

Analyzing Cryptography in Context:

The Case Study of Apple's CSAM Scanning Proposal

Ran Canetti (Boston University) Julie Ha (Boston University), and
Gabriel Kaptchuk (University of Maryland)

Cryptographic Applications Workshop 2024

“How *should* we analyze cryptographic *deployments*?”

The Apple PSI System

Abhishek Bhowmick
Apple Inc.

Dan Boneh
Stanford University

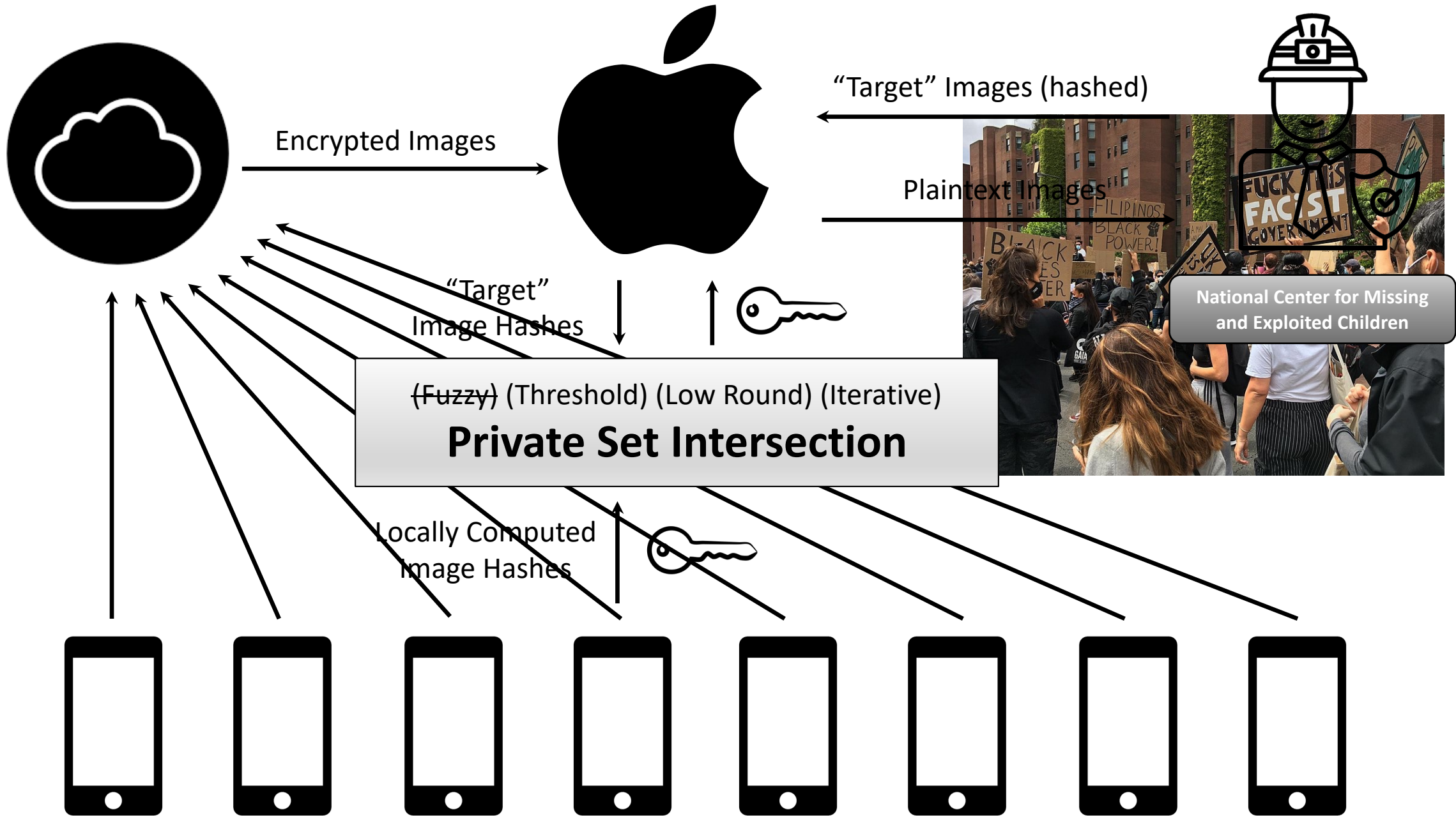
Steve Myers
Apple Inc.

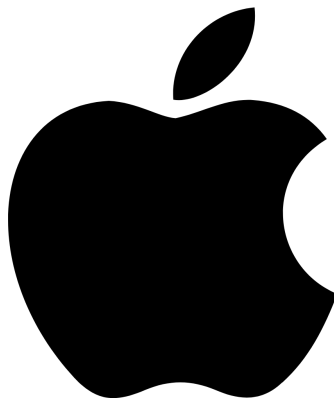
Kunal Talwar Karl Tarbe
Apple Inc. Apple Inc.

July 29, 2021

Abstract

This document describes the constraints that drove the design of the Apple *private set intersection* (PSI) protocol. Apple PSI makes use of a variant of PSI we call *private set intersection with associated data* (PSI-AD), and an extension called *threshold private set intersection with associated data* (tPSI-AD). We describe a protocol that satisfies the constraints, and analyze its security. The context and motivation for the Apple PSI system are described on the main project site.





$(img_1, img_2, \dots, img_n)$



$$pdata = (pk = g^\alpha, p_1 = H(img_1)^\alpha, p_2 = H(img_2)^\alpha, \dots, p_{n'} = H(img_{n'})^\alpha)$$

DDH tuple iff $img = img_i$

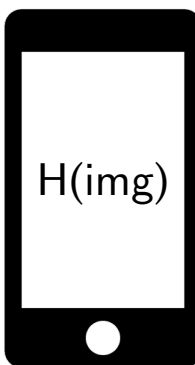
$$(g, pk, H(img), p_i)$$

$$(g, g^\alpha, H(img), H(img_i)^\alpha)$$

Sever Key Share

Shared Key

Client Key Share



(img_id, q, rct)

$$q \leftarrow H(img)^\beta \cdot g^\gamma$$

$$s \leftarrow p_i^\beta \cdot pk^\gamma \quad (g, pk, q, s)$$

$$rct \leftarrow Enc(KDF(s), (sh, adct))$$

$$adct \leftarrow Enc(adkey, key)$$

$$sh \leftarrow SecretShare(adkey, t)$$



Ma
@r

I've had
people f
CSAM s

1:59 AM · A

1,741 Retw



Ma
Re
Th
ma
ma



Ma
Ini
clo
sul



Ma
Th
be
let



E

Eric Rescorla



Matthew Green
@matthew_d_gree



Matthew Green
@matthew_d_green

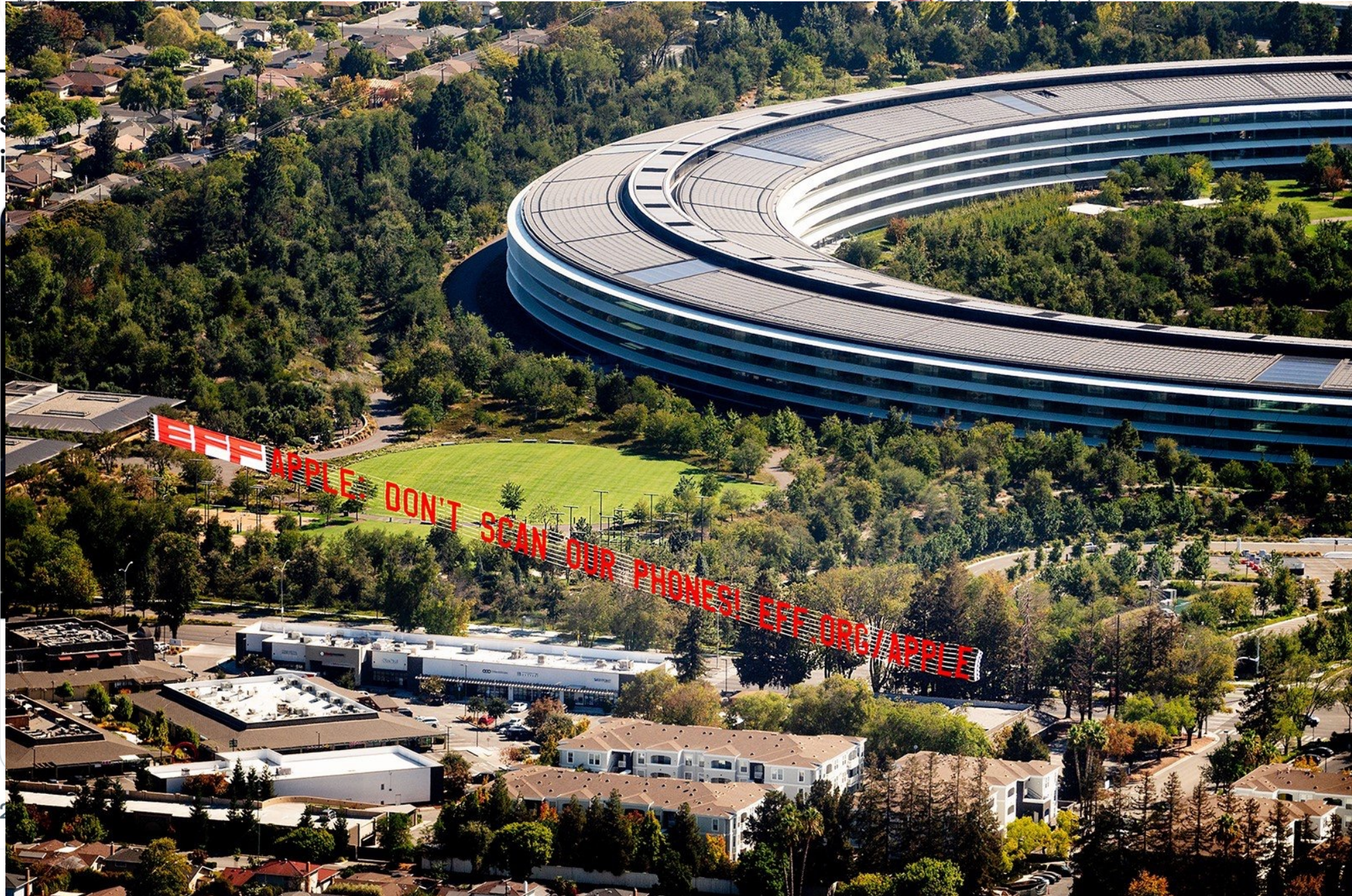


Sarah Jami
@SarahJarr



Sarah Jamie Lewis

...



technically,



If threshold of matches is exceeded, only then can Apple interpret voucher content for CSAM matches

Reviewed by Apple, report sent to NCMCEC

round new Photo scanning featu...
of new features coming later this
an. The new features include ...

Show replies

2) E
sec

that feature. It is reasonable to say the risk of abuse here is too f
use cloud-side scanning.



1

The system can be easily re-targeted

2

Malicious operators or hackers could abuse system

3

The system lacks robustness against malicious clients

Bugs in our Pockets: The Risks of Client-Side Scanning

Hal Abelson Ross Anderson Steven M. Bellovin
 Josh Benaloh Matt Blaze Jon Callas Whitfield Diffie
 Susan Landau Peter G. Neumann Ronald L. Rivest
 Jeffrey I. Schiller Bruce Schneier Vanessa Teague
 Carmela Troncoso
 October 15, 2021



Apple's Plan to "Think Different" About Encryption Opens a Backdoor to Your Private Life

BY INDIA MCKINNEY AND ERICA PORTNOY | AUGUST 5, 2021



YouTube

Search

Matthew Gre...

An evaluation of the risks of client-side scanning

Vanessa Teague, Bruce Schneier, Carmela Troncoso, Alex Stamos, Matthew Green

Presented by: Matthew Green, 4/14/2022 at RWC 2022

zoom

25:21 / 1:25:06 - Matt Green

“How *should* we analyze cryptographic *deployments*?”

“Well you should probably prove that your construction is secure...”

Parameters known to all parties:

- two parties: server and client,
- B is the maximum set size for the server and client,
- t is the threshold,
- s_{\max} is the maximum size of the set S of synthetics,
- all associated data values ad in \mathcal{D} have the same public length.

The functionality \mathcal{F} :

- Wait for input $X = \{x_1, x_2, \dots\}$ from the server;
abort if the server is corrupt and $|X| > B$.
- Send $|X|$ to the client; abort if the client is corrupt and aborts.
- Wait for input $\bar{Y} \in (\mathcal{U} \times \mathcal{ID} \times \mathcal{D})^m$ and $S \subseteq id(\bar{Y})$ from the client;
abort if the client is corrupt and $(m > B$ or $|S| > s_{\max})$.
- Send \bar{Y}_{id} to the server.
- If $|id(\bar{Y} \cap X) \setminus S| > t$:
send $\bar{Y}[id(\bar{Y} \cap X) \setminus S]_{\{id, ad\}}$ and S to the server,
otherwise:
send $id(\bar{Y} \cap X) \cup S$ to the server.

Figure 1: The ideal functionality \mathcal{F} for ftPSI-AD

Game G_{Π}^{SS}

INIT:

- 1 $c \leftarrow_{\$} \{0, 1\}$ // Random challenge bit
- 2 $H \leftarrow_{\$} \Pi.HS$ // Pick a function to be the random oracle

POST(X_0, X_1): // Adversary calls with sets X_0, X_1 .

- 3 **Require:** $|X_0| = |X_1|$ and $X_0, X_1 \subseteq \Pi.U$
- 4 $(pdata, skey) \leftarrow_{\$} \Pi.SePost^{HASH}(X_c)$
- 5 Return $pdata$ // Adversary is given this

HASH(W): // The random oracle

- 6 Return $H(W)$

FIN(c'): // Adversary provides guess $c' \in \{0, 1\}$ for c

- 7 Return $(c = c')$ // Result of game, true if $c' = c$ and false otherwise

Figure 2: Game defining server-security of protocol Π .

Downside: modularity is a double edged sword

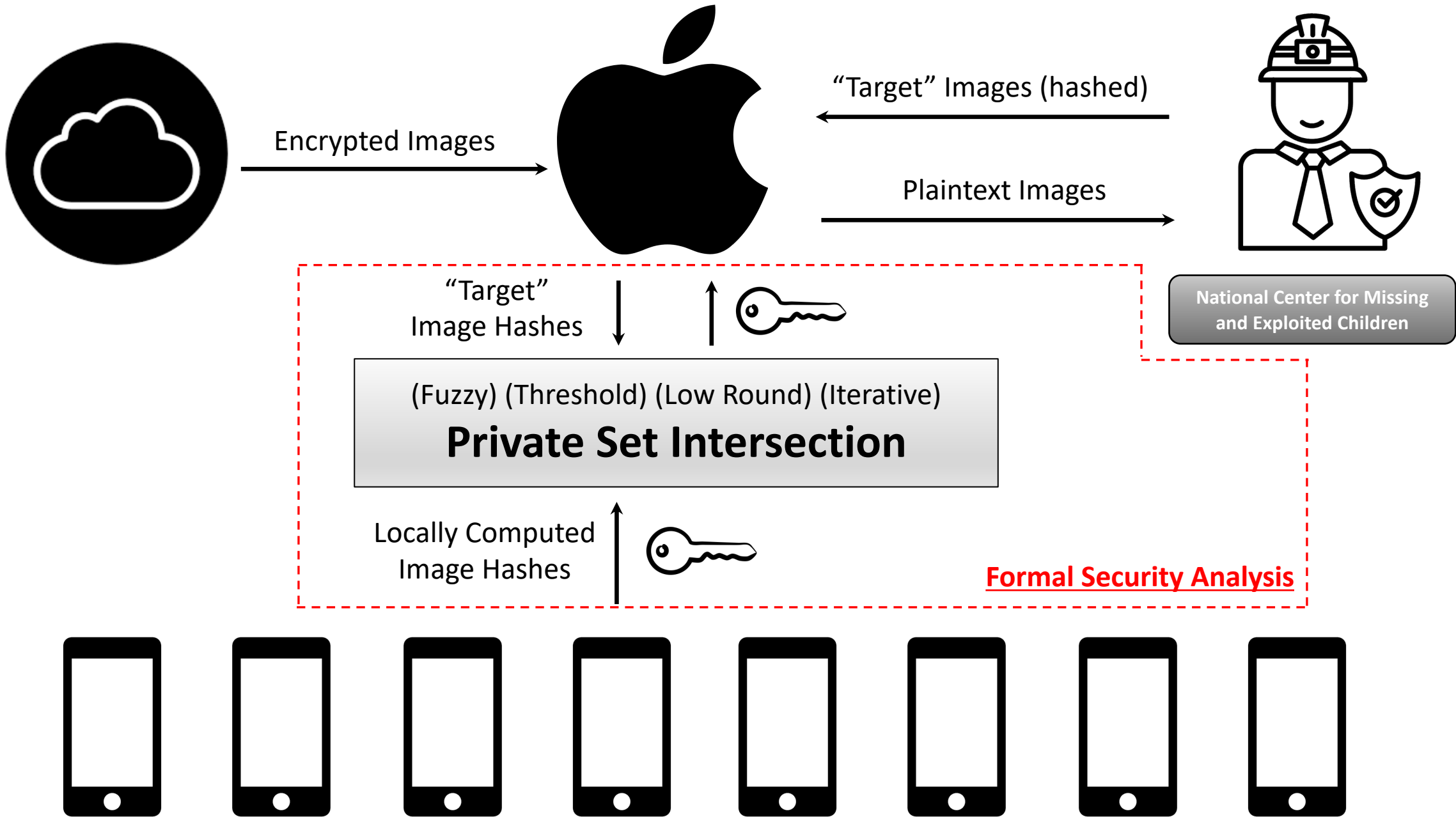
Downside: Game based proofs aren't very expressive

“How *should* we analyze cryptographic *deployments*?”

Our answer: “In context”

Option 1: Apple didn't know about the issues, and their analysis failed to highlight the problems

Option 2: Apple knew about all these issues, and our standard cryptographic analysis tools make it easy to suppress them



Convincingly Thorough

1 Broaden the analytical frame

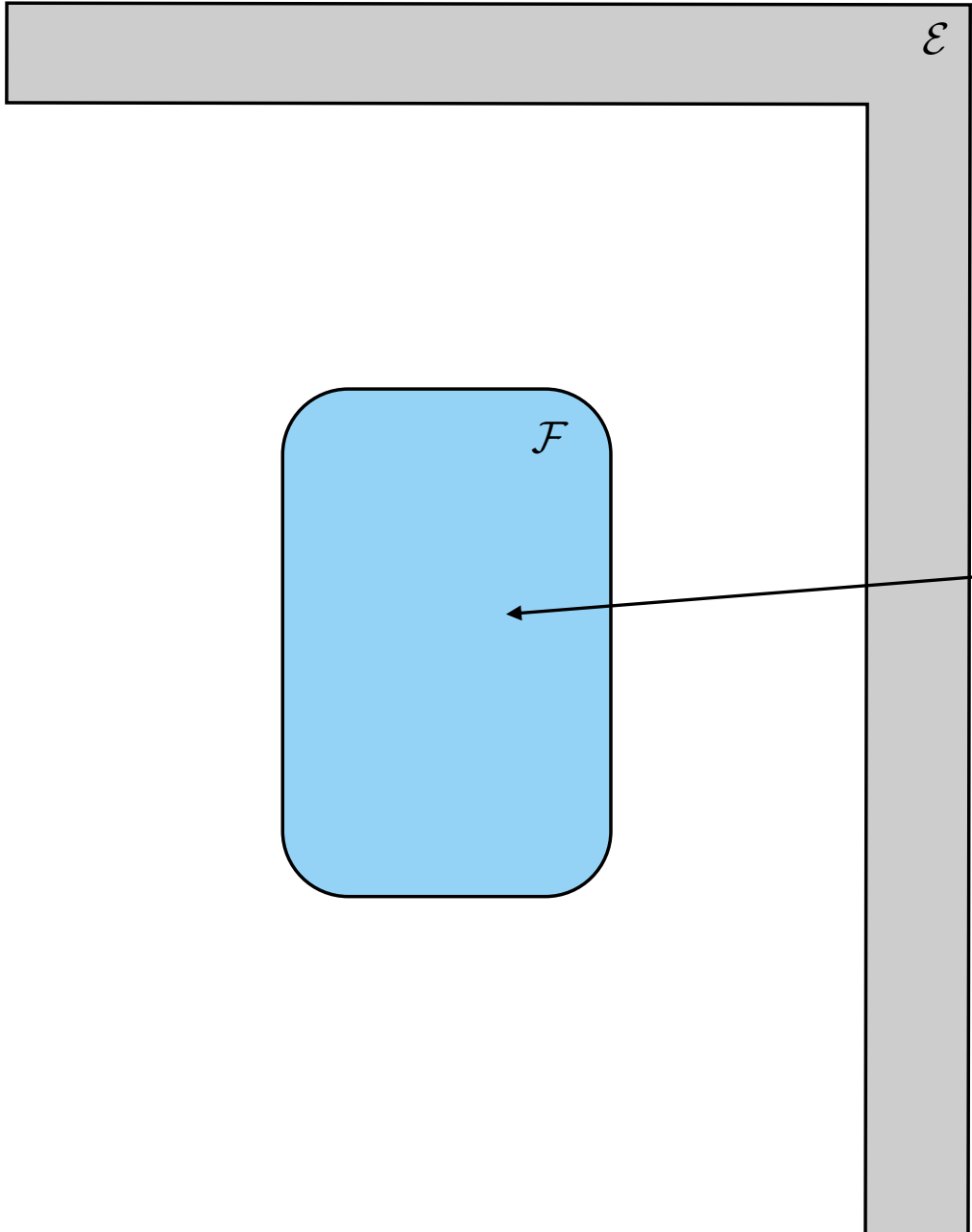
2 Systematic and enumerative

5 Low conceptual overhead

Surface questions of
social importance

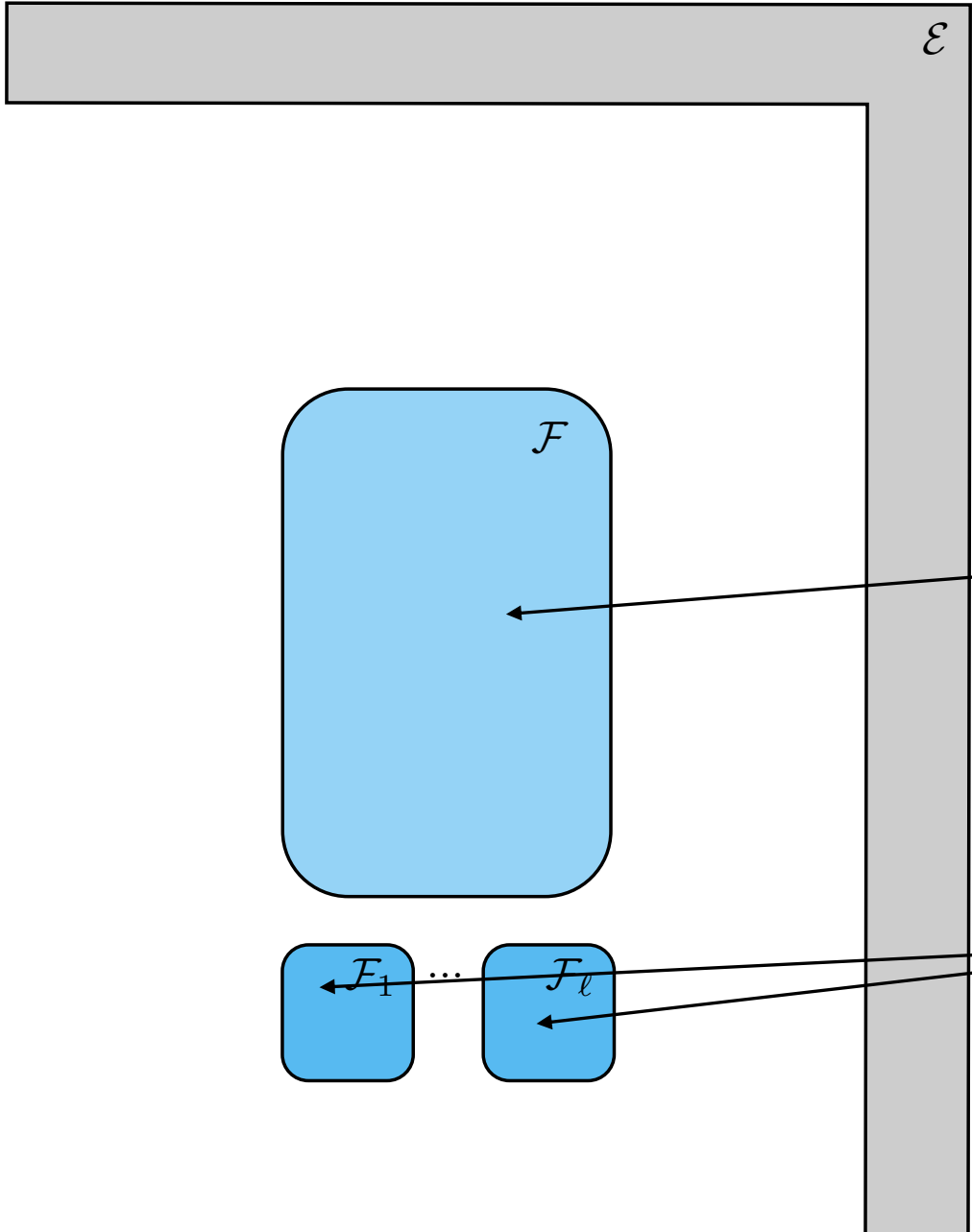
3 Highlights implications of trust

4 Provide a interdisciplinary bridge



New Ideal Functionality
Eg. Apple's PSI protocol

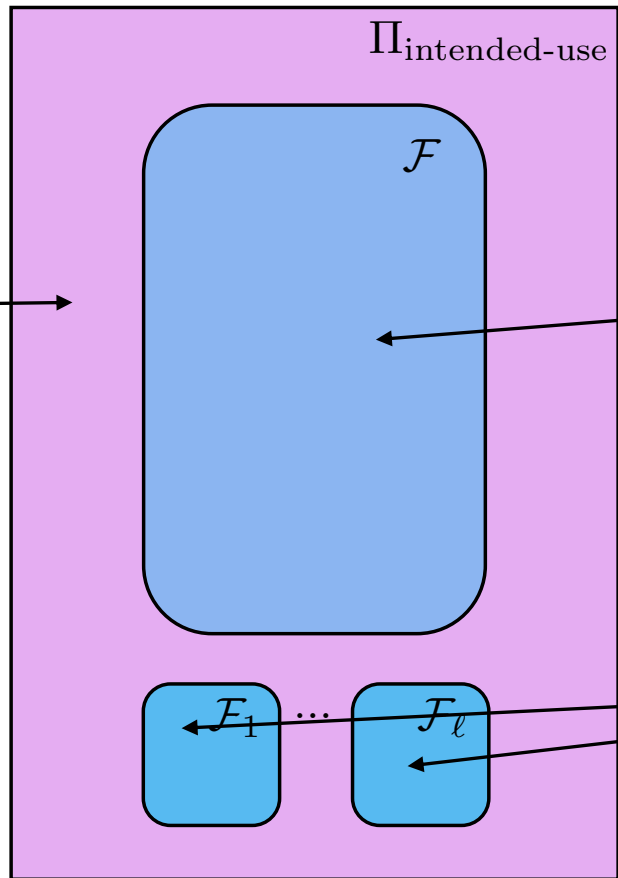




New Ideal Functionality
Eg. Apple's PSI protocol

Composed Ideal Functionalities
Other cryptographic components

Unproved Software
Describes “correct” behavior
of components for which
there is no formal proof

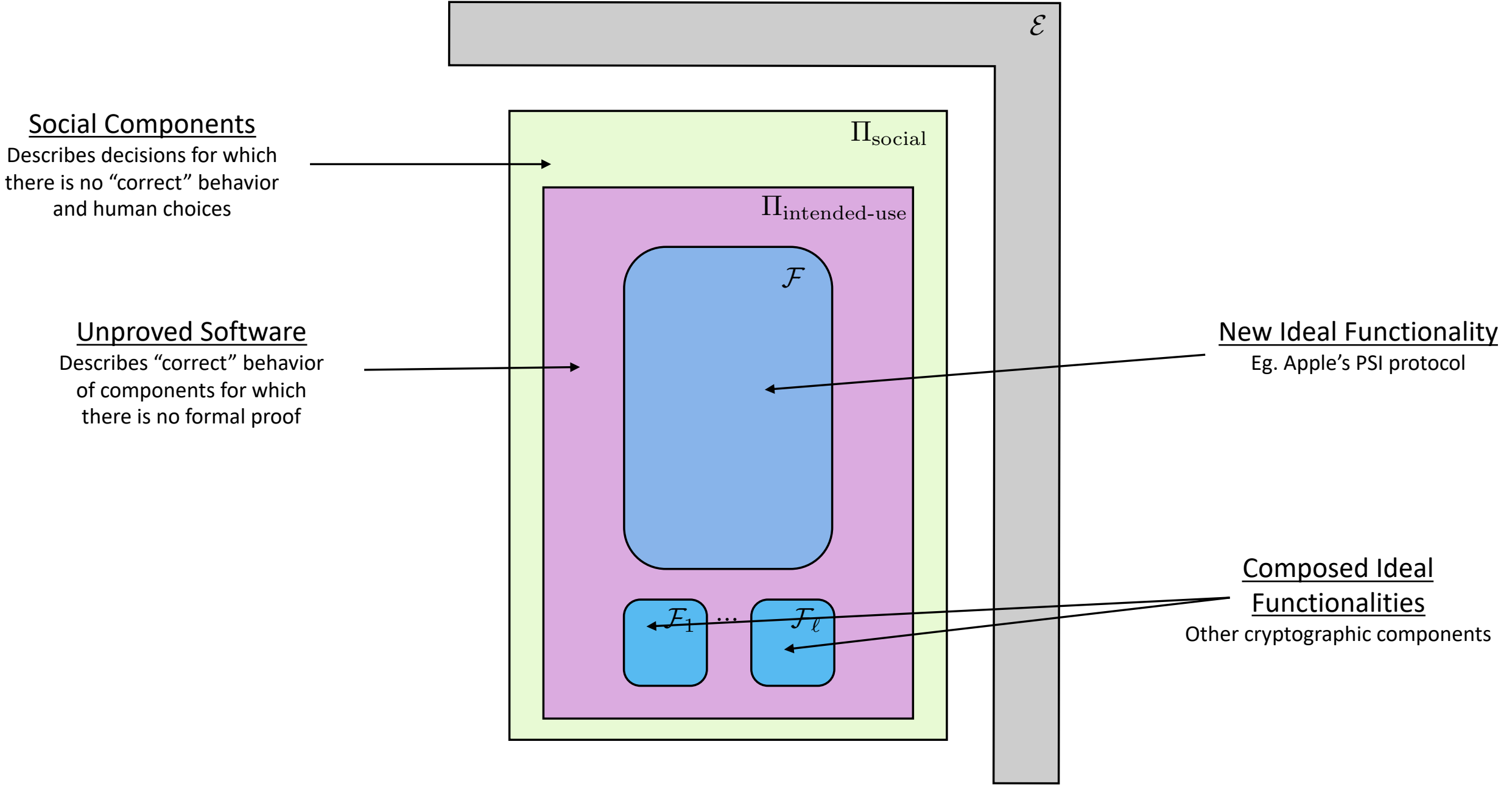


New Ideal Functionality
Eg. Apple’s PSI protocol



Composed Ideal
Functionalities
Other cryptographic components





Social Components

Describes decisions for which there is no "correct" behavior and human choices

Unproved Software

Describes "correct" behavior of components for which there is no formal proof

New Ideal Functionality

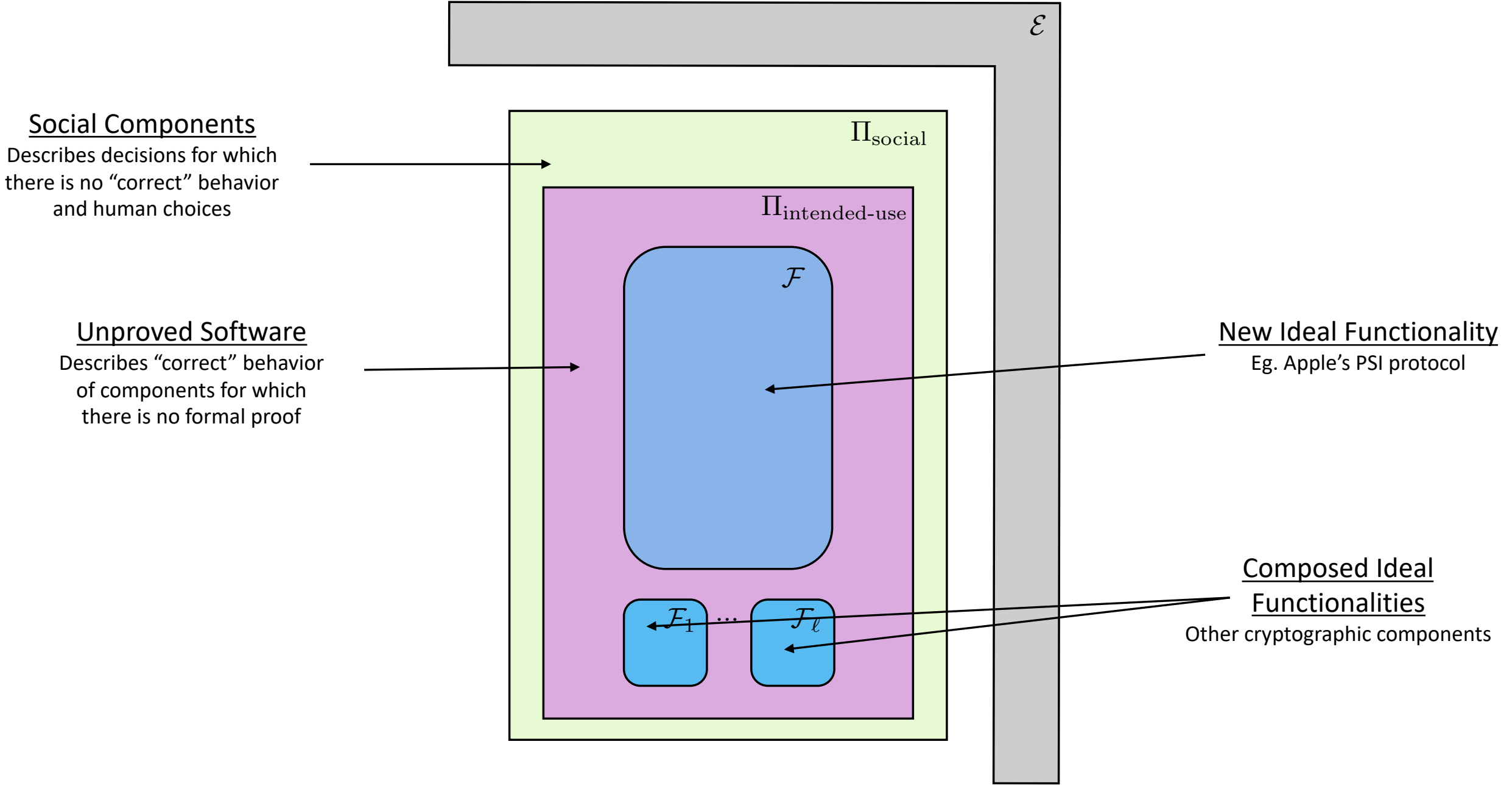
Eg. Apple's PSI protocol

Composed Ideal Functionalities

Other cryptographic components

Did we meet our goals?

- 1 Broaden the analytical frame – Trace full lifecycle of the data
- 2 Systematic and enumerative – Fixed process for generating protocols
- 3 Highlights implications of trust – Trust choices are explicit
- 4 Provide a interdisciplinary bridge – Highlights social components in “natural” language
- 5 Low conceptual overhead – *UC is super simple and easy to understand/work with*



Case Study: Apple's Proposed* CSAM Scanning System

$\mathcal{F}_{\text{Apple-CSAM-Scan}}$

The functionality is parameterized by a threshold \mathcal{T} . For all clients \mathcal{U} , initialize $\text{IsCounting}_{\mathcal{U}} = \text{False}$.

Initializing Clients: Upon input $(\text{InitScan}, \{\mathcal{U}_1, \dots, \mathcal{U}_n\}, \{X_{\mathcal{U}_1}, \dots, X_{\mathcal{U}_n}\})$ from *APPLE*:

1. For $\mathcal{U}_i \in \{\mathcal{U}_1, \dots, \mathcal{U}_n\}$, if $\text{IsInitiated}_{\mathcal{U}_i} = \text{True}$, continue. Otherwise record $X_{\mathcal{U}_i}$ and set $\text{IsInitiated}_{\mathcal{U}_i} = \text{True}$. Finally, send $(\text{InitScanComplete})$ to *APPLE*.

Process Initialization: Upon input (ProcessInit) from \mathcal{U} :

1. If $\text{IsInitiated}_{\mathcal{U}} \neq \text{True}$, set $\text{IsCounting}_{\mathcal{U}} = \text{True}$ and send $(\text{InitScanComplete}, |X_{\mathcal{U}}|)$ to \mathcal{U} . Otherwise, send $(\text{InitScanComplete}, \perp)$ to \mathcal{U} .

Scan Image: Upon input $(\text{ScanImage}, \text{img_id}, \text{img_hash}, k, \text{valid} \in \{\text{true}, \text{false}\})$ from client \mathcal{U} :

1. $\text{IsCounting}_{\mathcal{U}} = \text{True}$, return. If this is the first image received from \mathcal{U} , initialize variable $\text{counter}_{\mathcal{U}}$, list $\text{photolist}_{\mathcal{U}}$, and list $\text{keys}_{\mathcal{U}}$.

// \mathcal{F} alerts the Server that the user has uploaded a photo.

2. If $\text{img_hash} \notin X_{\mathcal{U}}$ or $\text{valid} = \text{false}$, send $(\text{ScanImageComplete}, \mathcal{U}, \text{false})$ to *APPLE*.
3. If $\text{img_hash} \in X_{\mathcal{U}}$, increment $\text{counter}_{\mathcal{U}}$, append img_id to $\text{photolist}_{\mathcal{U}}$ and append k to $\text{keys}_{\mathcal{U}}$. Then,
 - (a) If $\text{counter}_{\mathcal{U}} < \mathcal{T}$ send $(\text{ScanImageComplete}, \mathcal{U}, \text{true})$ to *APPLE*.
 - (b) If $\text{counter}_{\mathcal{U}} \geq \mathcal{T}$, send $(\text{ThresholdMet}, \mathcal{U}, \text{photolist}_{\mathcal{U}}, \text{keys}_{\mathcal{U}})$ to *APPLE*.

One Shot vs. Incremental

Parameters known to all parties:

- two parties: server and client,
- B is the maximum set size for the server and client,
- t is the threshold,
- s_{\max} is the maximum size of the set S of synthetics,
- all associated data values ad in \mathcal{D} have the same public length.

The functionality \mathcal{F} :

- Wait for input $X = \{x_1, x_2, \dots\}$ from the server; abort if the server is corrupt and $|X| > B$.
- Send $|X|$ to the client; abort if the client is corrupt and aborts.
- Wait for input $\bar{Y} \in (\mathcal{U} \times \mathcal{ID} \times \mathcal{D})^m$ and $S \subseteq id(\bar{Y})$ from the client; abort if the client is corrupt and $(m > B$ or $|S| > s_{\max})$.
- Send \bar{Y}_{id} to the server.
- If $|id(\bar{Y} \cap X) \setminus S| > t$:
 send $\bar{Y}_{id[id(\bar{Y} \cap X) \setminus S]_{\{id, ad\}}}$ and S to the server,
 otherwise:
 send $id(\bar{Y} \cap X) \cup S$ to the server.

$\mathcal{F}_{\text{Apple-CSAM-Scan}}$

The functionality is parameterized by a threshold \mathcal{T} . For all clients \mathcal{U} , initialize $IsCounting_{\mathcal{U}} = \text{False}$.

Initializing Clients: Upon input $(\text{InitScan}, \{\mathcal{U}_1, \dots, \mathcal{U}_n\}, \{X_{\mathcal{U}_1}, \dots, X_{\mathcal{U}_n}\})$ from $APPLE$:

1. For $\mathcal{U}_i \in \{\mathcal{U}_1, \dots, \mathcal{U}_n\}$, if $IsInitiated_{\mathcal{U}_i} = \text{True}$, continue. Otherwise record $X_{\mathcal{U}_i}$ and set $IsInitiated_{\mathcal{U}_i} = \text{True}$. Finally, send $(\text{InitScanComplete})$ to $APPLE$.

Process Initialization: Upon input (ProcessInit) from \mathcal{U} :

1. If $IsInitiated_{\mathcal{U}} \neq \text{True}$, set $IsCounting_{\mathcal{U}} = \text{True}$ and send $(\text{InitScanComplete}, |X_{\mathcal{U}}|)$ to \mathcal{U} . Otherwise, send $(\text{InitScanComplete}, \perp)$ to \mathcal{U} .

Scan Image: Upon input $(\text{ScanImage}, \text{img_id}, \text{img_hash}, k, \text{valid} \in \{\text{true}, \text{false}\})$ from client \mathcal{U} :

1. $IsCounting_{\mathcal{U}} = \text{True}$, return. If this is the first image received from \mathcal{U} , initialize variable $\text{counter}_{\mathcal{U}}$, list $\text{photolist}_{\mathcal{U}}$, and list $\text{keys}_{\mathcal{U}}$.
// \mathcal{F} alerts the Server that the user has uploaded a photo.
2. If $\text{img_hash} \notin X_{\mathcal{U}}$ or $\text{valid} = \text{false}$, send $(\text{ScanImageComplete}, \mathcal{U}, \text{false})$ to $APPLE$.
3. If $\text{img_hash} \in X_{\mathcal{U}}$, increment $\text{counter}_{\mathcal{U}}$, append img_id to $\text{photolist}_{\mathcal{U}}$ and append k to $\text{keys}_{\mathcal{U}}$. Then,
 - (a) If $\text{counter}_{\mathcal{U}} < \mathcal{T}$ send $(\text{ScanImageComplete}, \mathcal{U}, \text{true})$ to $APPLE$.
 - (b) If $\text{counter}_{\mathcal{U}} \geq \mathcal{T}$, send $(\text{ThresholdMet}, \mathcal{U}, \text{photolist}_{\mathcal{U}}, \text{keys}_{\mathcal{U}})$ to $APPLE$.

Figure 1: The ideal functionality \mathcal{F} for ftPSI-AD

Self Composition

Parameters known to all parties:

- two parties: server and client,
- B is the maximum set size for the server and client,
- t is the threshold,
- s_{\max} is the maximum size of the set S of synthetics,
- all associated data values ad in \mathcal{D} have the same public length.

The functionality \mathcal{F} :

- Wait for input $X = \{x_1, x_2, \dots\}$ from the server;
abort if the server is corrupt and $|X| > B$.
- Send $|X|$ to the client; abort if the client is corrupt and aborts.
- Wait for input $\bar{Y} \in (\mathcal{U} \times \mathcal{ID} \times \mathcal{D})^m$ and $S \subseteq id(\bar{Y})$ from the client;
abort if the client is corrupt and $(m > B$ or $|S| > s_{\max})$.
- Send \bar{Y}_{id} to the server.
- If $|id(\bar{Y} \cap X) \setminus S| > t$:
send $\bar{Y}[id(\bar{Y} \cap X) \setminus S]_{\{id, ad\}}$ and S to the server,
otherwise:
send $id(\bar{Y} \cap X) \cup S$ to the server.

Figure 1: The ideal functionality \mathcal{F} for ftPSI-AD

Extractability

4.4.1 Privacy for the server's dataset X against a malicious client

First, let's show that a malicious client \mathcal{A}_c^H that interacts with an honest server with input X , learns nothing about X other than its size.

For an adversary \mathcal{A}_c^H , a simulator Sim_c , and an environment \mathcal{Z} , define the following two random variables:

- $\text{REAL}_{\Pi, \mathcal{A}_c^H, \mathcal{Z}}$ is the random variable defined as the output of \mathcal{Z} in a real world execution after interacting with the malicious client \mathcal{A}_c^H as in Figure 4a.
- $\text{IDEAL}_{\mathcal{F}, \text{Sim}_c, \mathcal{Z}}$ is the random variable defined as the output of \mathcal{Z} in an ideal world execution after interacting with the simulator Sim_c as in Figure 4b.

Definition 2. We say that a protocol Π for \mathcal{F} is **server private** if for every efficient adversary \mathcal{A}_c , there exists an efficient adversary Sim_c , such that for every efficient environment \mathcal{Z} ,

$$\left| \Pr[\text{REAL}_{\Pi, \mathcal{A}_c^H, \mathcal{Z}} = 1] - \Pr[\text{IDEAL}_{\mathcal{F}, \text{Sim}_c, \mathcal{Z}} = 1] \right| \leq \epsilon$$

for some negligible ϵ .

4. Sim_c sends $\text{pdata} \leftarrow (L, P_1, \dots, P_{n'})$ and the nonces to \mathcal{A}_c^H .

This completes the description of Sim_c .

Case Study: Apple's Proposed* CSAM Scanning System

$\mathcal{F}_{\text{Apple-CSAM-Scan}}$

The functionality is parameterized by a threshold \mathcal{T} . For all clients \mathcal{U} , initialize $\text{IsCounting}_{\mathcal{U}} = \text{False}$.

Initializing Clients: Upon input $(\text{InitScan}, \{\mathcal{U}_1, \dots, \mathcal{U}_n\}, \{X_{\mathcal{U}_1}, \dots, X_{\mathcal{U}_n}\})$ from *APPLE*:

1. For $\mathcal{U}_i \in \{\mathcal{U}_1, \dots, \mathcal{U}_n\}$, if $\text{IsInitiated}_{\mathcal{U}_i} = \text{True}$, continue. Otherwise record $X_{\mathcal{U}_i}$ and set $\text{IsInitiated}_{\mathcal{U}_i} = \text{True}$. Finally, send $(\text{InitScanComplete})$ to *APPLE*.

Process Initialization: Upon input (ProcessInit) from \mathcal{U} :

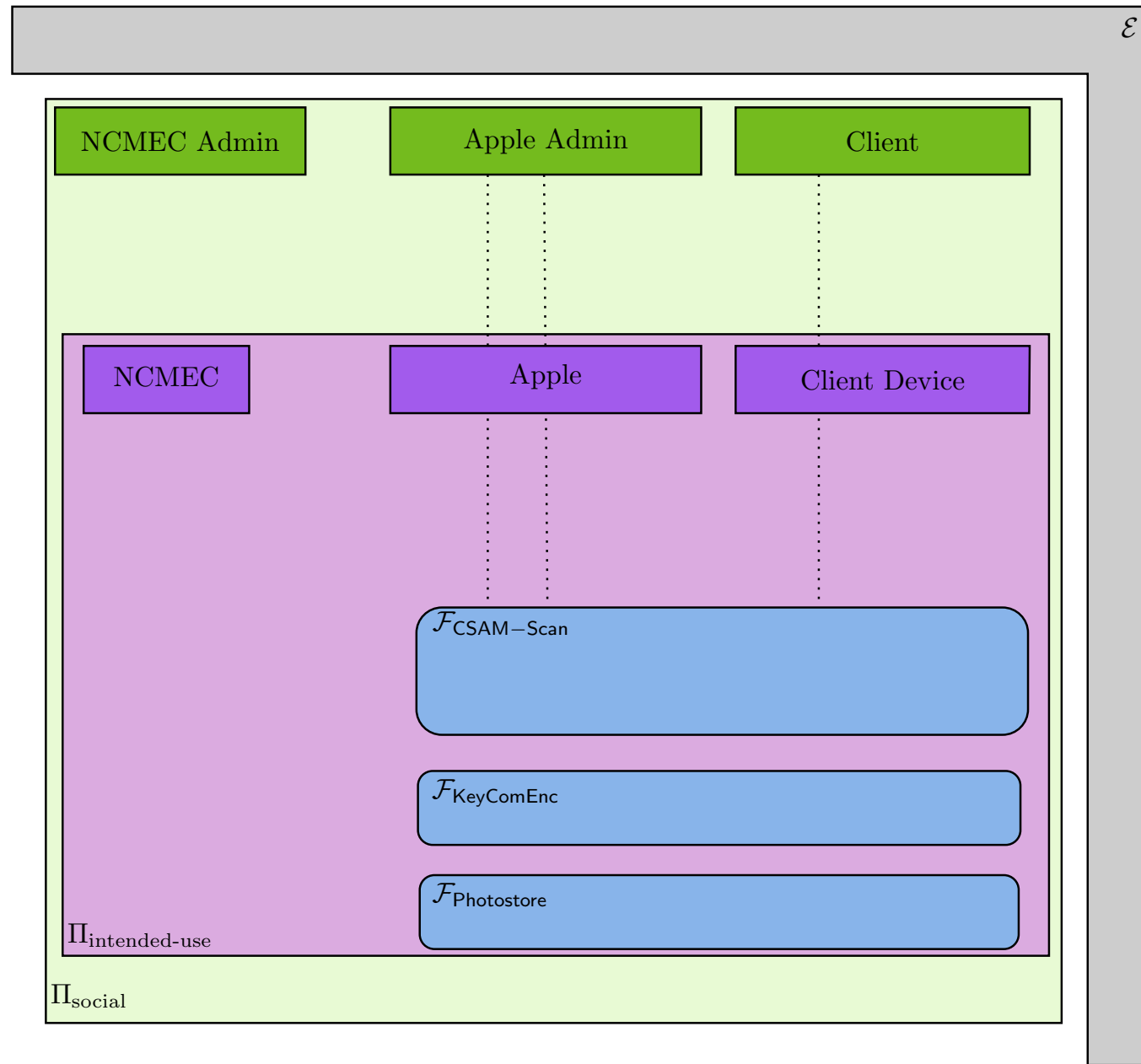
1. If $\text{IsInitiated}_{\mathcal{U}} \neq \text{True}$, set $\text{IsCounting}_{\mathcal{U}} = \text{True}$ and send $(\text{InitScanComplete}, |X_{\mathcal{U}}|)$ to \mathcal{U} . Otherwise, send $(\text{InitScanComplete}, \perp)$ to \mathcal{U} .

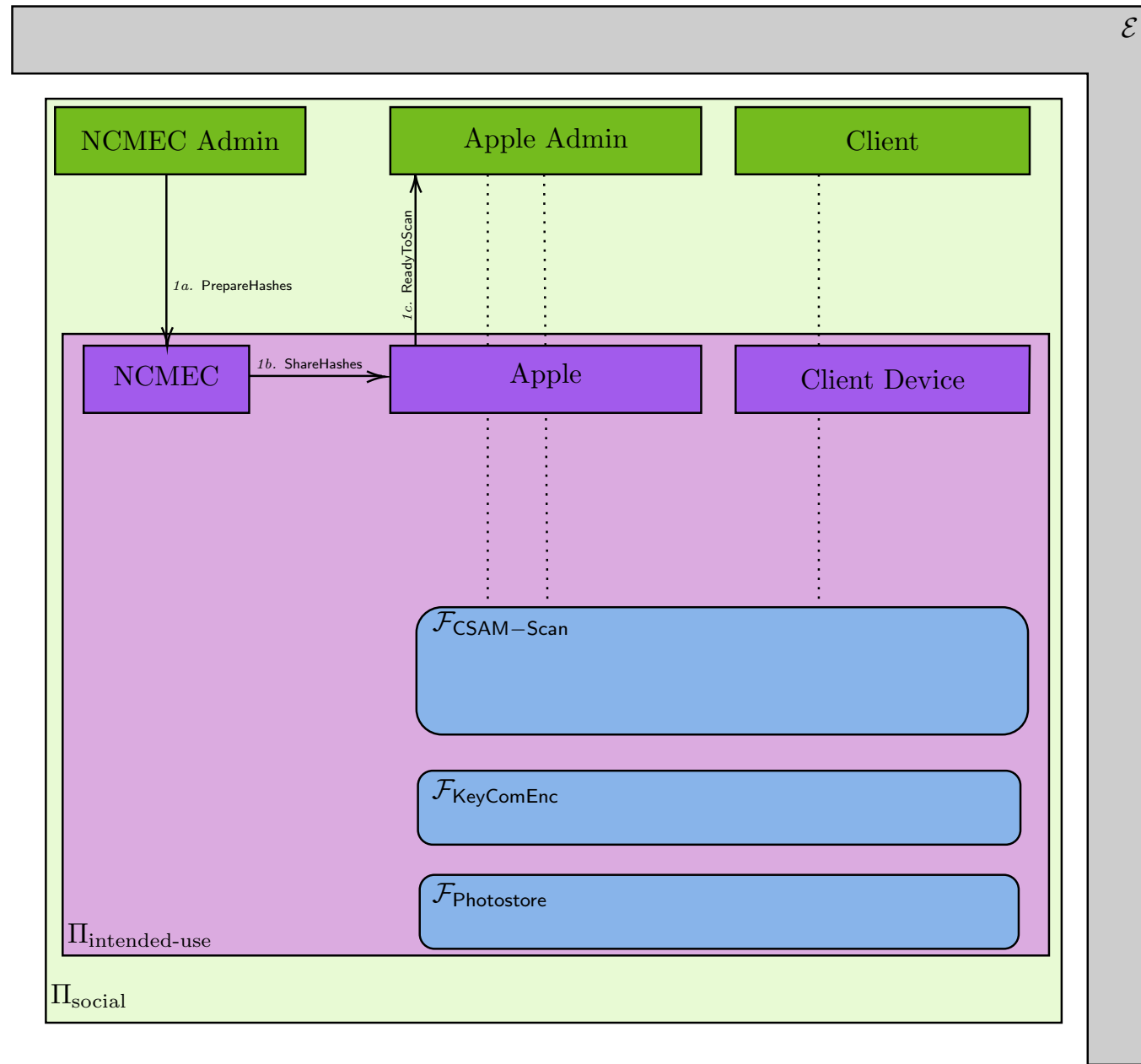
Scan Image: Upon input $(\text{ScanImage}, \text{img_id}, \text{img_hash}, k, \text{valid} \in \{\text{true}, \text{false}\})$ from client \mathcal{U} :

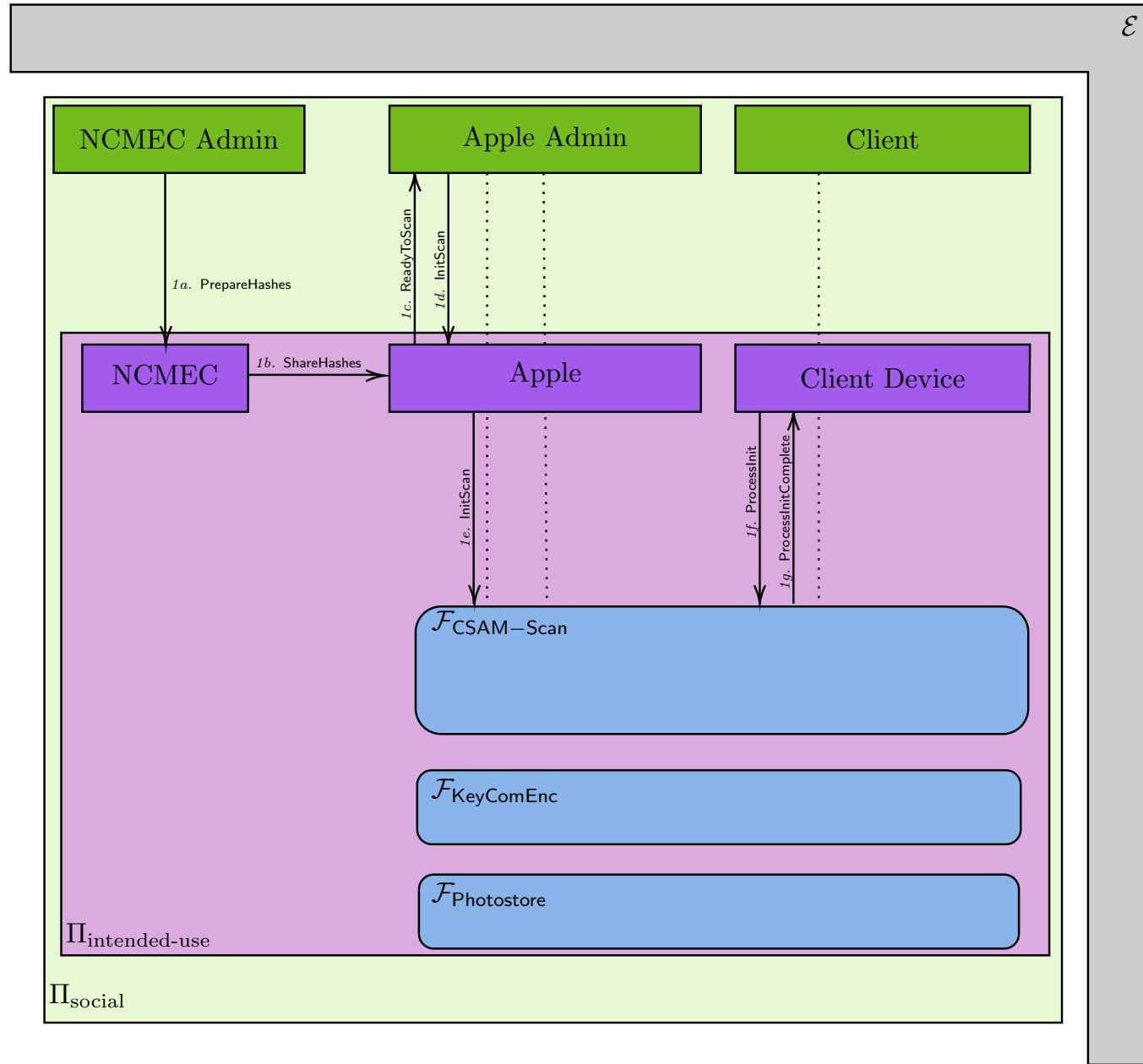
1. $\text{IsCounting}_{\mathcal{U}} = \text{True}$, return. If this is the first image received from \mathcal{U} , initialize variable $\text{counter}_{\mathcal{U}}$, list $\text{photolist}_{\mathcal{U}}$, and list $\text{keys}_{\mathcal{U}}$.

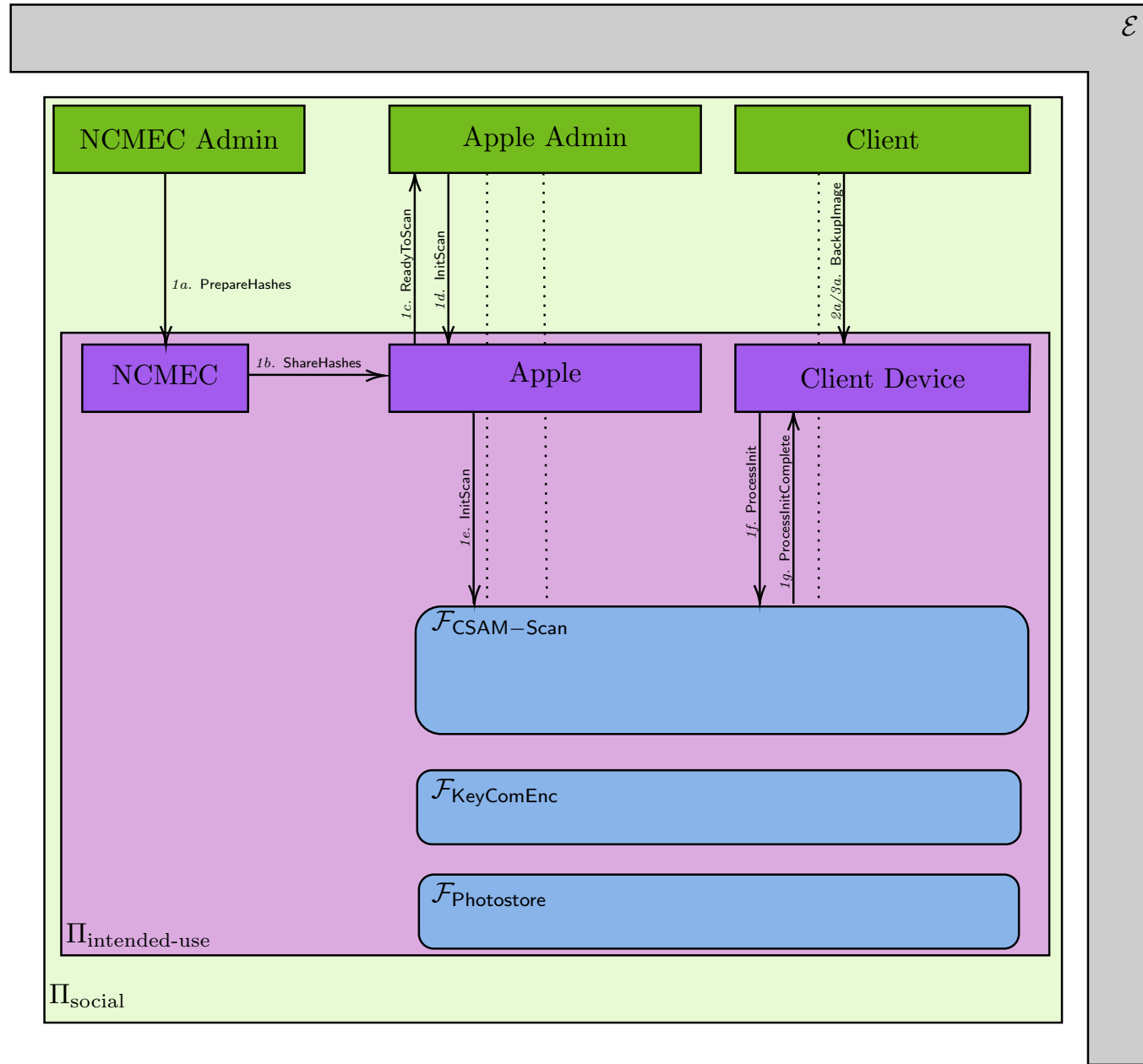
// \mathcal{F} alerts the Server that the user has uploaded a photo.

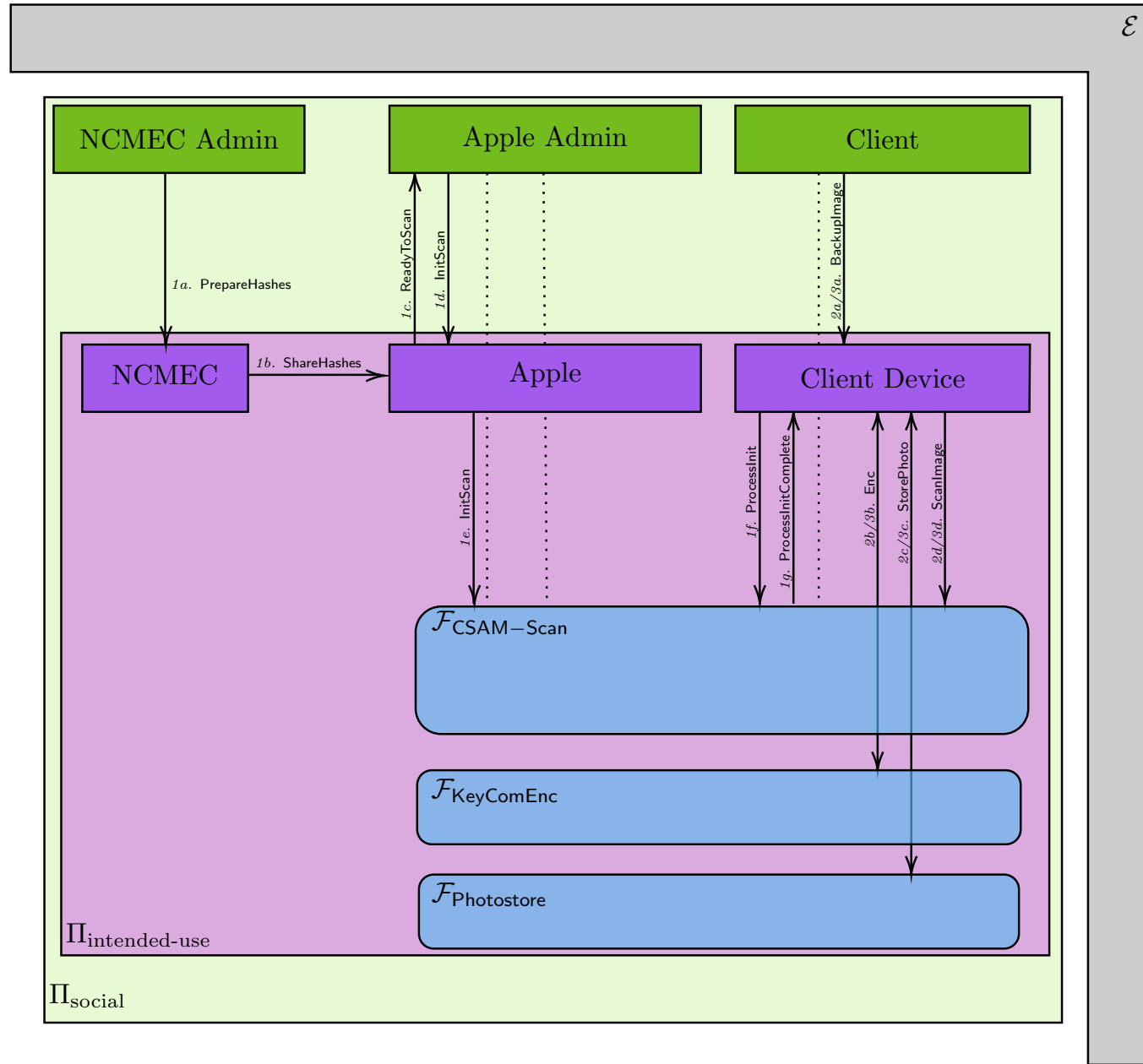
2. If $\text{img_hash} \notin X_{\mathcal{U}}$ or $\text{valid} = \text{false}$, send $(\text{ScanImageComplete}, \mathcal{U}, \text{false})$ to *APPLE*.
3. If $\text{img_hash} \in X_{\mathcal{U}}$, increment $\text{counter}_{\mathcal{U}}$, append img_id to $\text{photolist}_{\mathcal{U}}$ and append k to $\text{keys}_{\mathcal{U}}$. Then,
 - (a) If $\text{counter}_{\mathcal{U}} < \mathcal{T}$ send $(\text{ScanImageComplete}, \mathcal{U}, \text{true})$ to *APPLE*.
 - (b) If $\text{counter}_{\mathcal{U}} \geq \mathcal{T}$, send $(\text{ThresholdMet}, \mathcal{U}, \text{photolist}_{\mathcal{U}}, \text{keys}_{\mathcal{U}})$ to *APPLE*.

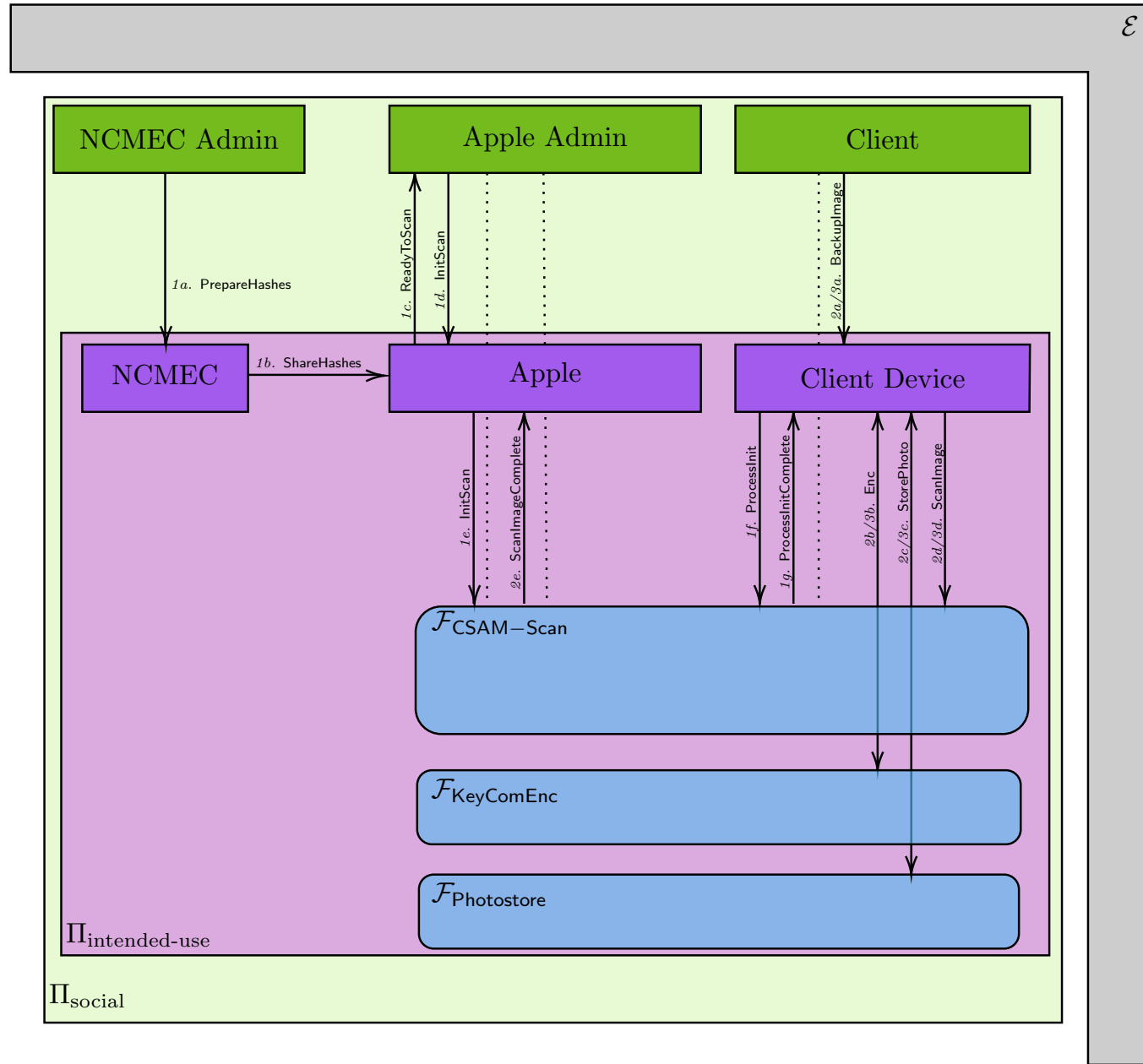


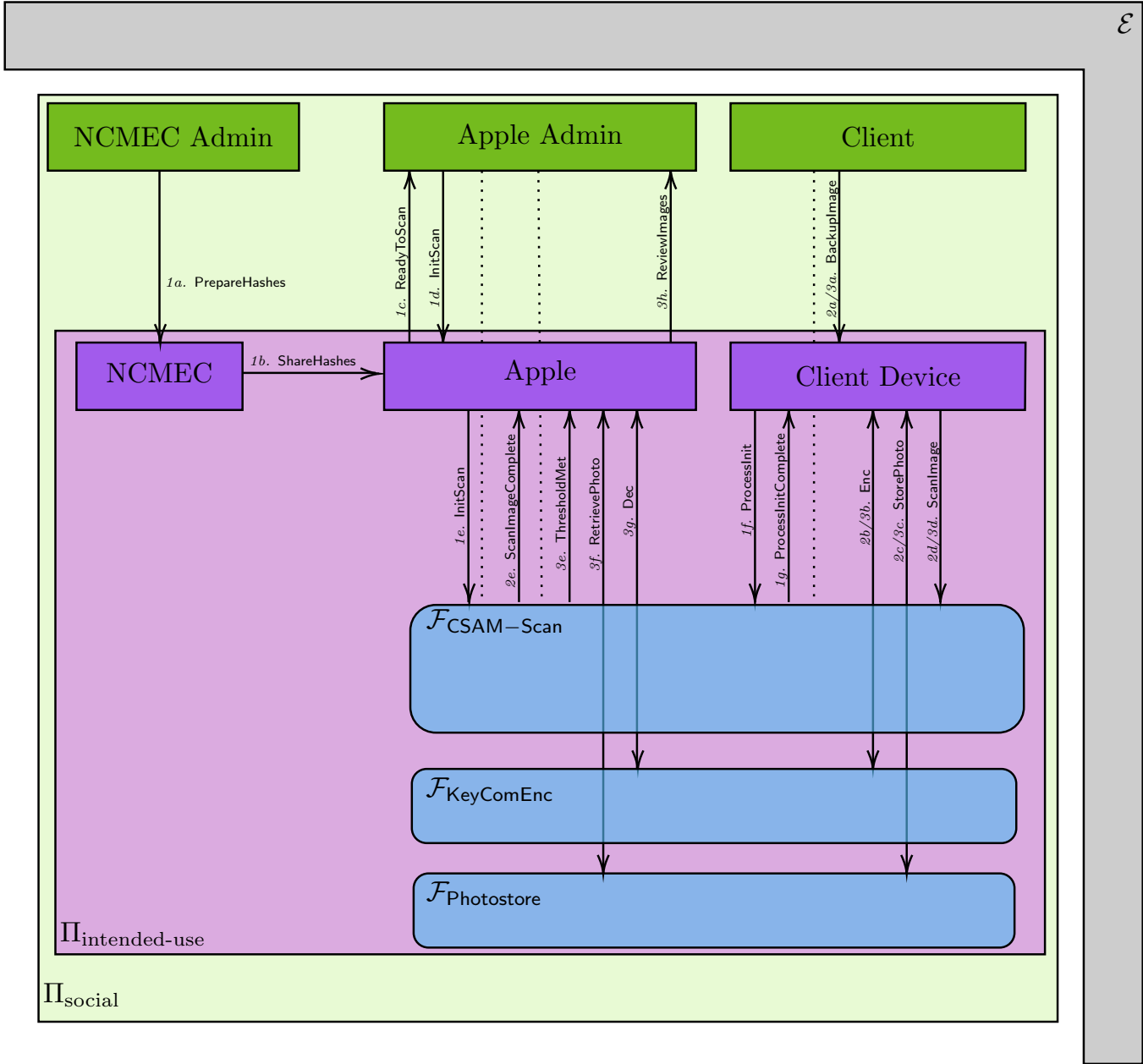








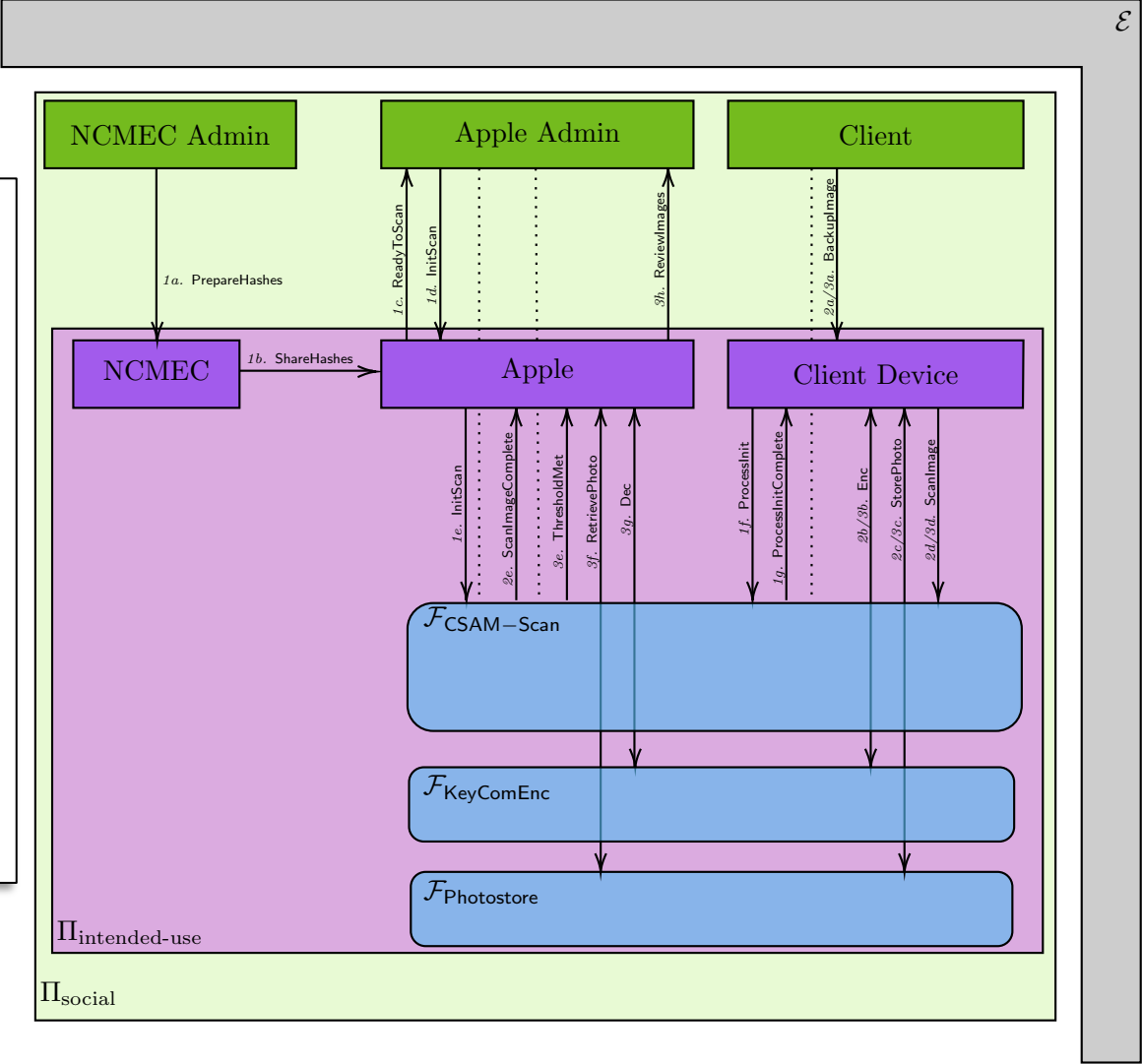




Does this *actually* help?

Top Down Analysis:

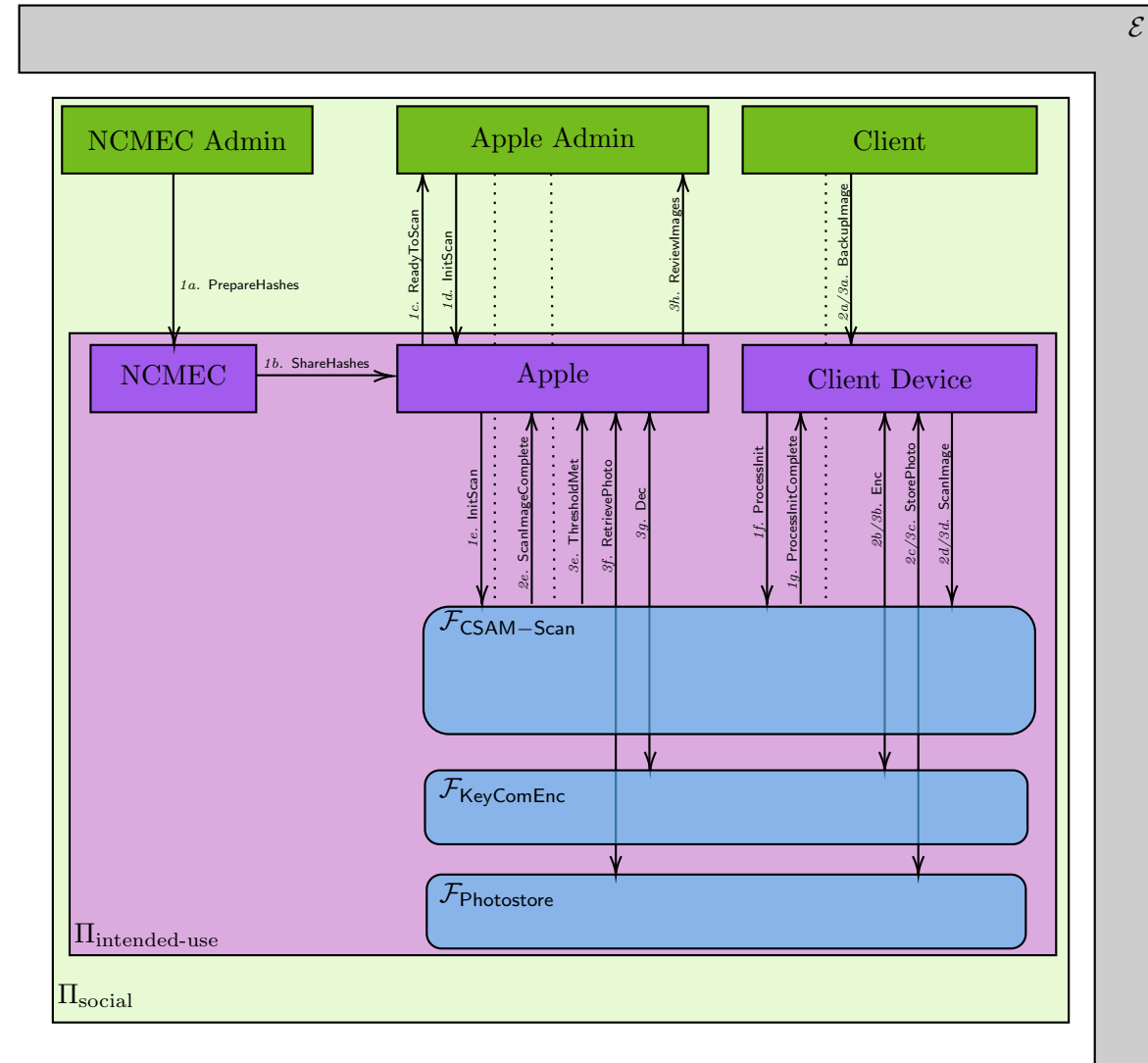
- CSAM Detection provides these privacy and security assurances:
- Apple does not learn anything about images that do not match the known CSAM database.
 - Apple can't access metadata or visual derivatives for matched CSAM images until a threshold of matches is exceeded for an iCloud Photos account.
 - 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.
 - Users can't identify which images were flagged as CSAM by the system.
- For detailed information about the cryptographic protocol and security proofs that the CSAM Detection process uses, see [The Apple PSI System](https://www.apple.com/child-safety/pdf/CSAM_Detection_Technical_Summary.pdf).



Does this *actually* help?

Bottom Up Analysis:

- What are the consequences of divergent/surprising/non-normative choices in Π_{social} ?
- What are the consequences of malicious/buggy software running $\Pi_{intended-use}$?
- What if software dependencies within $\Pi_{intended-use}$ do not have desired properties?



Where do we go from here?

More Social properties!

Potential Additional Properties:

1. Client devices can detect if apple is attempt to scan for any perceptual hashes not designated as CSAM by Apple [SchefflerKulshresthaMayer23]
2. A pre-determined judge must approve of any image decryptions before they are to happen
3. Even if the pre-determined judge becomes compromised, images can only be decrypted after the match threshold has been met
4. False positives should leave an indelible audit trail
5. The system should be publicly auditable. Namely, the public should be able to verify:
 - That the server is using the same set of hashes for all clients
 - That the perceptual hashes for which the server is scanning are those selected by NCMEC
 - That any images that the server claims were the decrypted as a result of the system were actually uploaded by the accused client
6. Clients can only use the photostore if they honestly use the scanning system

Key Take Aways:

1

We desperately need new analysis tools that prevent privacy theater

2

There is an opportunity to create Interdisciplinary boundary objects
[Star and Griesmemer 1989]

3

Apple's proposed system is a valuable case study and it remains understudied

Thanks!

