

Private Set Intersection (PSI)

Introduction, ORPF/Diffie-Hellman based PSI,
and the Apple CSAM detection system

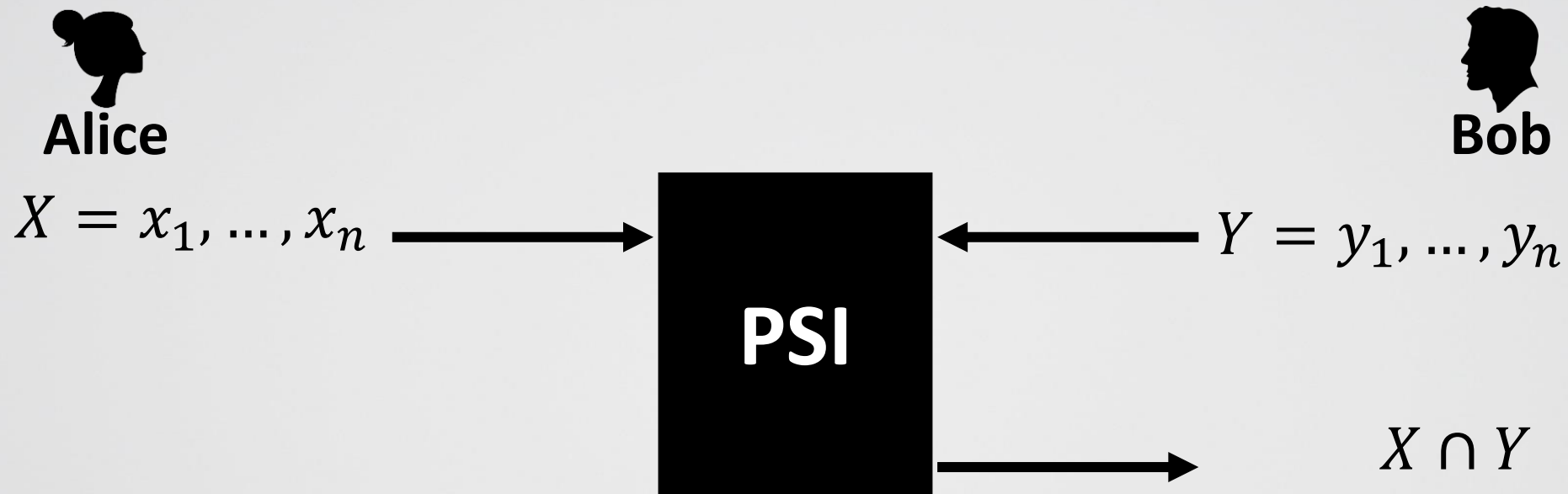
Benny Pinkas, Bar-Ilan University

In this talk

- PSI Introduction
- Diffie-Hellman based PSI, and PSI from OPRF
- The PSI system used by Apple for CSAM detection

Private Set Intersection (PSI)

Definition



Can also compute a function of the intersection, e.g. size (cardinality), or whether the intersection size is great than some threshold

Applications of PSI

- **Information sharing**, e.g., intersection of threat information
- **Matching**, e.g., testing compatibility of different properties (preferences, genomes...)
- **Join** DB operations
- **Analytics**: $\Pr(A / B) = \Pr(A \cap B) / \Pr(B)$
- **Identifying mutual contacts** (Signal app)
- **Computing ad conversion rates** (Google)

PSI Example – simple analytics for COVID vaccinations



University

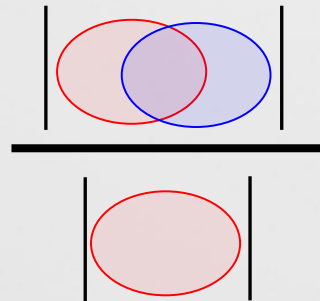
workers &
students ids

Ministry of
Health



vaccinated
people

Need to compute



Multi-party PSI



How can multiple fertility clinics compare their donor records?

Customer Risk Analysis

small merchant
(Alice)



Input: new customer

merchants, banks



customer databases

Output: do you know this guy?

Customer Risk Analysis

Here Alice's set only
contains one item!

small merchant
(Alice)



merchants, banks



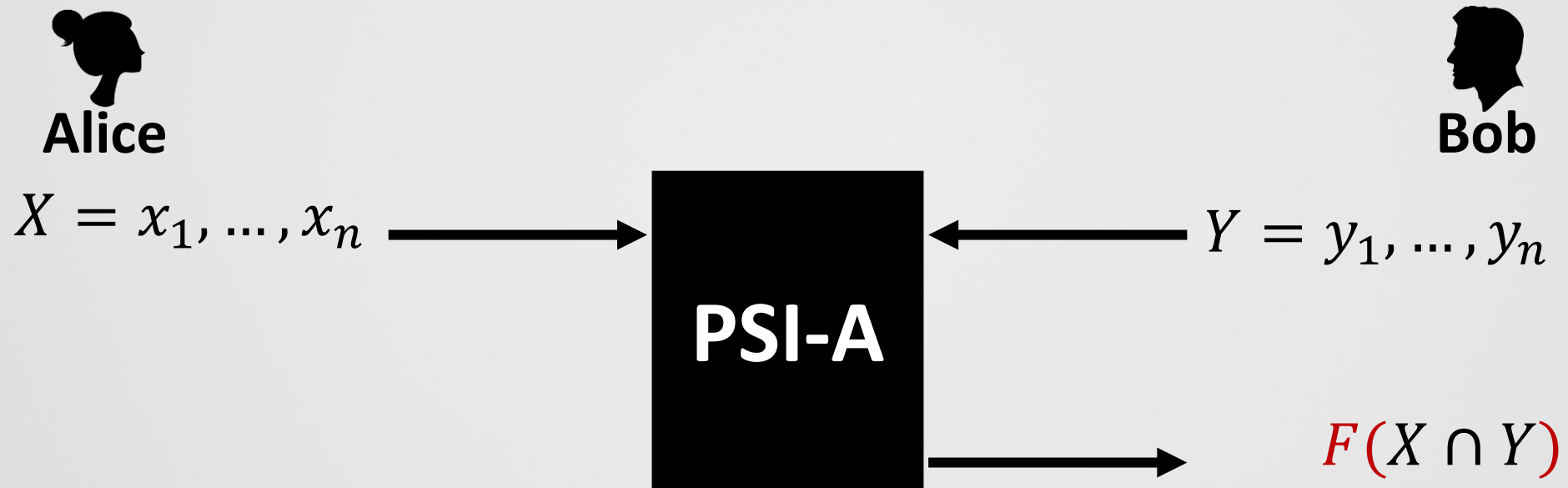
In more detail:

- Alice has a new customer x
 - Each other member has a set of known customers and their rating
 - Alice learns the rating of x by the other members
 - The other members do not learn x
-
- Or, Alice receives a purchase request to a certain address. She can check if purchases with of this credit card with the other members were sent to the same address.

Private Set Intersection (PSI) + Analytics

Definition

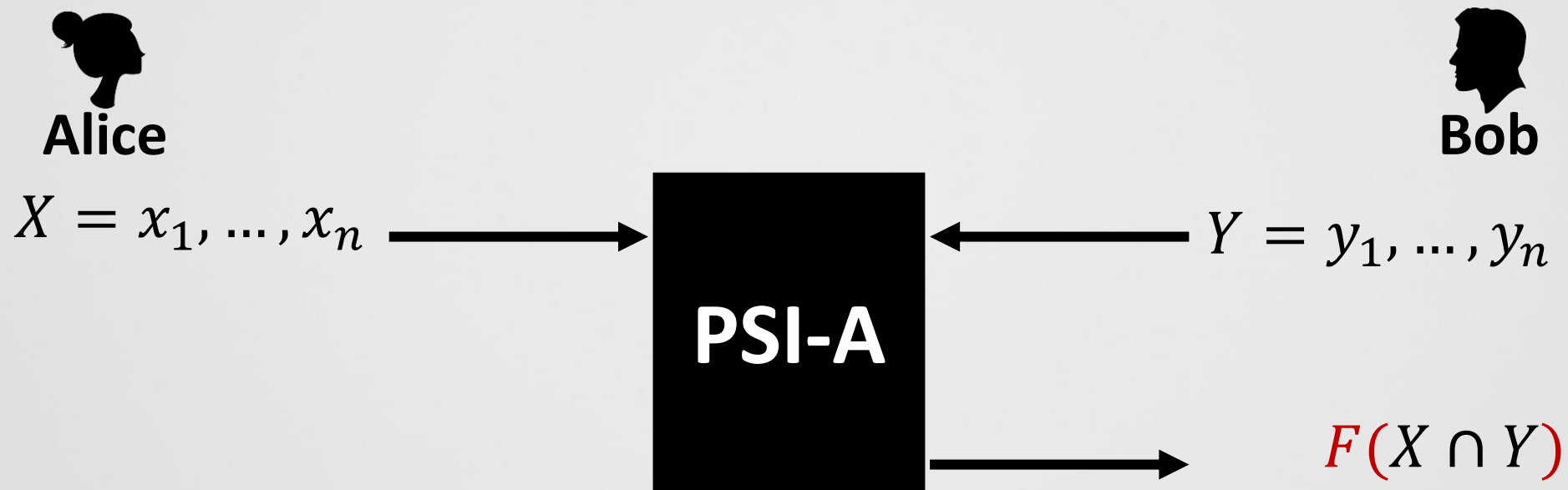
1. Post processing function F . E.g. $F = |X \cap Y|$



Private Set Intersection (PSI) + Payload Analytics

Definition

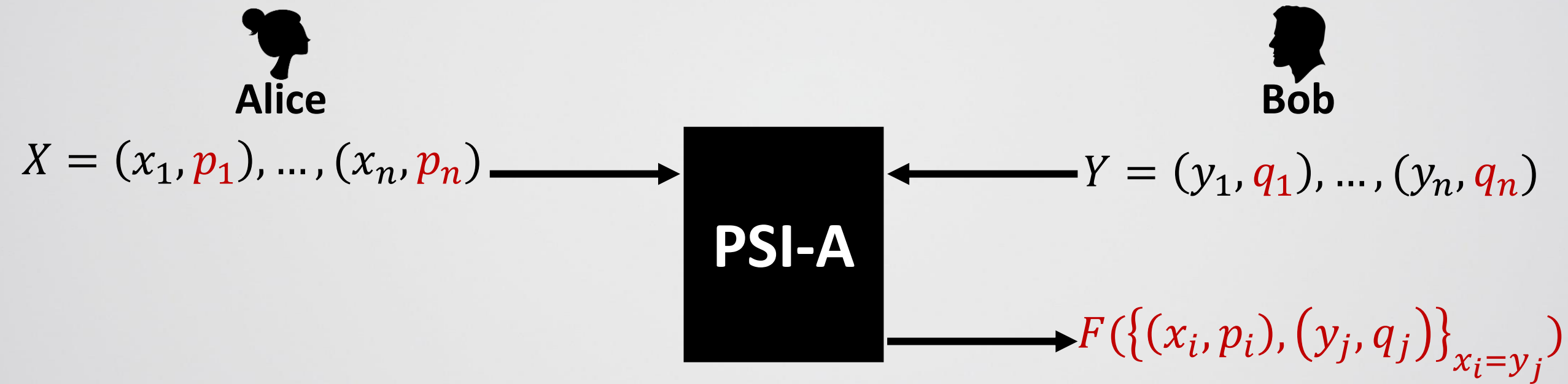
1. Post processing function F . E.g. $F = |X \cap Y|$
2. Items are associated with payloads (aka. PSI with data-transfer)



Private Set Intersection (PSI) + Payload Analytics

Definition

1. Post processing function F . E.g. $F = |X \cap Y|$
2. Items are associated with payloads (aka. PSI with data-transfer)



Private Set Intersection (PSI) - Analytics

Applications

Private Intersection-Sum Protocol with Applications to Attributing Aggregate Ad Conversions

Mihaela Ion[†], Ben Kreuter[†], Erhan Nergiz[†], Sarvar Patel[†],
Shobhit Saxena[†], Karn Seth[†], David Shanahan[†] and Moti Yung^{‡*}

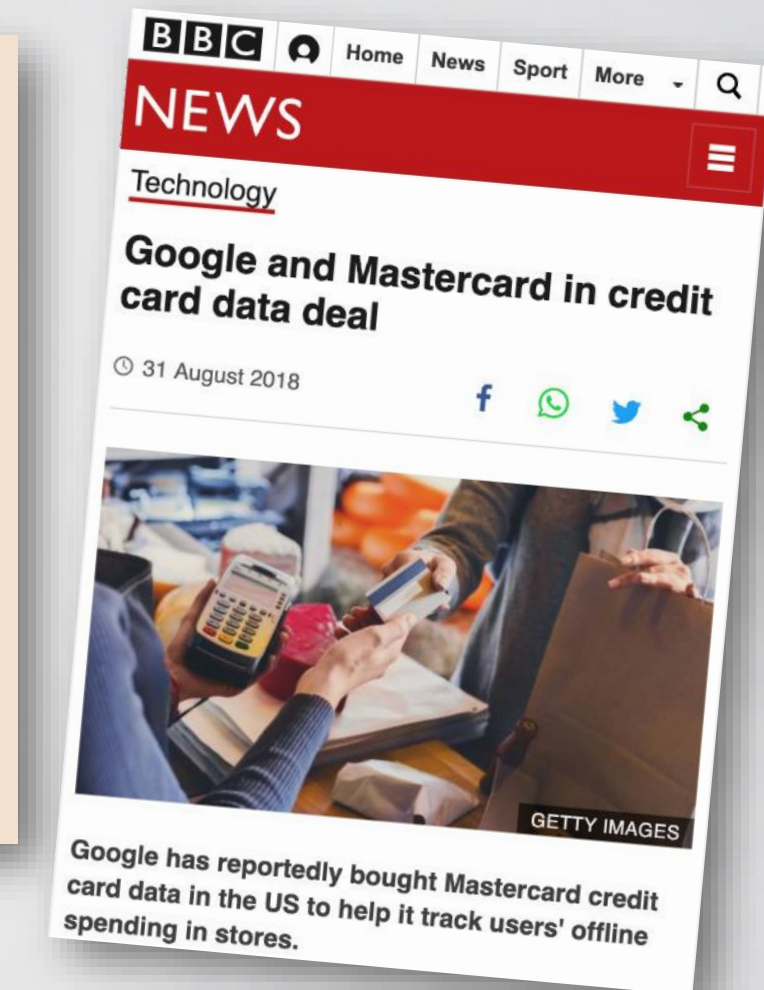
[†]{mion, benkreuter, anergiz, sarvar,
shobhitsaxena, karn, dshanahan}@google.com

Google Inc.

[‡]moti@cs.columbia.edu

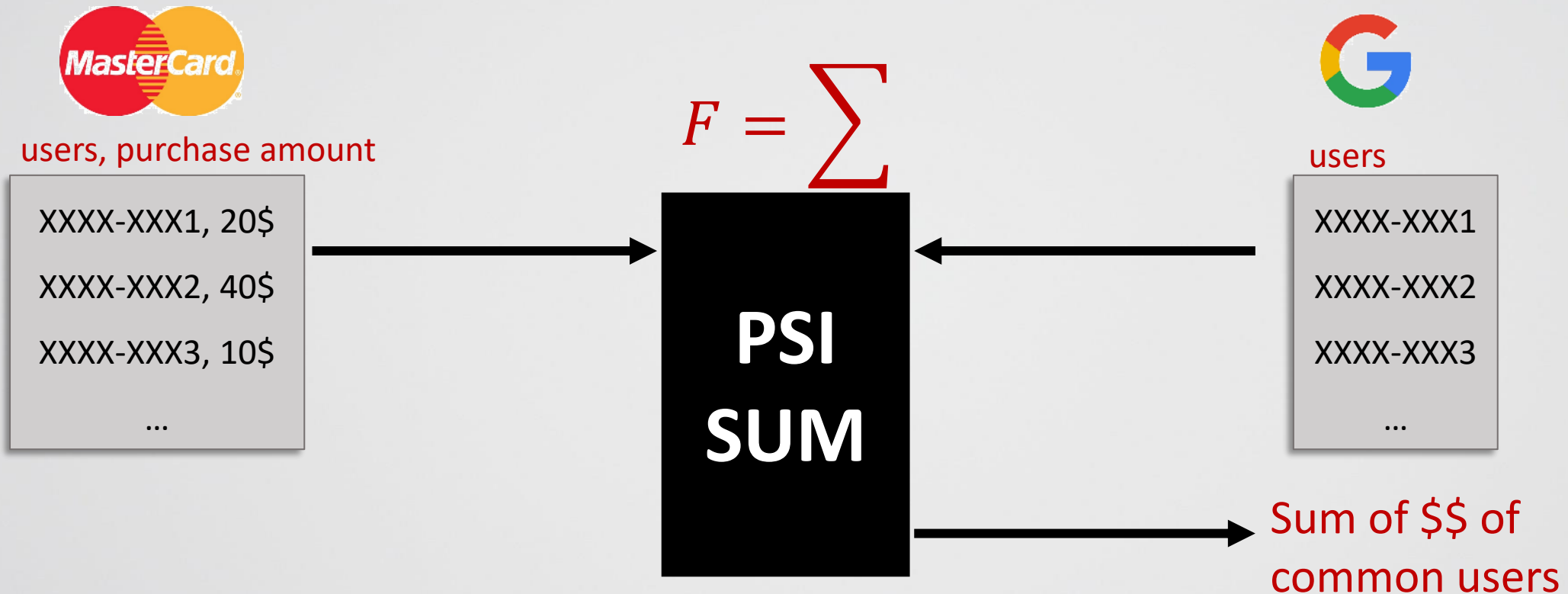
Columbia University and Snap Inc.

July 31, 2017



Private Set Intersection (PSI) - Analytics

Applications: Online Ads to Offline Purchase Conversion



Why did I investigate PSI?

- Relevant to many applications
- Cannot be efficiently solved using generic MPC (circuits)
- Solutions turned out to be based on interesting techniques

Pseudo-random function (PRF)

- Think of $AES_k(x)$
 - If we don't know the key k , the output looks random
- A PRF is a keyed function $F_k()$
- If k is chosen at random and kept secret, then a polynomial time program cannot distinguish between the output of $F_k()$ and a truly random function. (Even if it knows or chooses the inputs.)

Oblivious Pseudo-Random Function (OPRF)

A protocol for securely computing a PRF



Oblivious Pseudo-Random Function (OPRF)

Another OPRF version



Useful when there is no need for the server to choose the key K (as in PSI)

A pseudo-random function based on DDH

- Let G be a group in which the DDH assumption holds
- Let $H()$ be a publicly known random function mapping into G
 - $H()$ is modeled as a random oracle
- K is a key

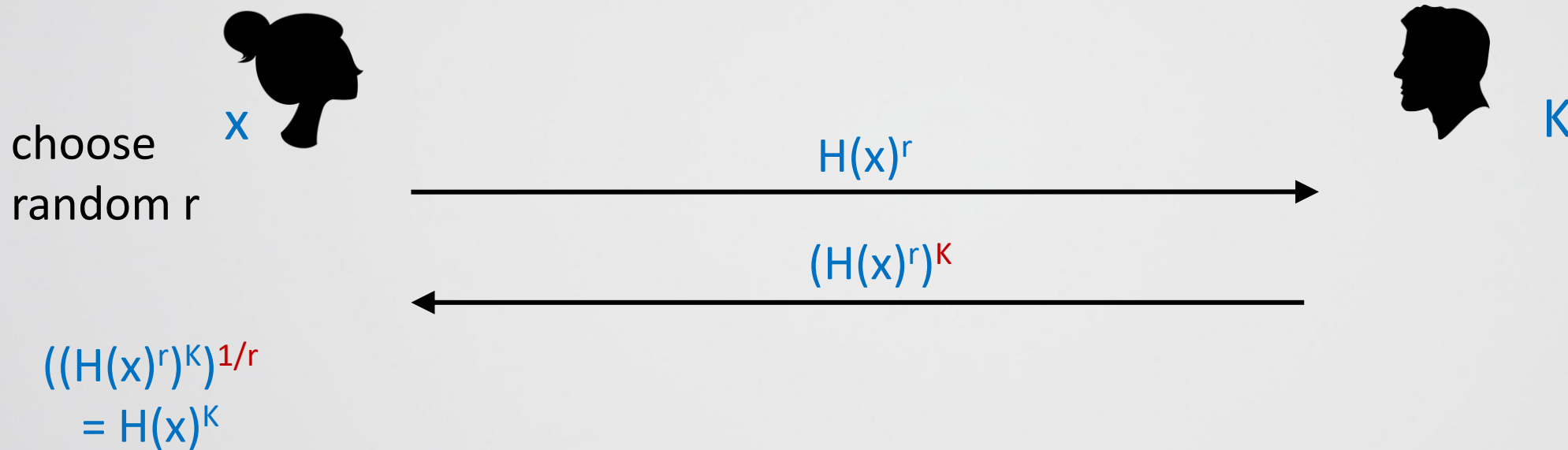
Typically
implemented
in ECC

DDH assumption: for generator g and random a, b, c cannot distinguish (g^a, g^b, g^{ab}) from (g^a, g^b, g^c)

The function $F_K(x) = (H(x))^K$ is pseudo-random

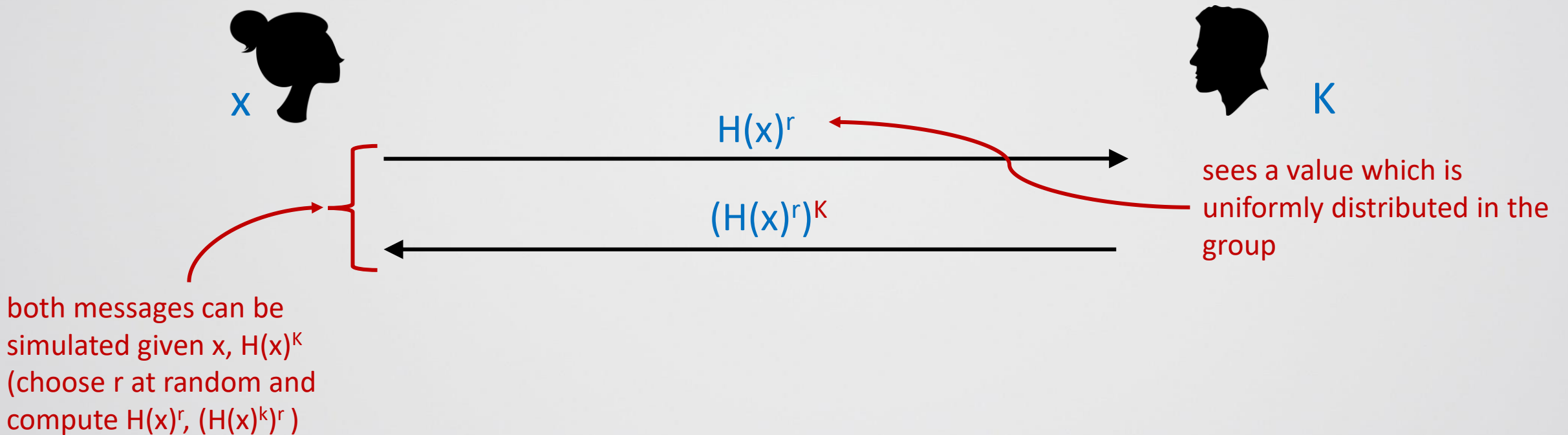
An OPRF construction based on DDH

A protocol for securely computing the PRF $F_K(x) = (H(x))^K$



An OPRF construction based on DDH

Security intuition:



An OPRF construction based on DDH

- Becoming an IETF standard
 - Oblivious Pseudorandom Functions (OPRFs) using Prime-Order Groups `draft-irtf-cfrg-voprf-01`
 - Implemented in elliptic curve groups
- Also a “verifiable OPRF” version
 - The server can prove that in all invocations it used the same key
 - The server publishes g^k
 - With each answer it proves that the discrete log of $(H(x)^r)^k$ to the base $(H(x)^r)$, is the same as the discrete log of g^k to the base g .

PSI based on DDH (described using OPRF)

Typically
implemented
in ECC

α
 x_1, \dots, x_n

β
 y_1, \dots, y_n

α is a randomization used by Alice
 β is a OPRF key known to Bob
PRF is $H(x)^\beta$

$(H(x_1))^\alpha, \dots, (H(x_n))^\alpha$

$(H(y_1))^\beta, \dots, (H(y_n))^\beta$

PRF of Bob's inputs

$((H(x_1))^\alpha)^\beta, \dots, ((H(x_n))^\alpha)^\beta$

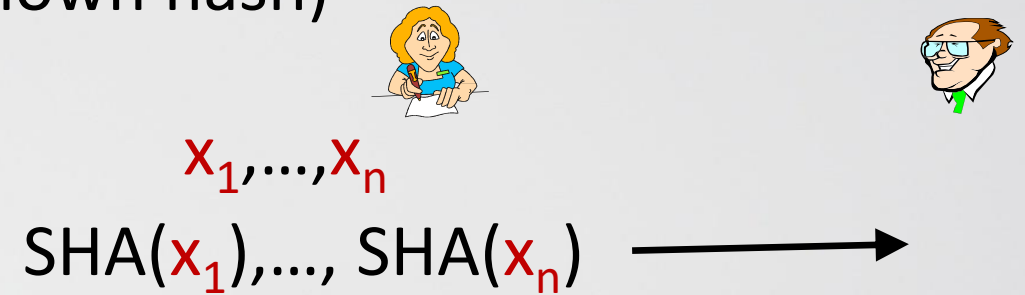
OPRF protocol

$((H(x_1))^\alpha)^\beta)^{1/\alpha}, \dots, ((H(x_n))^\alpha)^\beta)^{1/\alpha}$

Compares the two lists

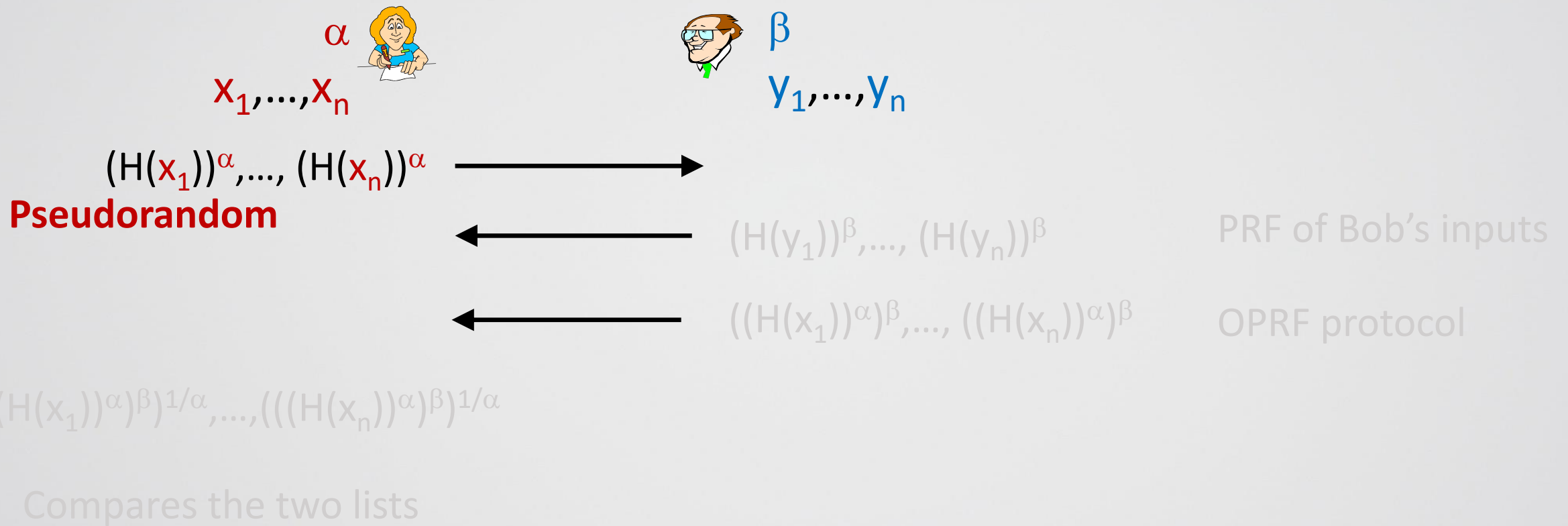
Security intuition

- Recall a trivial **insecure** PSI protocol (known hash)

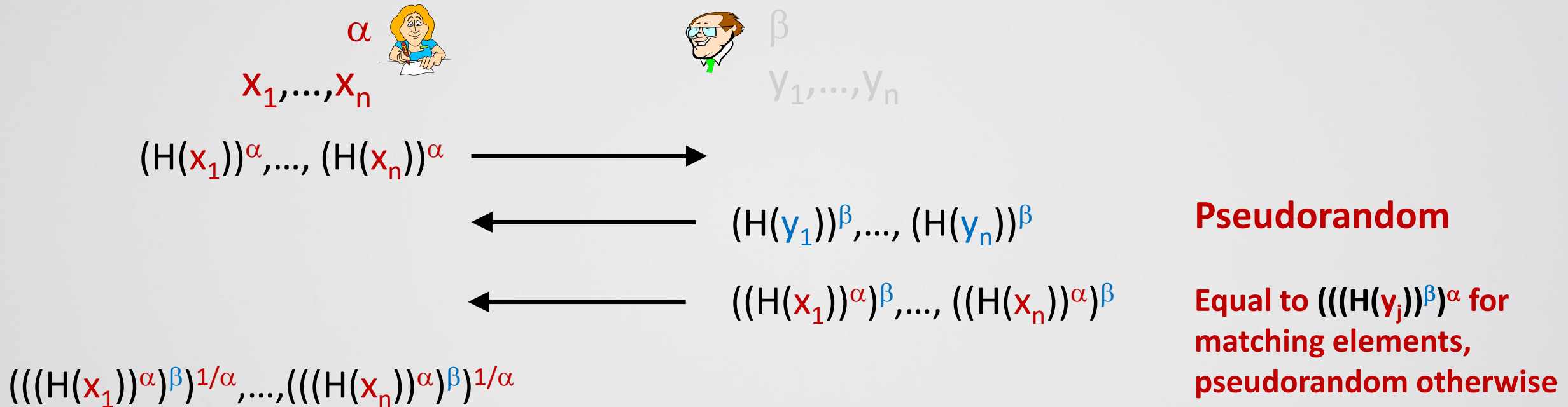


- This protocol is insecure against a dictionary attack
- Whereas in the OPRF-PSI protocol we described,
 - Bob sees the output of a hash that can only be computed by Alice $(H(x))^\alpha$
 - Alice sees the output of a hash that can only be computed by Bob $(H(y))^\beta$

Bob simulating his view

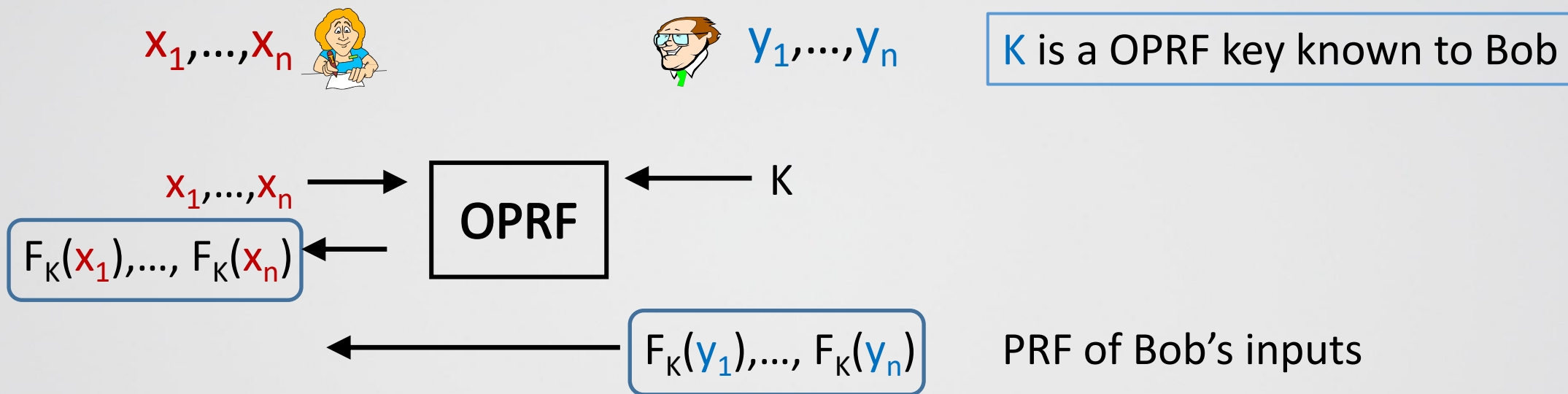


Alice simulating her view



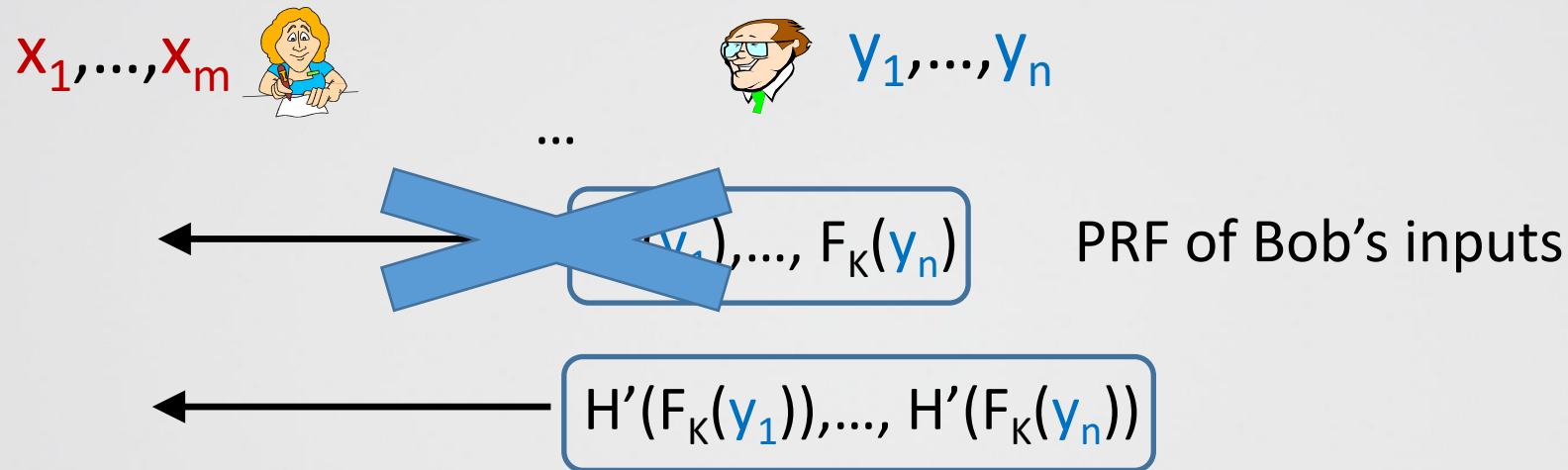
Compares the two lists

Template for PSI based on OPRF (basic version)



Compares the two lists

Sending the PRF values



- Instead, apply a hash function with a **shorter output** and send the results
- The output length of $H'()$ depends on desired false positive probability, set sizes, and the data structure.
- Can also encode the results in an efficient data structure supporting searches (dictionary), such as a Bloom filter or a cuckoo filter

PSI of sets of unequal sizes

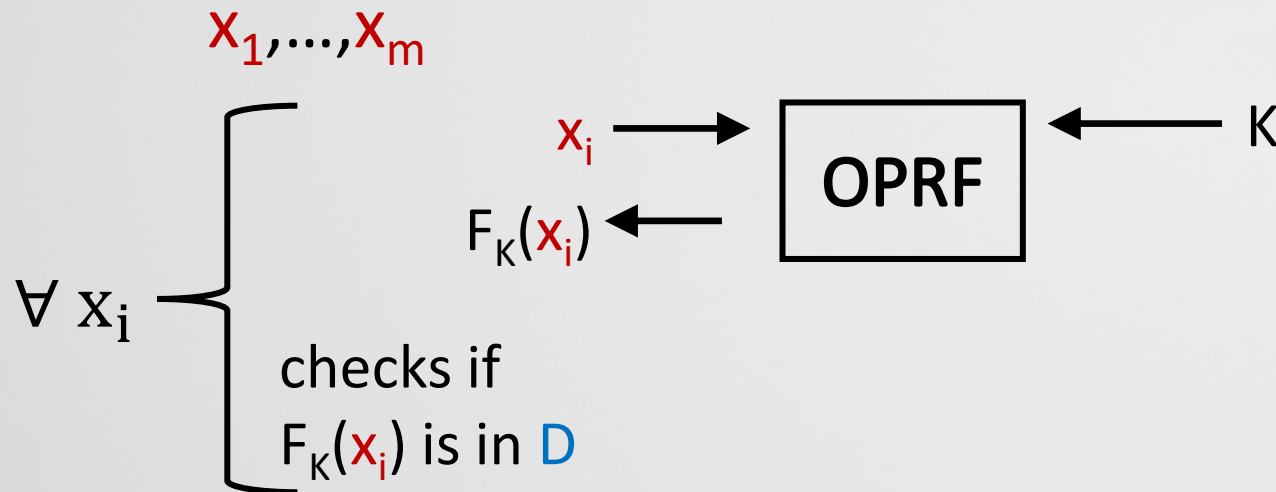


y_1, \dots, y_n

K is a OPRF key known to Bob

$\longleftarrow D(F_K(y_1), \dots, F_K(y_n))$

Dictionary encoding of the
PRF outputs of Bob's inputs



PSI of sets of unequal sizes



y_1, \dots, y_n

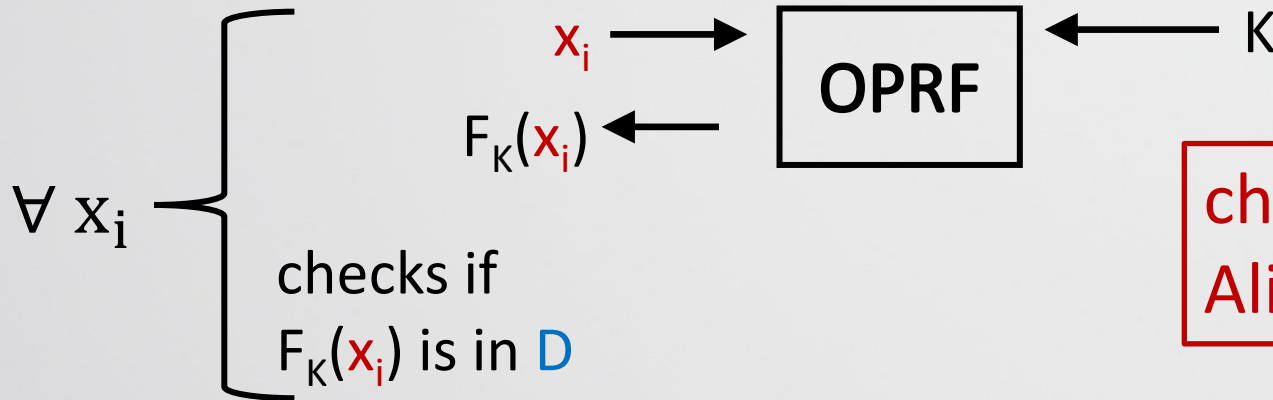
K is a OPRF key known to Bob

(before Alice learns her inputs)

$\leftarrow D(F_K(y_1), \dots, F_K(y_n))$

Dictionary encoding of the PRF outputs of Bob's inputs

x_1, \dots, x_m



checks for a match whenever Alice has a new x_i

Implementations of OPRFs

- Based on DH and random oracle $H()$ (described earlier)
- Based on DH alone (Naor-Reingold)
- Based on RSA
- Based on 2-party MPC of AES or other ciphers (communication depends on circuit size)
- Post quantum secure:
 - Based on the ring learning-with-errors problem and the short-integer-solution problem in one dimension
 - Based on isogenies of supersingular elliptic curves

Some applications need

- Verifiable OPRF
- Malicious security

Apple's CSAM Detection System (my interpretation of it)

"May your research area be relevant to current affairs" is the
"May you live in interesting times" curse-disguised-as-a-
blessing for the academic set" - Riana Pfefferkorn

Background

- CSAM = Child Sexual Abuse Material
- All companies, including Apple, allow law enforcement to access cloud data, if they are provided the right warrant
 - But this is different than scanning **all** data
- Apple wanted / was required to scan uploaded photos (not on phones or in messaging)

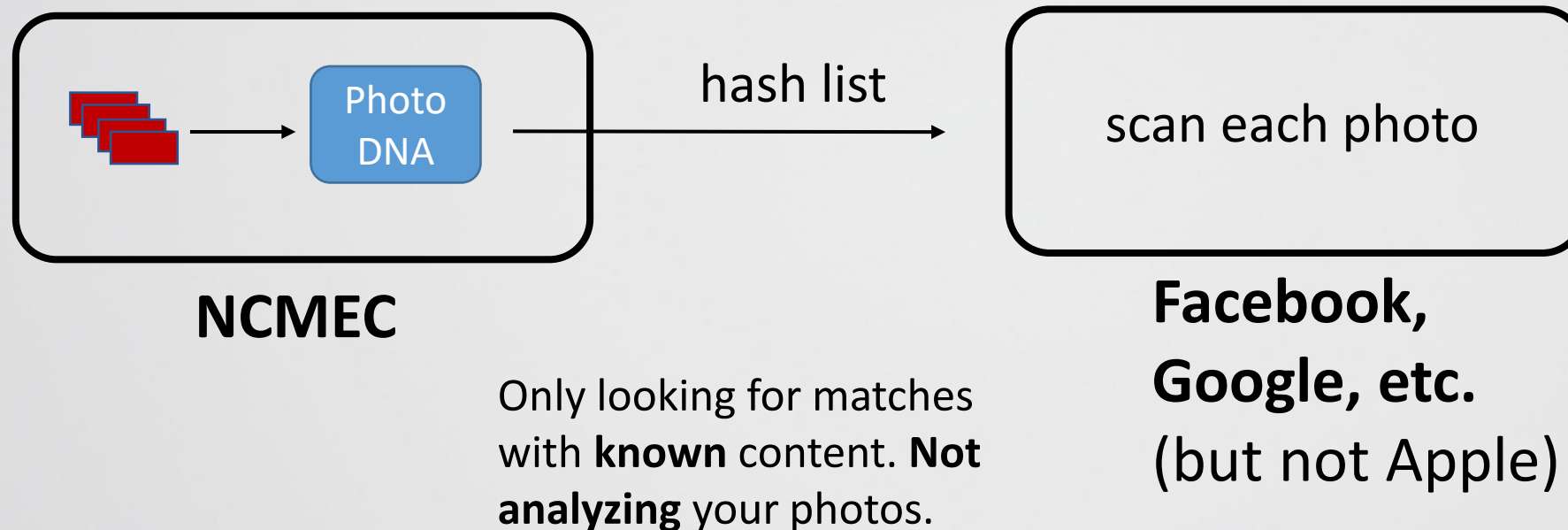
The Internet Is Overrun With Images of Child Sexual Abuse. What Went Wrong?

Online predators create and share the illegal material, which is increasingly cloaked by technology. Tech companies, the government and the authorities are no match.

Last year, tech companies reported over 45 million online photos and videos of children being sexually abused — more than double what they found the previous year.

Current industry practices for CSAM detection

Scan data (photos) in the clear



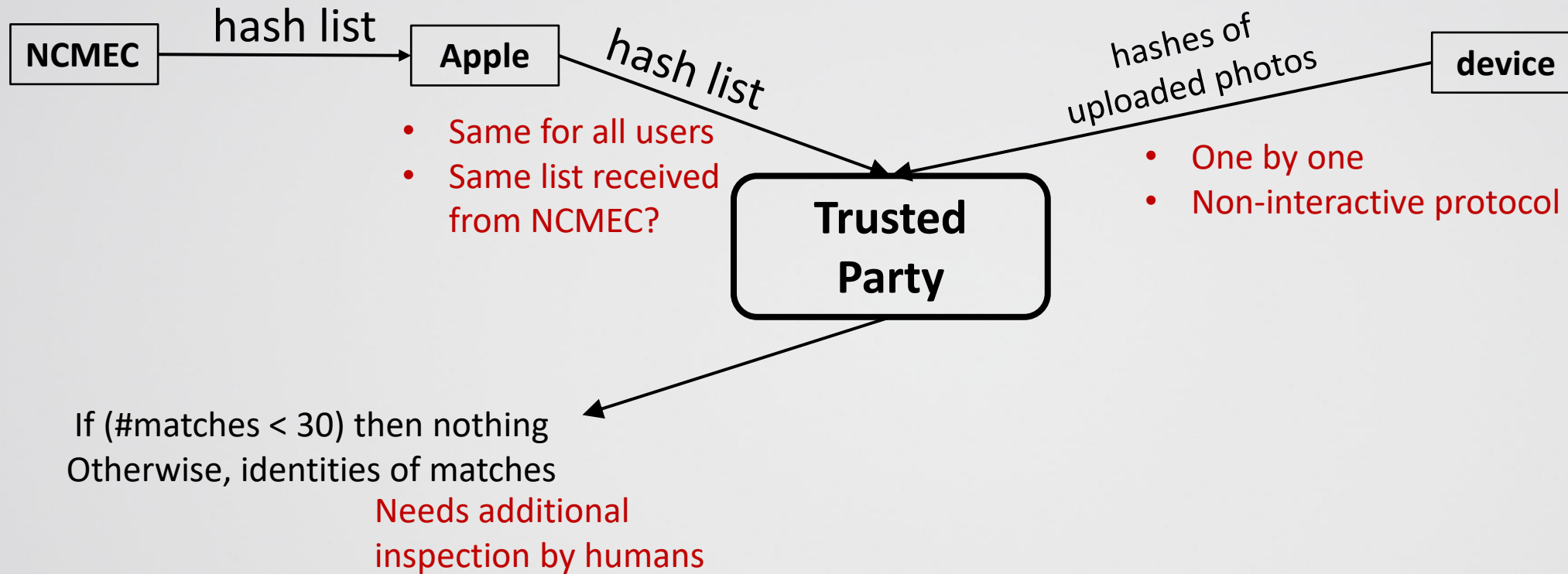
Apple's design objectives

- Apple learns nothing about images that do not match **known** CSAM
 - Apple can't access any data of matched CSAM images until a **threshold** of 30 matches is exceeded for an iCloud account
 - The risk of the system incorrectly flagging an account is extremely low
 - In addition, Apple **manually** reviews all reports made to NCMEC to ensure reporting accuracy
-
- Users can't access or view the database of known CSAM images/hashes
 - Users can't identify which images were flagged as CSAM by the system

Protecting users

Protecting society

Ideal world view

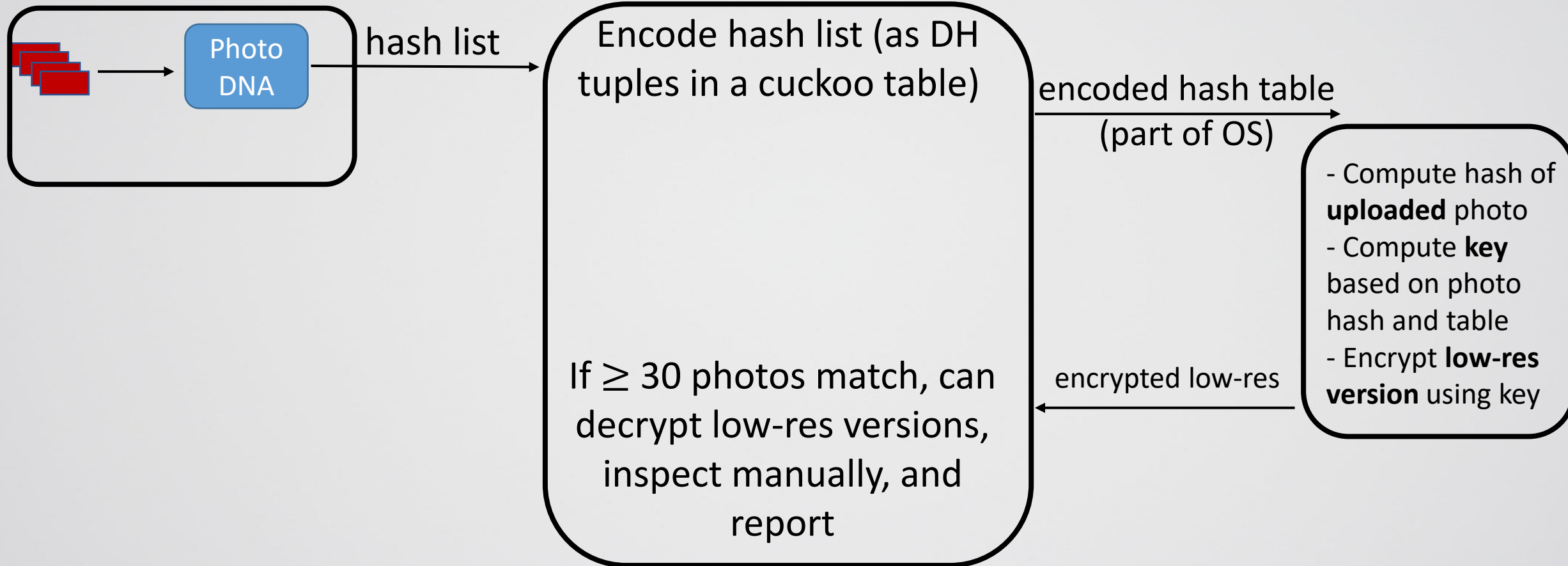


Apple's CSAM scanning system

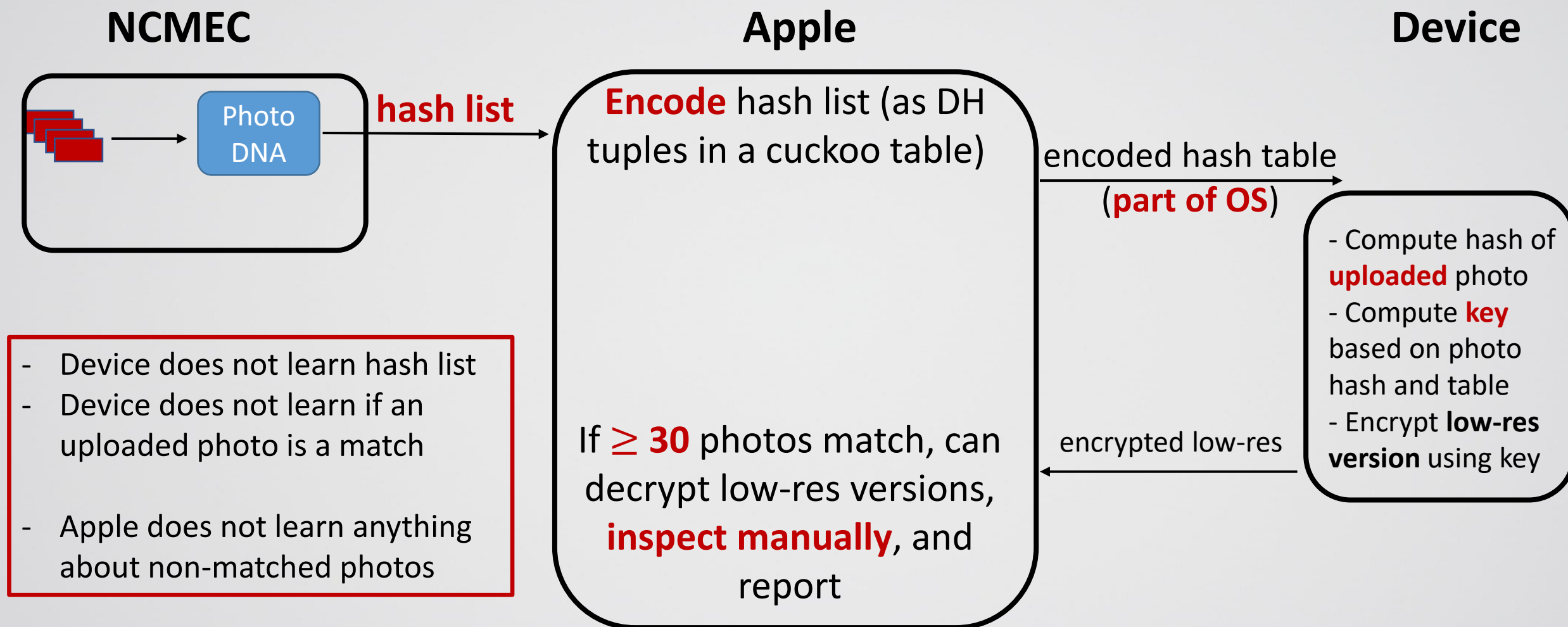
NCMEC

Apple

Device

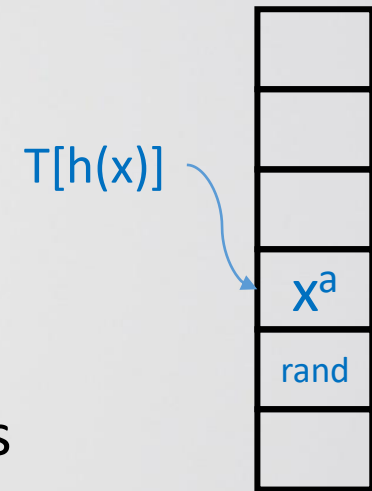


Apple's CSAM scanning system



Server preprocessing

- Choose a random secret key (exponent) a
- Construct a (cuckoo hash) table storing the CSAM hash list
- For every **hash** x in the CSAM list, store x^a in location $T[h(x)]$
 - (This is actually a cuckoo table with two hash functions)
 - Elements that do not find a place are dropped
- Put random values in empty entries
 - Since $F(x)=x^a$ is a **prf** (x is already the output of a hash function), this is **indistinguishable** without knowledge of the key a
- Send **table** and g^a to client (as part of the operating system)



DH self reduction [Naor-Reingold]

- $(g, L=g^a, T=g^b, P=g^{ab})$ is a DH tuple

- Pick random c, d

- $Q = g^d T^c, \quad S = L^d P^c$

- If (g, L, T, P) is a DH tuple then, since $L=g^a$ and $P=T^a$, we get that $S=Q^a$

- If (g, L, T, P) is **not** a DH tuple then (Q, S) is a uniformly random pair

Client side

- Client has photo it wants to upload (\mathbf{x} = hash, \mathbf{id} , \mathbf{data} = associated data)
 - Picks random \mathbf{c}, \mathbf{d}
 - Computes $\mathbf{Q} = \mathbf{x}^{\mathbf{c}} \mathbf{g}^{\mathbf{d}}$ and $\mathbf{S} = \mathbf{T}[\mathbf{h}(\mathbf{x})]^{\mathbf{c}} (\mathbf{g}^{\mathbf{a}})^{\mathbf{d}}$
 - If $\mathbf{T}[\mathbf{h}(\mathbf{x})] = \mathbf{x}^{\mathbf{a}}$ then $\mathbf{S} = \mathbf{Q}^{\mathbf{a}}$. Otherwise, based on DH self reducibility, \mathbf{Q} and \mathbf{S} are a uniformly random pair of values.
 - Idea: Send \mathbf{Q} and use \mathbf{S} as a **key**. If this is a DH tuple then the server, which knows \mathbf{a} , can recover key \mathbf{S} as $\mathbf{Q}^{\mathbf{a}}$. Otherwise, \mathbf{S} is uniformly random.
- Security:
 - Client sees the table, which looks pseudo-random without knowledge of key \mathbf{a}
 - If \mathbf{x} is **not** in table, client encrypts with a key \mathbf{S} which is random.

Threshold PSI

- Client Picks a single random master key **adkey** for all items.
- Client has input item (**x** = hash, **id**, **data** = associated data)
 - Computes **Enc(adkey, data)**
 - Generates a Shamir **share** of **adkey** with threshold **t = 30**
 - Computes **Q, S** as described before
 - Sends **Q** and authenticated encryption with key **S** of (**share**, **Enc(adkey, data)**)
- If **x** is in hash table, server knows **S** and can decrypt (**share**, **Enc(adkey, data)**)
- After 30 such shares, server can recover **adkey** and decrypt data
 - But we do not want to the server to identify which users have 29 shares, etc.

Threshold PSI

- For users who have < 30 matches, need to hide $\#matches$ from server
- Solution: users send at random dummy matches – messages for which server can decrypt (**dummy-share**, **dummy-Enc**)
- The server will have some real shares and some dummy shares. It must be able to decrypt whenever $\#real\ shares \geq 30$.

Threshold PSI

- The server will now have some real shares and some dummy shares. It must be able to decrypt whenever $\# \text{real shares} \geq 30$.
- A solution based on a decoding algorithm of Coppersmith-Sudan
 - Given t correct “Shamir+” shares and s random shares, can decode secret
 - Caveat: The size of each share is larger by a factor of $s+1$
 - E.g., for $t=30$ and $s=100$, each share contains 101 field elements...
 - Doesn't matter for this application, since the share size is “negligible” compared to the size of the uploaded photo

Associated data

- The associated data includes a low resolution “visual derivative” of the image.
- If 30 matches are found, human operators examine these visual derivatives, and if they look like CSAM then the user is reported.
- (associated data size \gg data sent in PSI protocol)

Potential problems beyond the technical solution

- False positives
- Apple is scanning on my device
- Who controls which items are searched for? Governments coerce Apple into searching for specific photos
- “A first crack in the dam”

Relevant papers

- <https://www.apple.com/child-safety/>
- In particular
 - Technical description [https://www.apple.com/child-safety/pdf/Apple PSI System Security Protocol and Analysis.pdf](https://www.apple.com/child-safety/pdf/Apple_PSI_System_Security_Protocol_and_Analysis.pdf)
 - Another technical description <https://decentralizedthoughts.github.io/2021-08-29-the-private-set-intersection-psi-protocol-of-the-apple-csam-detection-system/>
 - Threat model [https://www.apple.com/child-safety/pdf/Security Threat Model Review of Apple Child Safety Features.pdf](https://www.apple.com/child-safety/pdf/Security_Threat_Model_Review_of_Apple_Child_Safety_Features.pdf)