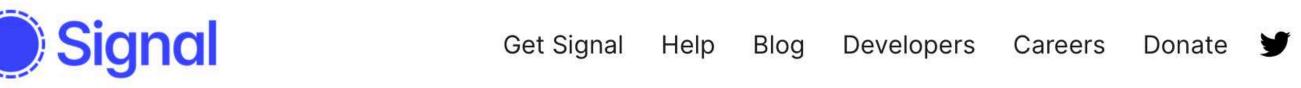
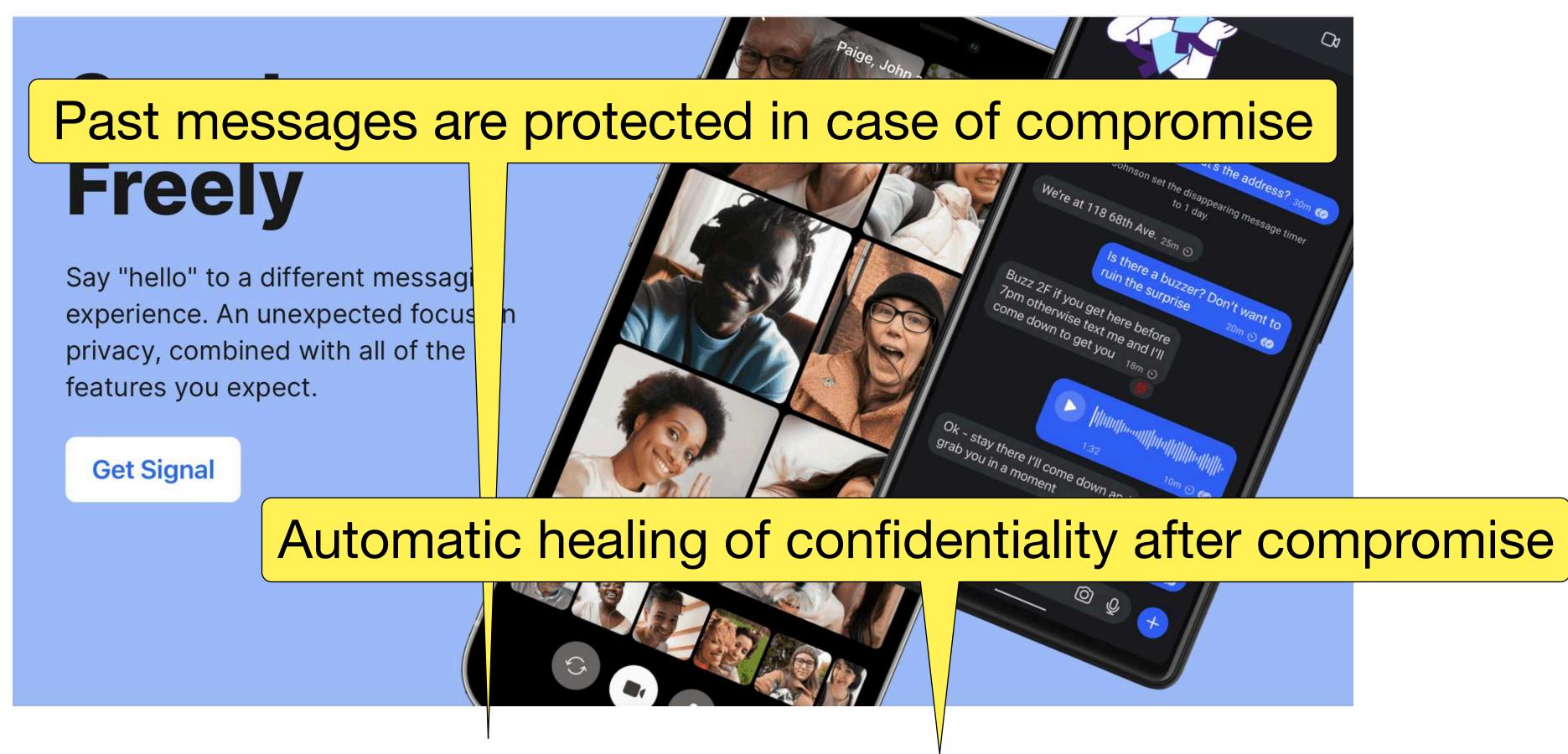
# Beyond the ratchet: practical challenges in secure messaging

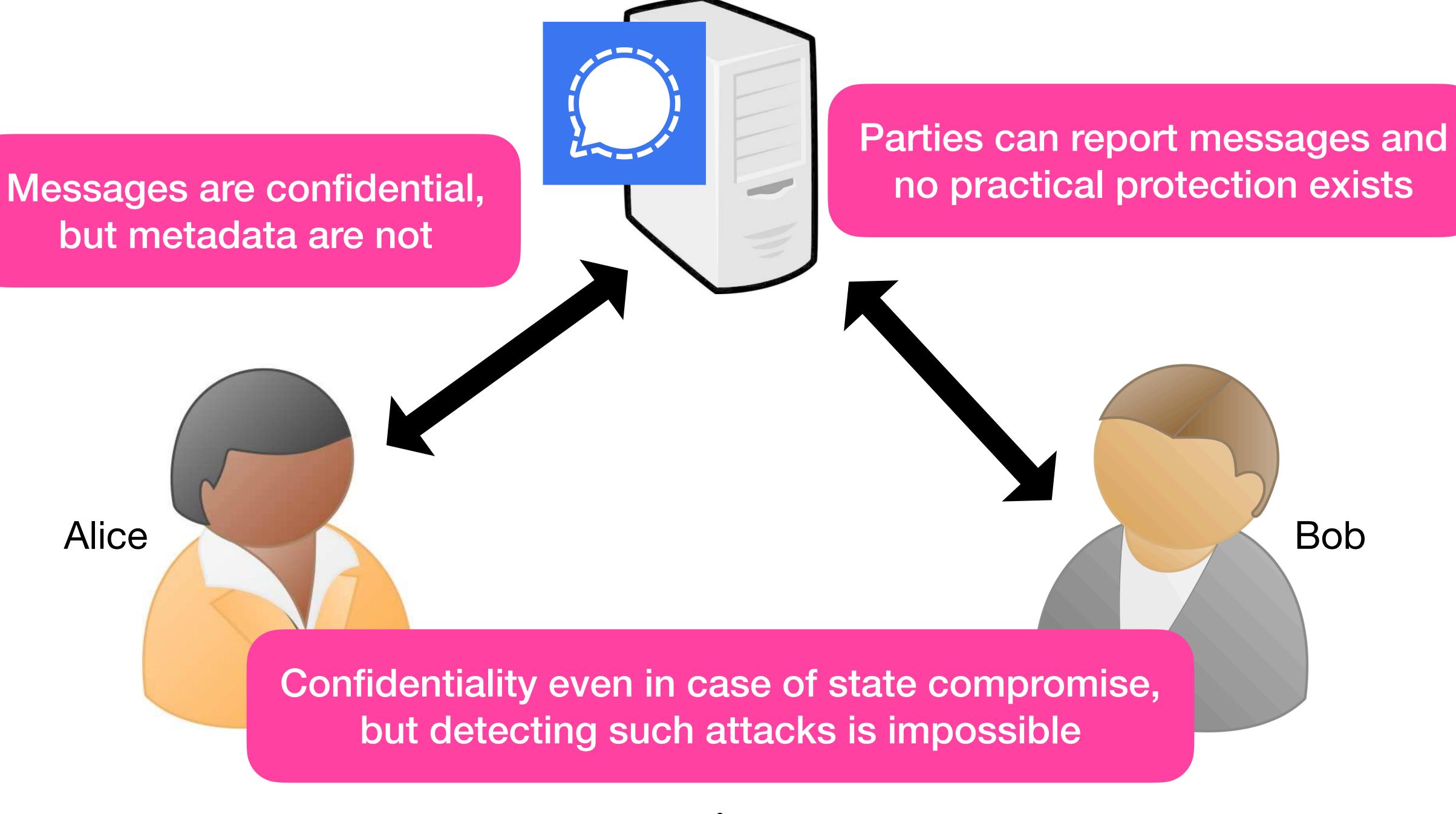
PhD oral exam of Simone Colombo, 11th of December 2024

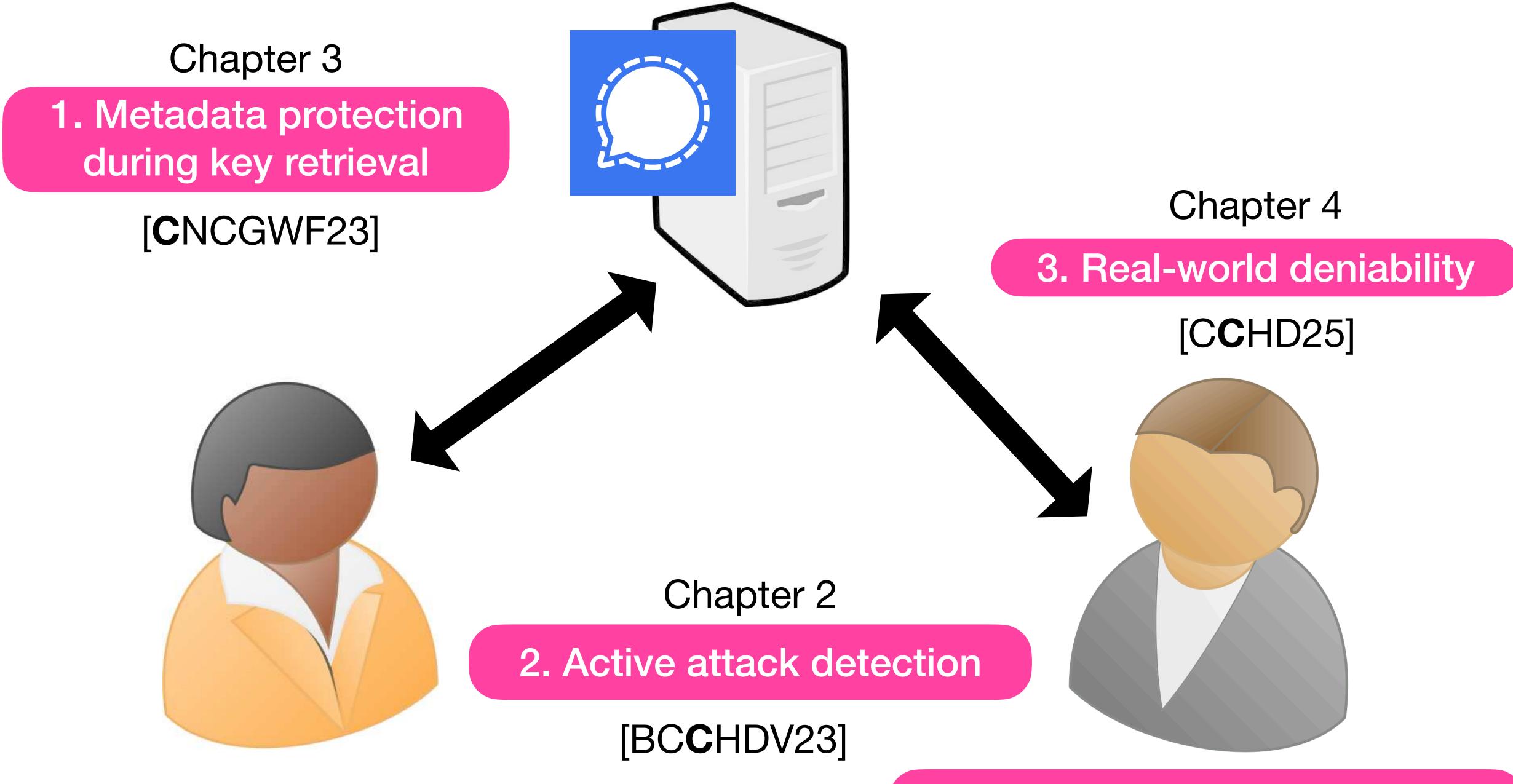




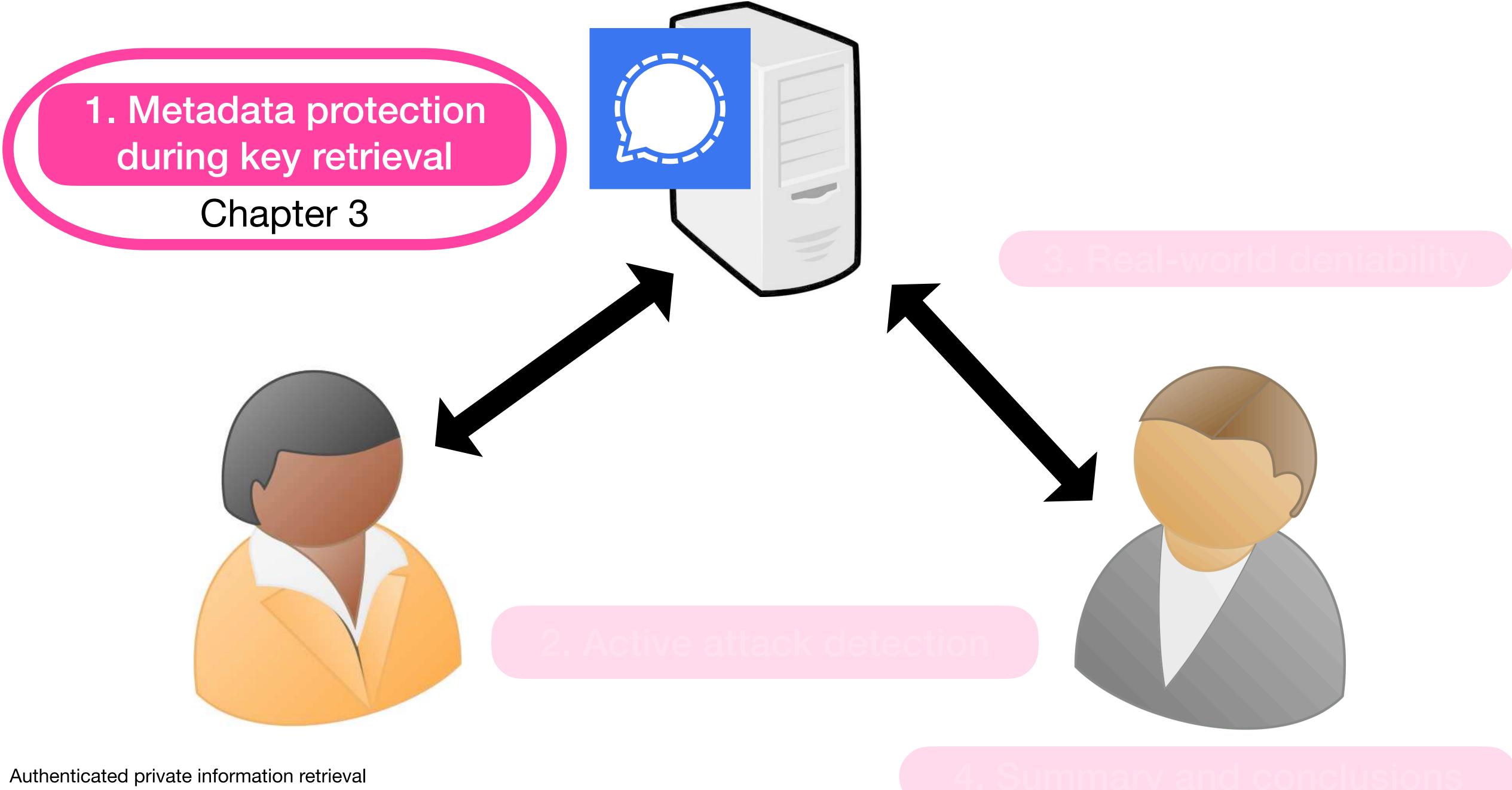
end-to-end encryption, forward secrecy, post-compromise security

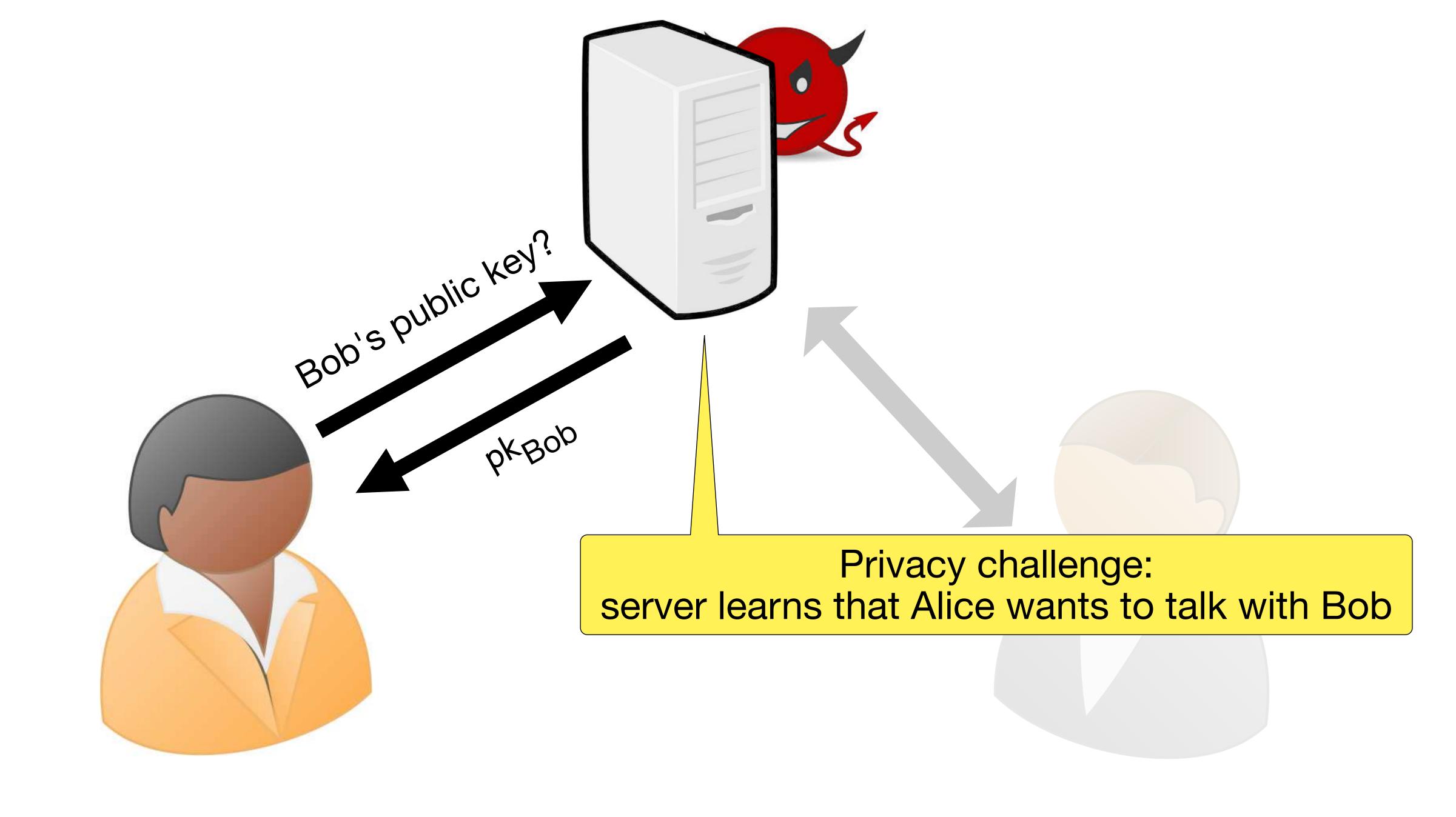
Despite these important improvements, several challenges remain





4. Summary and conclusions





# Private information retrieval (PIR) [CGKS95]

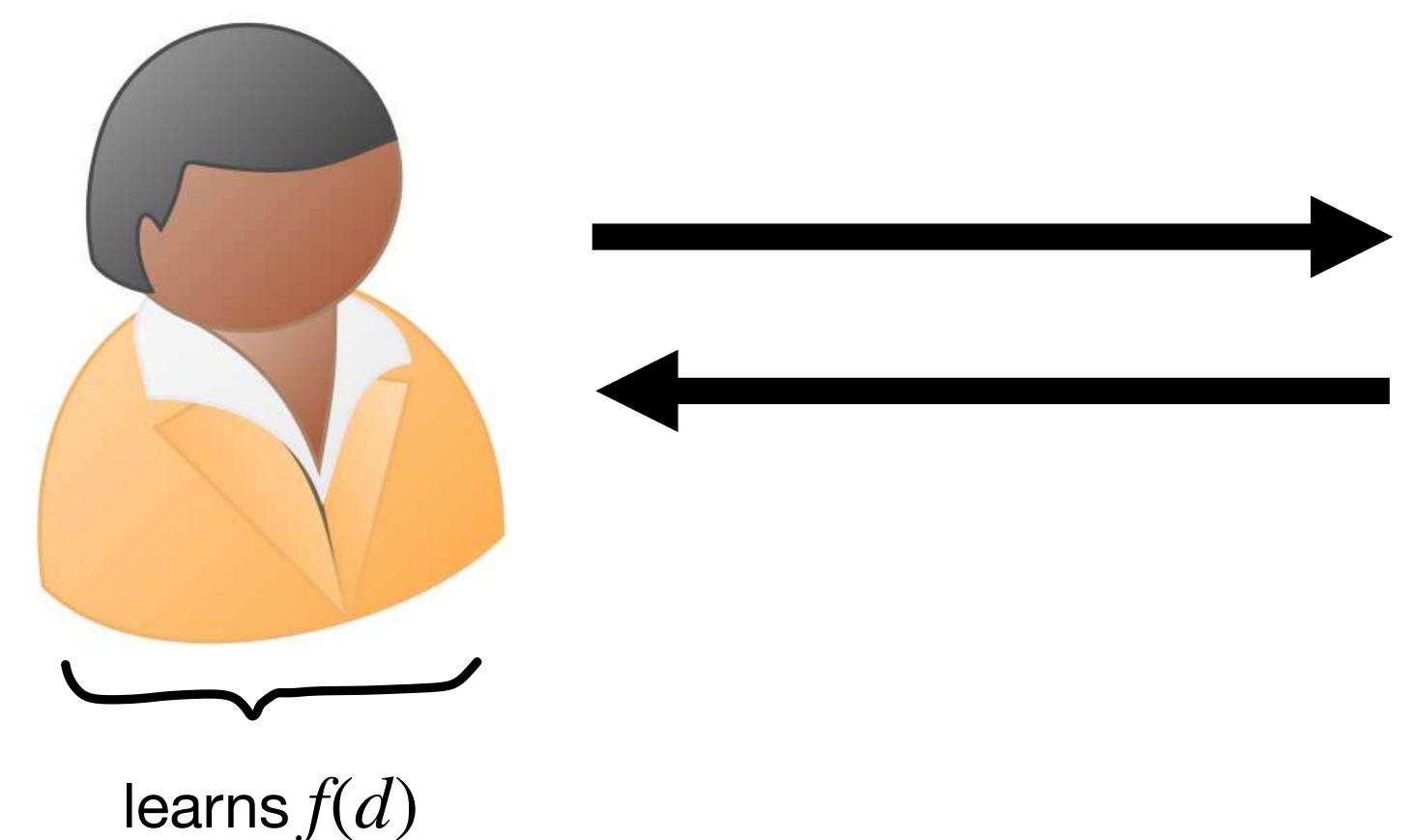
holds database  $d \in \mathbb{F}^N$ holds index  $i \in \{1,...,N\}$ 

learns  $d_i$ 

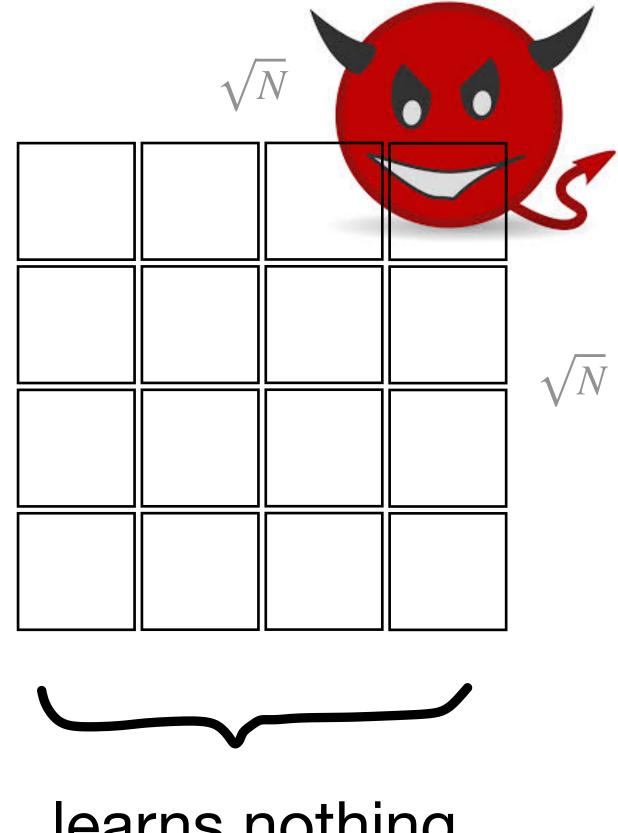
learns nothing

# Private information retrieval (PIR) [CGKS95,WYGVZ17]

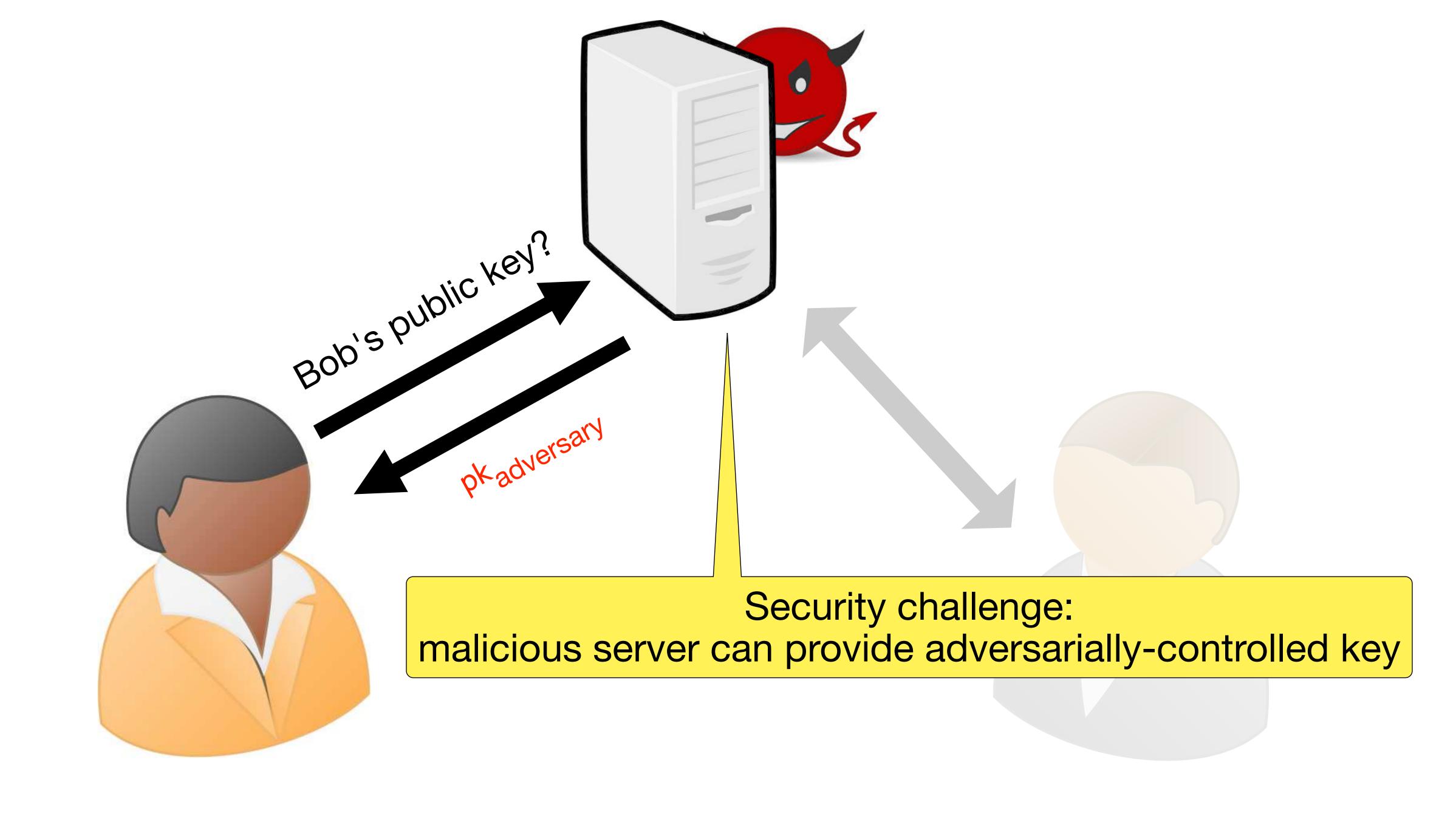
holds function  $f: \mathbb{F}^N \to \mathbb{F}$ 



holds database  $d \in \mathbb{F}^N$ 



learns nothing



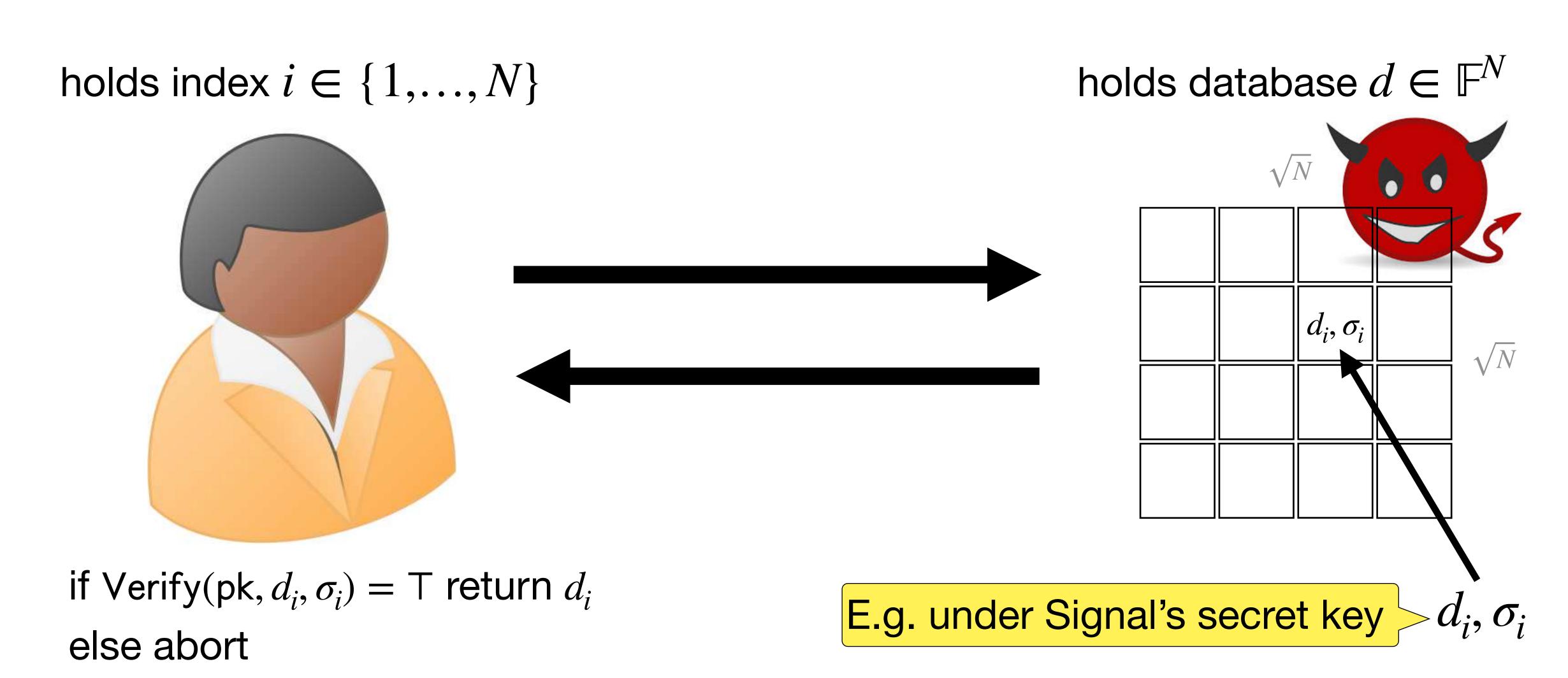
## PIR does not consider integrity

holds database  $d \in \mathbb{F}^N$ holds index  $i \in \{1, ..., N\}$ learns nothing learns wrong  $d_i'$ 

## PIR does not consider integrity

holds database  $d \in \mathbb{F}^N$ holds index  $i \in \{1,...,N\}$ learns wrong pkadversary learns nothing

# PIR and authentication are not enough



## PIR and authentication are not enough

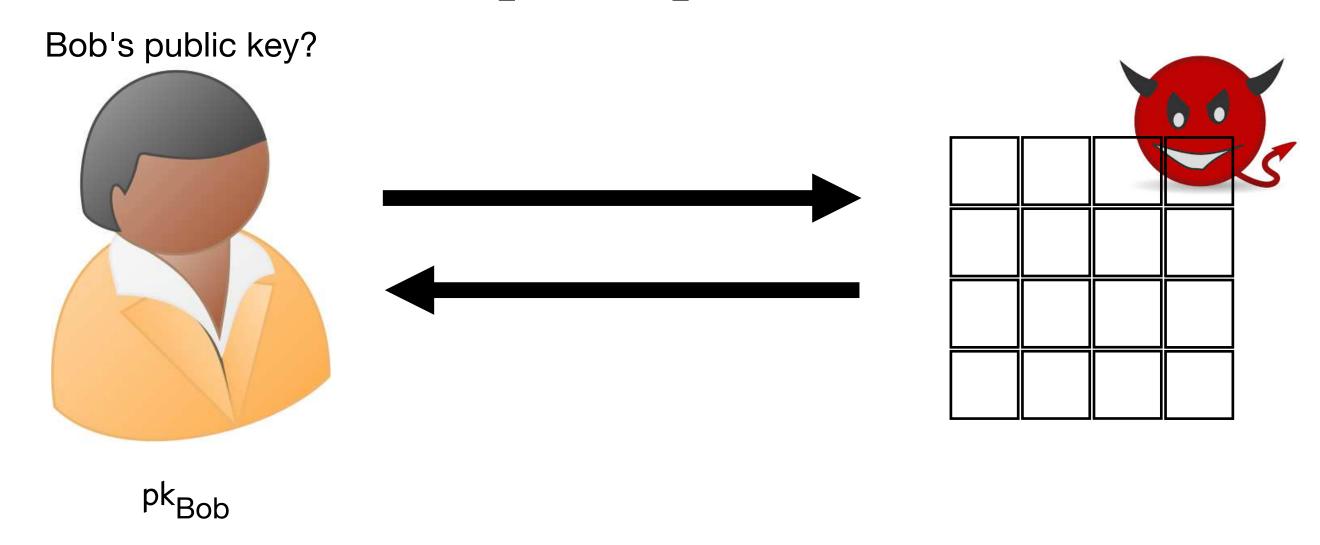
holds index  $i \in \{1,...,N\}$ 

holds database  $d \in \mathbb{F}^N$ 

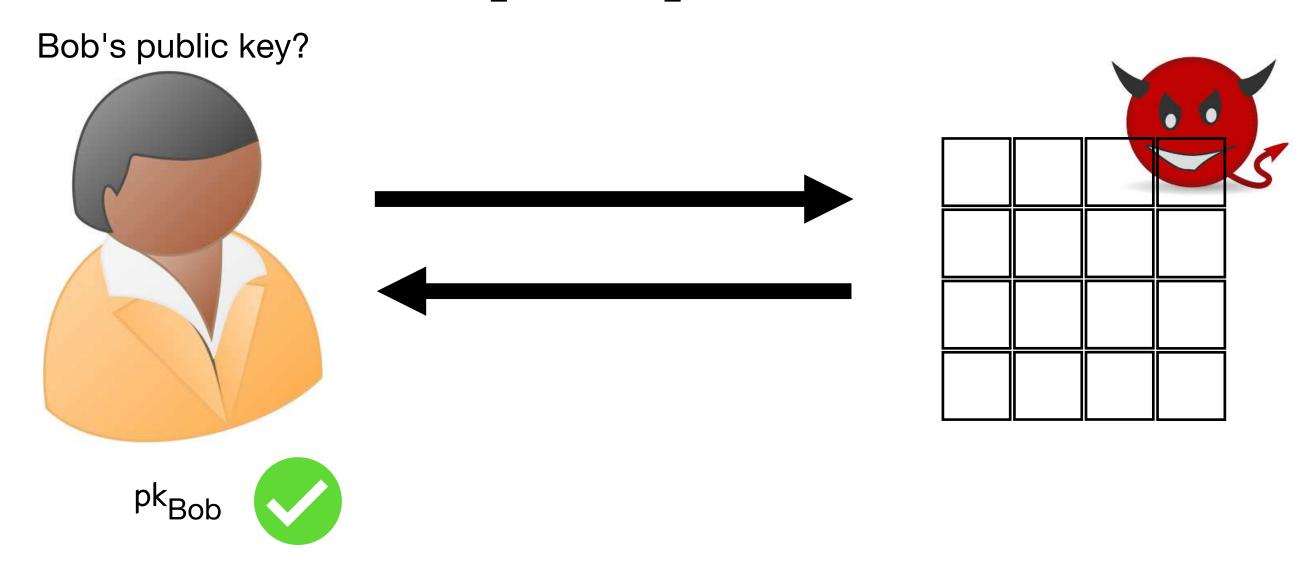
#### Our contribution: authenticated private information retrieval

- First definition of authenticated PIR in multi-server and single-server settings.
- Multi-server schemes to fetch records and evaluate functions on database.
- Two single-server schemes to fetch single-bit records.
- Implementation and evaluation of all the schemes that we propose.
- Keyd, a PGP public-key directory service that builds on authenticated PIR.

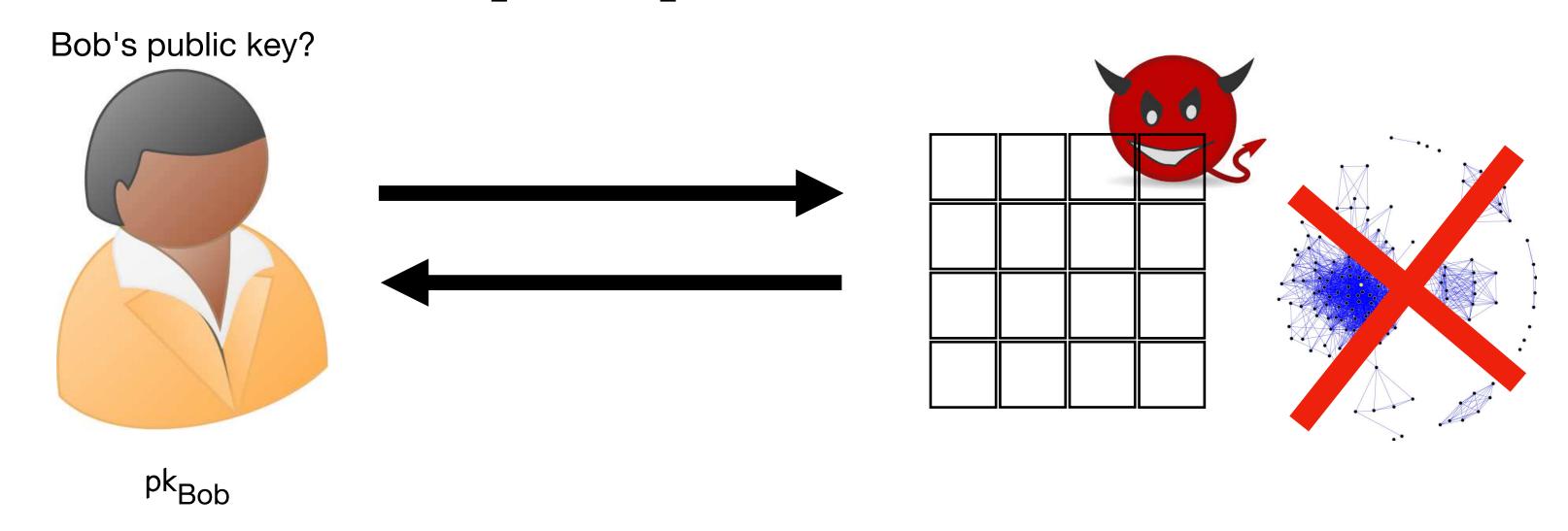
The verify  $(pk, a_i, \sigma_i) = 1$  return  $a_i$  by revealing she queries the i<sup>th</sup> entry: else abort selective-failure attack [KS06].



• Efficiency: Total communication is sublinear in the size of the database.

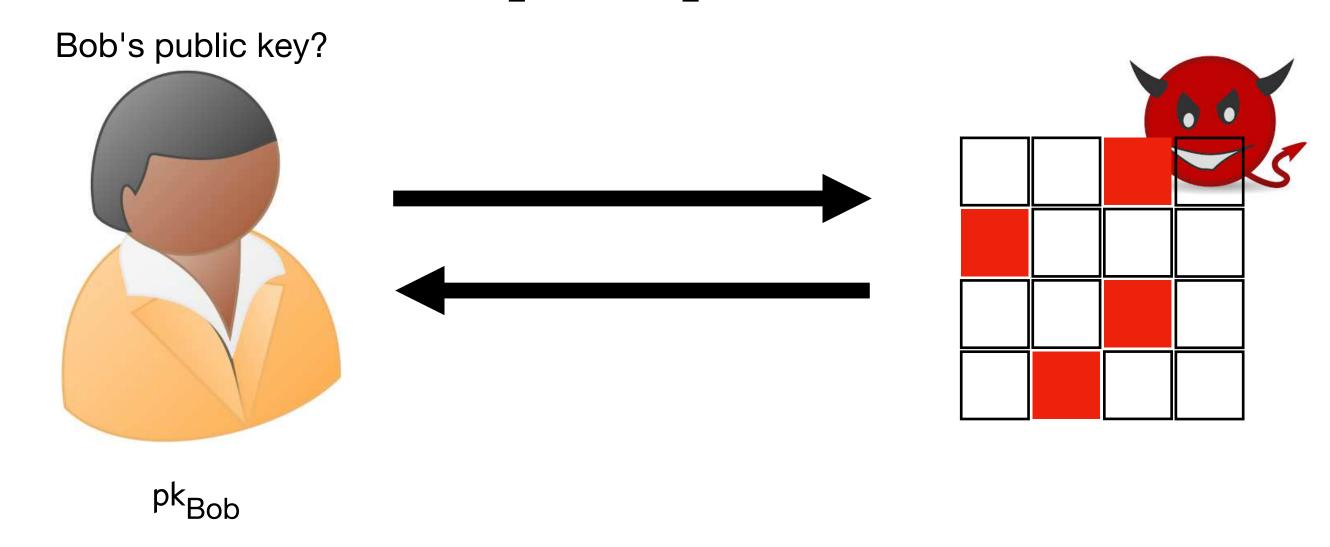


- Efficiency: Total communication is sublinear in the size of the database.
- Correctness: If client and server are honest, the client recovers pk<sub>Bob</sub>.



- Efficiency: Total communication is sublinear in the size of the database.
- Correctness: If client and server are honest, the client recovers pk<sub>Bob</sub>.
- Privacy: The server(s) learns nothing about the content of the client's query, even if the server(s) learns whether the client aborted during reconstruction.

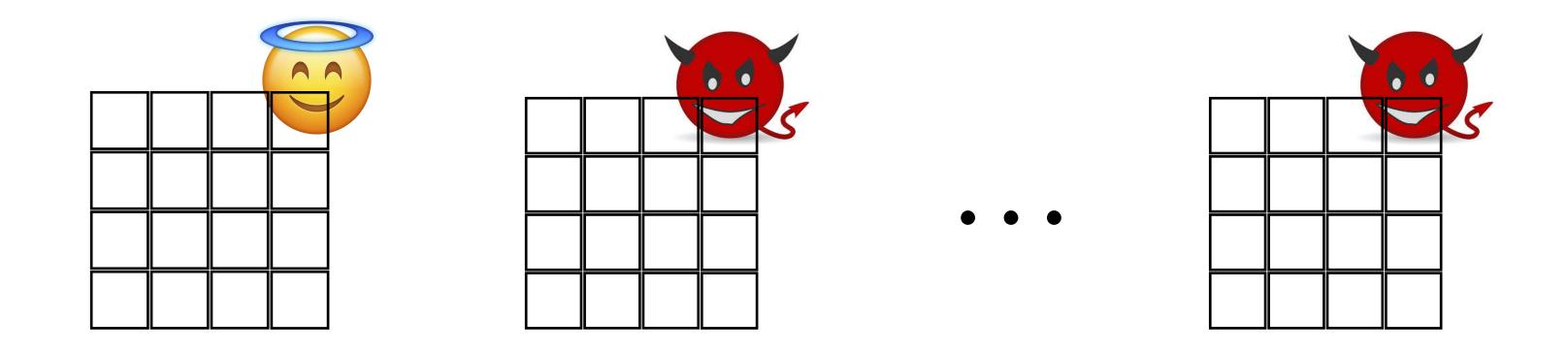
Selective-failure attacks



- Efficiency: Total communication is sublinear in the size of the database.
- Correctness: If client and server are honest, the client recovers pk<sub>Bob</sub>.
- Privacy: The server(s) learn nothing about the content of the client's query, even if the server(s) learn whether the client aborted during reconstruction.
- Integrity: The client either outputs the authentic pk<sub>Bob</sub> or aborts, except with negligible probability.

### How to define authentic data?

Honest server's view of the database.



## Multi-server schemes

#### (1) Multi-servers, single-record query

Given a Merkle-tree scheme, on a database of size N

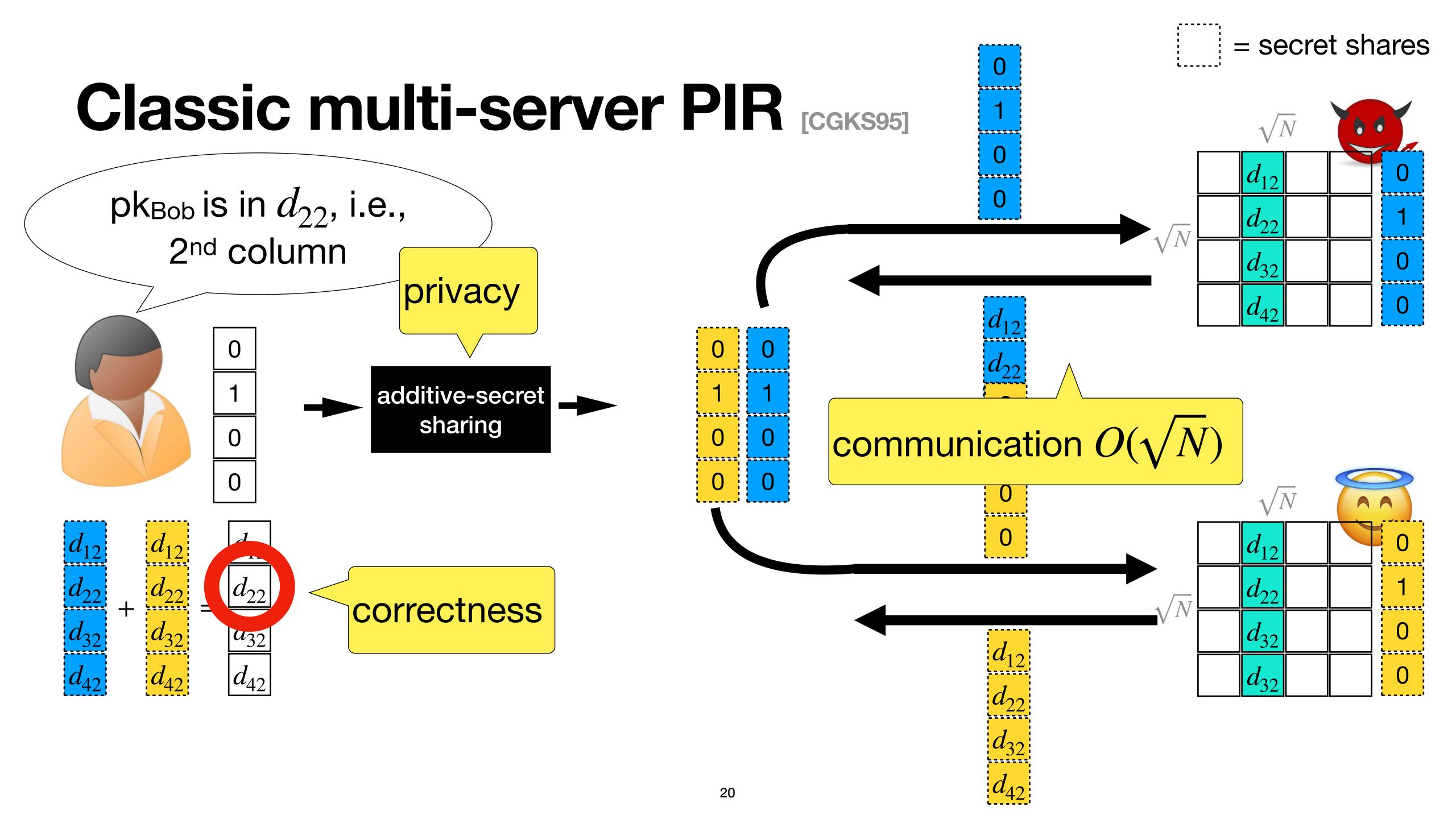
- the per-query communication is  $O(\log N)$ , same as unauthenticated PIR,
- the integrity error is negligible.

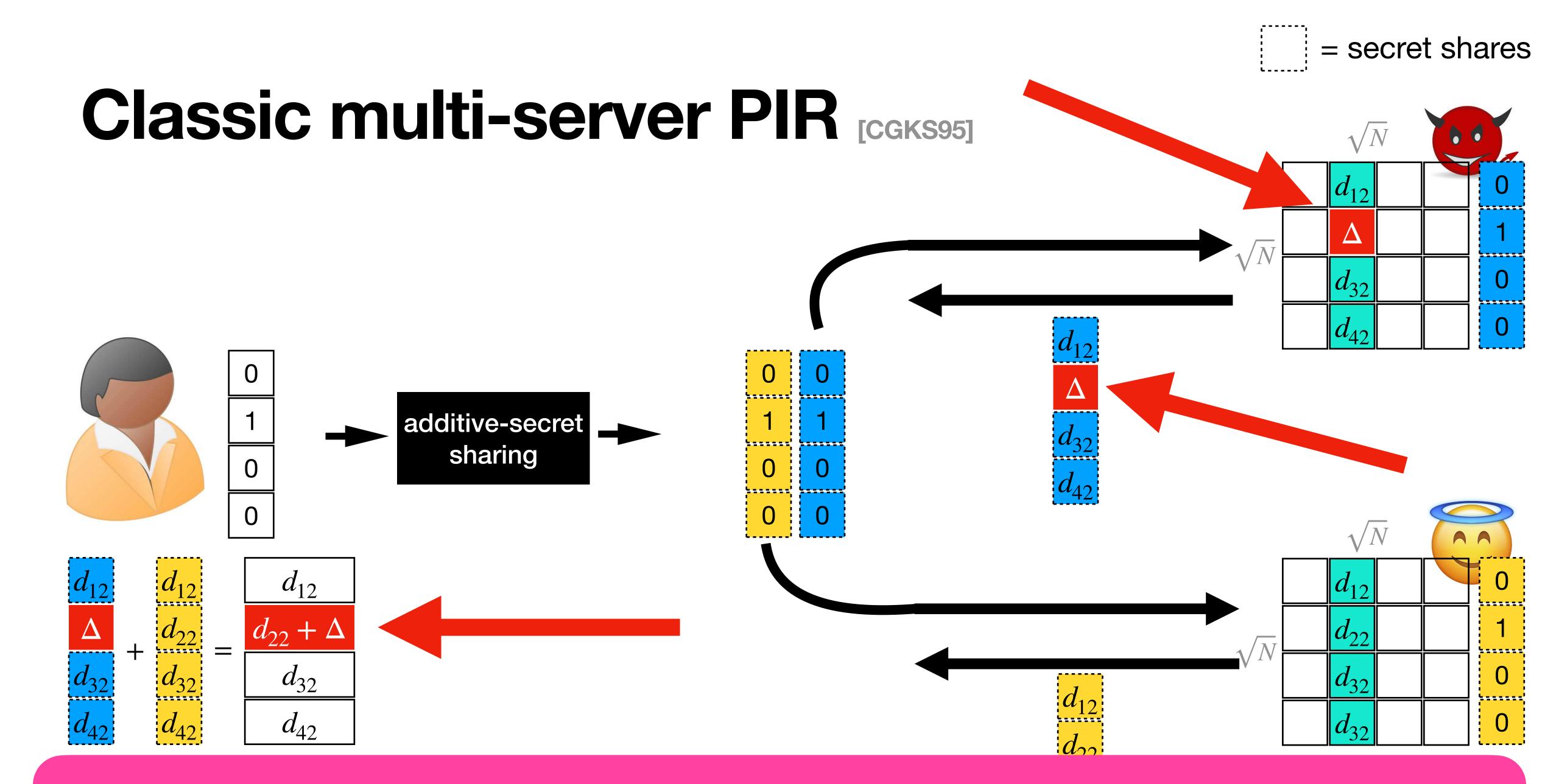
#### (2) Two-servers, single-record and aggregate queries

Given PRG and a field  $\mathbb F$ , on a database of size N

- the per-query communication is  $O(\log N)$ , same as unauthenticated PIR,
- the integrity error is  $1/|\mathbb{F}|$

This talk (roughly)



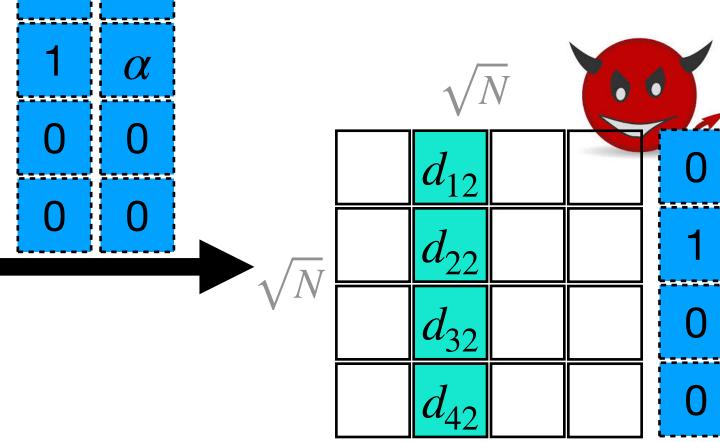


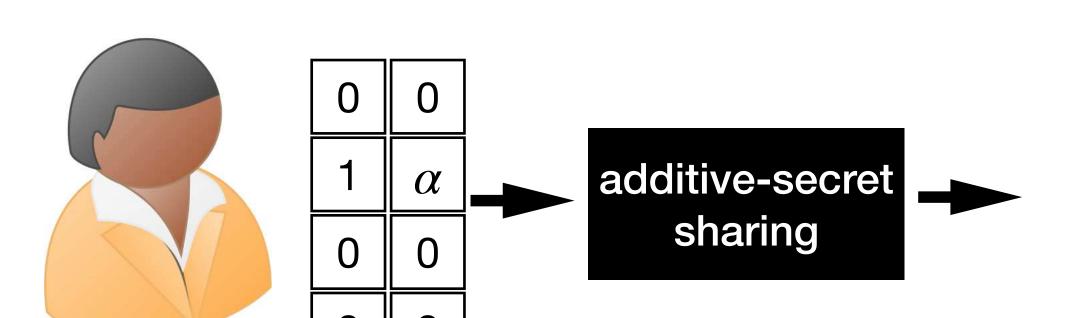
Key idea: two correlated queries, one for data and one to authenticate

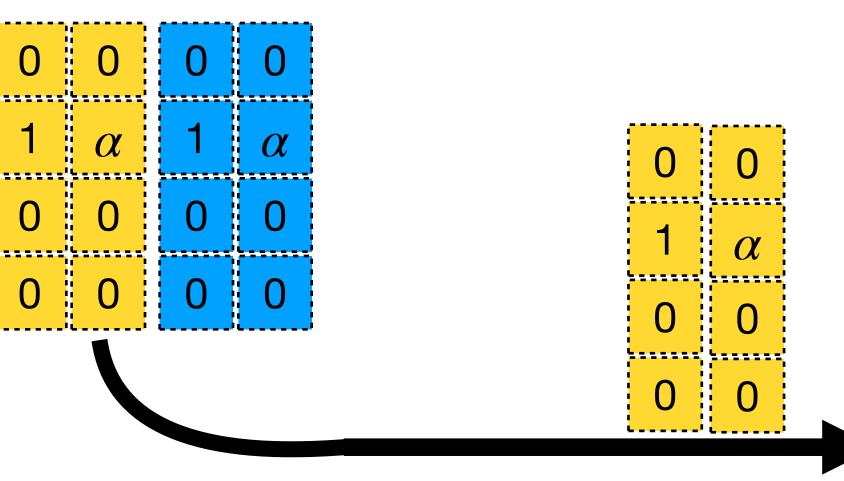
# Our contribution Authentica

samples random  $\alpha \in_R \mathbb{F}$ 

## Authenticated multi-server PIR

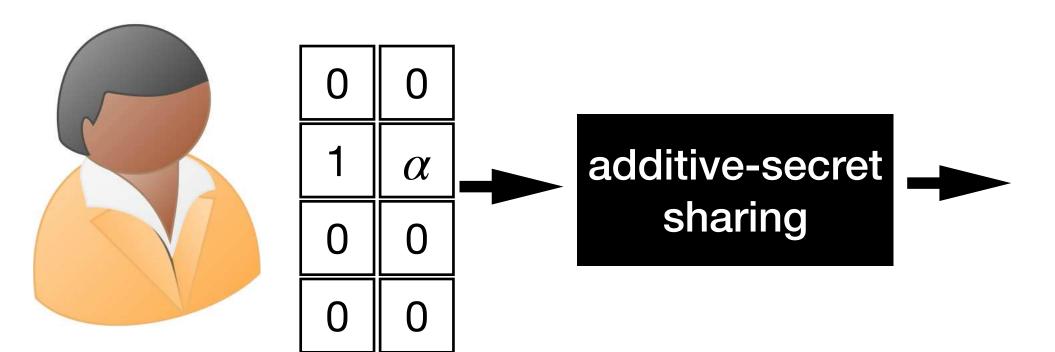


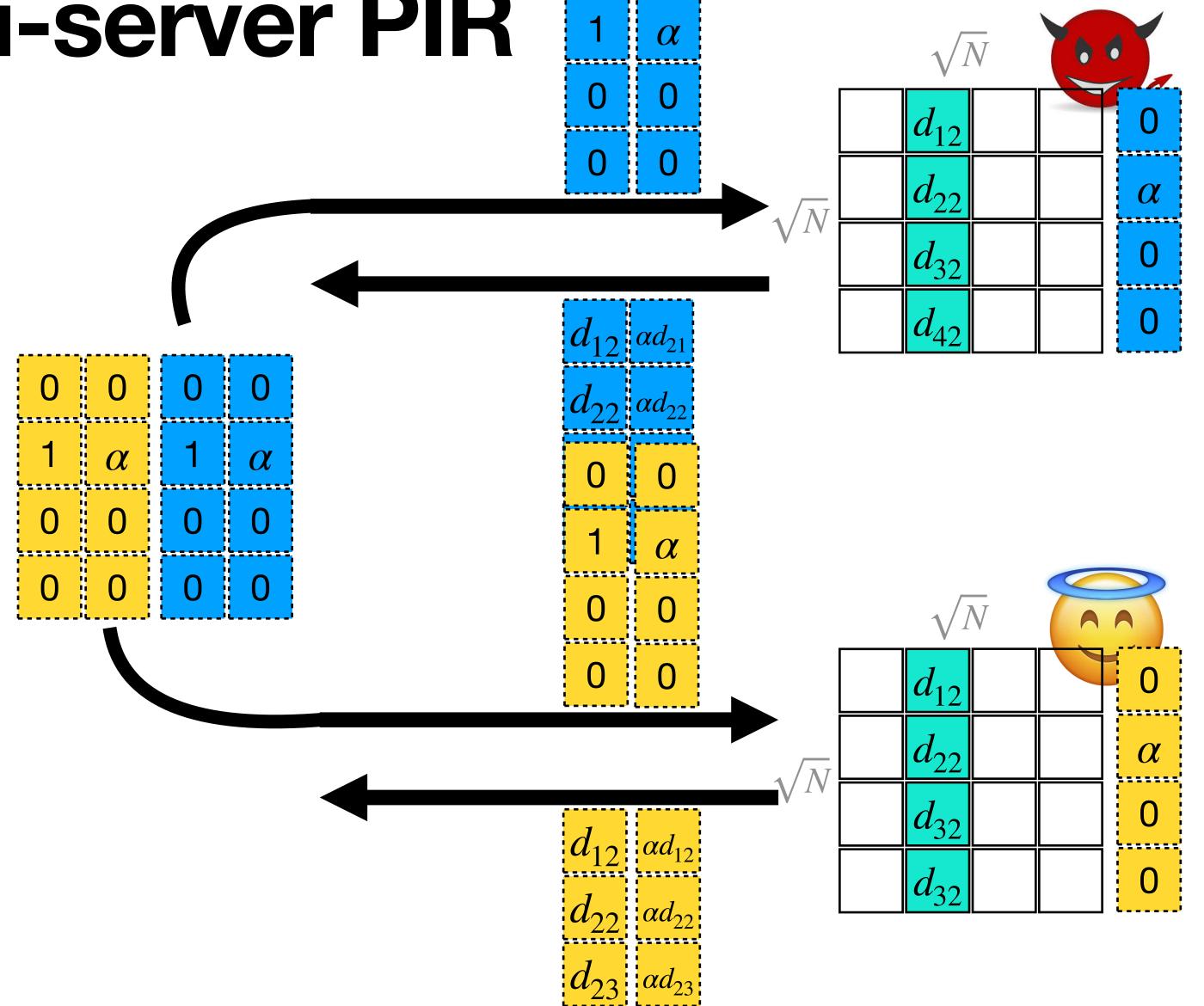




## Authenticated multi-server PIR

samples random  $\alpha \in_R \mathbb{F}$ 





# Authenticated multi-server PIR integrity

$$if \ lpha \cdot \left( egin{array}{c|c} d_{12} & d_{12} \ d_{22} \ d_{32} \ d_{42} \ \end{array} 
ight) = egin{array}{c|c} lpha d_{12} & lpha d_{12} \ lpha d_{22} \ lpha d_{32} \ \end{array} + egin{array}{c|c} lpha d_{22} \ ad_{32} \ d_{42} \ \end{array} 
ight) = egin{array}{c|c} lpha d_{22} \ lpha d_{32} \ \end{array} + egin{array}{c|c} lpha d_{22} \ ad_{32} \ \end{array} 
ight)$$

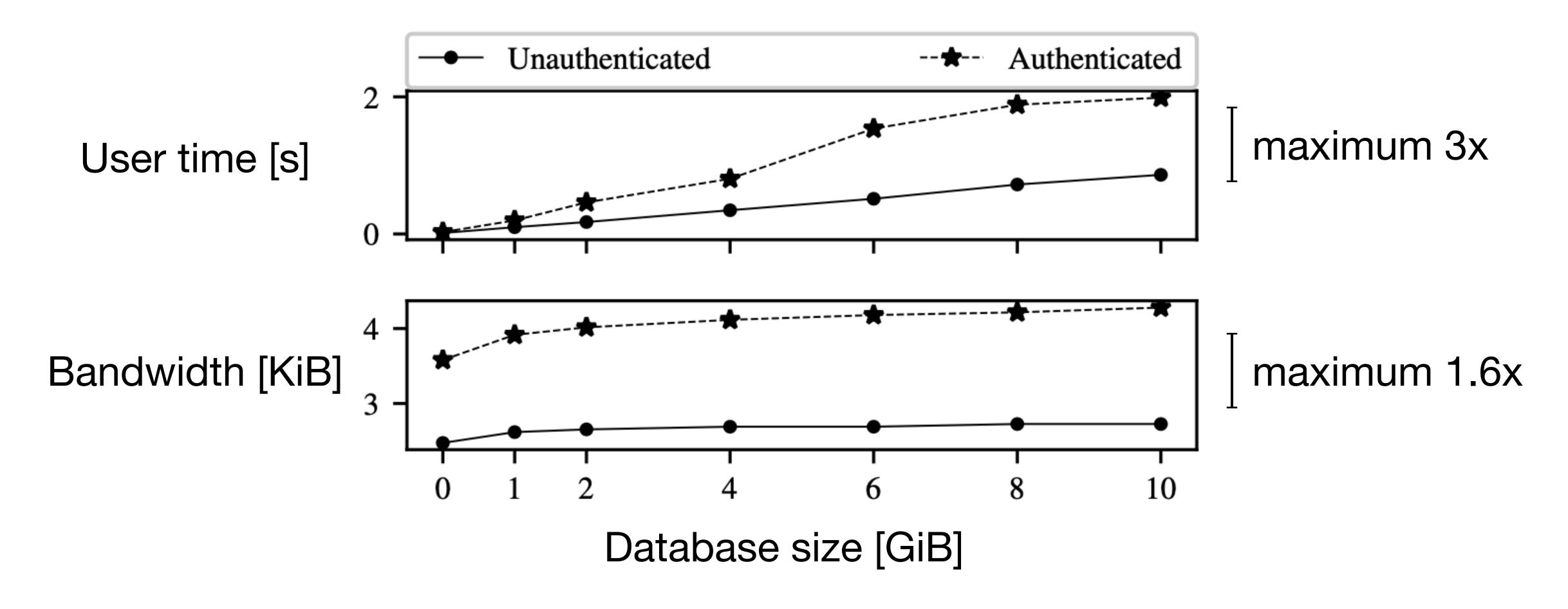


return second element of

else abort

communication  $O(\sqrt{N})$ , function secret sharing reduces to  $O(\lambda \log N)$  [BGI16]

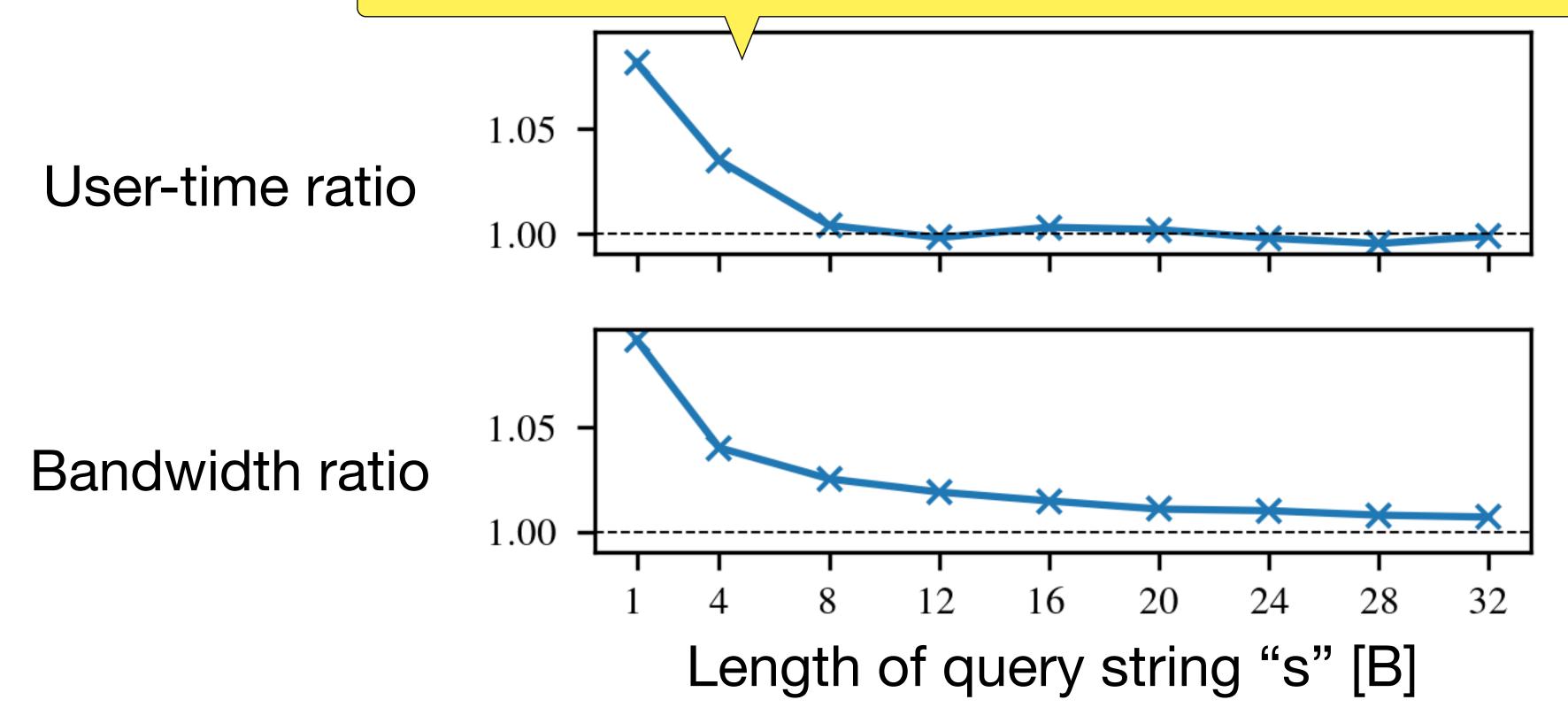
# Evaluation: single-record queries (Merkle)



Cost of retrieving a 1KiB record

## Evaluation: aggregate queries

ratio of authenticated and classic unauthenticated PIR



SELECT COUNT(\*) FROM keys WHERE email LIKE "%s"

Count emails that end with string "s"

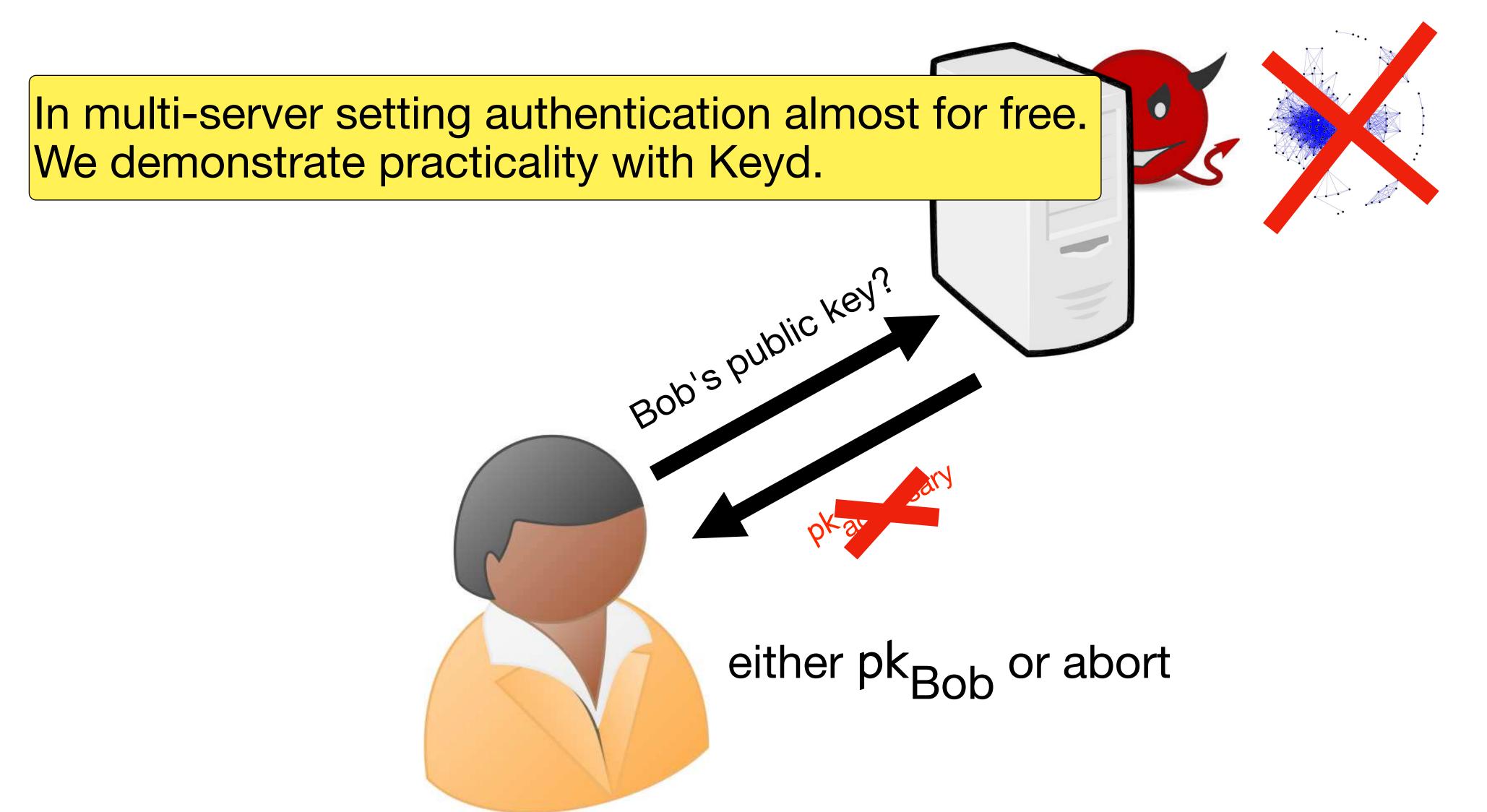


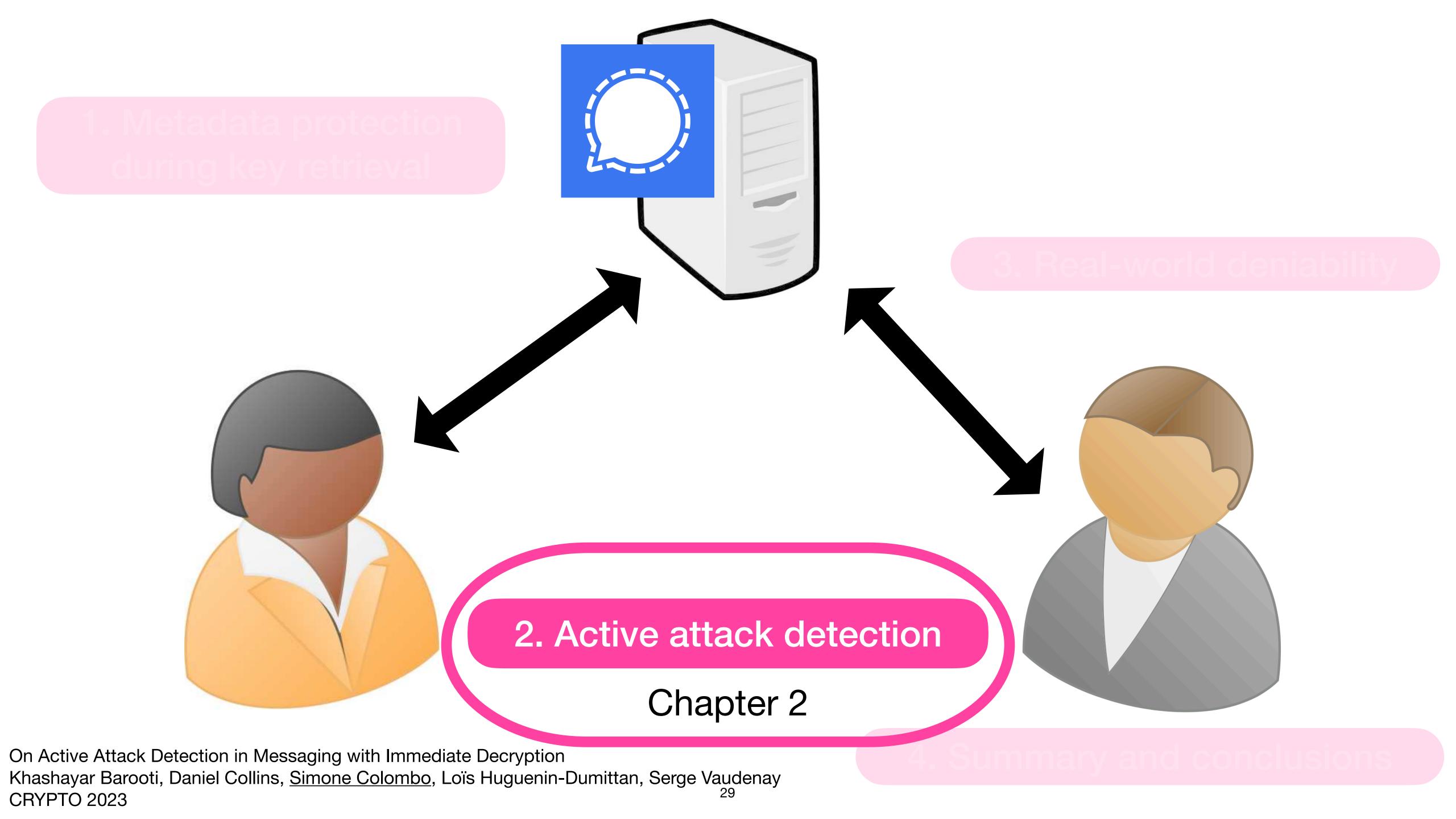


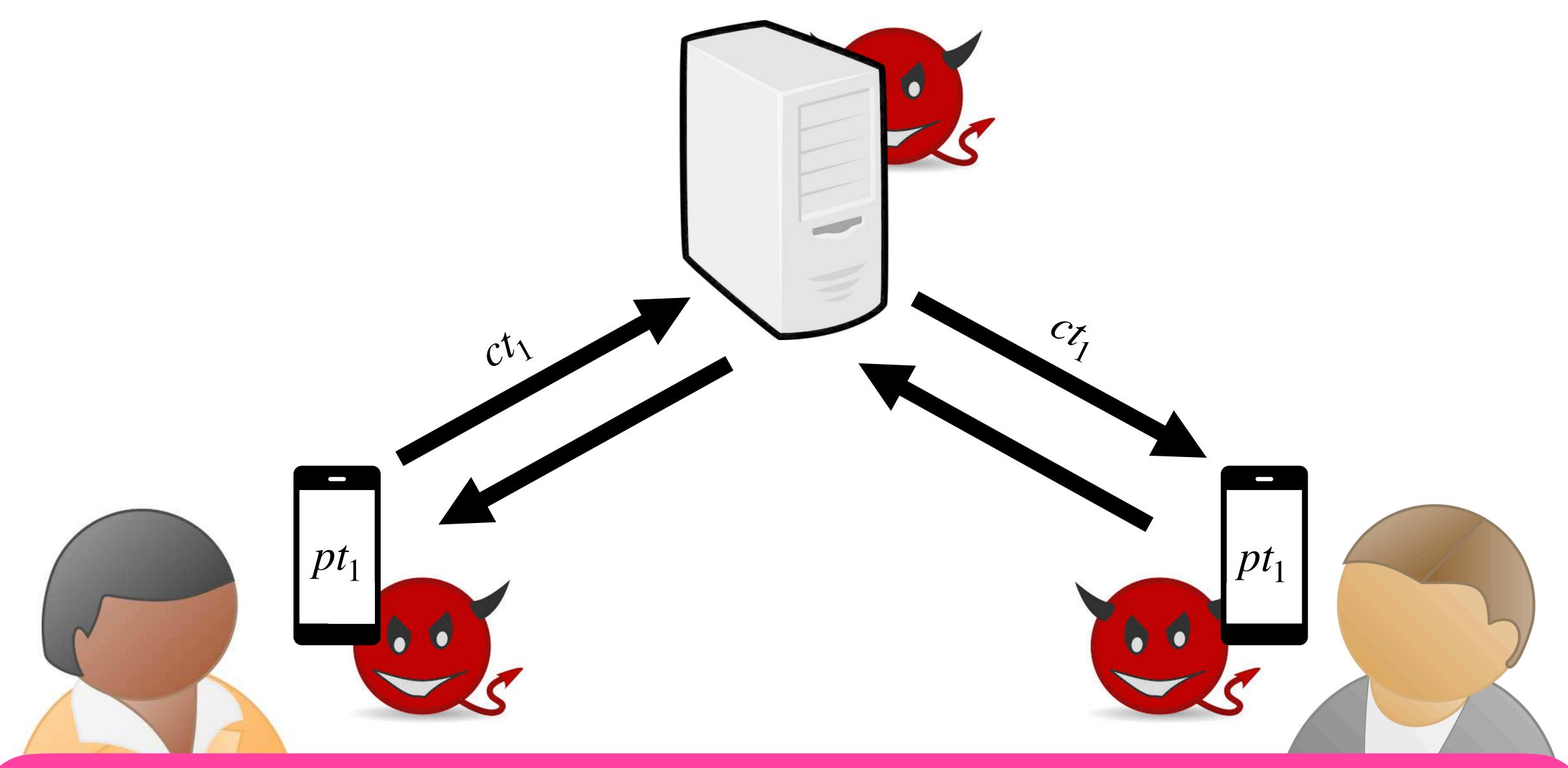




## Summary of metadata protection during key retrieval



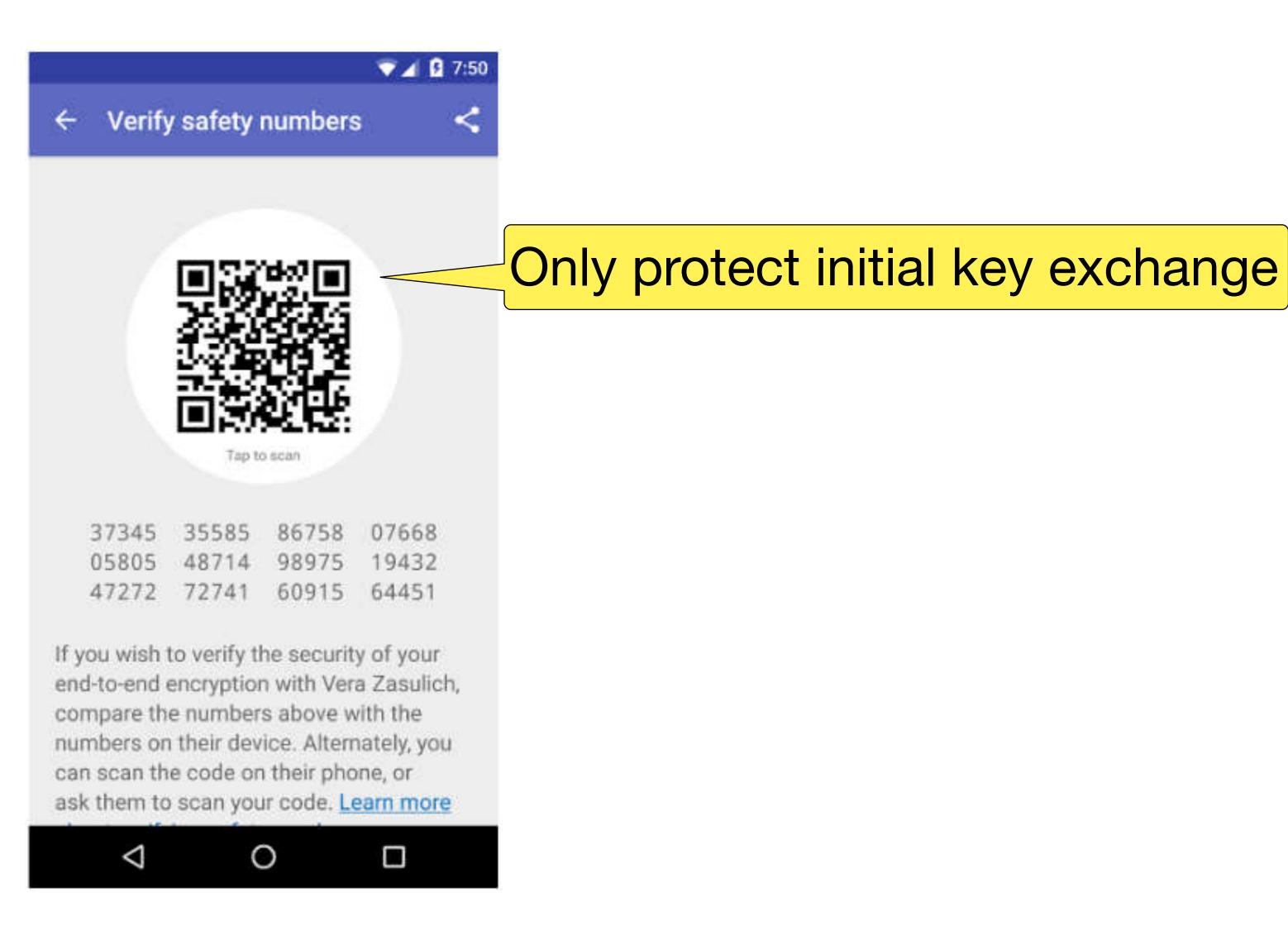




Forward secrecy and post-compromise security protect against active attacks.

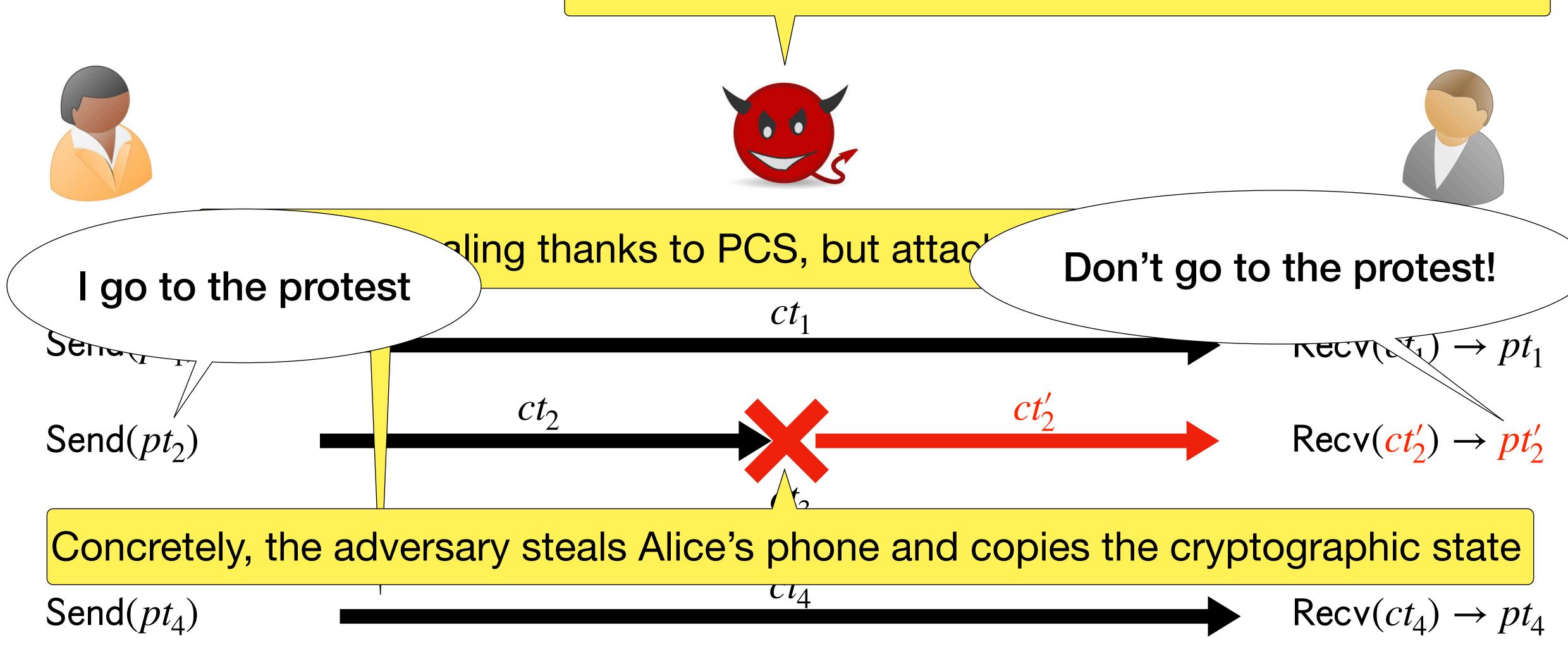
What about detection of these attacks?

# Signal's safety numbers



## Active attacks

In our model the adversary can expose states, control randomness and invoke algorithms via oracles



## Immediate decryption (ID) [ACD19]

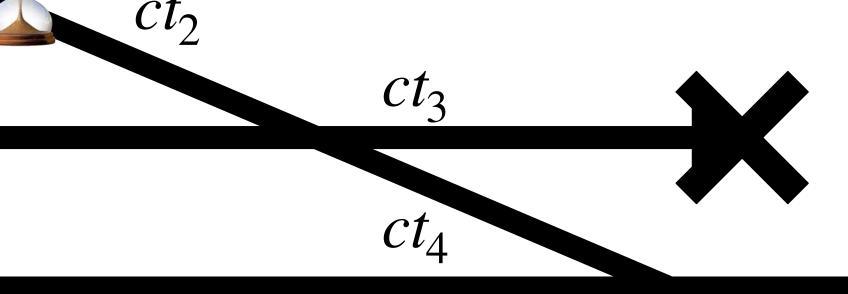
Schemes support out-of-order delivery and message loss at the protocol level.

Immediate decryption requires an ordinal for each sent and received message  $Send(pt_1) \rightarrow (num_1, ct_1)$ ightharpoonup Recv $(ct_1) \rightarrow (num_1, pt_1)$ 

 $Send(pt_2) \rightarrow (num_2, ct_2)$ 

 $Send(pt_3) \rightarrow (num_3, ct_3)$ 

 $Send(pt_4) \rightarrow (num_4, ct_4)$ 



 $Recv(ct_4) \rightarrow (num_4, pt_4)$ 

 $Recv(ct_2) \rightarrow (num_2, pt_2)$ 

# Immediate decryption (ID) [ACD19]

Schemes support out-of-order delivery and message loss at the protocol level.

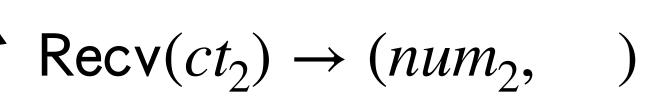


Notation simplified for the rest of the talk

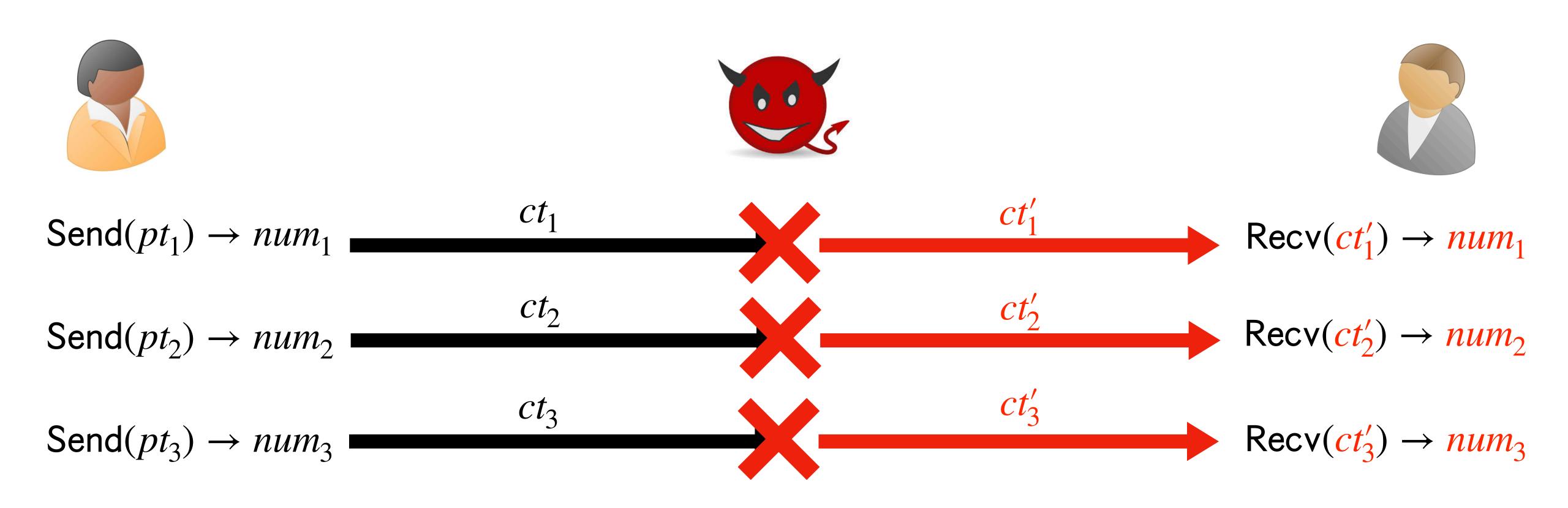


Our contribution: active attack detection with immediate decryption

- New authenticated ratcheted communication primitive.
- In-band detection with immediate decryption with r-RID and s-RID notions.
- Out-of-band detection with ID with new r-UNF and s-UNF security notions.
- Optimisations for s-RID and s-UNF security towards practicality.

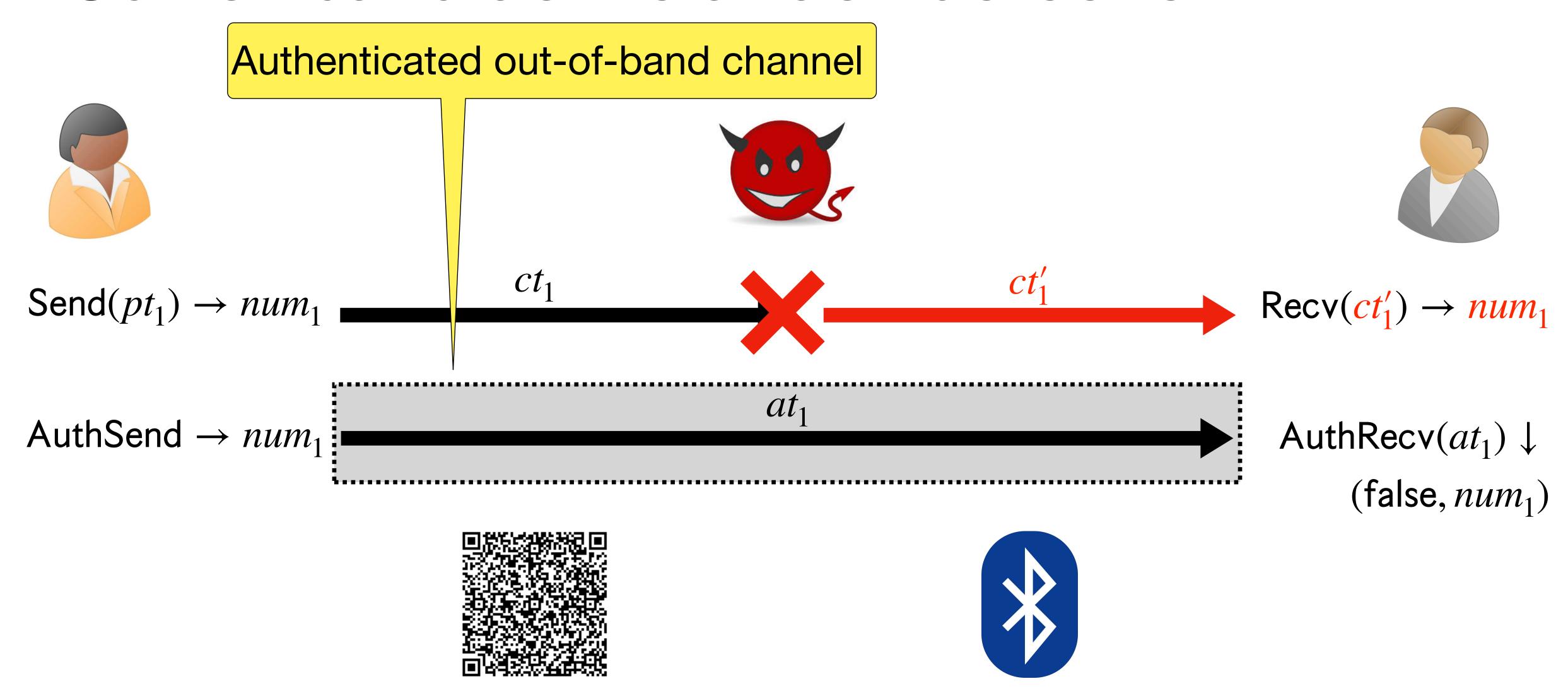


## Out-of-band active attack detection

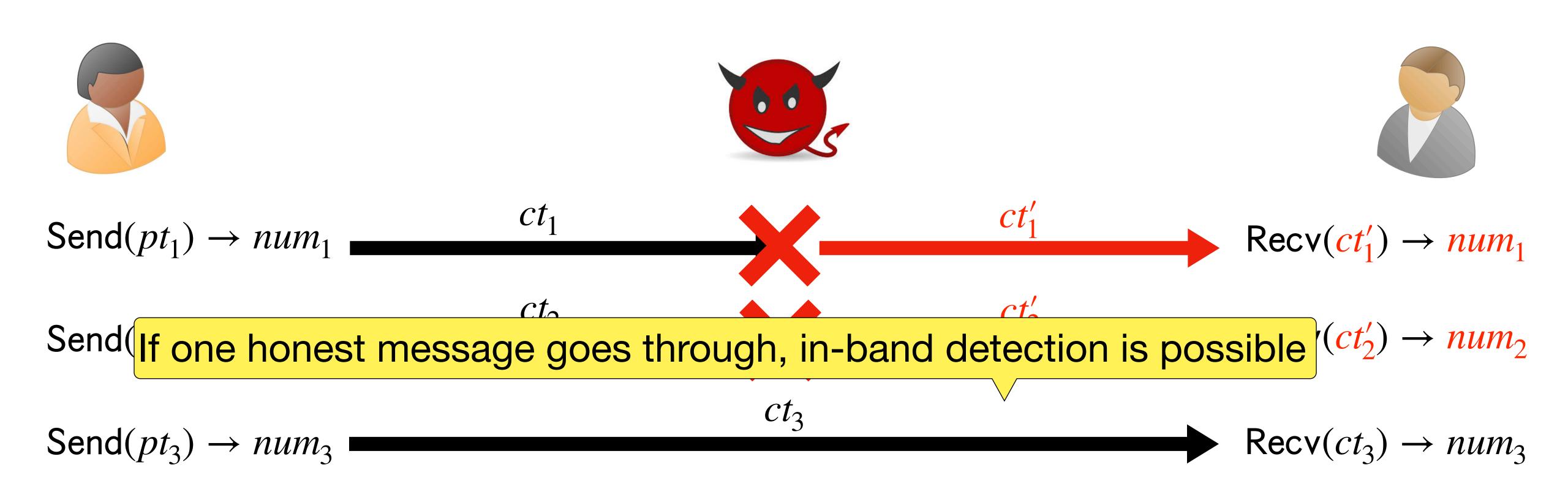


If the adversary blocks all messages, we must use an out-of-band channel

## Out-of-band active attack detection



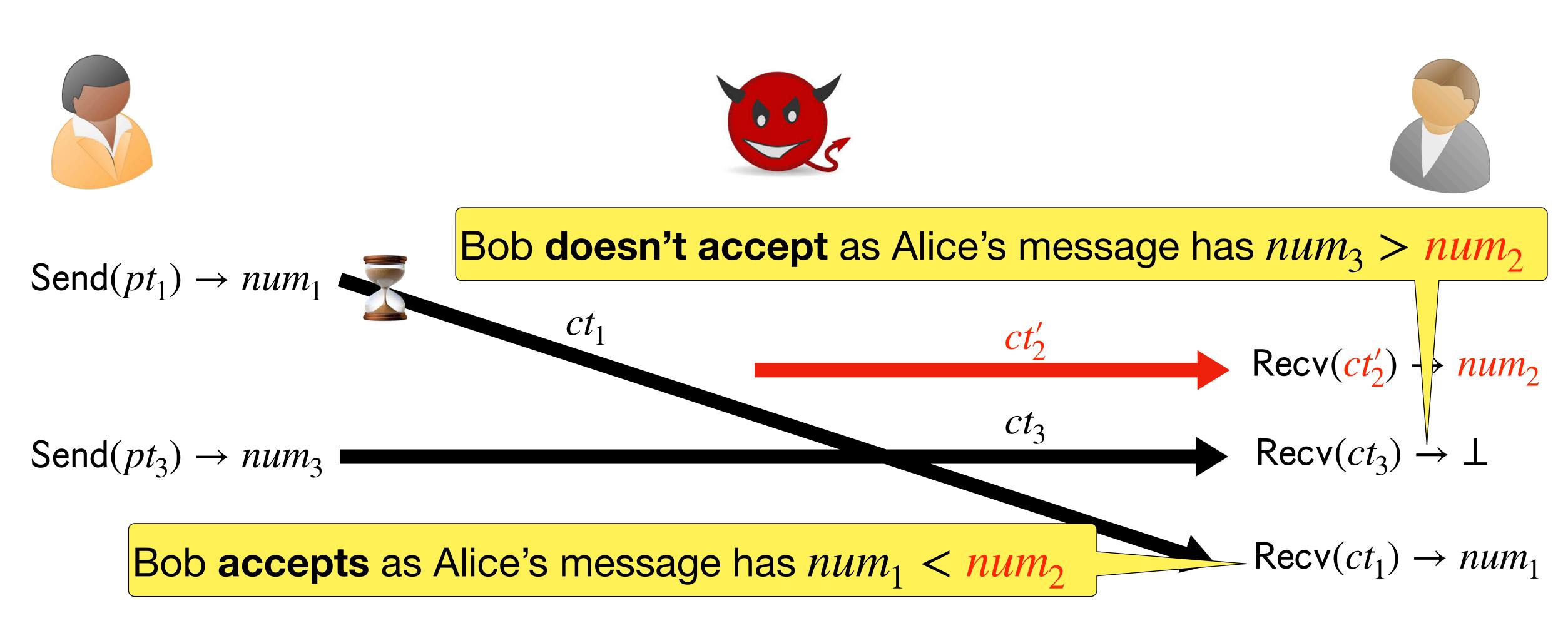
#### In-band active attack detection



We define in-band active attack detection through r-RID and s-RID security

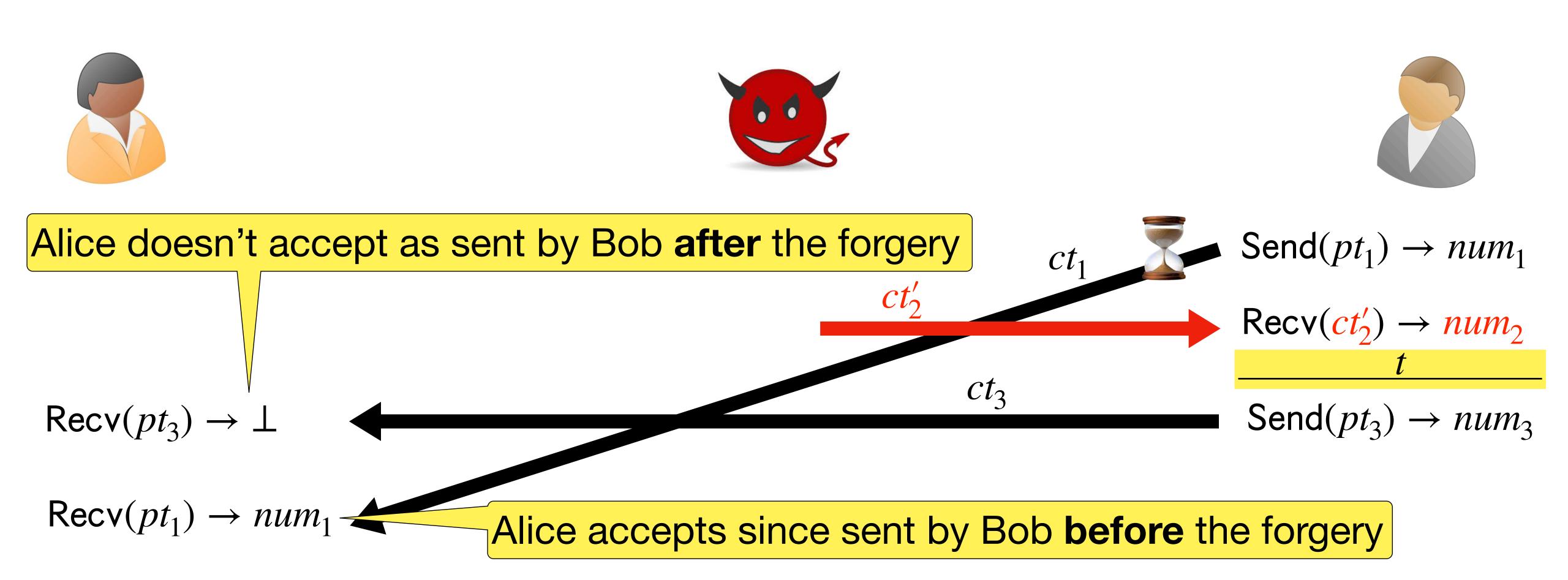
#### r-RID security

If Bob receives a forgery with num, then Bob rejects all Alice's messages with num' > num.



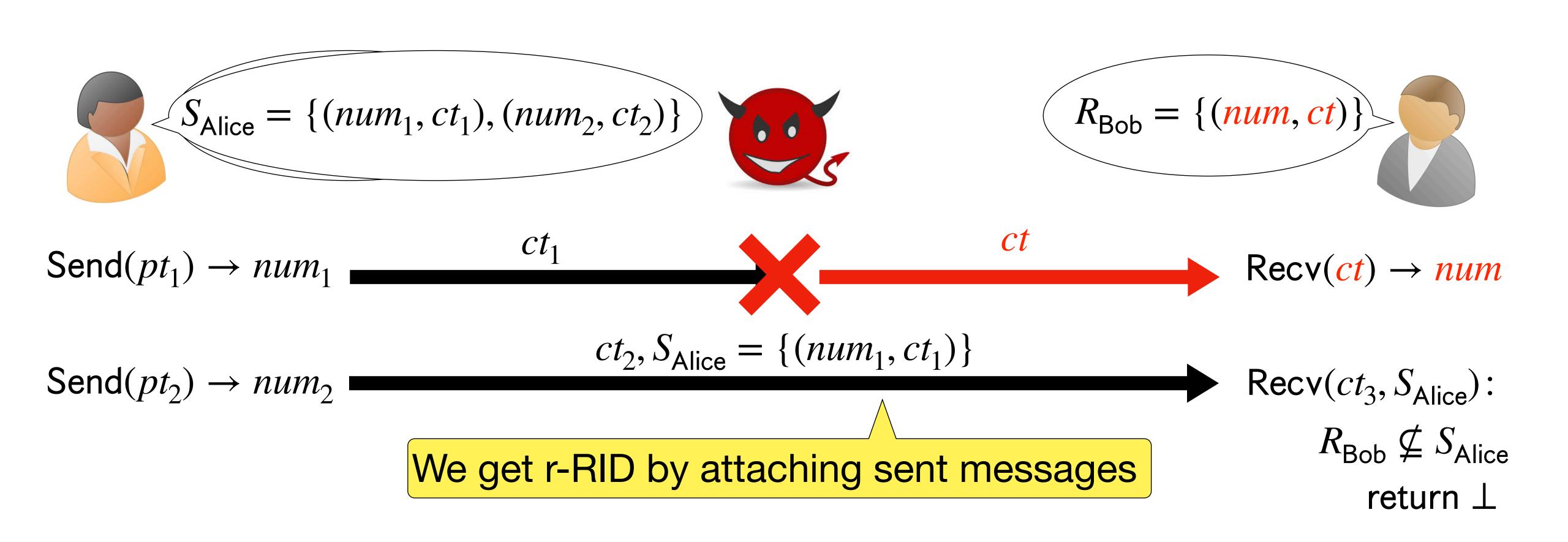
#### s-RID security

If Bob receives a forgery with num at time t, then Alice rejects all Bob's messages sent after time t.



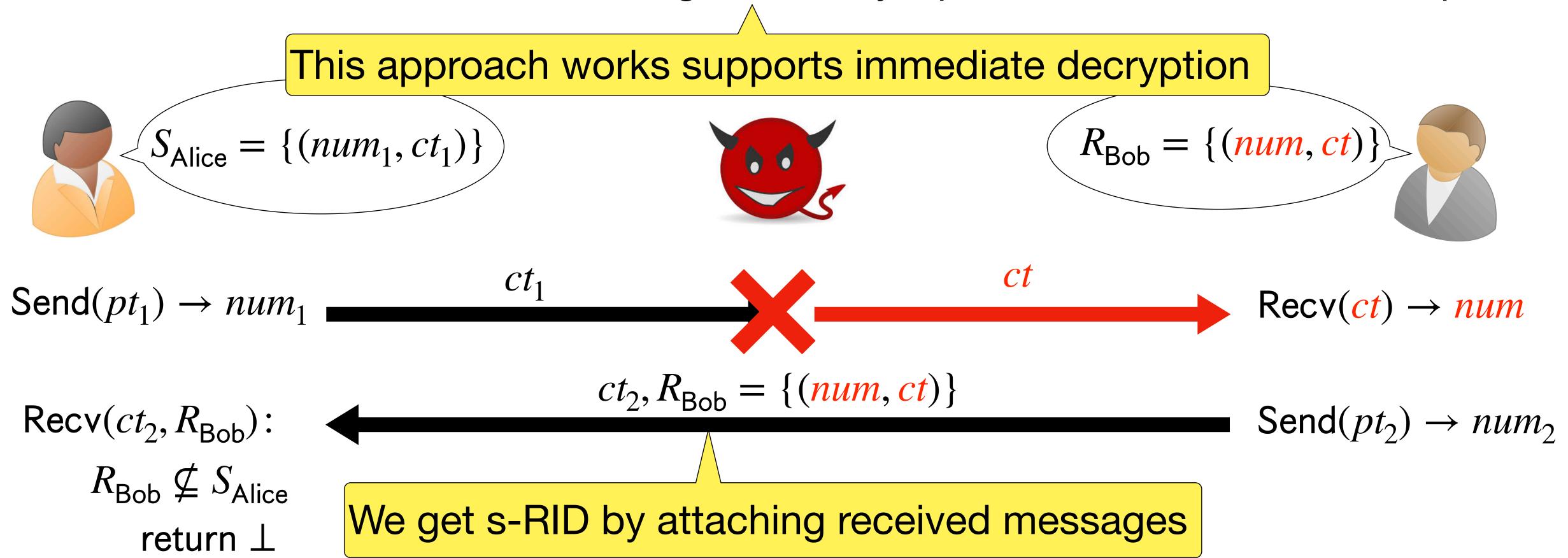
### A simple RID construction (r-RID + s-RID)

Attach all sent and received ciphertexts to every ciphertext and check at reception.



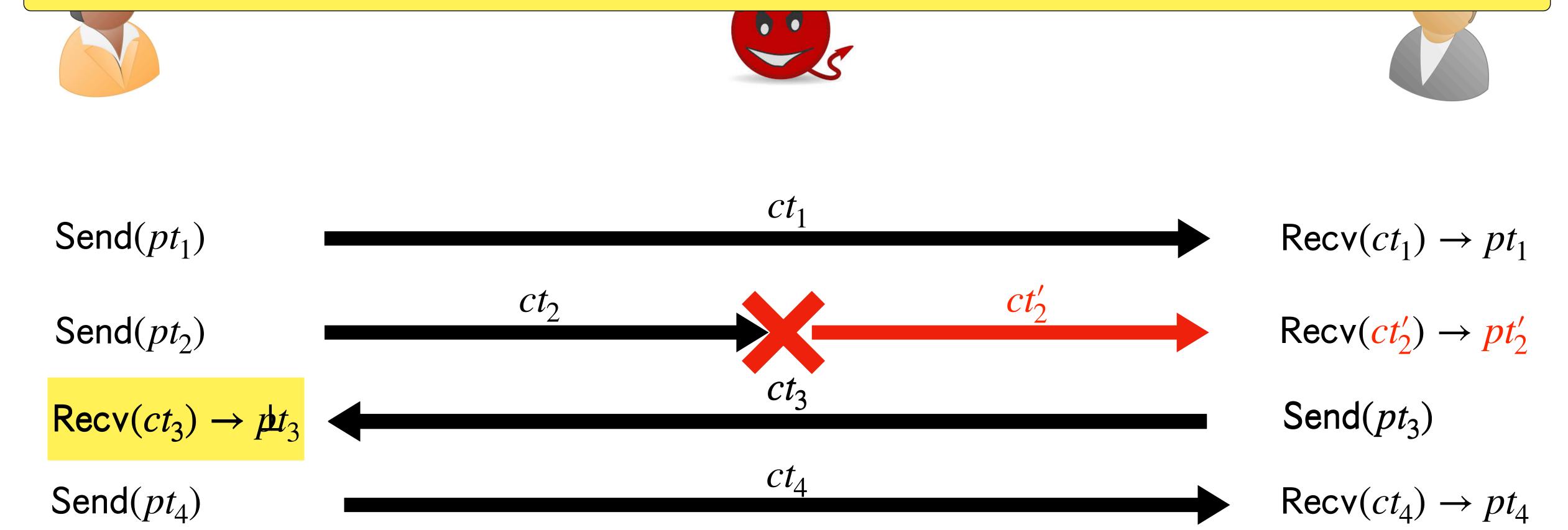
## A simple RID construction (r-RID + s-RID)

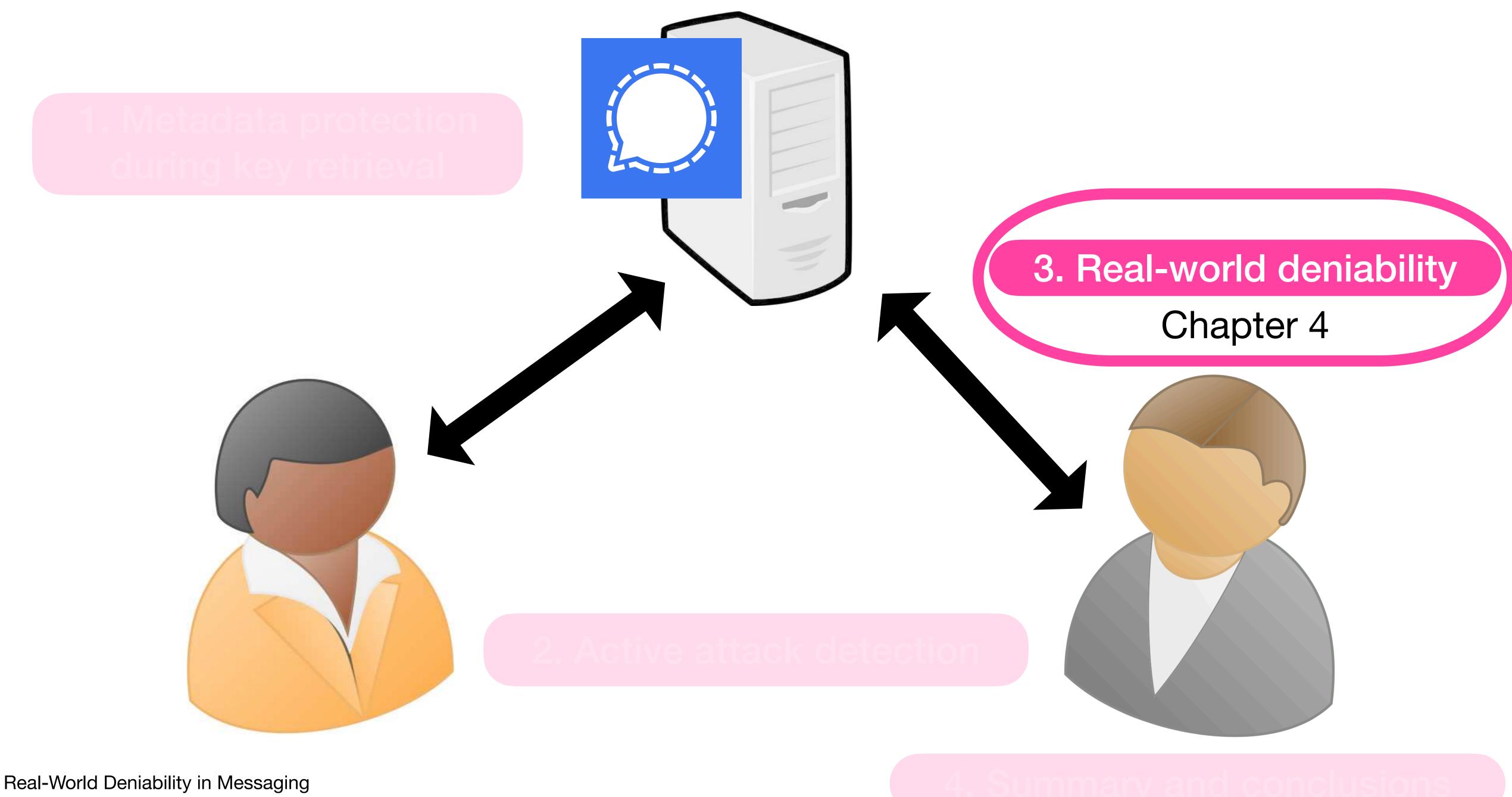
Attach all sent and received messages to every ciphertexts and check at reception.



#### Summary of active attack detection

If adversary blocks in-band channel, parties can use out-of-band authenticated one. s-RID is likely practical and we propose optimisations to reduce overhead.





Real-World Deniability in Messaging Daniel Collins, <u>Simone Colombo</u>, Loïs Huguenin-Dumittan PETS 2025, RWC 2023





Let's go to the protest!

Our contribution: analysis of deniability's impracticality, and solutions to achieve real-world deniability.

- Model to analyze real-world deniability in messaging.
- Technical case studies: Signal application and DKIM protected email.
- Legal case study: 140 Swiss court cases that use WhatsApp as evidence.
- Discussion on deniability and on how to achieve real-world deniability.

#### Technical case study: Signal

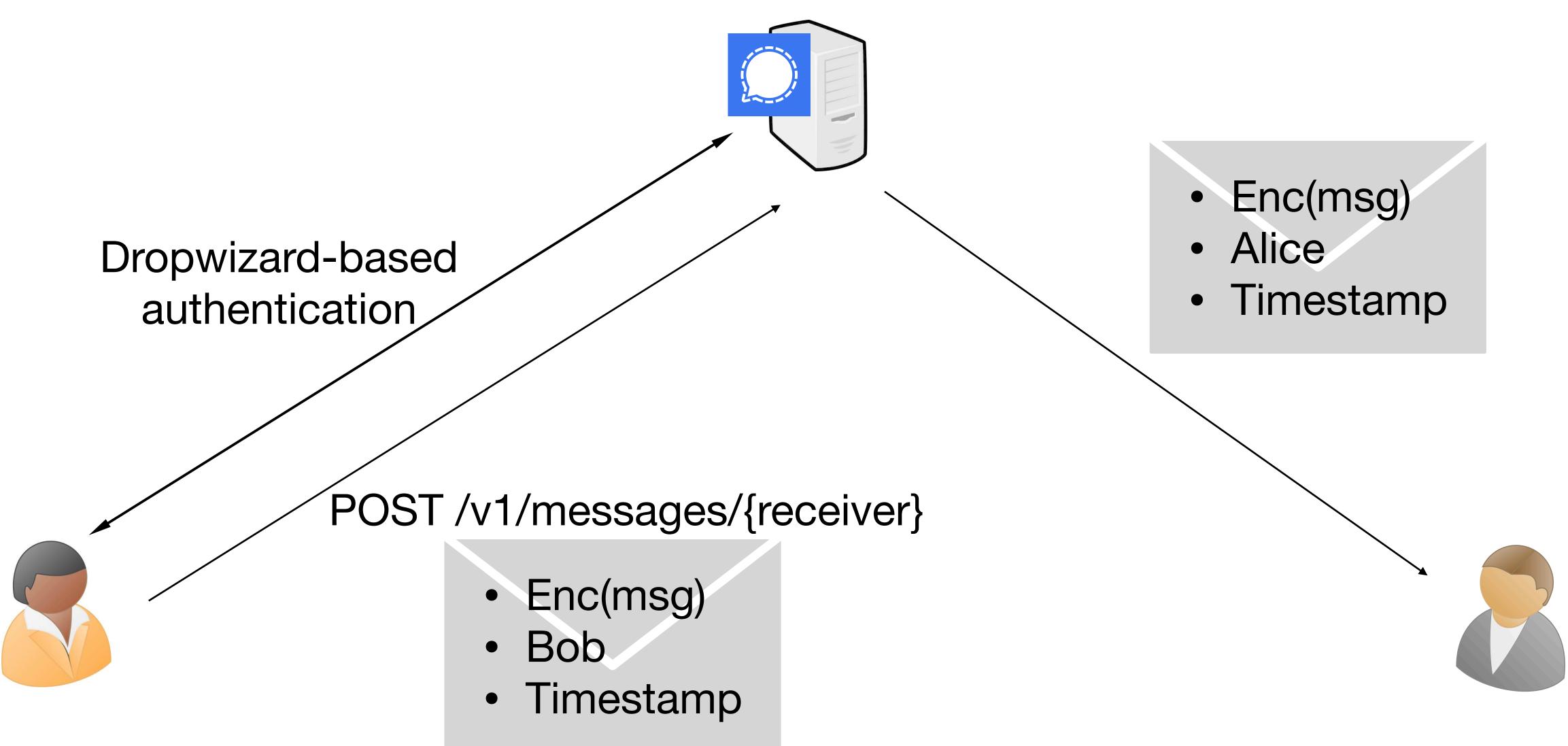
Signal claims to provide deniability and recent works show it achieves some form of cryptographic deniability [VGIK20, FJ24].

# On the Cryptographic Deniability of the Signal Protocol

Nihal Vatandas<sup>1</sup>, Rosario Gennaro<sup>1</sup>, Bertrand Ithurburn<sup>1</sup>, and Hugo Krawczyk<sup>2</sup>

#### Is this sufficient in practice?

#### Signal with classic authentication



#### Classic authentication hinders deniability



If Bob's phone contains Alice's message, then < If server logs, even worse

- either Alice really sent it after authenticating with the server, or
- Bob modified the local message database.

Signal is technically undeniable unless Bob knows how to tamper with the device.

What about the legal impact of deniability?

#### Legal case study methodology

- Manual analysis of 341 penal cases in Switzerland that mention "WhatsApp".
- Research questions:

  No mention of Signal in cases
  - Do judges in Swiss courts use WhatsApp as evidence?
  - When they do, is their usage contested by any of the parties involved?
  - What are the reasons used to dispute the legal validity of such messages?
  - How do judges respond to these disputes?

#### Legal case study results

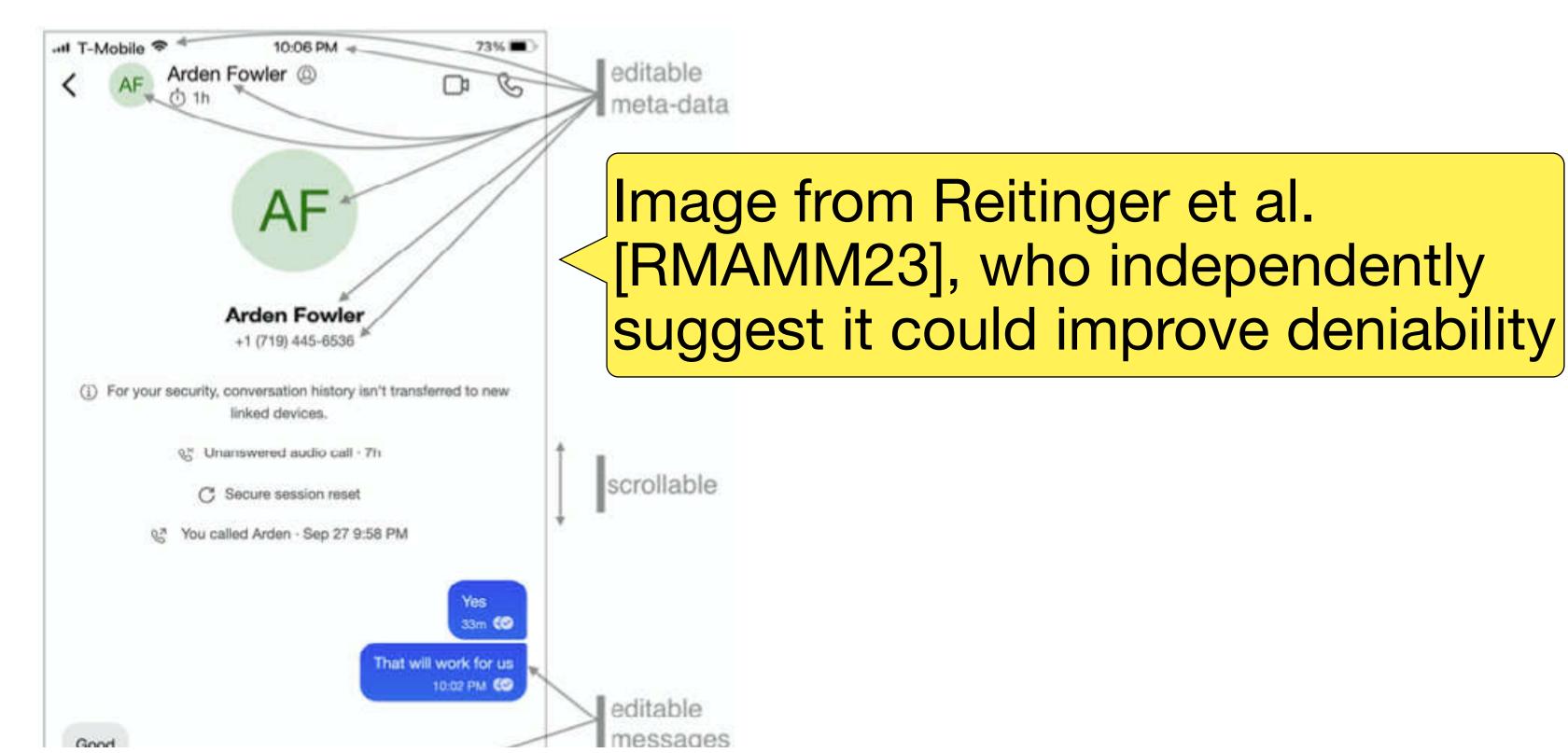
Deniability is not invoked in these two cases

Total cases	N/A	Evidence	Contested	Rejected
341	201 (59%)	140 (41%)	2 (0.6%)	0
341	201 (3970)	140 (4170)	2 (0.070)	U

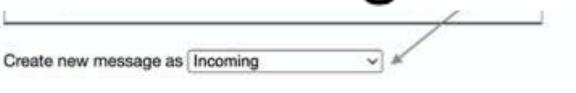
Yadav et al. [YGS23] reach similar results in an analysis of US court cases

Cryptographic deniability fails technically and (likely) legally: what to do?

#### A possible solution



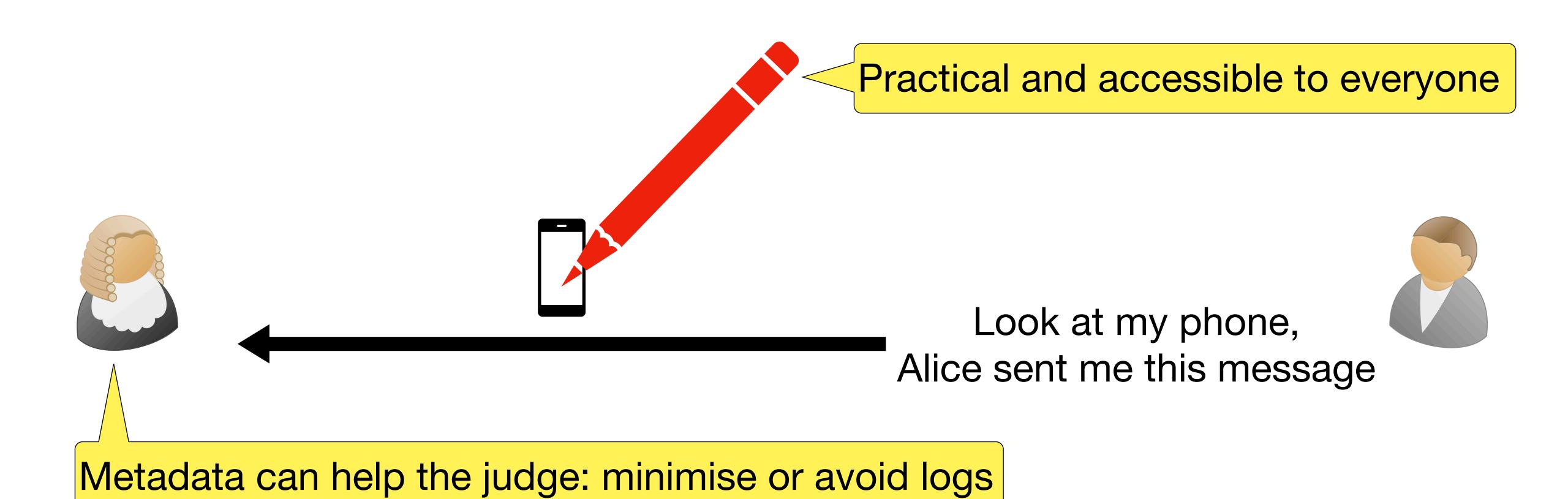
- either Alice really sent it after authenticating with the server, or
- Bob modified the local message database.

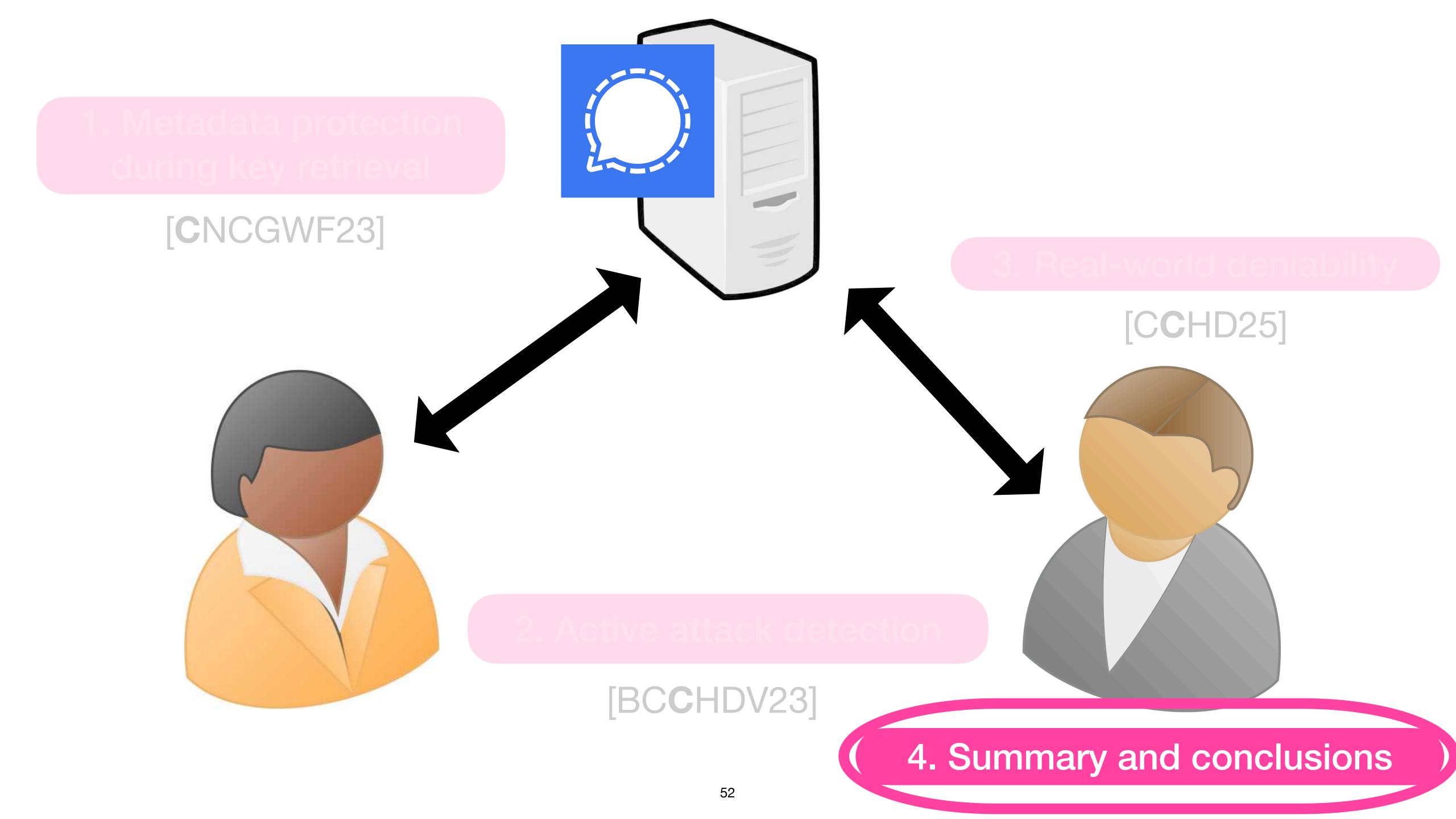




#### Summary of real-world deniability

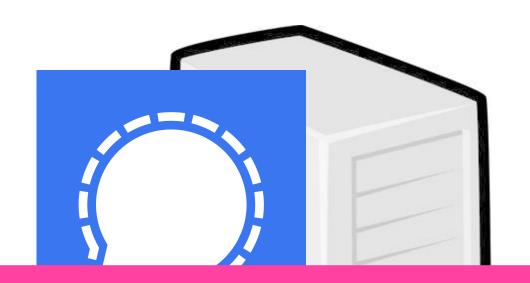
As well as available and functional





### Summary of contributions

	Before	After
Metadata protection	Classic PIR without integrity protection	Authenticated PIR provides privacy and integrity
Active attack detection	Detection without immediate decryption, partial detection or additional rounds.	Emcient detection while supporting immediate decryption
Deniability	Only cryptographic deniability	Cryptographic deniability fails technically and legally: local messages modification can work



#### Technical real-world integration

