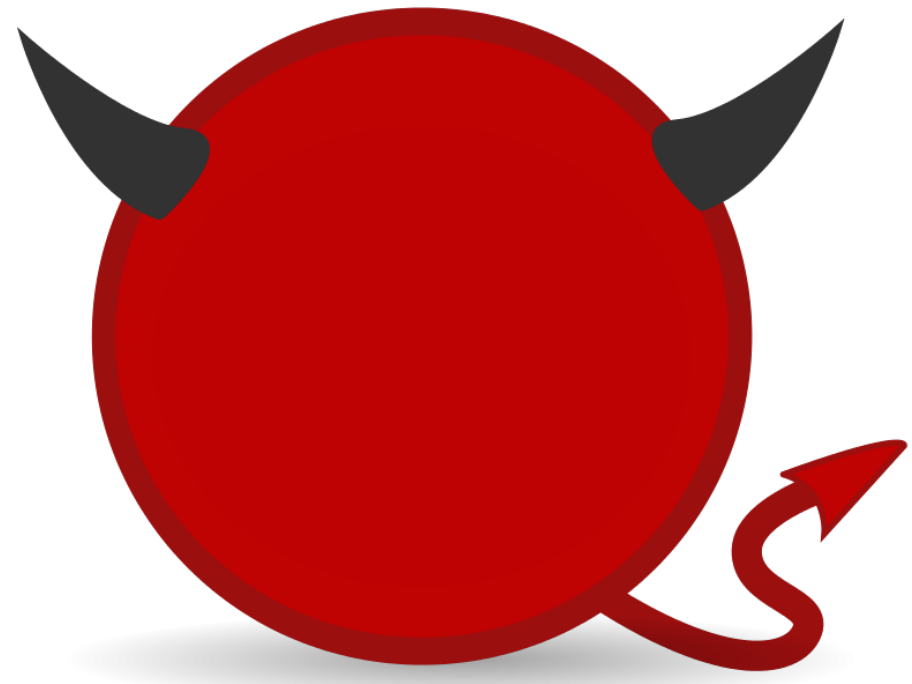
Simple idea: P can choose the challenge itself



Cryptographic hash function
(e.g. SHA 256)
Formally, a random oracle



challenge = $H(\text{commitment})$



Garbler

$$\text{Enc}(K_a^0, \text{Enc}(K_b^0, K_c^0))$$

$$\text{Enc}(K_a^0, \text{Enc}(K_b^1, K_c^0))$$

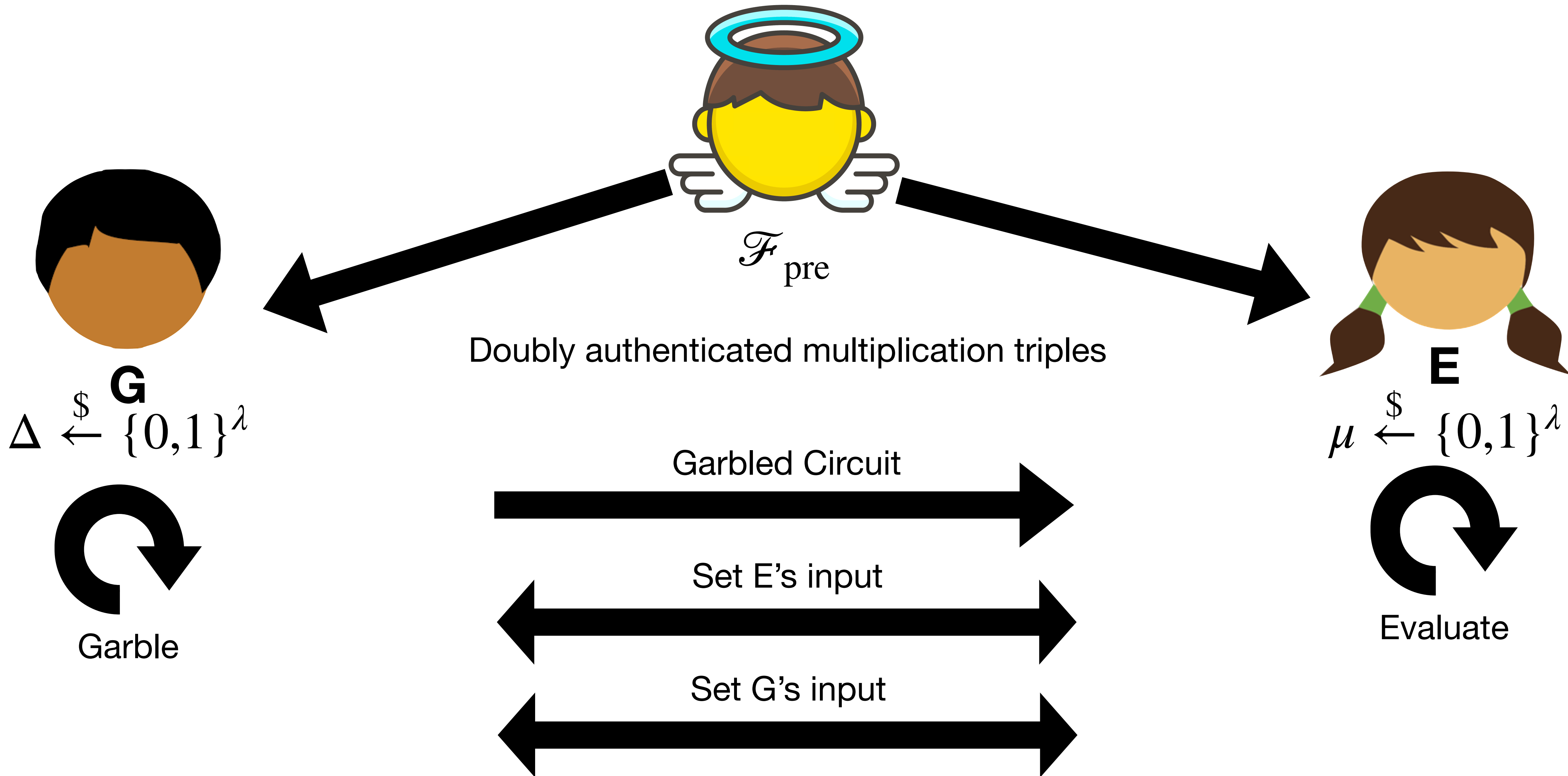
$$\text{Enc}(K_a^1, \text{Enc}(K_b^0, K_c^0))$$

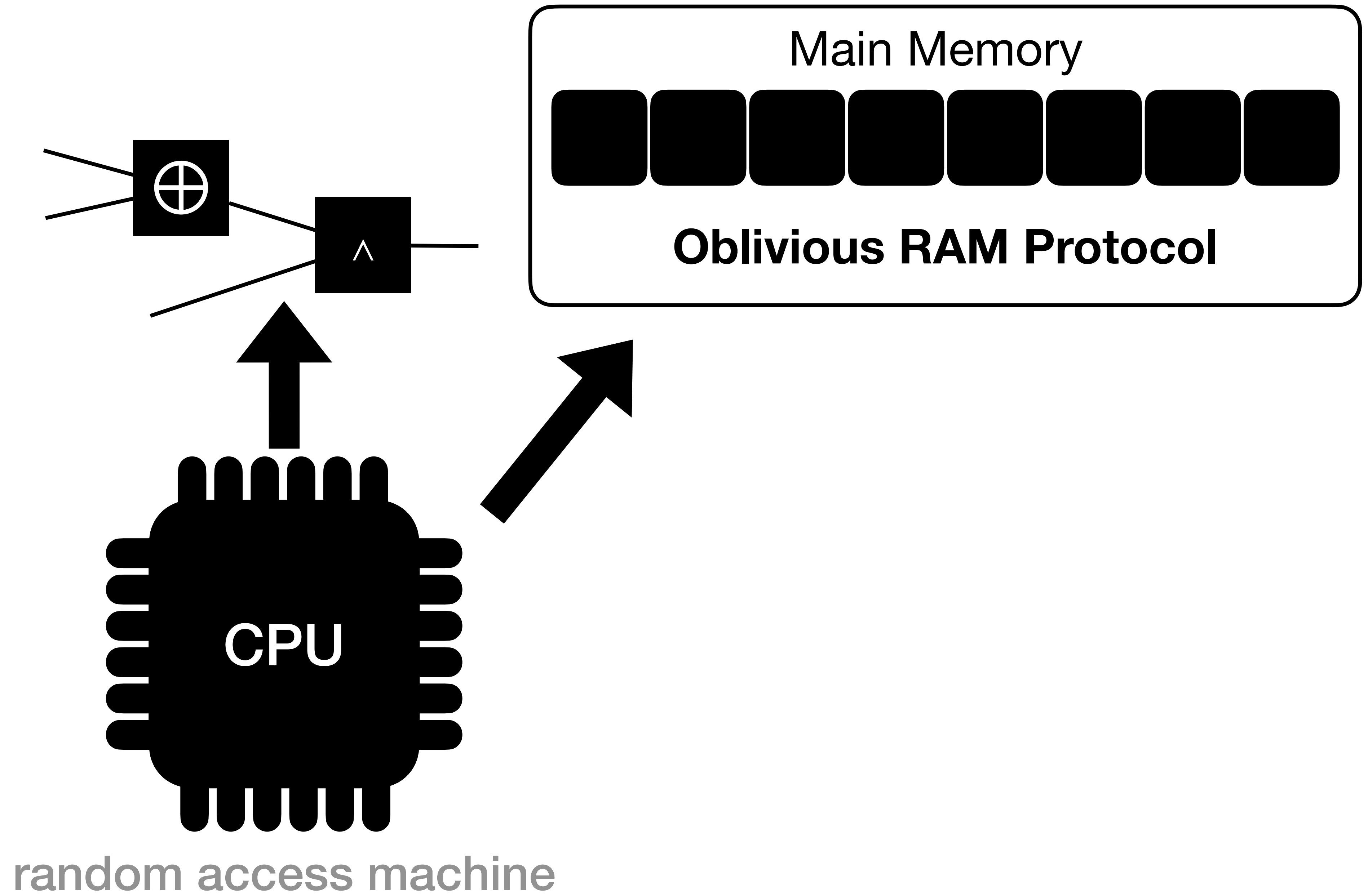
$$\text{Enc}(K_a^1, \text{Enc}(K_b^1, K_c^1))$$

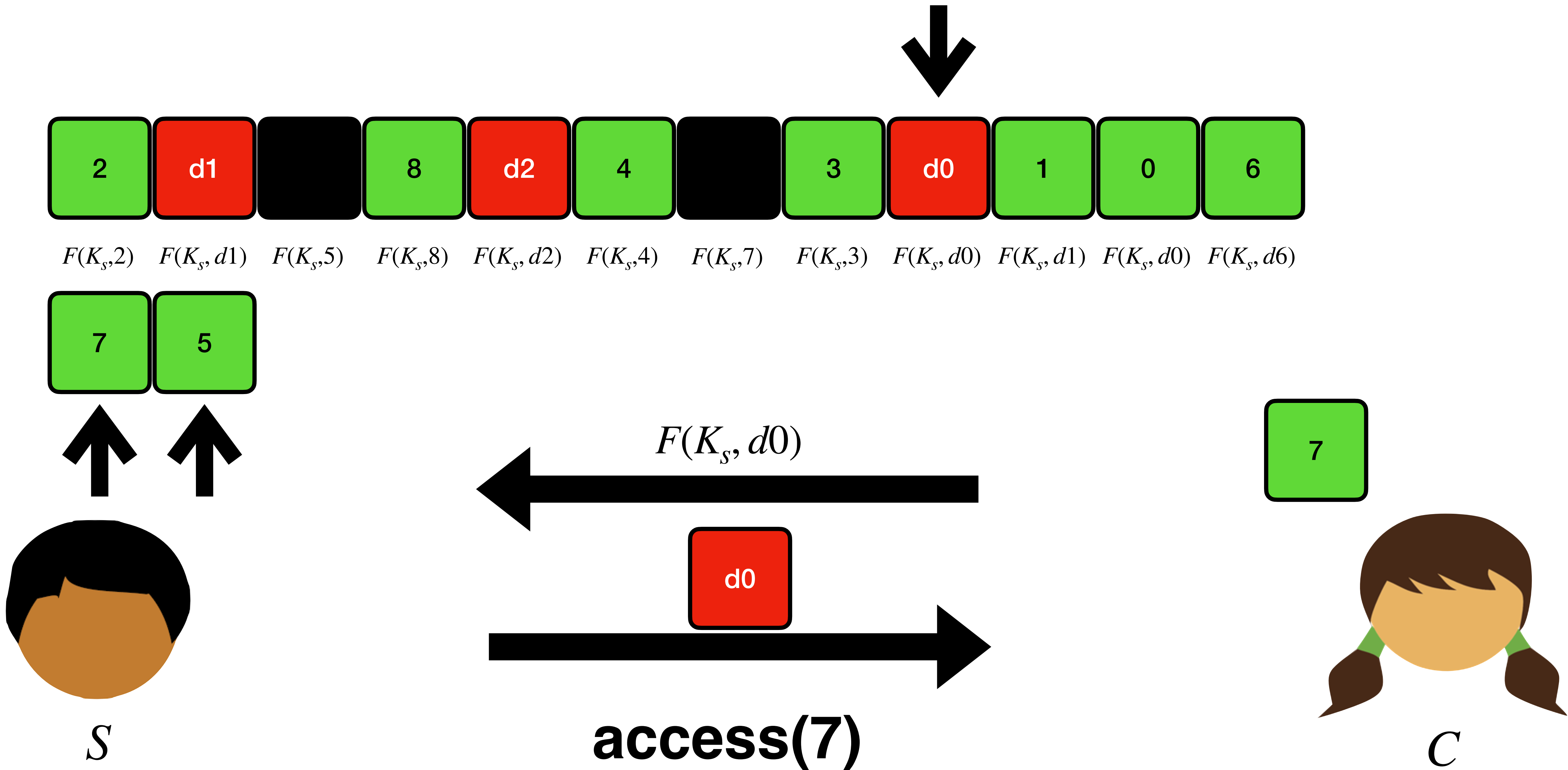
Why can't we simulate G?

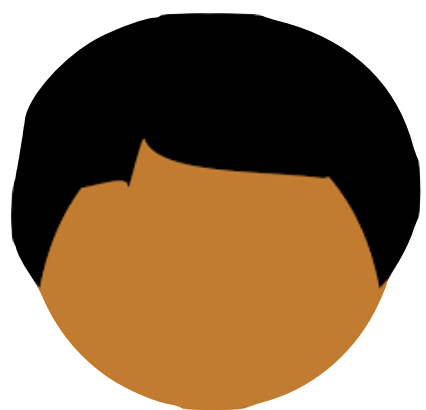
G can encrypt each gate *freely*

E has no way to tell if gate it
correctly garbled

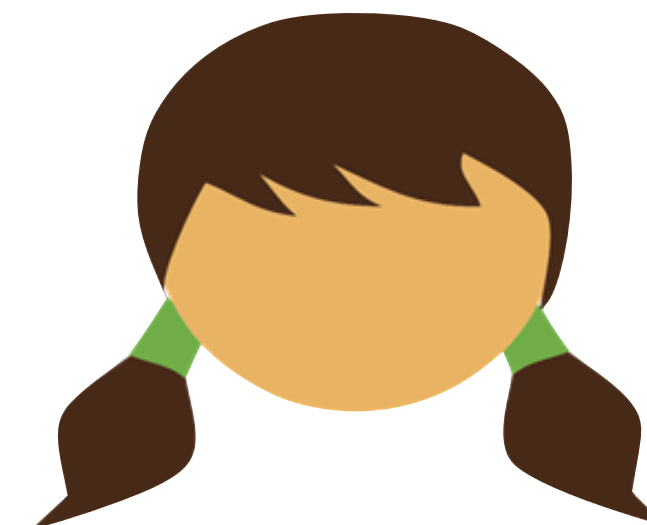
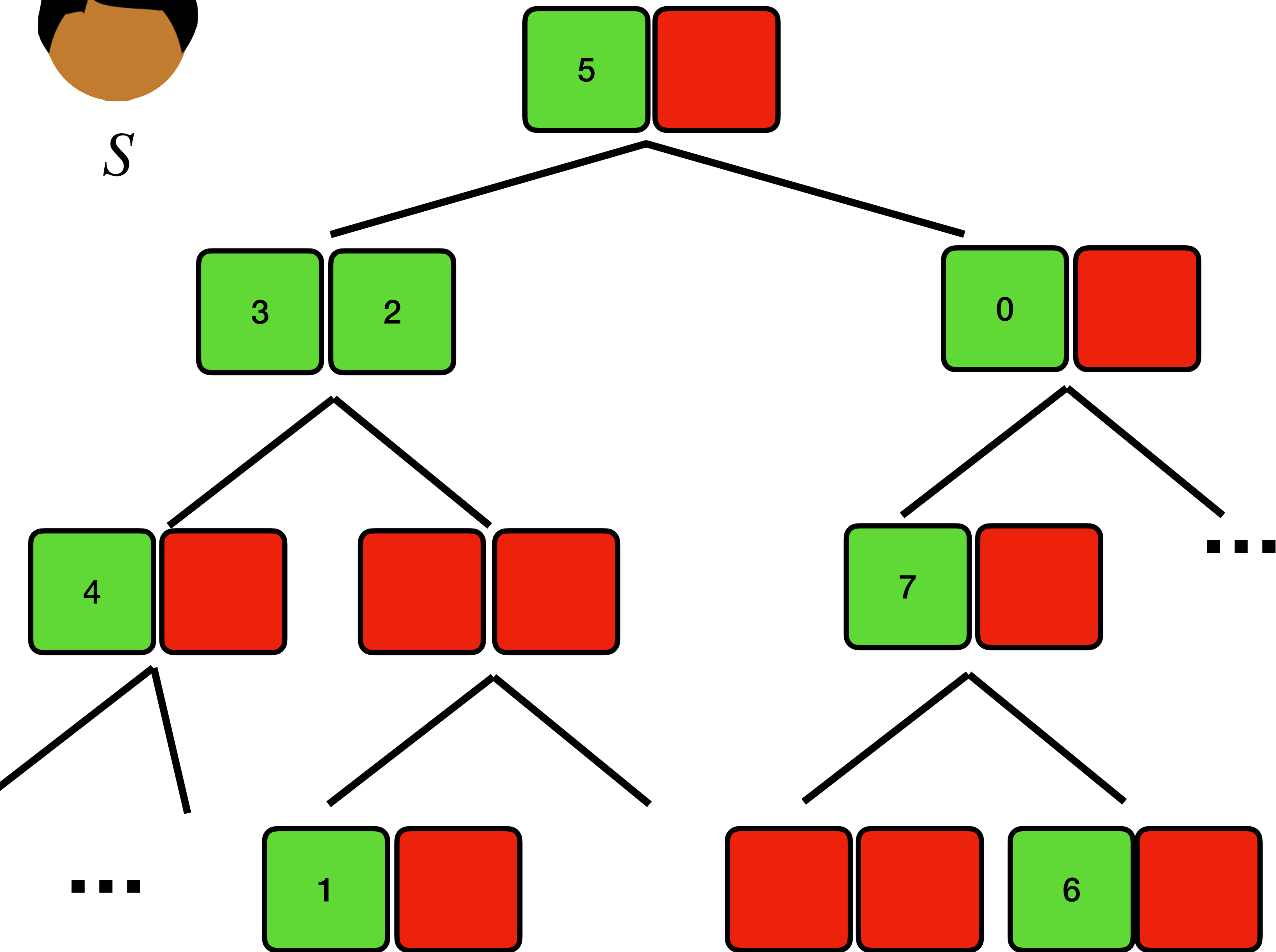








S



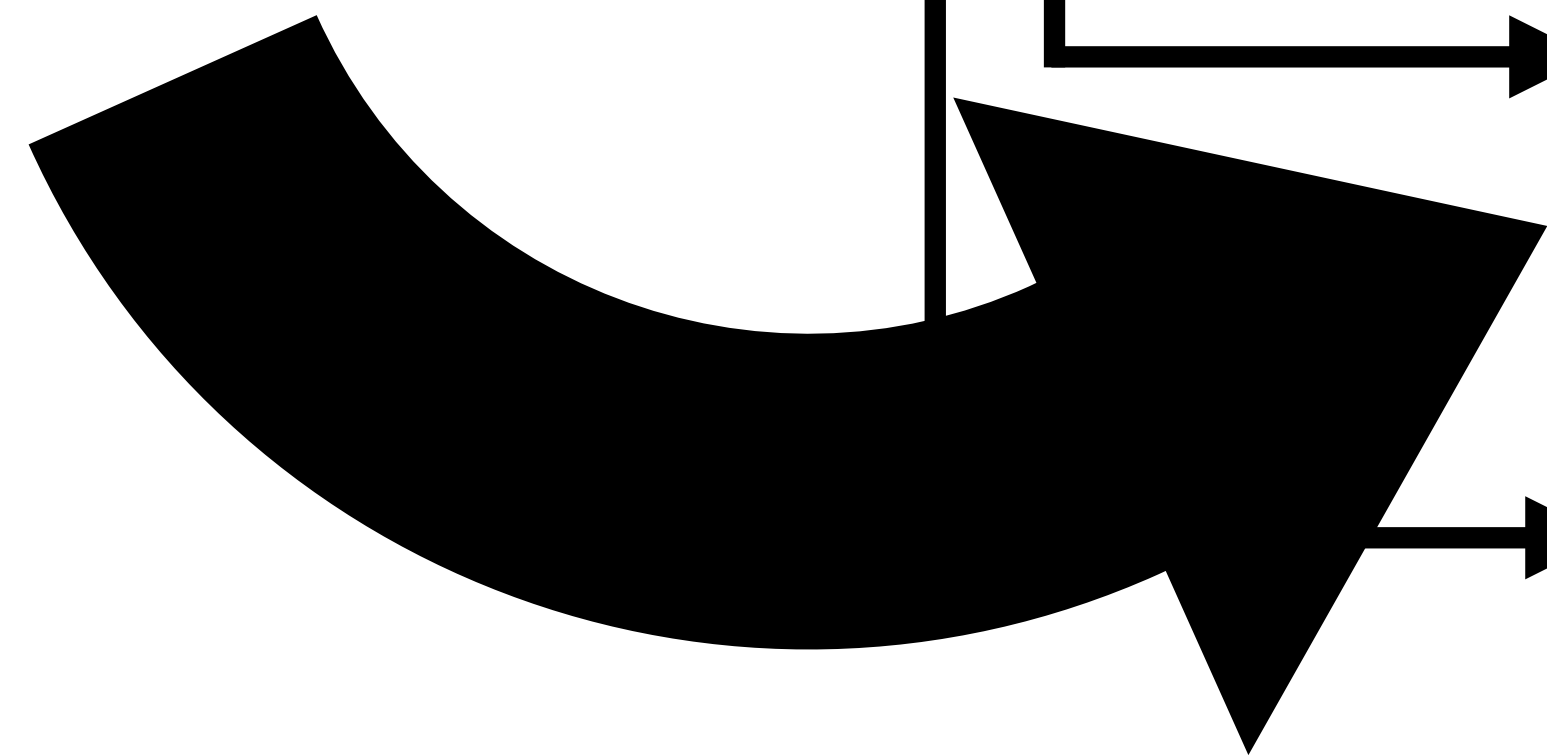
C

Path Invariant: Each node is assigned a uniformly random leaf

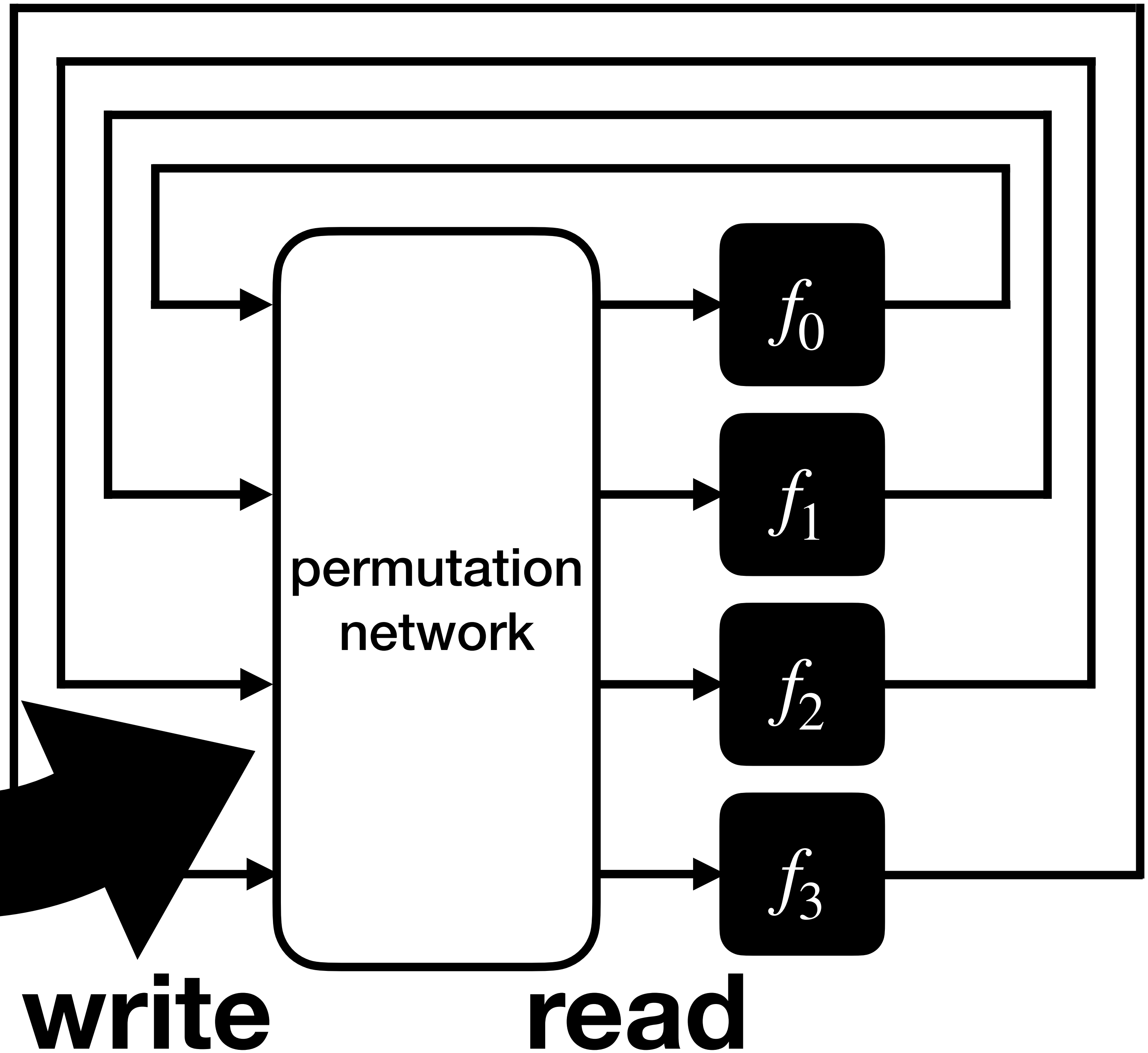
Logical address	Leaf
0	10
1	5
2	7
...	...

Position Map

**Quasilinear
Size!**

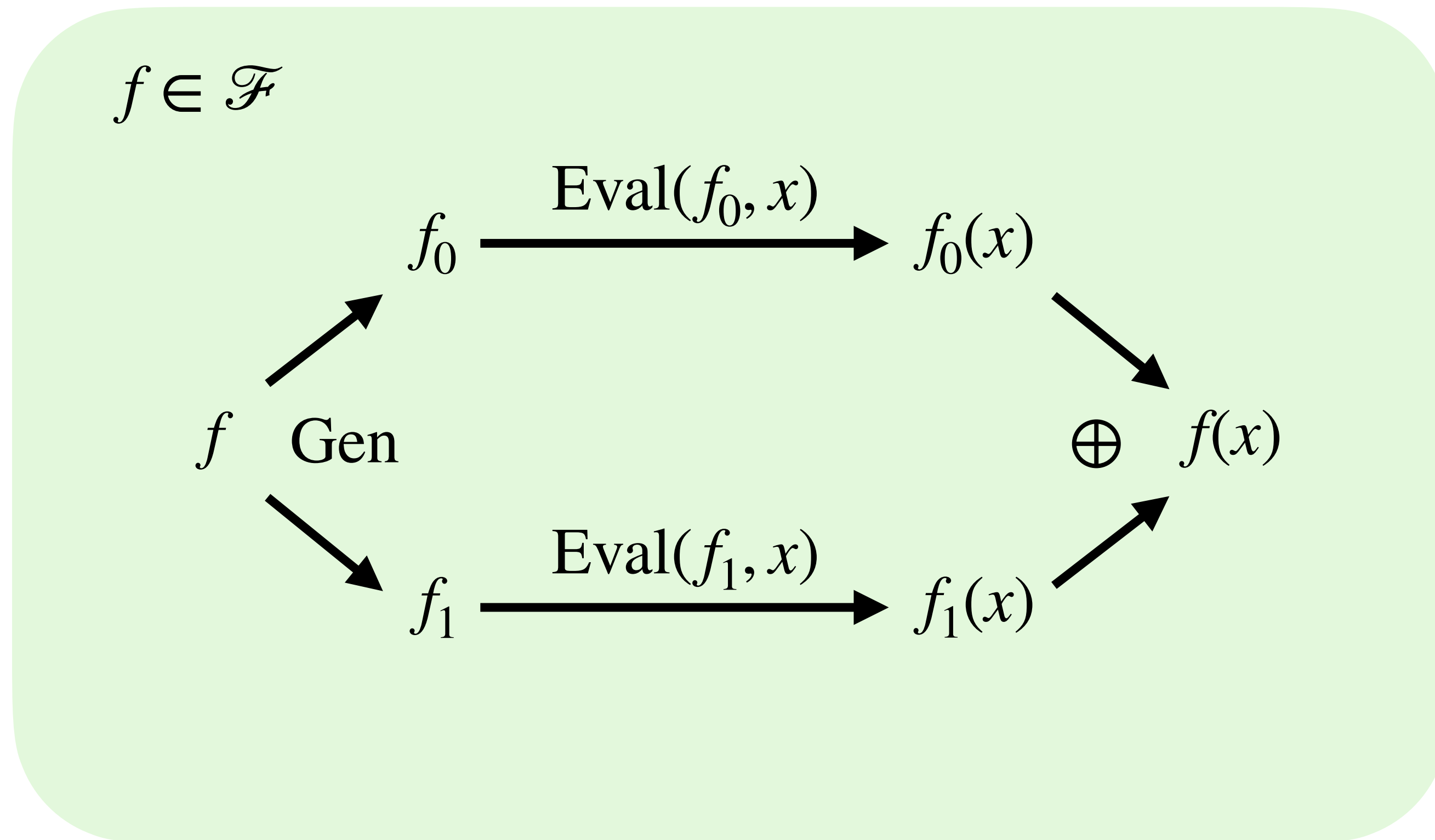


write



read

Distributed Point Function



$\text{point}_0(i, \cdot)$

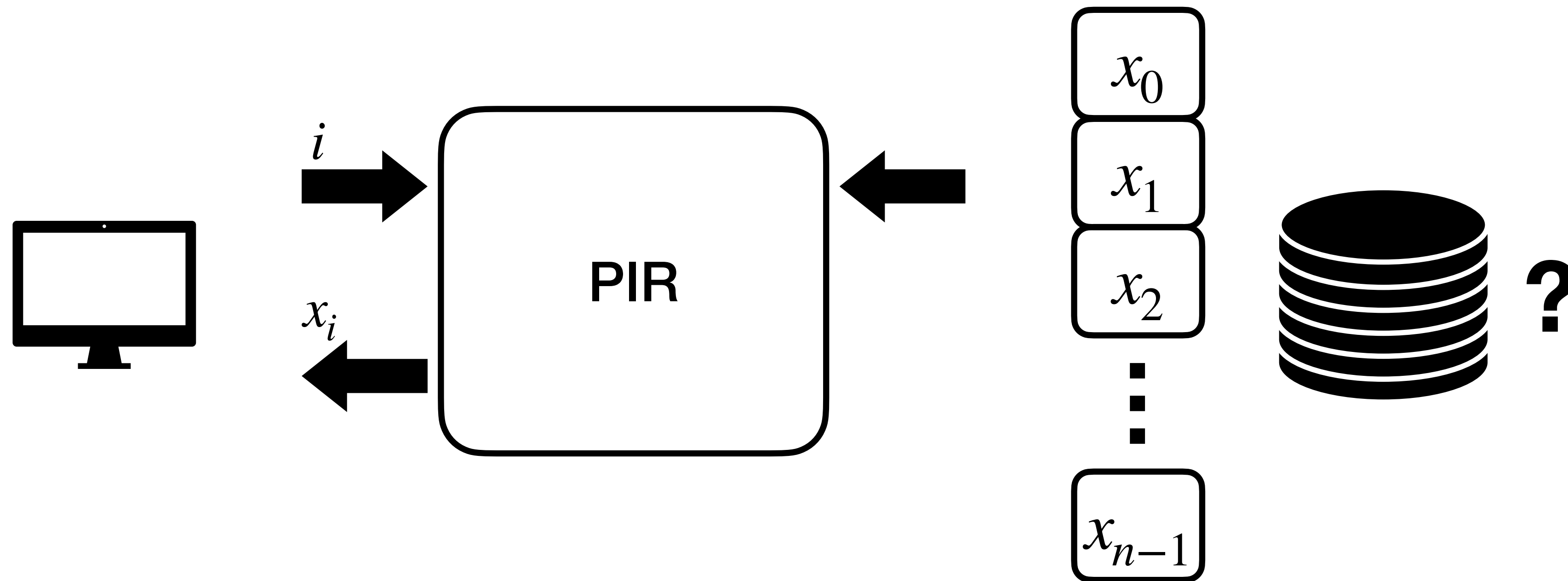


$\text{point}_1(i, \cdot)$

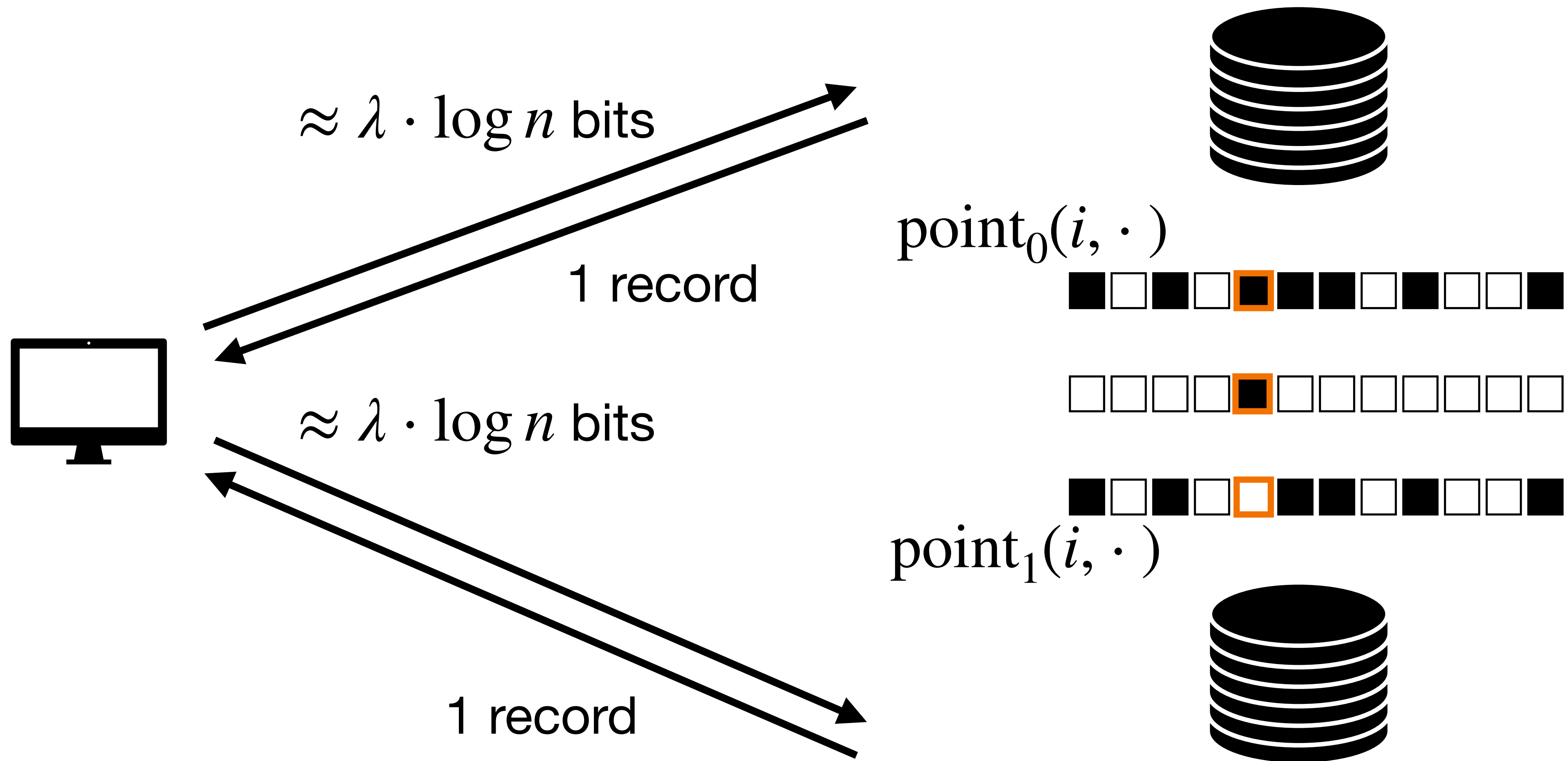


$$\text{point}(i, x) = \begin{cases} 1 & \text{if } x = i \\ 0 & \text{otherwise} \end{cases}$$

Private Information Retrieval

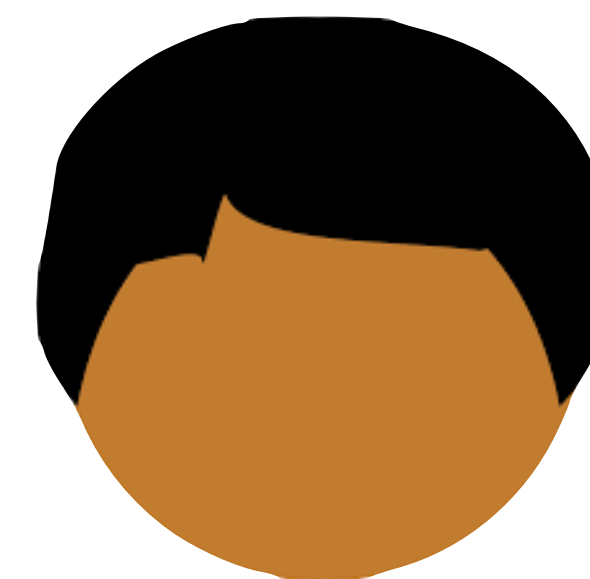


**Client wishes to privately
query one element from a
large database**





PSI



{13, 17, 25, 45, 52, 101}

{1, 4, 17, 19, 21, 45, 100}

Efficient Circuit-based PSI via Cuckoo Hashing

Benny Pinkas¹, Thomas Schneider², Christian Weinert², and Udi Wieder³

¹ Bar-Ilan University
benny@pinkas.net

² TU Darmstadt
{thomas.schneider,christian.weinert}@crisp-da.de

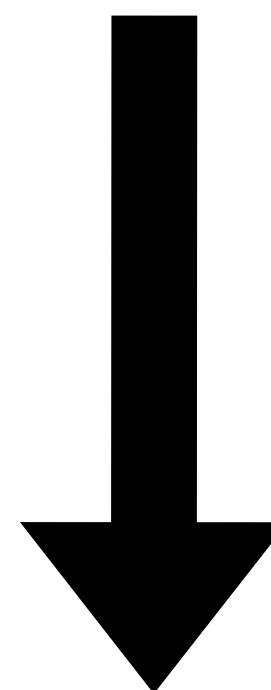
³ VMware Research
udi.wieder@gmail.com

Abstract. While there has been a lot of progress in designing efficient custom protocols for computing Private Set Intersection (PSI), there has been less research on using generic Multi-Party Computation (MPC) protocols for this task. However, there are many variants of the set intersection functionality that are not addressed by the existing custom PSI solutions and are easy to compute with generic MPC protocols (e.g., comparing the cardinality of the intersection with a threshold or measuring ad conversion rates).

Generic PSI protocols work over circuits that compute the intersection. For sets of size n , the best known circuit constructions conduct $O(n \log n)$ or $O(n \log n / \log \log n)$ comparisons (Huang et al., NDSS'12 and Pinkas et al., USENIX Security'15). In this work, we propose new circuit-based protocols for computing *variants of the intersection* with an almost linear number of comparisons. Our constructions are based on new variants of Cuckoo hashing in two dimensions.

We present an asymptotically efficient protocol as well as a protocol with better concrete efficiency. For the latter protocol, we determine the required sizes of tables and circuits experimentally, and show that the run-time is concretely better than that of existing constructions.

The protocol can be extended to a larger number of parties. The proof technique presented in the full version for analyzing Cuckoo hashing in



{17,45}

Special case of MPC

“Just use MPC”

Because it is a special case, we can hope for much more efficiency

Via OPRF

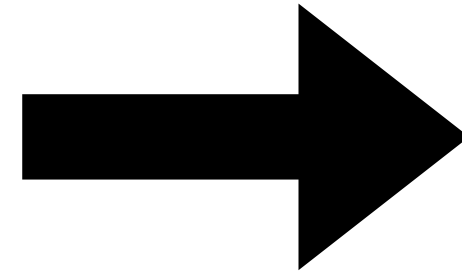
Semi-honest Protocols

GMW Protocol

Multi-party

Multi-round

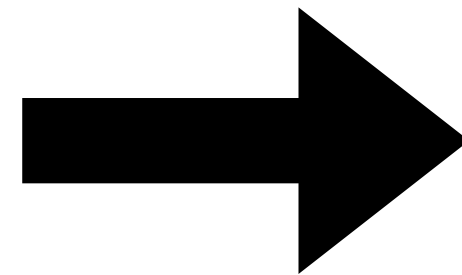
N-1 corruptions



Garbled Circuit

Constant Round

Two Party (multiparty via BMR)



Malicious Protocols

BDOZ Protocol

Authenticated Garbling

BGW Protocol

Multi-party

Guaranteed output delivery

Primitives

Oblivious Transfer/OT Extension

Secret Sharing

Authenticated Secret Sharing

Distributed Point Functions

Related Problems

PIR

PSI

OPRF