# Impact Tracing: Identifying the Culprit of Misinformation in Encrypted Messaging Systems

**Zhongming Wang**, Tao Xiang, Xiaoguo Li, Biwen Chen, Guomin Yang, Chuan Ma, and Robert H. Deng
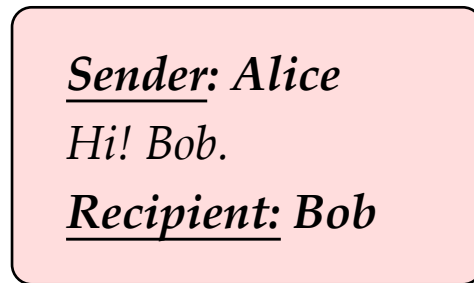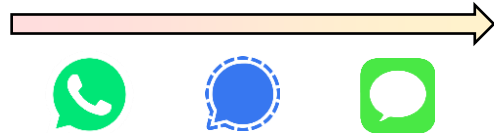
NDSS Symposium 2025

CHONGQING UNIVERSITY

SMU
SINGAPORE MANAGEMENT
UNIVERSITY

# Content Moderation for EEMSs

End-to-end encrypted messaging systems (EEMS):

Only the END users can read the messages.

# Content Moderation for EEMSs

End-to-end encrypted messaging systems (EEMS):

Only the END users can read the messages.

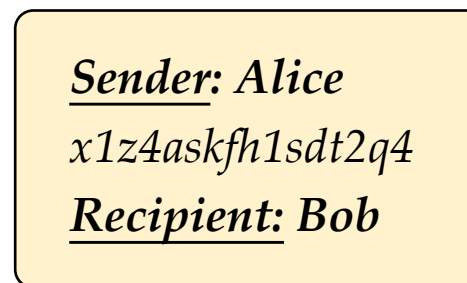**Transport level security**

*Sender: Alice*

*Hi! Bob.*

*Recipient: Bob*

**Platform**

**End-to-end security**

*Sender: Alice*

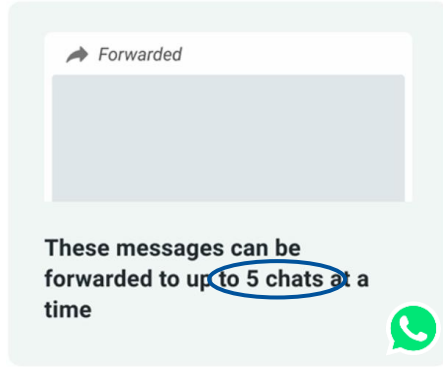*x1z4askfh1sdt2q4*

*Recipient: Bob*

?

Problematic messages proliferates in EEMSs.

End-to-end encryption obstructs content moderation :(

# Message Forwarding & Tracing

Misinformation propagate rapidly through forwarding.

# Message Forwarding & Tracing

Misinformation propagate rapidly through forwarding.



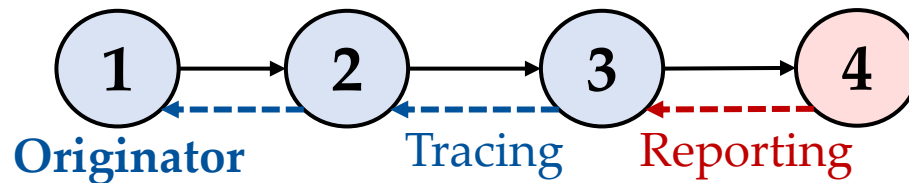Forward limit

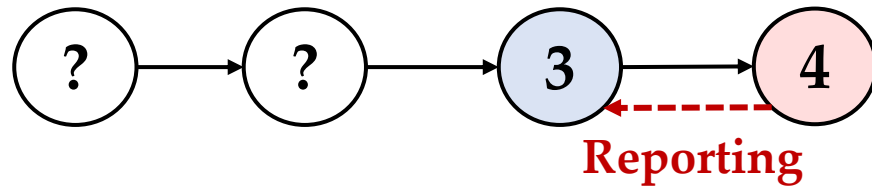# Message Forwarding & Tracing

Misinformation propagate rapidly through forwarding.



Forward limit

**Traceability** enables <u>tracing</u> after user <u>reporting</u>.

- The platform can disclose the dissemination path.



Originator    Tracing    Reporting

# Tracing Policies



**Message Franking**
[Facebook17,RWC], [GLR17,CRYPTO],
[TGL+19,CRYPTO], [GPE25,NDSS], etc.

**Message Traceback**
[TMR19,CCS], [KTW22,ESORICS]

**Source Tracing**
[PEB21,CCS], [IAV22,Usenix Sec.], [LRTY22,NDSS],
[BGJP23,EUROCRYPT], [BE24,PETS], etc.

Which part of the dissemination path is tracked during tracing?

# Traceability vs. Privacy

- **Traceability:** Reveal the culprits of spreading misinformation.
- **Privacy:** Reveal nothing about forwarding path.

# Traceability vs. Privacy

- **Traceability:** Reveal the culprits of spreading misinformation.

- **Privacy:** Reveal nothing about forwarding path.

# Traceability vs. Privacy

- **Traceability:** Reveal the culprits of spreading misinformation.

- **Privacy:** Reveal nothing about forwarding path.

# Traceability vs. Privacy

- **Traceability:** Reveal the culprits of spreading misinformation.

- **Privacy:** Reveal nothing about forwarding path.

# Traceability vs. Privacy

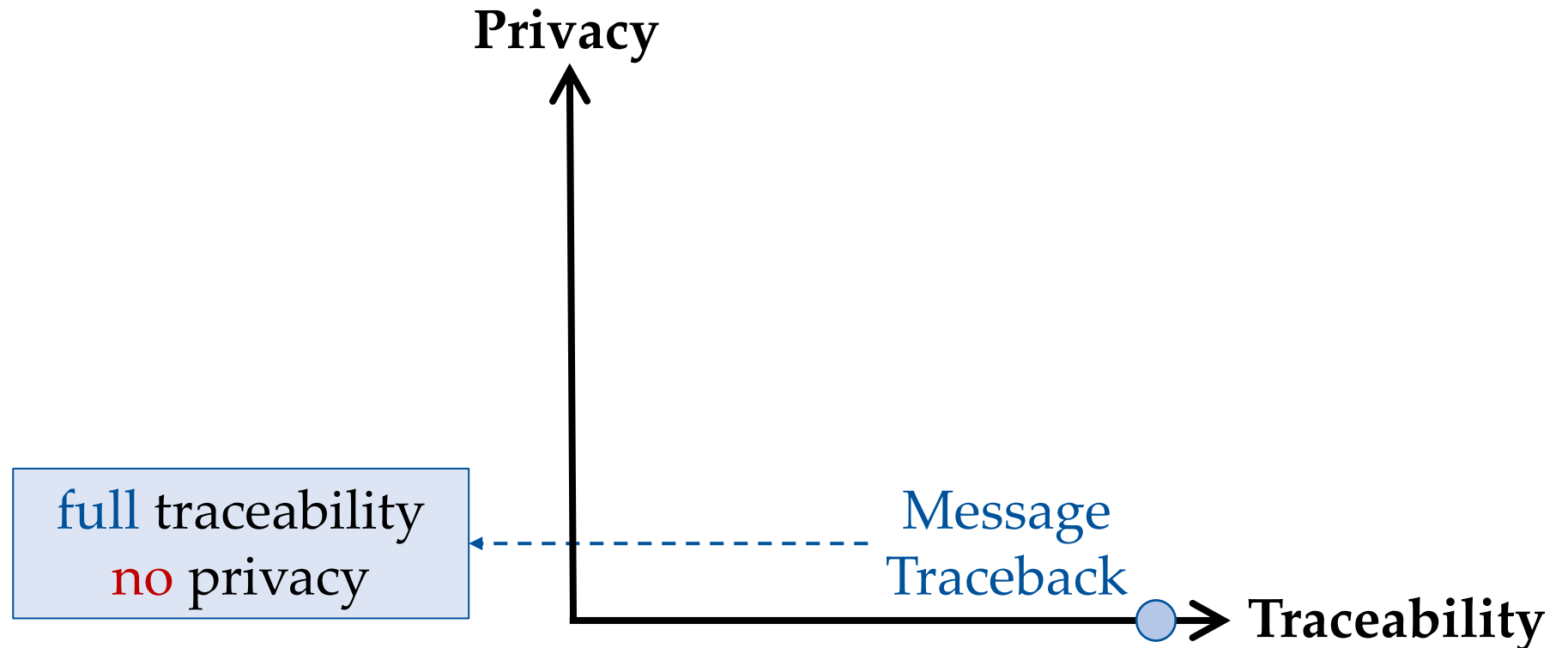Is there a <u>tracing policy</u> that <span style="color:darkred">balances traceability and privacy</span>, but also provides <span style="color:darkblue"><u>practical values</u></span> to EEMSs?

# Our Solution: Impact Tracing!

A small group of users (called influential spreaders) *significantly* contribute to spreading misinformation.



**Traceability:** The platform *can* identify <u>influential spreaders</u>.

**Privacy:** The platform *cannot* uncover <u>non-influential users</u>.

# The Design: Enabling Reporting



**Sender**

**Platform**

**Recipient**

Messaging phase

Reporting phase

# The Design: Enabling Reporting

**Sender**                                **Platform**                              **Recipient**

Compute $k \leftarrow \$ \, \mathbb{Z}_q$    $\xrightarrow{\quad E2EE(m,k),\, tag \quad}$    Store tag    $\xrightarrow{\quad E2EE(m,k),\, tag \quad}$    Check tag $=?\, F_k(m)$

$tag \leftarrow F_k(m)$

**Messaging phase**

**Reporting phase**

# The Design: Enabling Reporting



**Sender**

**Platform**

Confidentiality: Learns nothing about m.

cipient

Compute $k \leftarrow \$ \, \mathbb{Z}_q$

tag $\leftarrow F_k (m)$

$E2EE(m, k), tag$

Store tag

$E2EE(m, k), tag$

Check tag =? $F_k(m)$

**Messaging phase**

**Reporting phase**

# The Design: Enabling Reporting



**Sender**     **Platform**     **Recipient**

Compute $k \leftarrow \$ \, \mathbb{Z}_q$

$\xrightarrow{\text{E2EE}(m, k), \text{tag}}$

Store tag

$\xrightarrow{\text{E2EE}(m, k), \text{tag}}$

$\text{tag} \leftarrow F_k(m)$

Check tag $=? \, F_k(m)$

**Messaging phase**

$(m, k)$

Compute tag $\leftarrow F_k(m)$ $\longleftarrow$ Submit$(m, k)$ as report

Check whether tag is in the storage

**Reporting phase**

# The Design: Enabling Reporting

**Sender**

**Platform**

**Recipient**

Compute $k \leftarrow \$ \ \mathbb{Z}_q$

$\xrightarrow{\text{E2EE}(m, k), \text{tag}}$

Store tag

$\xrightarrow{\text{E2EE}(m, k), \text{tag}}$

Check tag $=?\ F_k(m)$

tag $\leftarrow F_k(m)$

**Messaging phase**

$(m, k)$

Compute tag $\leftarrow F_k(m)$ $\xleftarrow{\hspace{2cm}}$ Submit$(m, k)$ as report

Check whether tag is in the storage

**Accountability:**
Cannot smear honest user.
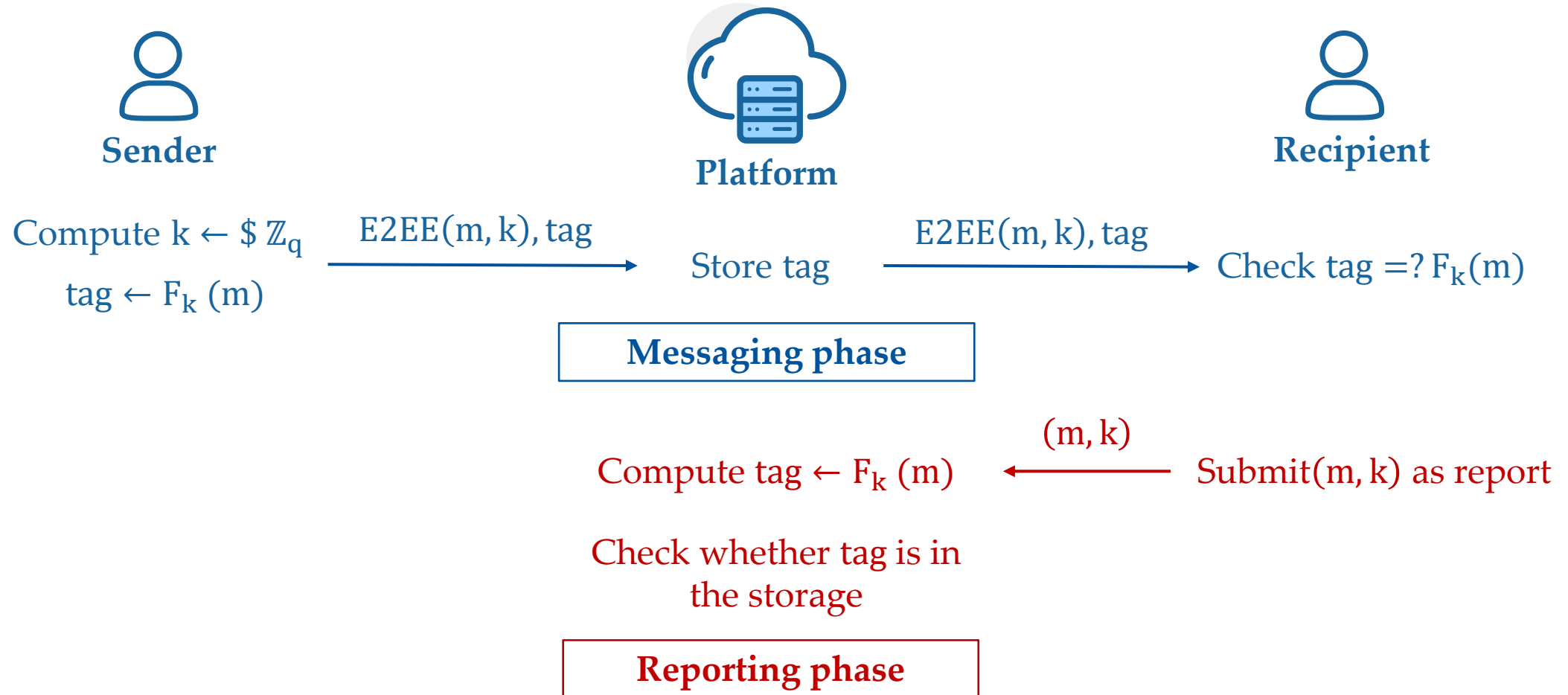
**Reporting phase**

# The Design: Enabling Reporting

# The Design: Enabling Reporting



**Messaging phase**

Originator $\xrightarrow[k_1 \leftarrow \$ \mathbb{Z}_q]{(m, k_1, tag_1)}$ Forwarder 1 $\xrightarrow[k_2 \leftarrow \$ \mathbb{Z}_q]{(m, k_2, tag_2)}$ Forwarder 2 $\xrightarrow[k_3 \leftarrow \$ \mathbb{Z}_q]{(m, k_3, tag_3)}$ Reporter

$(m, k_3)$

Compute $tag_3 \leftarrow F_{k_3}(m)$

Check whether $tag_3$ exists

How to obtain $k_2$ from $k_3$?

Forwarder 2

Platform

**Reporting phase**

# The Design: Enabling Reporting

Originator $\xrightarrow[k_1 \leftarrow \$\, \mathbb{Z}_q]{(m, k_1, \text{tag}_1)}$ Forwarder 1 $\xrightarrow[k_2 \leftarrow \$\, \mathbb{Z}_q]{(m, k_2, \text{tag}_2)}$ Forwarder 2 $\xrightarrow[k_3 \leftarrow \$\, \mathbb{Z}_q]{(m, k_3, \text{tag}_3)}$ Reporter

How to link $k_1$, $k_2$, and $k_3$ implicitly?

How to obtain $k_2$ from $k_3$?

$(m, k_3)$

Compute $\text{tag}_3 \leftarrow F_{k_3}(m)$

Check whether $\text{tag}_3$ exists

Forwarder 2

Platform

Reporting phase

8 / 14

# The Design: Reporting then Tracing

# The Design: Reporting then Tracing

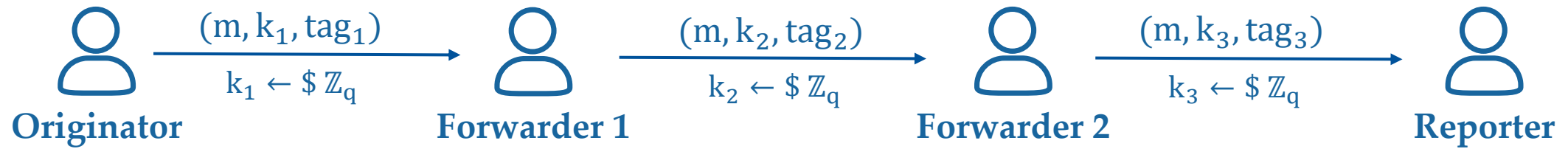$(m, k_1, tag_1)$     $(m, k_2, tag_2)$     $(m, k_3, tag_3)$

**Originator**     **Forwarder 1**     **Forwarder 2**     **Reporter**

$tk_{01} \leftarrow H(ik_0 || U_1)$     $tk_{12} \leftarrow H(ik_1 || U_2)$     $tk_{23} \leftarrow H(ik_2 || U_3)$

$k_1 \leftarrow Enc_{tk_{01}}(k_0)$     $k_2 \leftarrow Enc_{tk_{12}}(k_1)$     $k_3 \leftarrow Enc_{tk_{23}}(k_2)$

$k_0 \leftarrow \$ \, \mathbb{Z}_q$

# The Design: Reporting then Tracing

**ik:** Users' secret identity key; **U:** Users' public identity; **Enc:** Block cipher



**Originator**

$(m, k_1, tag_1)$

**Forwarder 1**
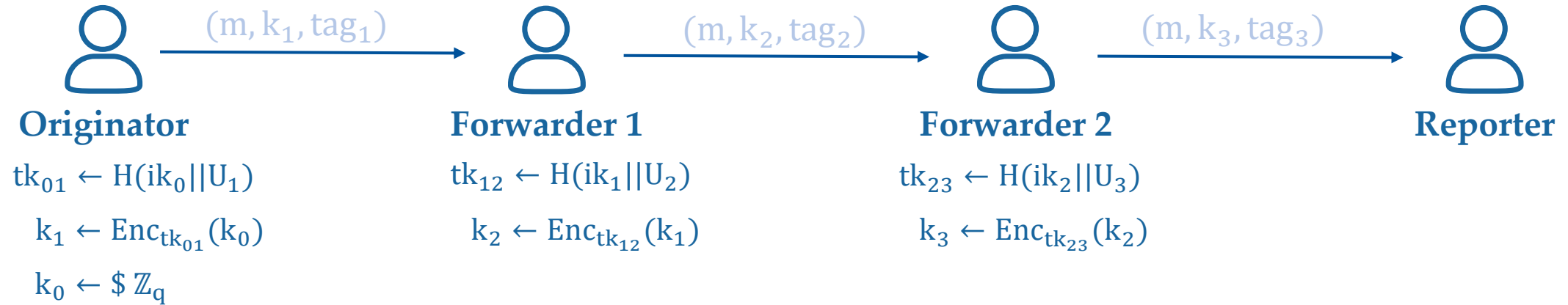
$(m, k_2, tag_2)$

**Forwarder 2**

$(m, k_3, tag_3)$

**Reporter**

$tk_{01} \leftarrow H(ik_0 || U_1)$

$tk_{12} \leftarrow H(ik_1 || U_2)$

$tk_{23} \leftarrow H(ik_2 || U_3)$

$k_1 \leftarrow Enc_{tk_{01}}(k_0)$  $\dashrightarrow$  $k_2 \leftarrow Enc_{tk_{12}}(k_1)$  $\dashrightarrow$  $k_3 \leftarrow Enc_{tk_{23}}(k_2)$

$k_0$   Implicit encryption chain

# The Design: Reporting then Tracing

**Originator**      $(m, k_1, tag_1)$      **Forwarder 1**      $(m, k_2, tag_2)$      **Forwarder 2**      $(m, k_3, tag_3)$      **Reporter**

$tk_{01} \leftarrow H(ik_0 || U_1)$      $tk_{12} \leftarrow H(ik_1 || U_2)$      $tk_{23} \leftarrow H(ik_2 || U_3)$

$k_1 \leftarrow Enc_{tk_{01}}(k_0)$  ⇢  $k_2 \leftarrow Enc_{tk_{12}}(k_1)$  ⇢  $k_3 \leftarrow Enc_{tk_{23}}(k_2)$

$(m, k_3)$

$k_0$

**Implicit encryption chain**

**Start at $k_3$**

Comp. $tag_3 \leftarrow F_{k_3}(m)$

Check whether $tag_3$ exists

**Platform**

**Forwarder 2**

**Reporting-then-tracing phase**

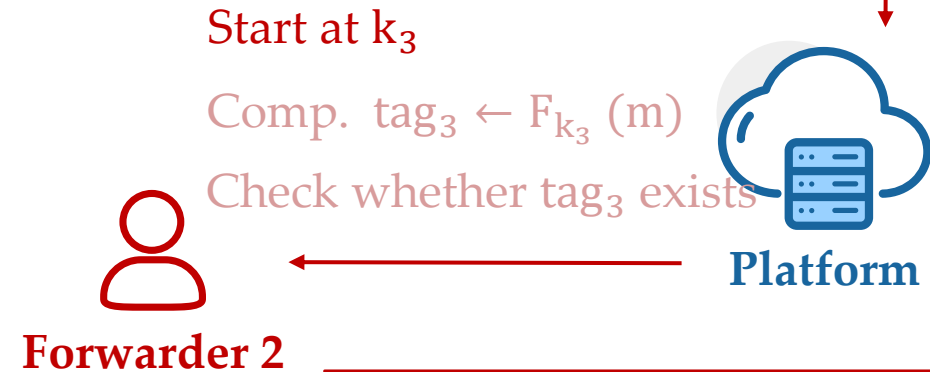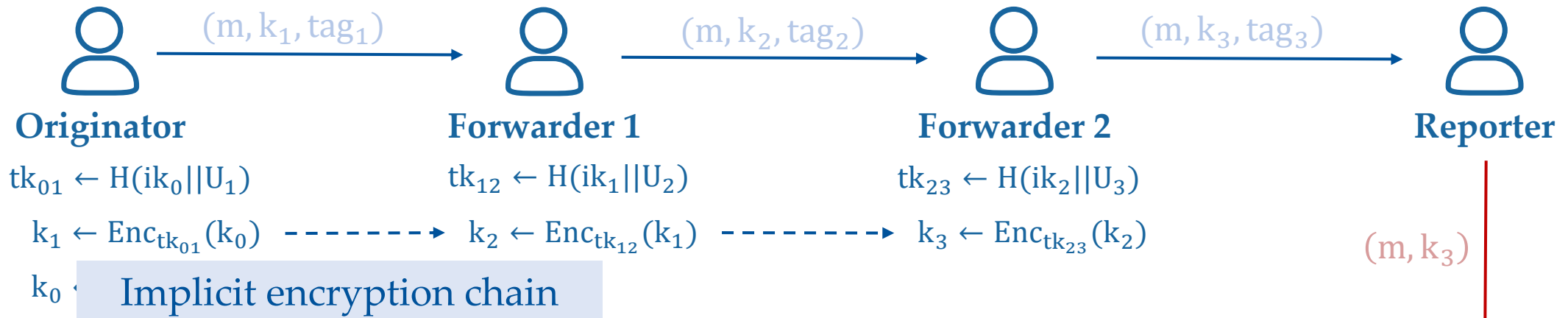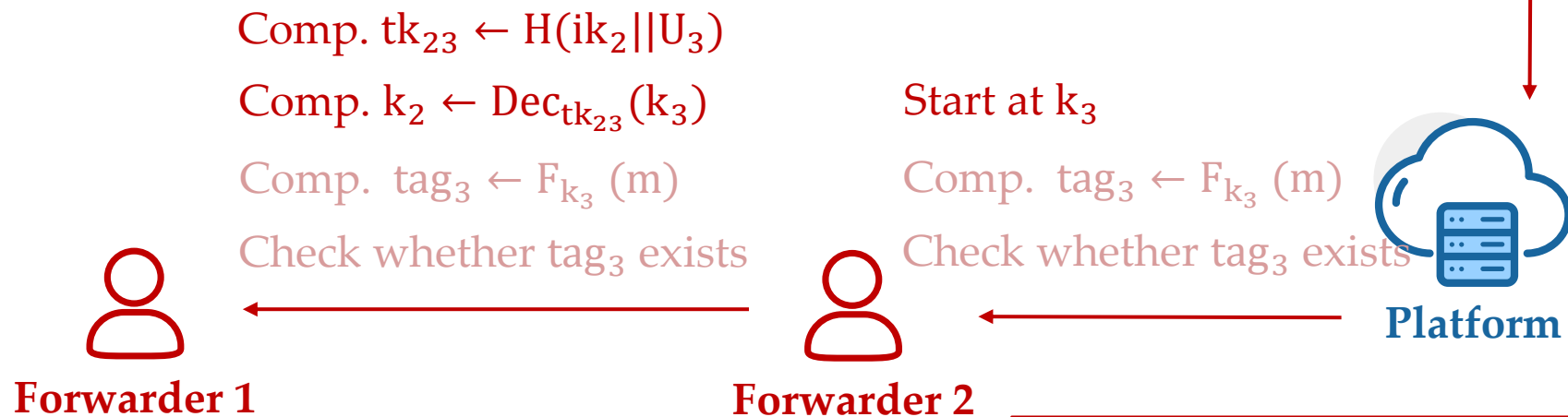# The Design: Reporting then Tracing

**ik:** Users' secret identity key; **U:** Users' public identity; **Enc:** Block cipher



**Originator**

$tk_{01} \leftarrow H(ik_0 || U_1)$

$k_1 \leftarrow Enc_{tk_{01}}(k_0)$

$k_0$

**Implicit encryption chain**

**Forwarder 1**

$tk_{12} \leftarrow H(ik_1 || U_2)$

$k_2 \leftarrow Enc_{tk_{12}}(k_1)$

**Forwarder 2**

$tk_{23} \leftarrow H(ik_2 || U_3)$

$k_3 \leftarrow Enc_{tk_{23}}(k_2)$

**Reporter**

$(m, k_1, tag_1)$ $(m, k_2, tag_2)$ $(m, k_3, tag_3)$

$(m, k_3)$

Comp. $tk_{23} \leftarrow H(ik_2 || U_3)$

Comp. $k_2 \leftarrow Dec_{tk_{23}}(k_3)$

Comp. $tag_3 \leftarrow F_{k_3}(m)$

Check whether $tag_3$ exists

Start at $k_3$

Comp. $tag_3 \leftarrow F_{k_3}(m)$

Check whether $tag_3$ exists

**Platform**

**Forwarder 1**

**Forwarder 2**

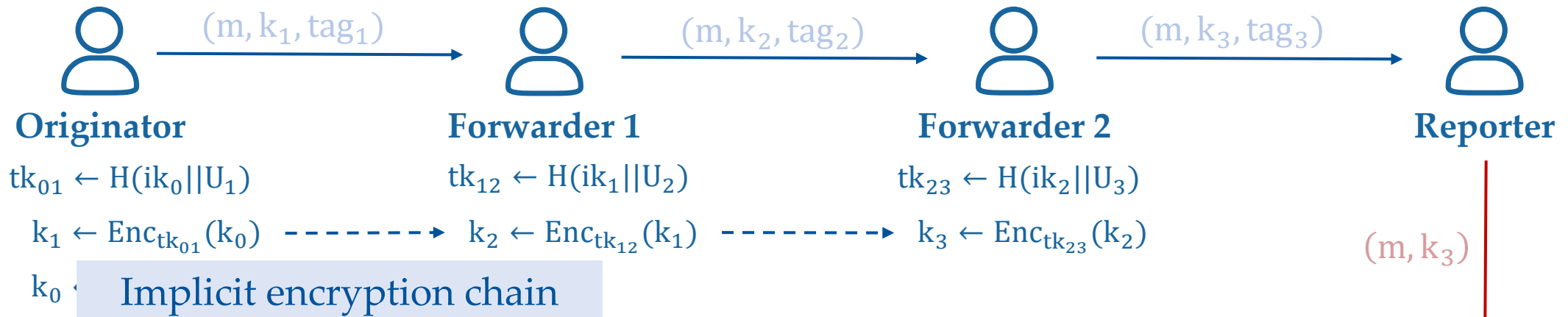**Reporting-then-tracing phase**

# The Design: Reporting then Tracing

**ik:** Users' secret identity key;  **U:** Users' public identity;  **Enc:** Block cipher



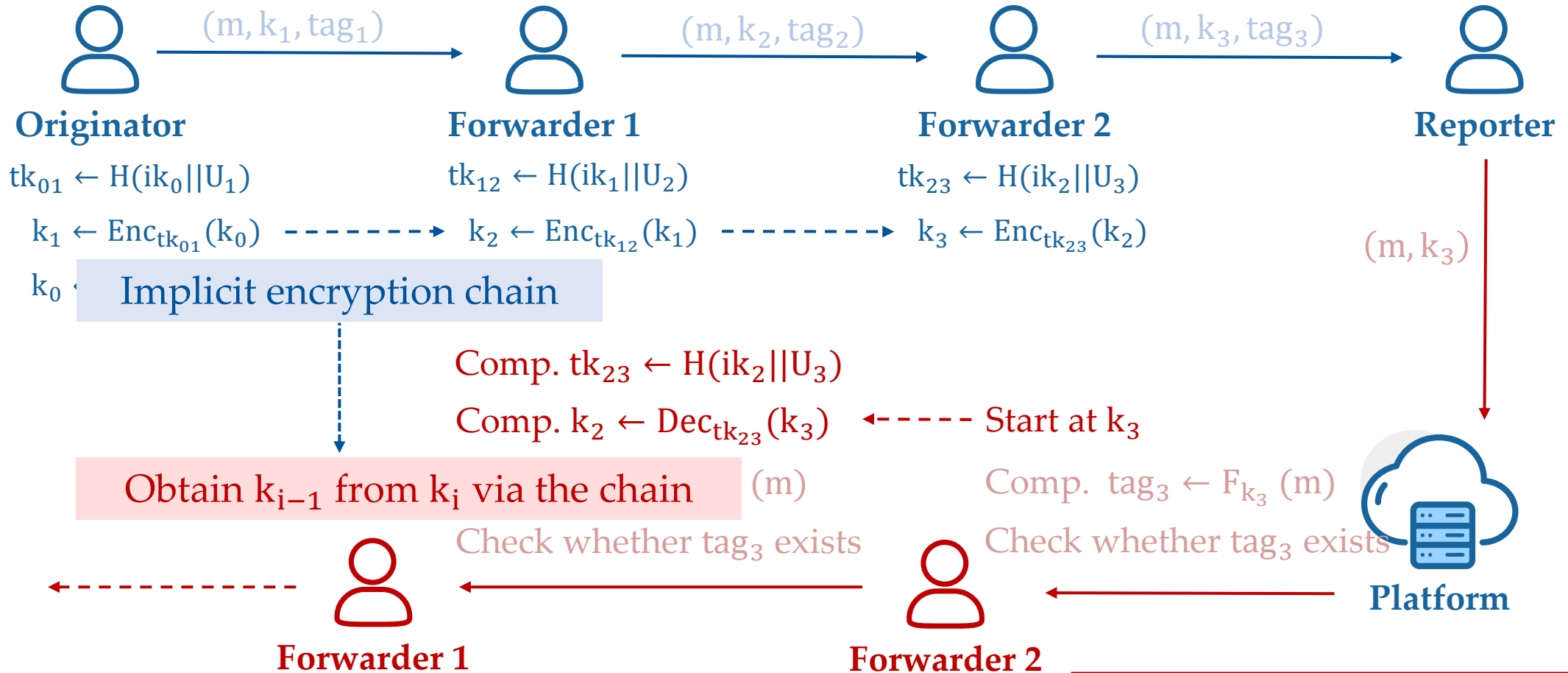**Originator** $\xrightarrow{(m, k_1, tag_1)}$ **Forwarder 1** $\xrightarrow{(m, k_2, tag_2)}$ **Forwarder 2** $\xrightarrow{(m, k_3, tag_3)}$ **Reporter**

$tk_{01} \leftarrow H(ik_0 || U_1)$       $tk_{12} \leftarrow H(ik_1 || U_2)$       $tk_{23} \leftarrow H(ik_2 || U_3)$

$k_1 \leftarrow Enc_{tk_{01}}(k_0)$  - - - - →  $k_2 \leftarrow Enc_{tk_{12}}(k_1)$  - - - - →  $k_3 \leftarrow Enc_{tk_{23}}(k_2)$

$(m, k_3)$

$k_0$

**Implicit encryption chain**

Comp. $tk_{23} \leftarrow H(ik_2 || U_3)$

Comp. $k_2 \leftarrow Dec_{tk_{23}}(k_3)$  ← - - - - Start at $k_3$

Obtain $k_{i-1}$ from $k_i$ via the chain   (m)   Comp. $tag_3 \leftarrow F_{k_3}(m)$

Check whether $tag_3$ exists        Check whether $tag_3$ exists

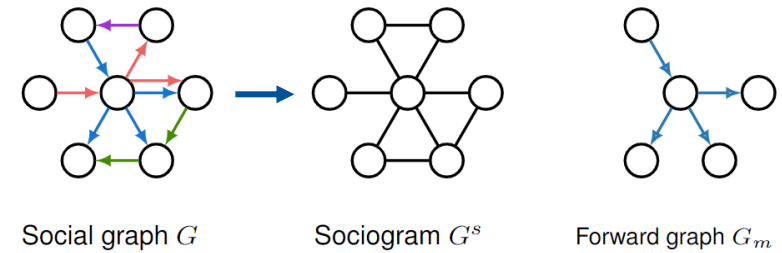**Forwarder 1**          **Forwarder 2**          **Platform**

**Reporting-then-tracing phase**

# The Design: Reporting then Tracing

Messaging phase:

- Users forward messages.
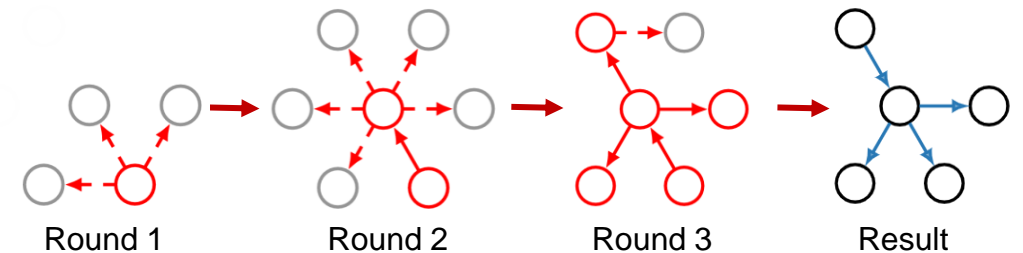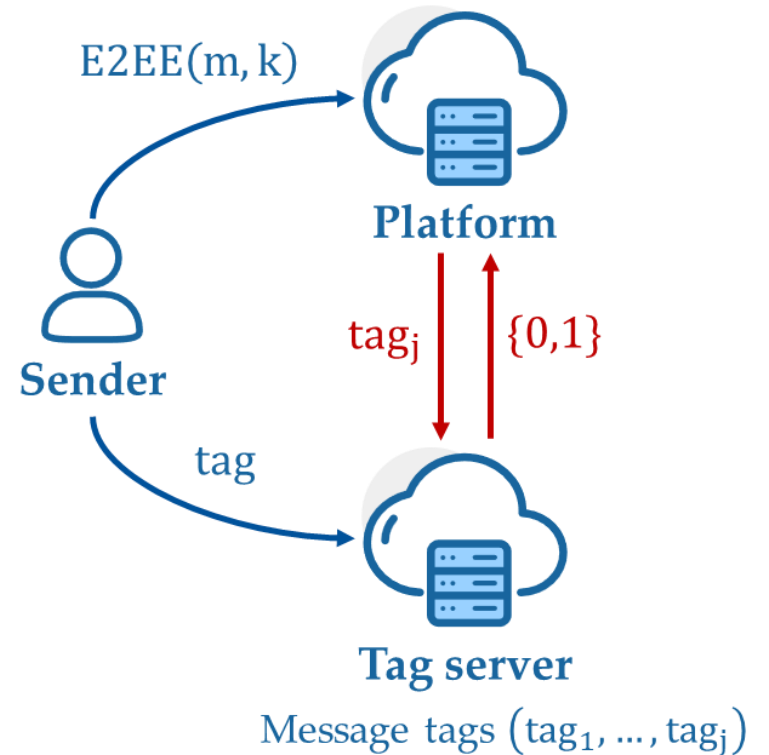
- Platform collects sociogram & stores tags.



Social graph $G$     Sociogram $G^s$     Forward graph $G_m$

# The Design: Reporting then Tracing

Messaging phase:

- Users forward messages.

- Platform collects sociogram & stores tags.



Social graph $G$     Sociogram $G^s$     Forward graph $G_m$

Reporting-then-Tracing phase:

1) Compute tag key and tag.

2) Check the existence of the tag.

3) Start next round at existed vertices.



Round 1     Round 2     Round 3     Result

# The Design: Noising the Traceback

Introduce random response to add noise.

- Two servers: Platform & Tag server



$E2EE(m, k)$

**Platform**

**Sender**

$tag_j$ | $\{0,1\}$

tag

**Tag server**

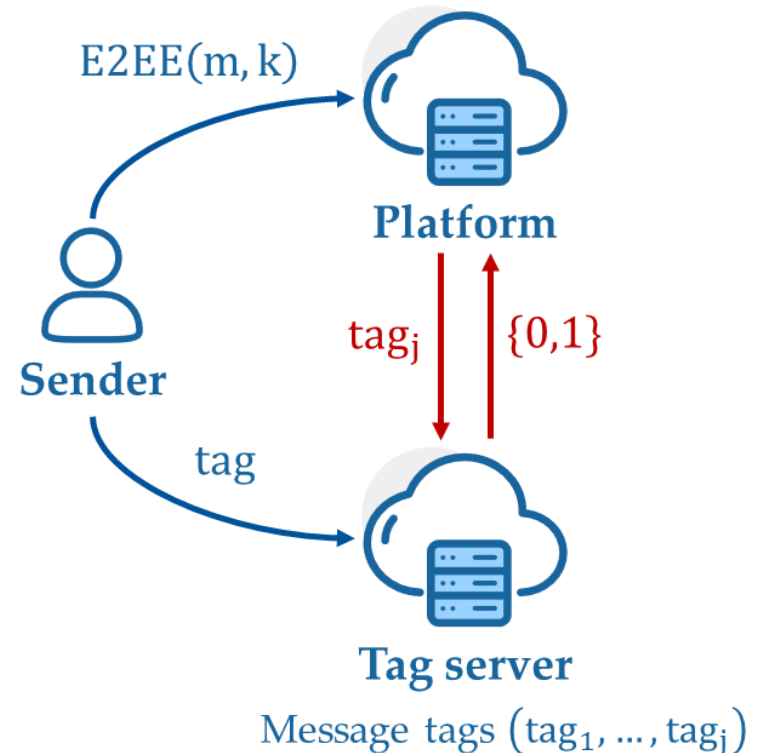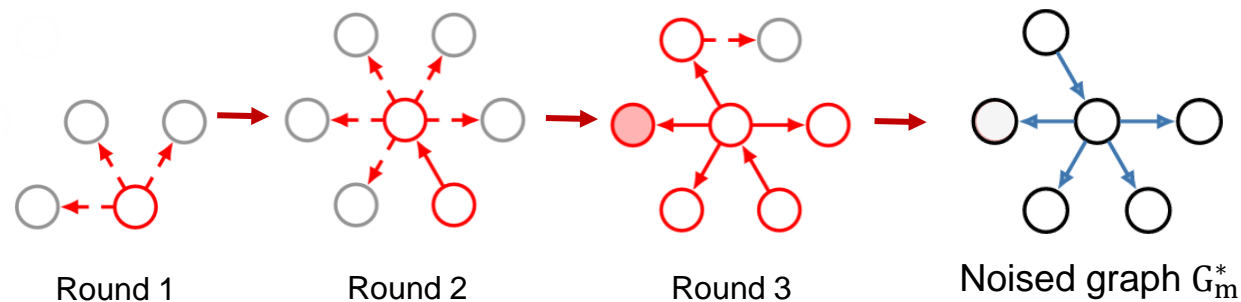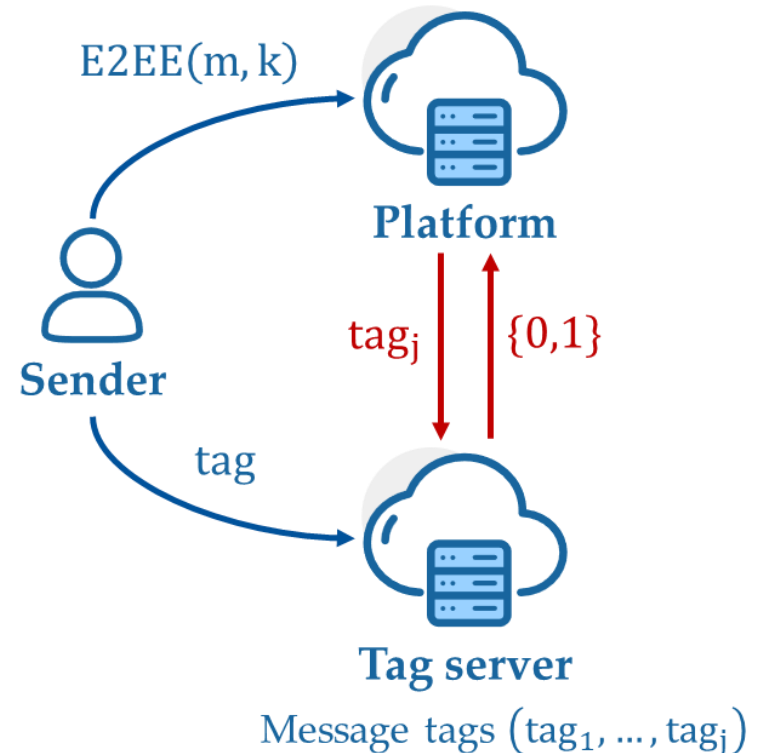Message tags $(tag_1, ..., tag_j)$

# The Design: Noising the Traceback

Introduce random response to add noise.

- Two servers: Platform & Tag server

- Tag server responses a query as:

  if $tag_j$ exists, return 1

  else, return 0 w.p. $1 - \psi$

  return 1 w.p. $\psi$



E2EE(m, k)

Platform

Sender

$tag_j$   {0,1}

tag

Tag server

Message tags $(tag_1, \ldots, tag_j)$

# The Design: Noising the Traceback

Introduce <span style="color:red">random response</span> to add noise.

- Two servers: Platform & Tag server

- Tag server responses a query as:

  if $tag_j$ exists, return 1

  else, return 0 w.p. $1 - \psi$

  return 1 w.p. $\psi$



$E2EE(m, k)$

Platform

Sender

$tag_j$    $\{0,1\}$

tag

Tag server

Message tags $(tag_1, ..., tag_j)$



Round 1     Round 2     Round 3     Noised graph $G_m^*$
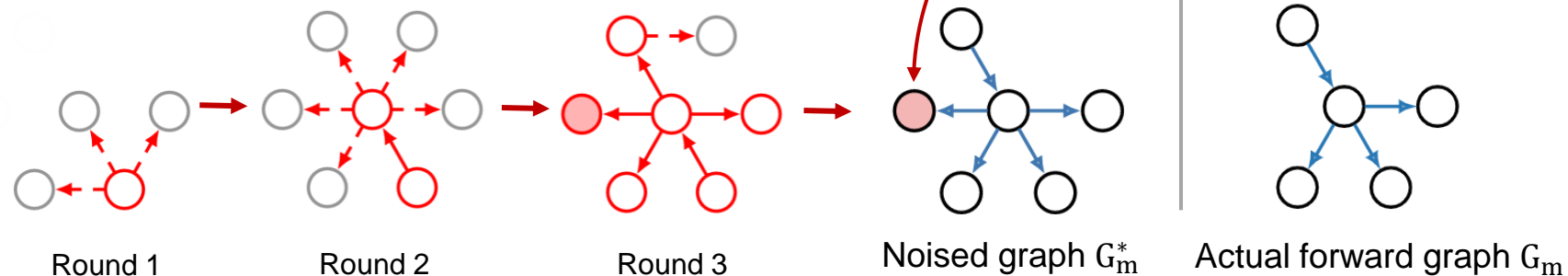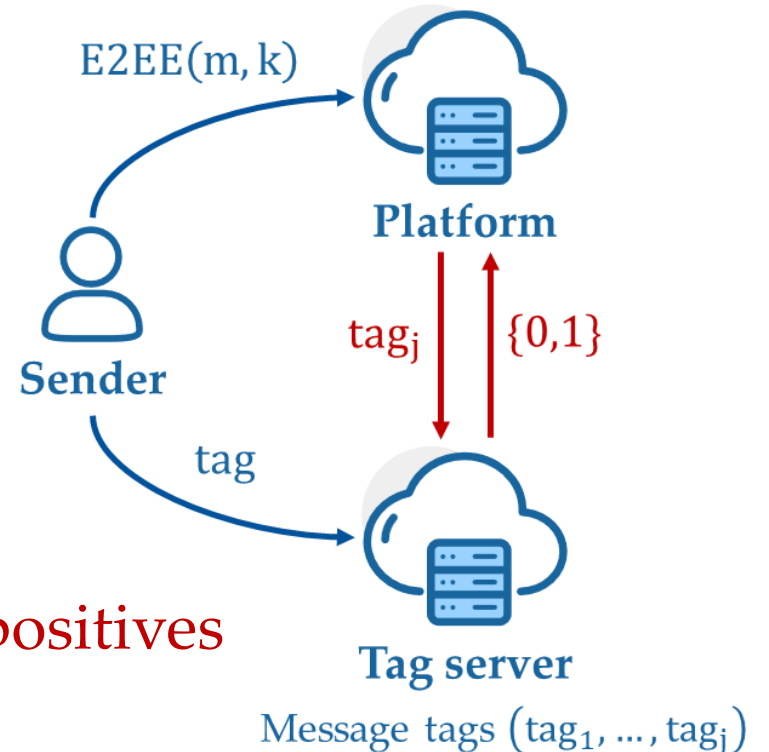
# The Design: Noising the Traceback

Introduce random response to add noise.

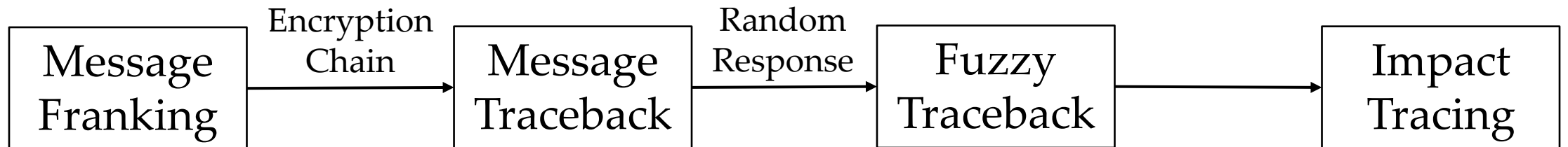- Two servers: Platform & Tag server

- Tag server responses a query as:

    if $\text{tag}_j$ exists, return 1

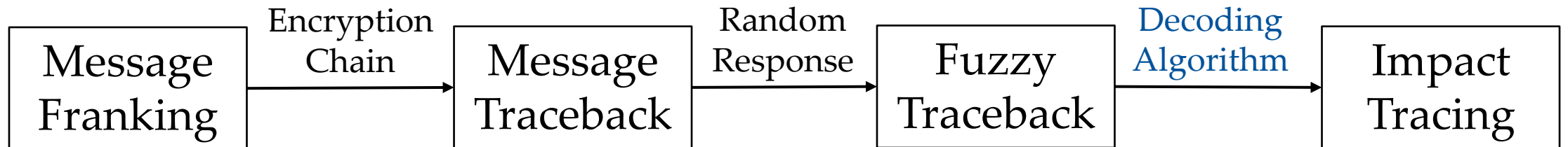    else, return 0 w.p. $1 - \psi$

    return 1 w.p. $\psi$



$E2EE(m, k)$

Platform

Sender

$\text{tag}_j$  $\{0,1\}$

tag

Tag server

Message tags $(\text{tag}_1, ..., \text{tag}_j)$

False positives

Round 1     Round 2     Round 3     Noised graph $\mathbb{G}_m^*$     Actual forward graph $\mathbb{G}_m$

# Remaining Concerns

```
┌─────────────┐  Encryption  ┌─────────────┐  Random    ┌─────────────┐            ┌─────────────┐
│  Message    │    Chain     │  Message    │  Response  │   Fuzzy     │            │  Impact     │
│  Franking   │ ───────────▶ │  Traceback  │ ─────────▶ │  Traceback  │ ─────────▶ │  Tracing    │
└─────────────┘              └─────────────┘            └─────────────┘            └─────────────┘
```

# Remaining Concerns

How to reveal influential spreaders from noised graph?

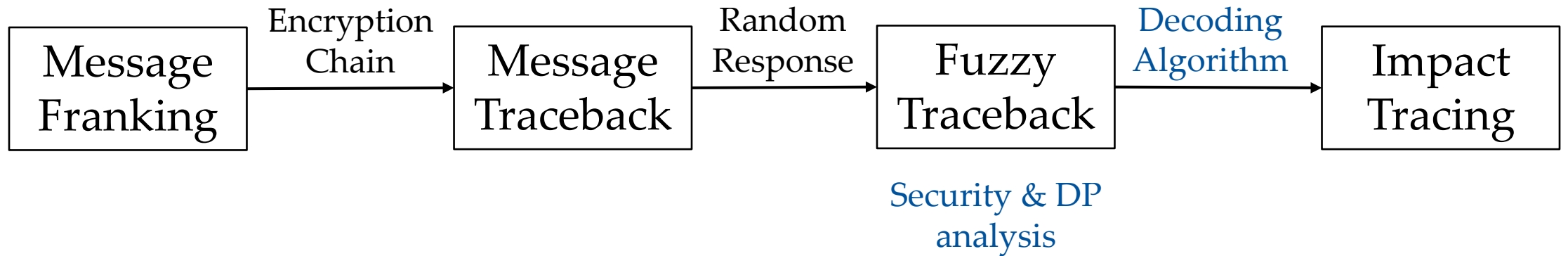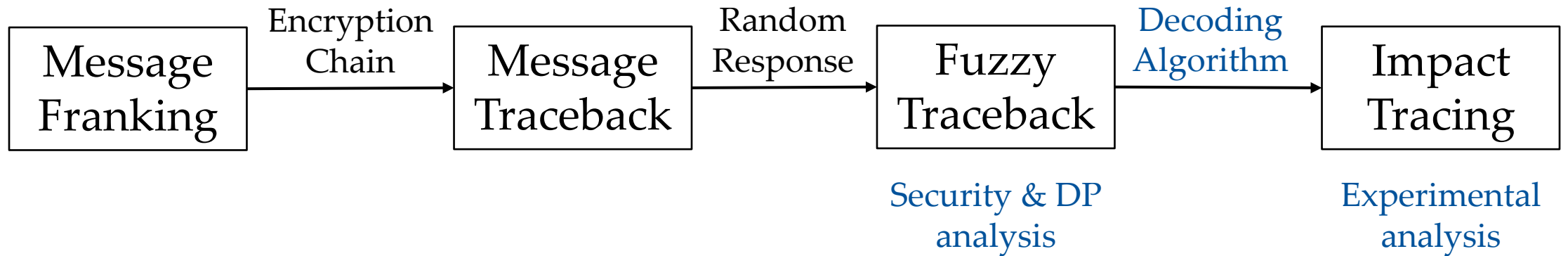- Design a decoding algorithm to identify influential spreaders.

# Remaining Concerns

How to reveal influential spreaders from noised graph?

- Design a decoding algorithm to identify influential spreaders.

Does the decoding algorithm violate privacy?
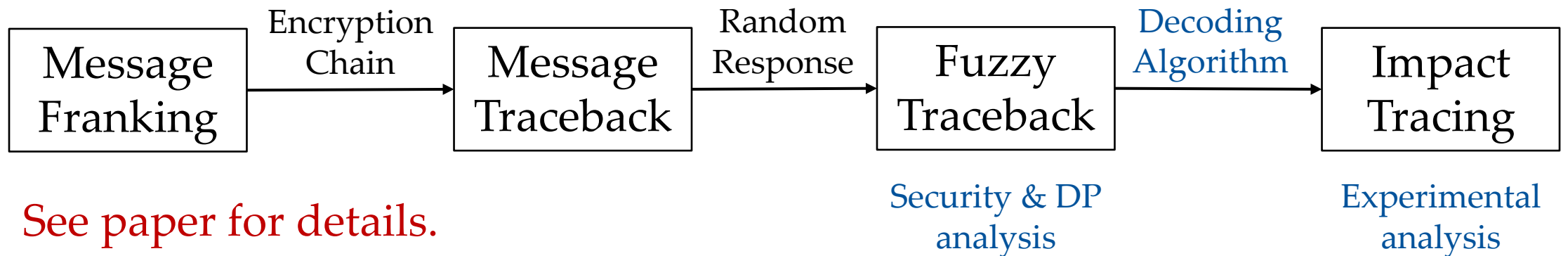
- Prove the noised graph satisfies differential privacy.

# Remaining Concerns

How to reveal influential spreaders from noised graph?

- Design a decoding algorithm to identify influential spreaders.

Does the decoding algorithm violate privacy?

- Prove the noised graph satisfies differential privacy.

# Remaining Concerns

How to reveal influential spreaders from noised graph?

- Design a decoding algorithm to identify influential spreaders.

Does the decoding algorithm violate privacy?

- Prove the noised graph satisfies differential privacy.

| Message Franking | Encryption Chain → | Message Traceback | Random Response → | Fuzzy Traceback | Decoding Algorithm → | Impact Tracing |
|---|---|---|---|---|---|---|

See paper for details.

Security & DP analysis

Experimental analysis

# Performance

## Bandwidth & Storage (Byte)

| Tracing Policy | Schemes[†] | Bandwidth | | Storage | |
|---|---|---|---|---|---|
| | | $\mathcal{S}$ - $\mathcal{P}$ | $\mathcal{P}$ - $\mathcal{R}$ | $\mathcal{C}$ | $\mathcal{P}$ |
| Source Tracing | PEB21 [39] | 256 | 320 | 160 | - |
| | LRTY22 [30] | 243 | 243 | 243 | - |
| | IAV22 [21] | 380 | 484 | 380 | - |
| Message Traceback | TMR19 [50] | 96 | 80 | 34 | 104 |
| | KTW22 [27] | 203 | 203 | 16 | 136 |
| Impact Tracing | Ours | 96[‡] | 72 | 16 | 6 |

$\mathcal{S}$, $\mathcal{R}$, $\mathcal{P}$, $\mathcal{C}$: Sender, recipient, platform, and client.

Store 1 billion messages: 5.6 GB

# Performance

## Bandwidth & Storage (Byte)

| Tracing Policy | Schemes[†] | Bandwidth | | Storage | |
|---|---|---|---|---|---|
| | | $\mathcal{S}$ - $\mathcal{P}$ | $\mathcal{P}$ - $\mathcal{R}$ | $\mathcal{C}$ | $\mathcal{P}$ |
| Source Tracing | PEB21 [39] | 256 | 320 | 160 | - |
| | LRTY22 [30] | 243 | | 243 | - |
| | IAV22 [21] | 380 | | 380 | - |
| Message Traceback | TMR19 [50] | 96 | 80 | 34 | 104 |
| | KTW22 [27] | 203 | 203 | 16 | 136 |
| Impact Tracing | Ours | 96[‡] | 72 | 16 | 6 |

71%

$\mathcal{S}$, $\mathcal{R}$, $\mathcal{P}$, $\mathcal{C}$: Sender, recipient, platform, and client.
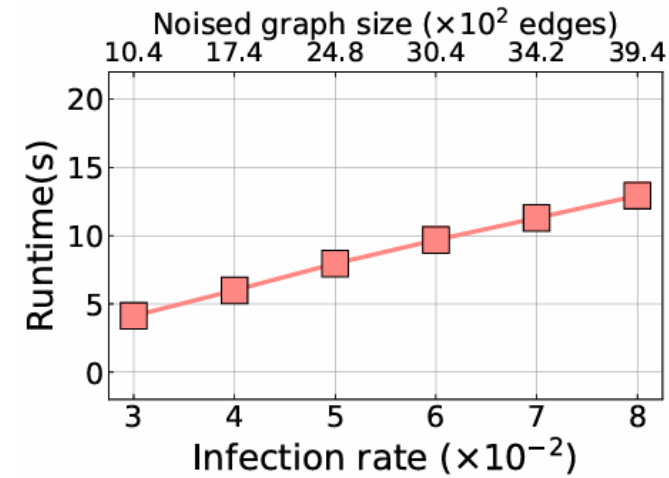
Store 1 billion messages: 5.6 GB

# Performance

## Bandwidth & Storage (Byte)

| Tracing Policy | Schemes[†] | Bandwidth | | Storage | |
|---|---|---|---|---|---|
| | | $\mathcal{S}$ - $\mathcal{P}$ | $\mathcal{P}$ - $\mathcal{R}$ | $\mathcal{C}$ | $\mathcal{P}$ |
| Source Tracing | PEB21 [39] | 256 | 320 | 160 | - |
| | LRTY22 [30] | 243 | | 243 | - |
| | | | 71% | | |
| | IAV22 [21] | 380 | | 380 | - |
| Message Traceback | TMR19 [50] | 96 | 80 | 34 | 104 |
| | KTW22 [27] | 203 | 5% | 94% | |
| Impact Tracing | Ours | 96[‡] | 72 | 16 | 6 |

$\mathcal{S}$, $\mathcal{R}$, $\mathcal{P}$, $\mathcal{C}$: Sender, recipient, platform, and client.

Store 1 billion messages: 5.6 GB

# Performance

## Bandwidth & Storage (Byte)

| Tracing Policy | Schemes[†] | Bandwidth | | Storage | |
|---|---|---|---|---|---|
| | | $\mathcal{S}$ - $\mathcal{P}$ | $\mathcal{P}$ - $\mathcal{R}$ | $\mathcal{C}$ | $\mathcal{P}$ |
| Source Tracing | PEB21 [39] | 256 | 320 | 160 | - |
| | LRTY22 [30] | 243 | | 243 | - |
| | IAV22 [21] | 380 | 71% | 380 | - |
| Message Traceback | TMR19 [50] | 96 | 80 | 34 | 104 |
| | KTW22 [27] | 203 | 5% | 94% | 6 |
| Impact Tracing | Ours | 96[‡] | 72 | 16 | 6 |

$\mathcal{S}$, $\mathcal{R}$, $\mathcal{P}$, $\mathcal{C}$: Sender, recipient, platform, and client.

Store 1 billion messages: 5.6 GB

## Runtime of tracing (s)



Trace graph with 4,000 edges: 15 s

Transmit 1-KB message: 0.3 ms

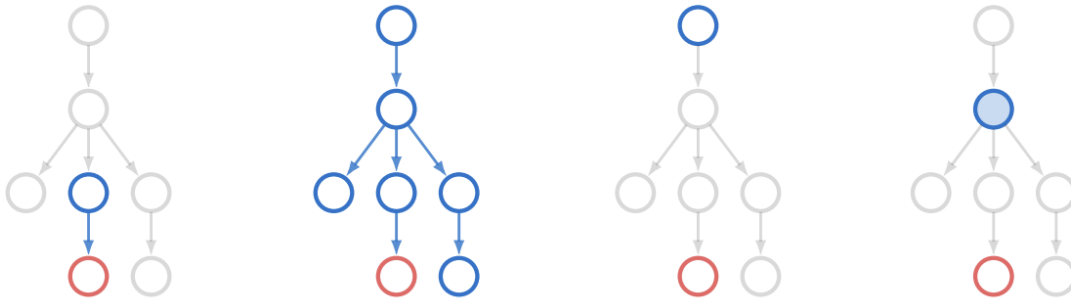# Performance

## Bandwidth & Storage (Byte)

| Tracing Policy | Schemes[†] | Bandwidth | | Storage | |
|---|---|---|---|---|---|
| | | $\mathcal{S}$ - $\mathcal{P}$ | $\mathcal{P}$ - $\mathcal{R}$ | $\mathcal{C}$ | $\mathcal{P}$ |
| Source Tracing | PEB21 [39] | 256 | 320 | 160 | - |
| | LRTY22 [30] | 243 | 71% | 243 | - |
| | IAV22 [21] | 380 | | 380 | - |
| Message Traceback | TMR19 [50] | 96 | 80 | 34 | 104 |
| | KTW22 [27] | 203 | 5% | 94% | |
| Impact Tracing | Ours | 96[‡] | 72 | 16 | 6 |

$\mathcal{S}, \mathcal{R}, \mathcal{P}, \mathcal{C}$: Sender, recipient, platform, and client.

Store 1 billion messages: 5.6 GB

## Runtime of tracing (s)



Trace graph with 4,000 edges: 15 s

Transmit 1-KB message: 0.3 ms

# Summary

Introduce impact tracing: balances traceability and privacy.

Design fuzzy message traceback and decoding algorithm.

Analyze security, privacy, and utility formally.



Message franking     Message traceback     Source tracing     **Impact tracing**
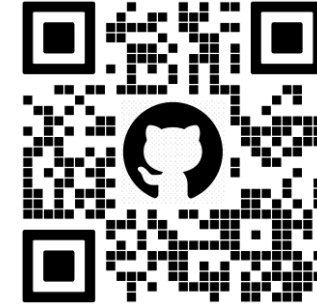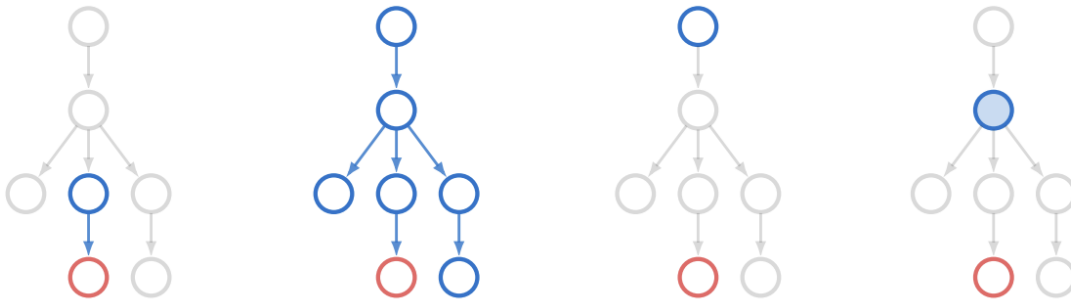
# Summary

Introduce impact tracing: balances traceability and privacy.

Design fuzzy message traceback and decoding algorithm.

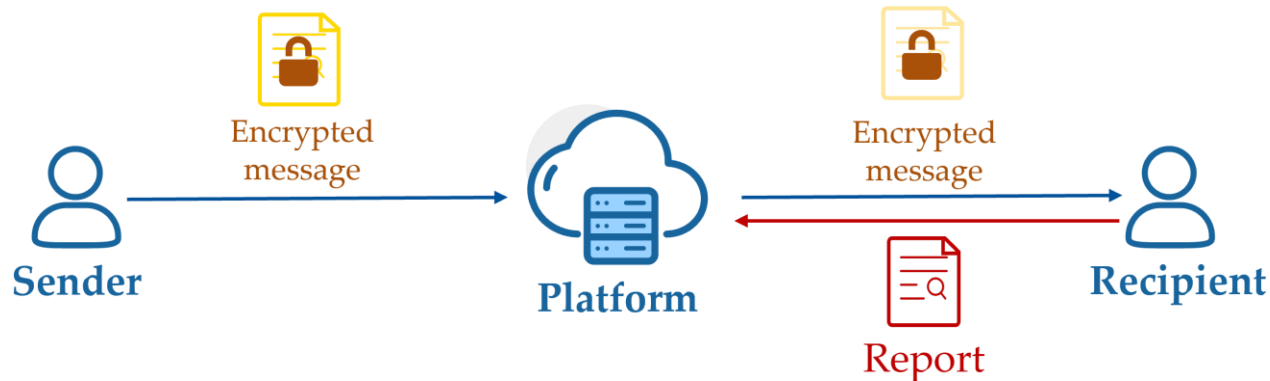Analyze security, privacy, and utility formally.



Scan for code.

Thank for listening :) Any questions?

# Bonus Slides

# Message Franking

Message franking enables a recipient to <u>report a message</u>.

And, the platform can <u>authenticate</u> that the sender actually sent it.



**Accountability:** 1) Recipients cannot smear honest senders.

2) Senders cannot evade reporting.

**Confidentiality:** The platform learn nothing about unreported messages.
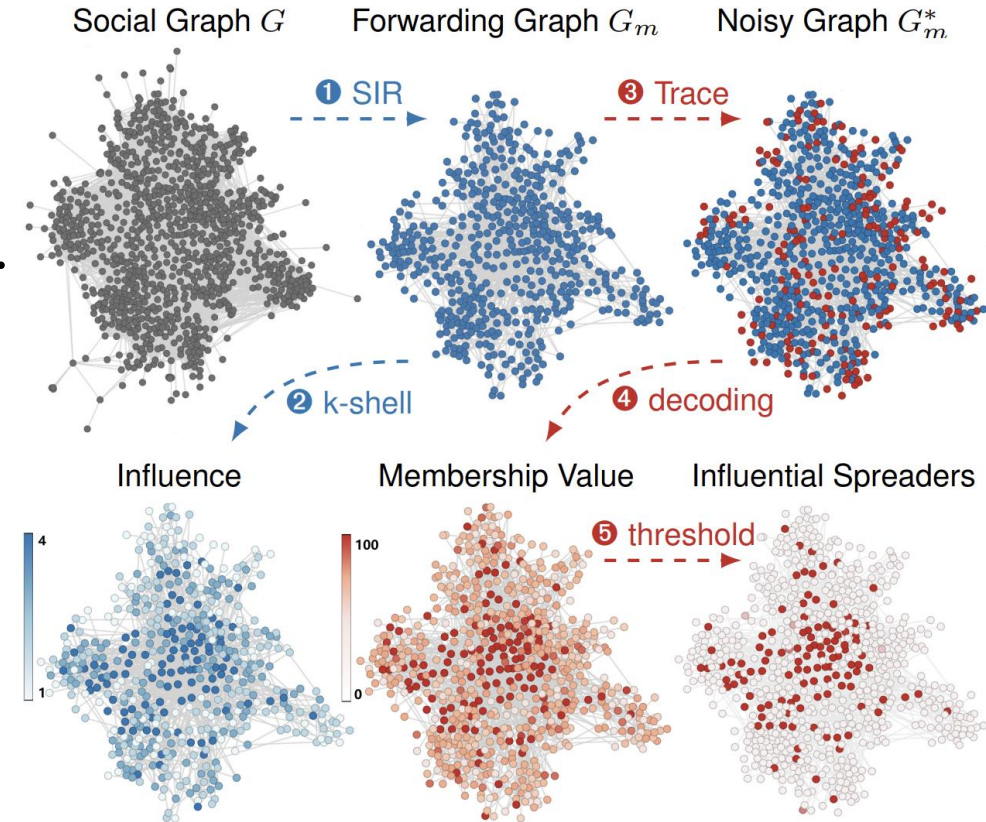
# Simulation on Real-World Datasets

## Influence evaluation

1. Simulate forwarding using SIR model.

2. Evaluate vertices' impact using k-shell.

## Impact tracing

3. Trace with fuzzy message traceback.

4. Decode the traced noisy graph.
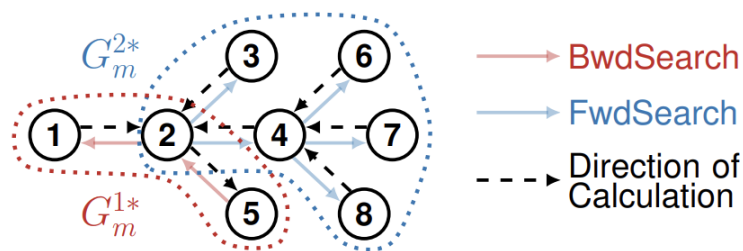
5. Output a set of influential spreaders.

# The Design: Decoding the Result

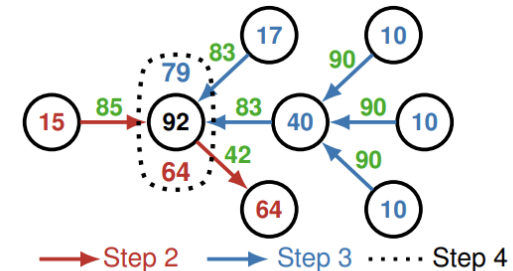Decoding algorithm: Computes a fuzzy value for each vertex.

- The false positives of a vertex satisfy binomial distribution.

- A vertex has *only* one true precursor.

- A vertex is true positive if one of its descendants is true.



(a) llustration of $G_m^{1*}$ and $G_m^{2*}$

| Graph | Vertex | $n(v)$ | $n^*(v)$ | $a_v(\%)$ |
|-------|--------|--------|----------|-----------|
| $G_m^{1*}$ | v(2) | 50 | 1 | 85 |
| | v(2) | 50 | 2 | 83 |
| $G_m^{2*}$ | v(4) | 100 | 3 | 90 |
| | v(5) | 10 | 1 | 42 |

(b) Step 1: caculate $\alpha_v$

(c) Step 2-4: caculate $\beta_v^1$, $\beta_v^2$, and $\mu_v$
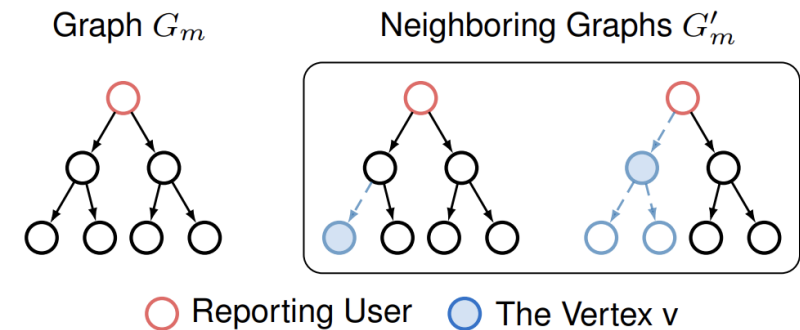
# IAS-DP: Defining Privacy

## Individualized Asymmetric Subtree Differential Privacy

- **Individualized**: Privacy budgets vary per user.

- **Asymmetric**: Traceability adopts one-side noise.

- **Subtree**: All impacts caused by one user on one forwarding graph.

**Definition 2 ($\varepsilon_v$-IAS-DP).** A randomized algorithm $\mathcal{M}$ is $\varepsilon_v$-IAS-DP if given a graph G the following equation holds for any $S \in Range(\mathcal{M})$ and neighboring subgraph pair $(G_m, G'_m)$, where $G'_m$ is obtained by removing a subtree $tree(v)$ in $G_m$.

$$\Pr\left[\mathcal{M}(G_m) = S\right] \le e^{\varepsilon_v} \cdot \Pr[\mathcal{M}(G'_m) = S],$$

where $\varepsilon_v$ is the privacy budget for vertex $v$.

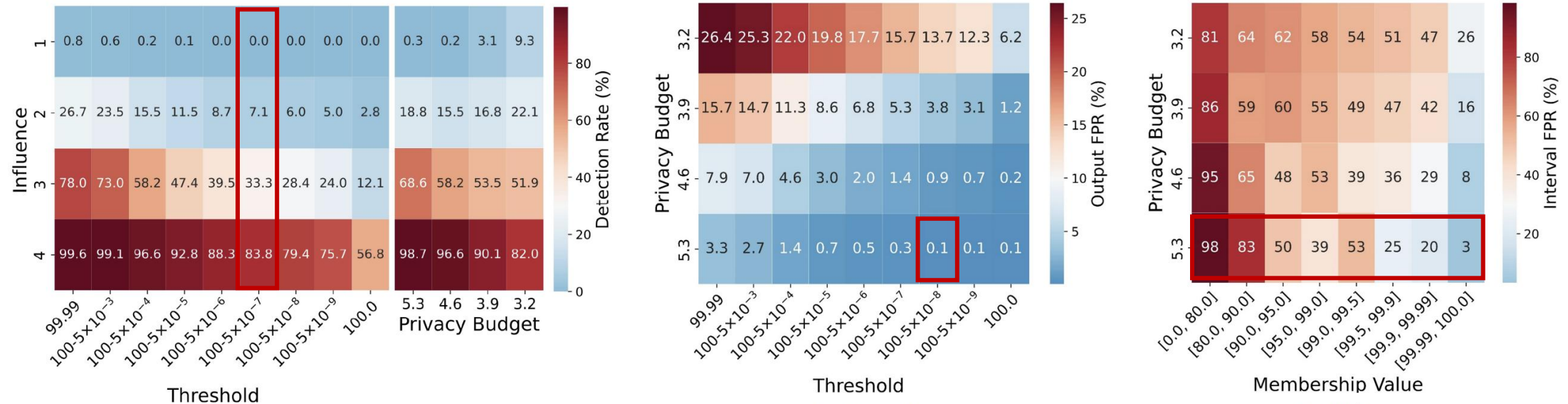Graph $G_m$        Neighboring Graphs $G'_m$

○ Reporting User    ○ The Vertex v

**Theorem 1.** The fuzzy message traceback scheme satisfies $\varepsilon_v$-IAS-DP, where $\varepsilon_v = \ln(1/\psi^n)$, $n$ is the number of edges in $tree(v)$, and $\psi$ is the FPR of random response.

# Evaluating Utility & Privacy

**Detection rate:** The output contains as <u>most influential users</u> as possible.

**Output FPR:** The output should contain as <u>few false positives</u> as possible.

**Interval FPR:** Non-influential users should be hidden by <u>sufficient noise</u>.



Identifies <u>84% of the most influential spreaders</u> and <u>no the least influential users</u> with <u>99.9% correctness</u> (i.e., less than 0.25 false positives on average).