

# Asymmetric Message Franking

Content Moderation for Metadata-Private End-to-End Encryption

Nirvan Tyagi

Paul Grubbs

Julia Len

Ian Miers

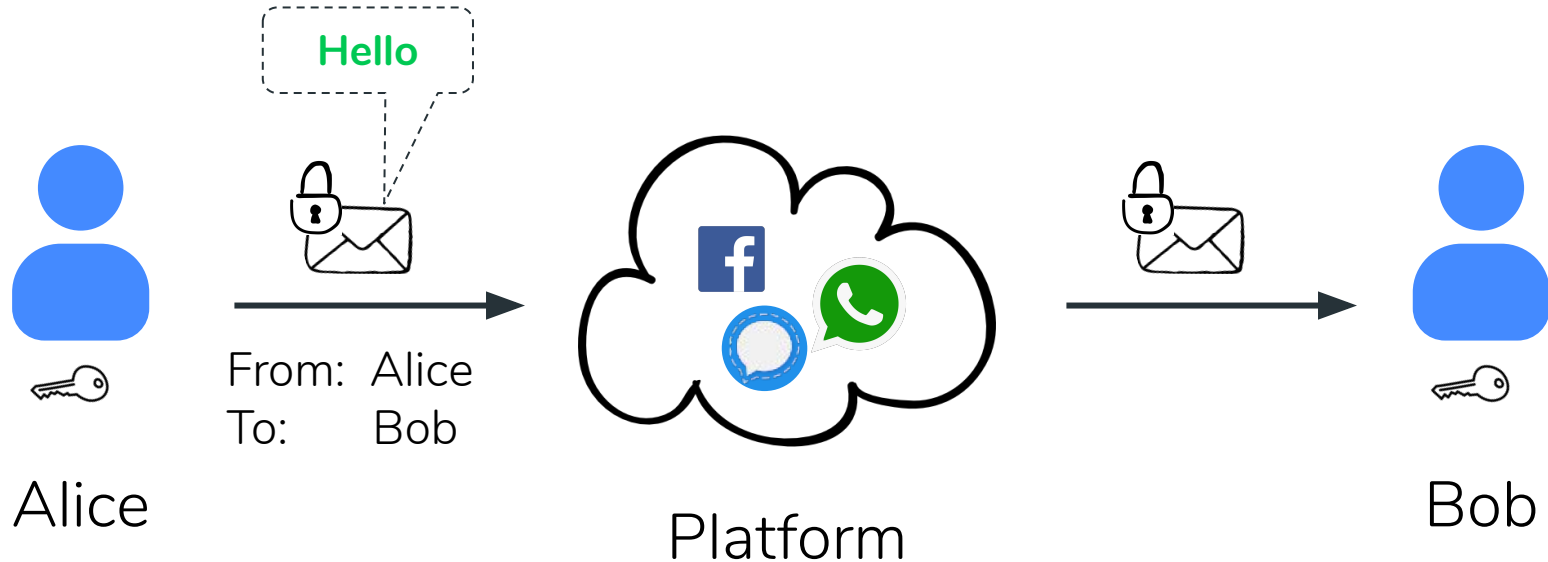
Tom Ristenpart

# Setting: End-to-end encrypted messaging



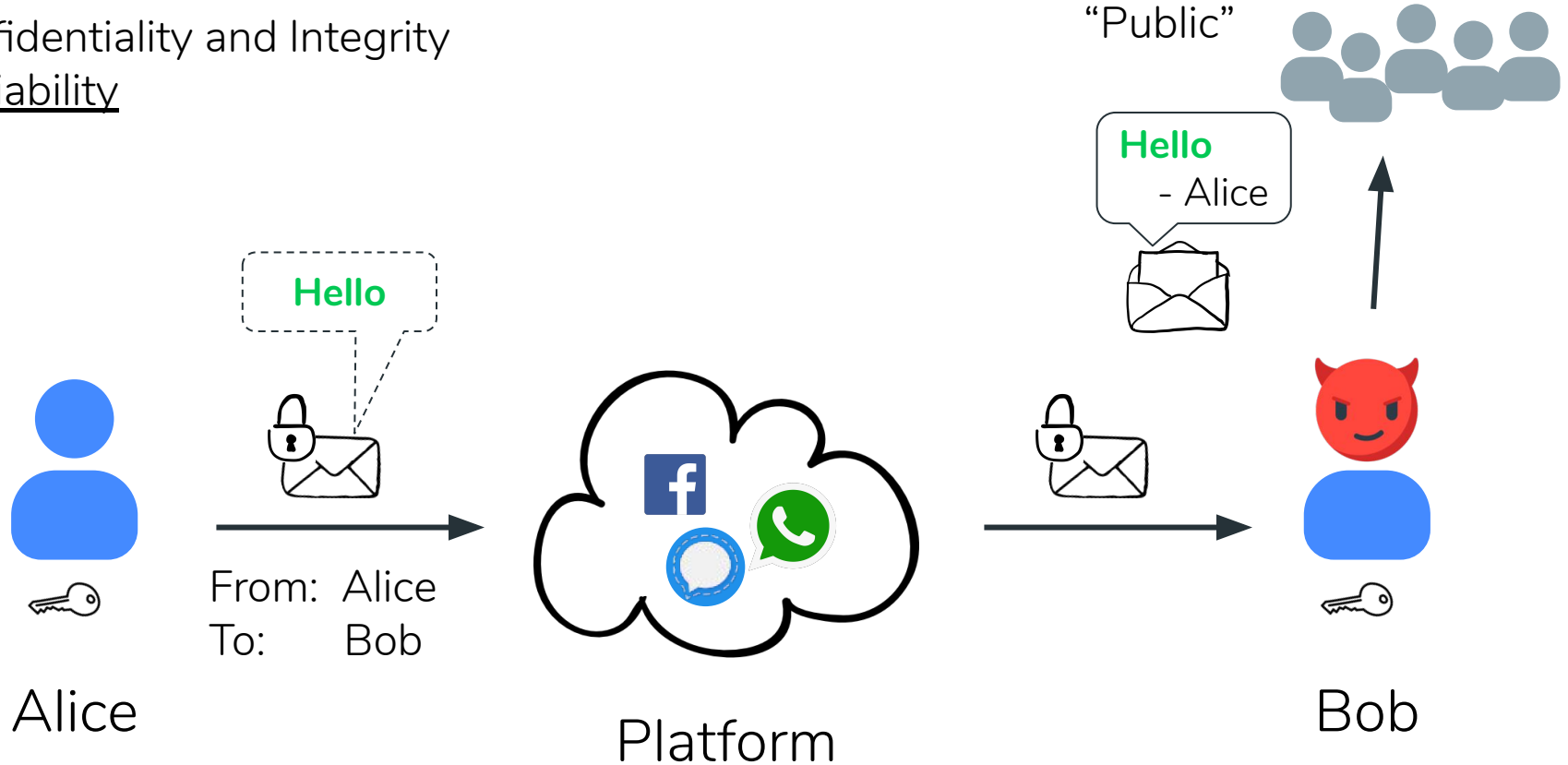
# Setting: End-to-end encrypted messaging

- Confidentiality and Integrity



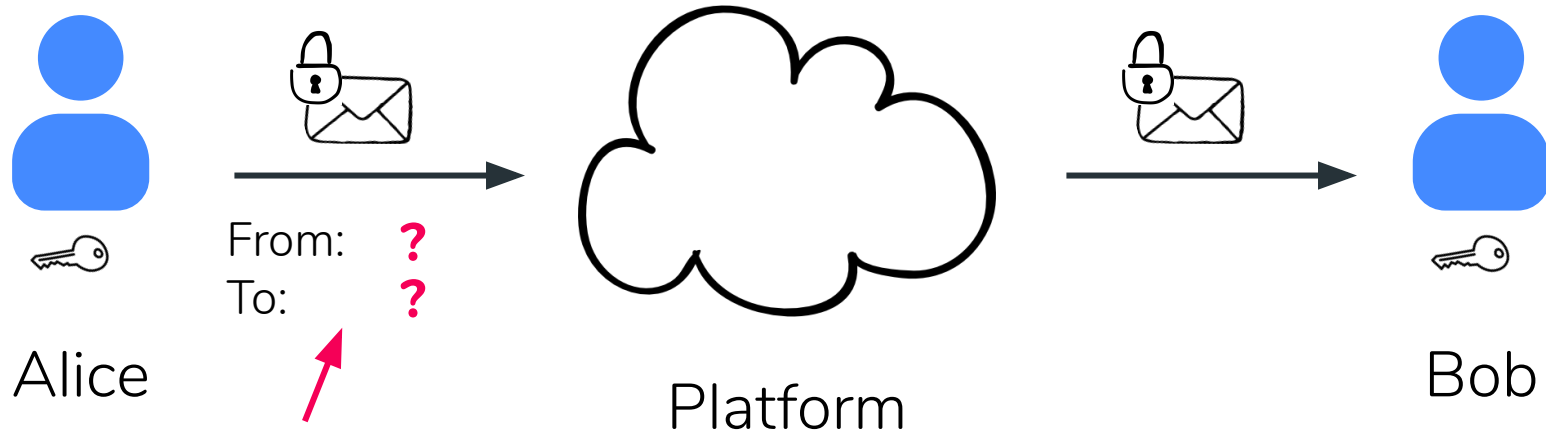
# Setting: End-to-end encrypted messaging

- Confidentiality and Integrity
- Deniability



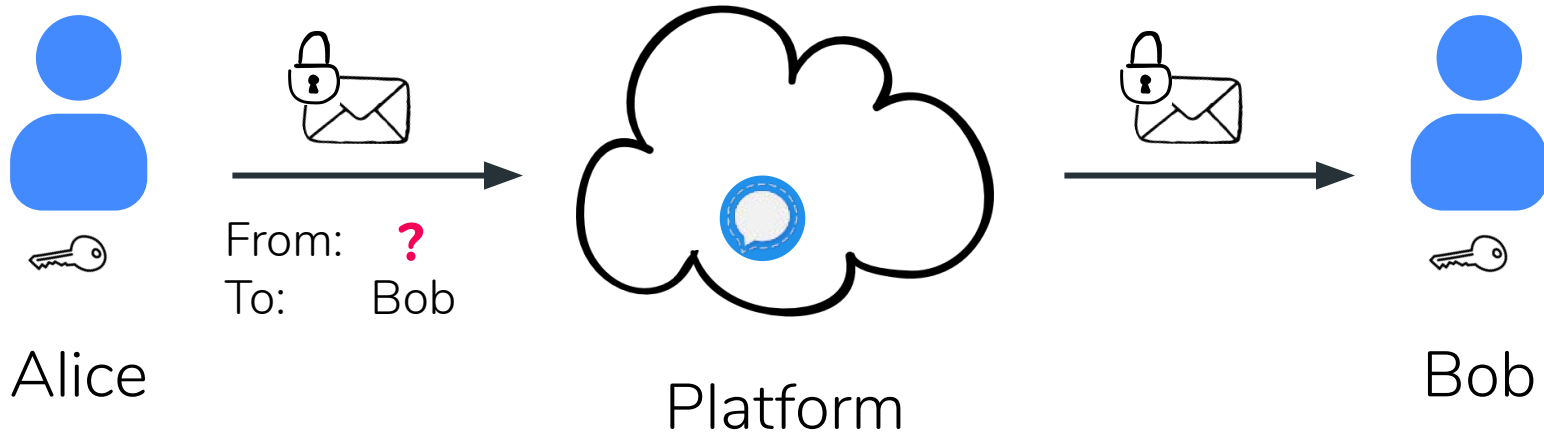
# Setting: End-to-end encrypted messaging

- Confidentiality and Integrity
- Deniability
- Metadata privacy

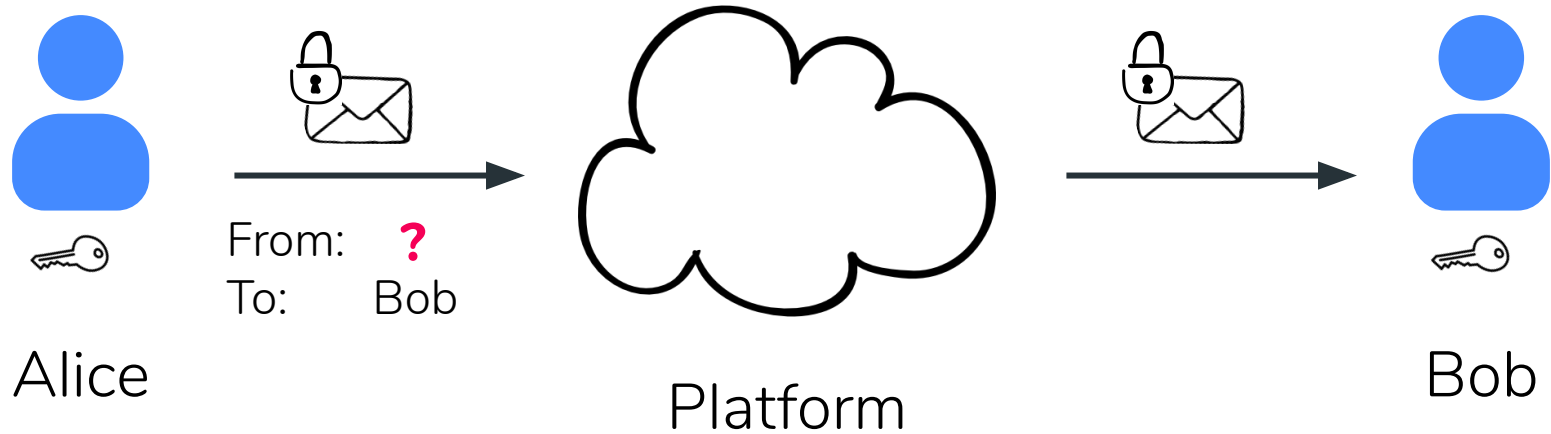


# Setting: End-to-end encrypted messaging

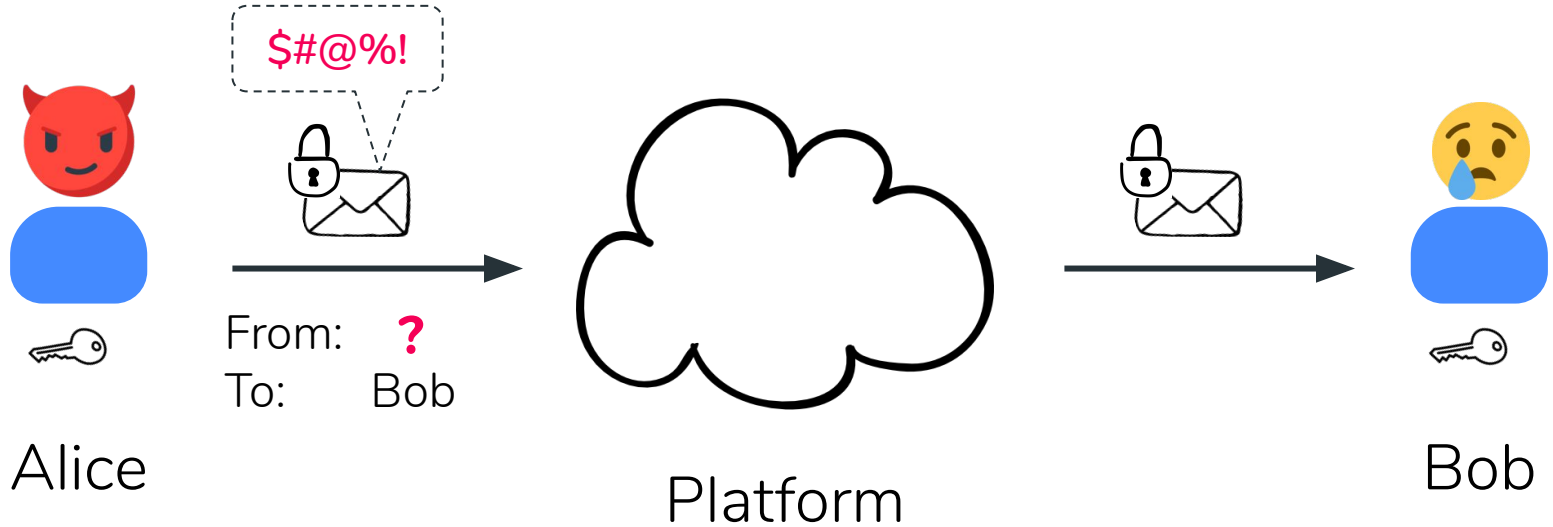
- Confidentiality and Integrity
- Deniability
- Metadata privacy



# What about abuse?

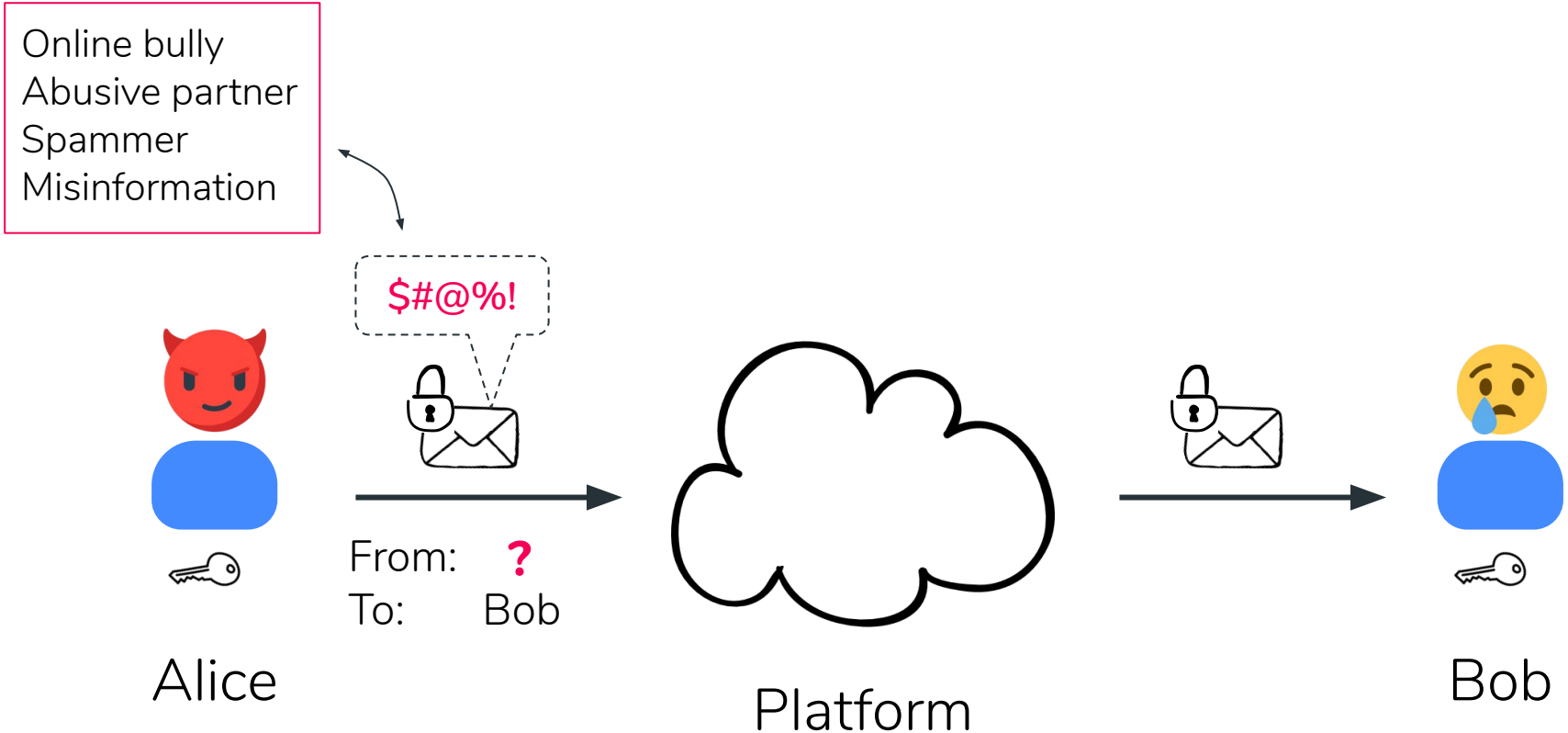


# What about abuse?

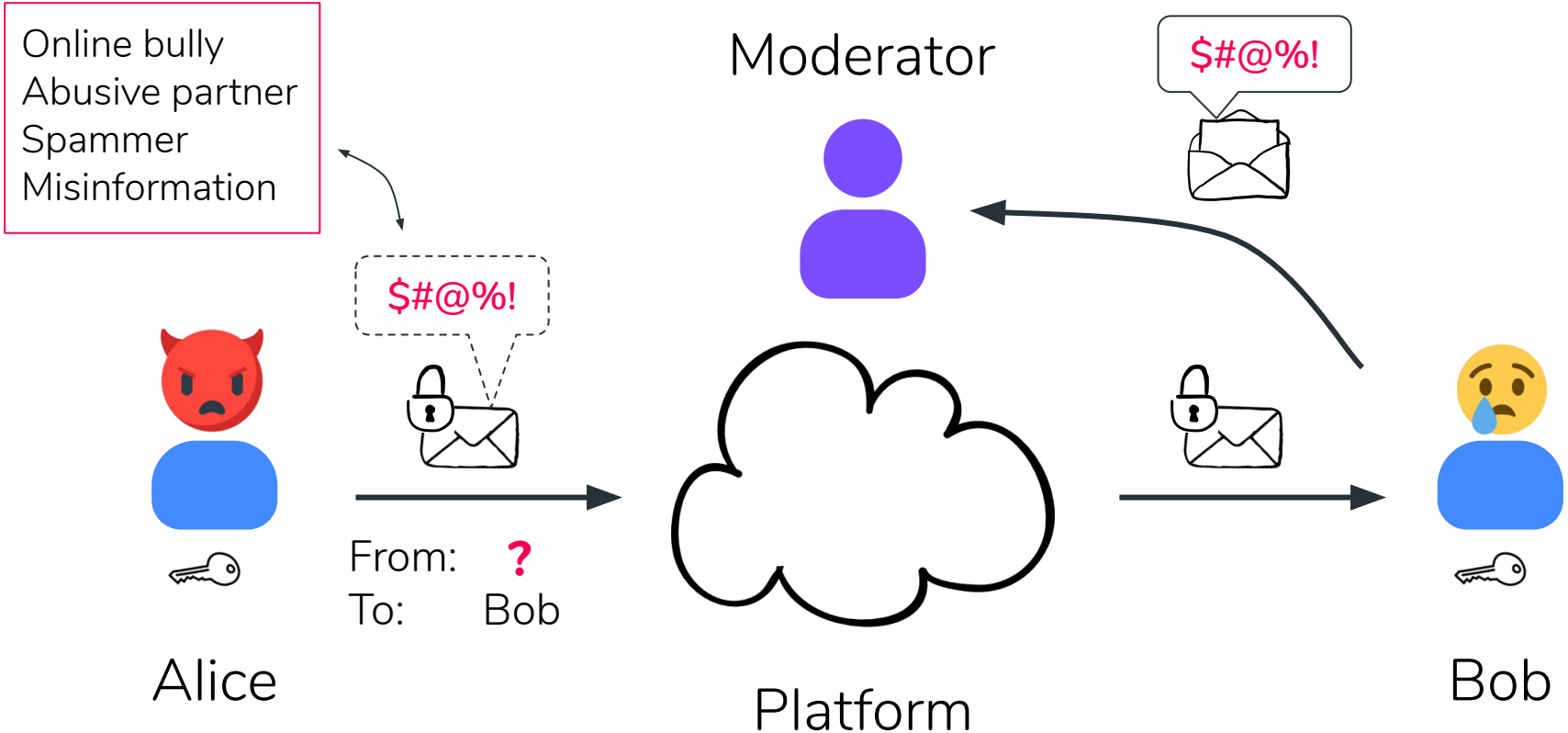




# What about abuse?



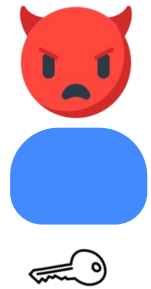
# What about abuse?



# What about abuse?

Moderation is a big priority:  
Facebook employs ≈15K content moderators\*

- Online bully
- Abusive partner
- Spammer
- Misinformation

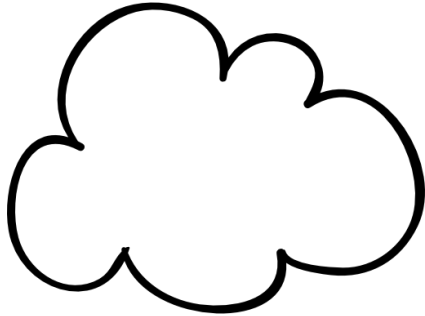
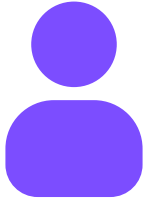


Alice

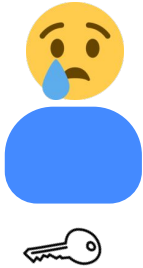
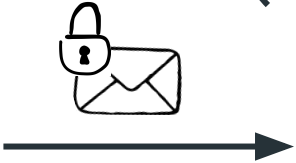


From: ?  
To: Bob

Moderator



Platform



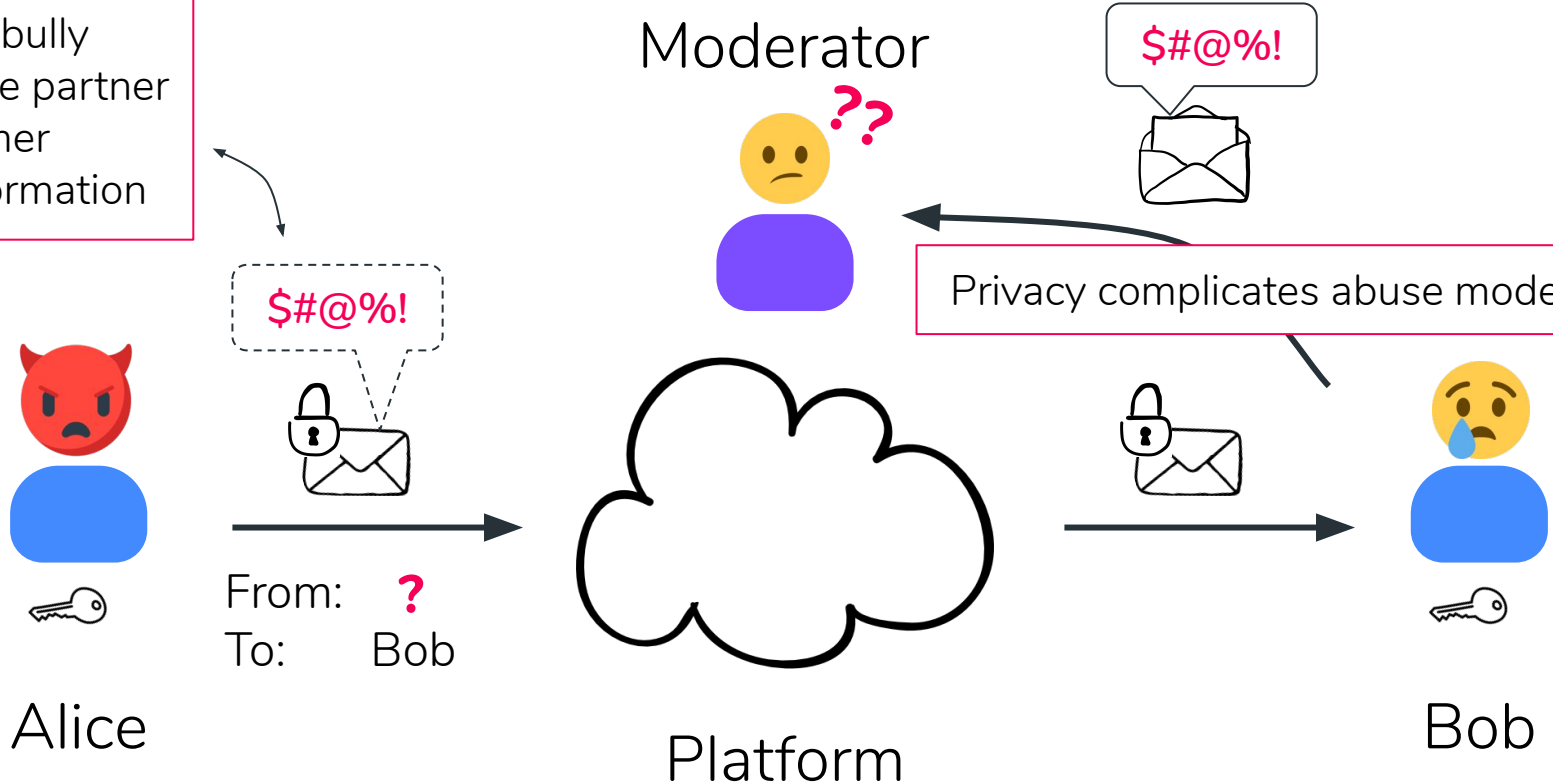
Bob

\* "The secret lives of Facebook moderators in America" [The Verge 2019]

# What about abuse?

Moderation is a big priority:  
Facebook employs ≈15K content moderators\*

- Online bully
- Abusive partner
- Spammer
- Misinformation



Privacy complicates abuse moderation!

\* “The secret lives of Facebook moderators in America” [The Verge 2019]

# What about abuse?

Moderation is a big priority:  
Facebook employs ≈15K content moderators\*

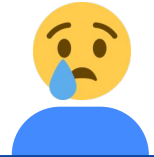
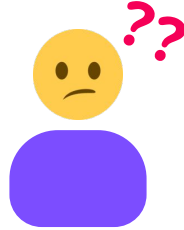
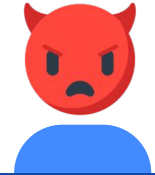
Online bully  
Abusive partner  
Spammer  
Misinformation

Moderator

\$#@%!

\$#@%!

Privacy complicates abuse moderation!



Can we balance need for **accountability** via moderation  
with **privacy** goals?

Platform

\* "The secret lives of Facebook moderators in America" [The Verge 2019]

# Our contributions

- **Asymmetric Message Franking (AMF)**: a new cryptographic primitive for content moderation
  - **Metadata-privacy**: message sender and/or recipient identities hidden
  - **Third-party moderation**: moderator decoupled from message-delivery platform
- Formal accountability and deniability security notions for content moderation
- Construction inspired by “designated-verifier” signatures
- Implementation and proof-of-concept deployment

# Prior work on moderation in E2E encryption

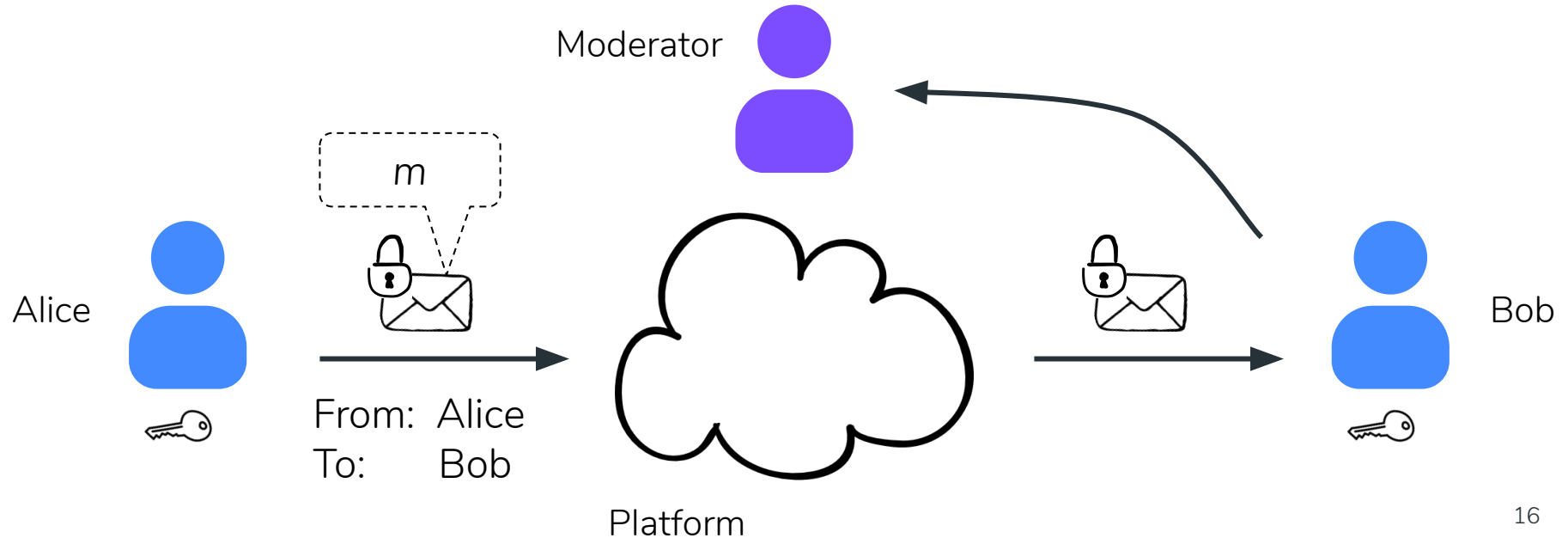
Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption

# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption

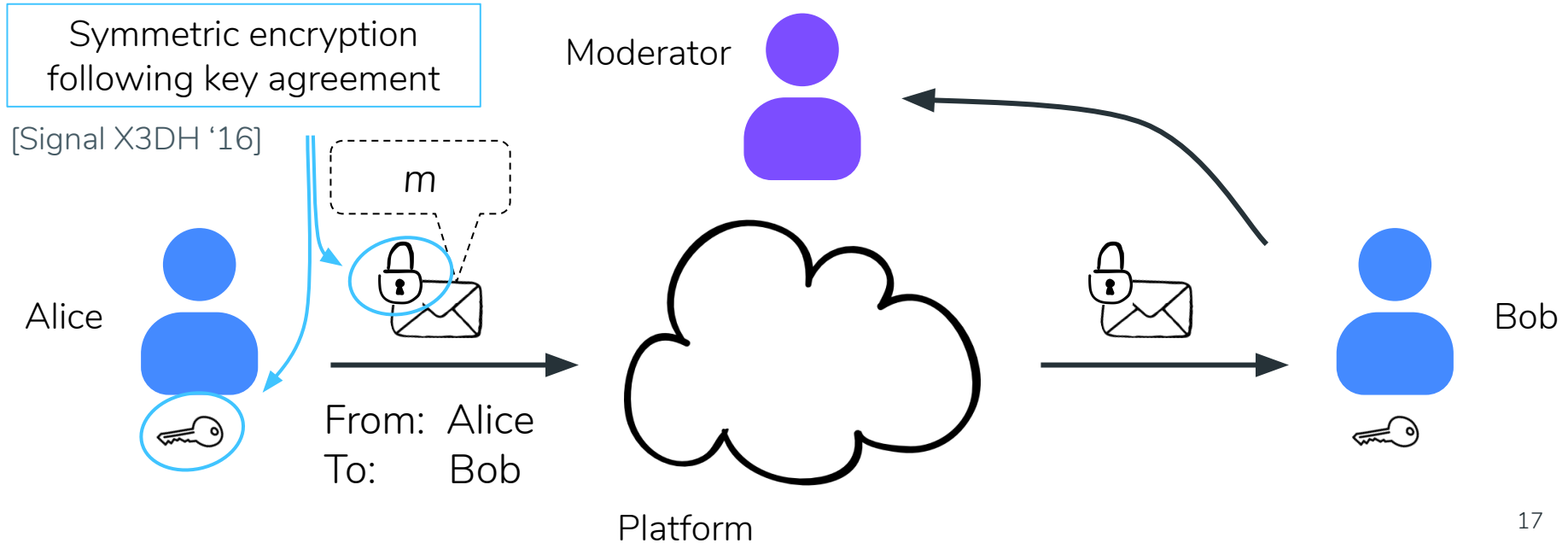




# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

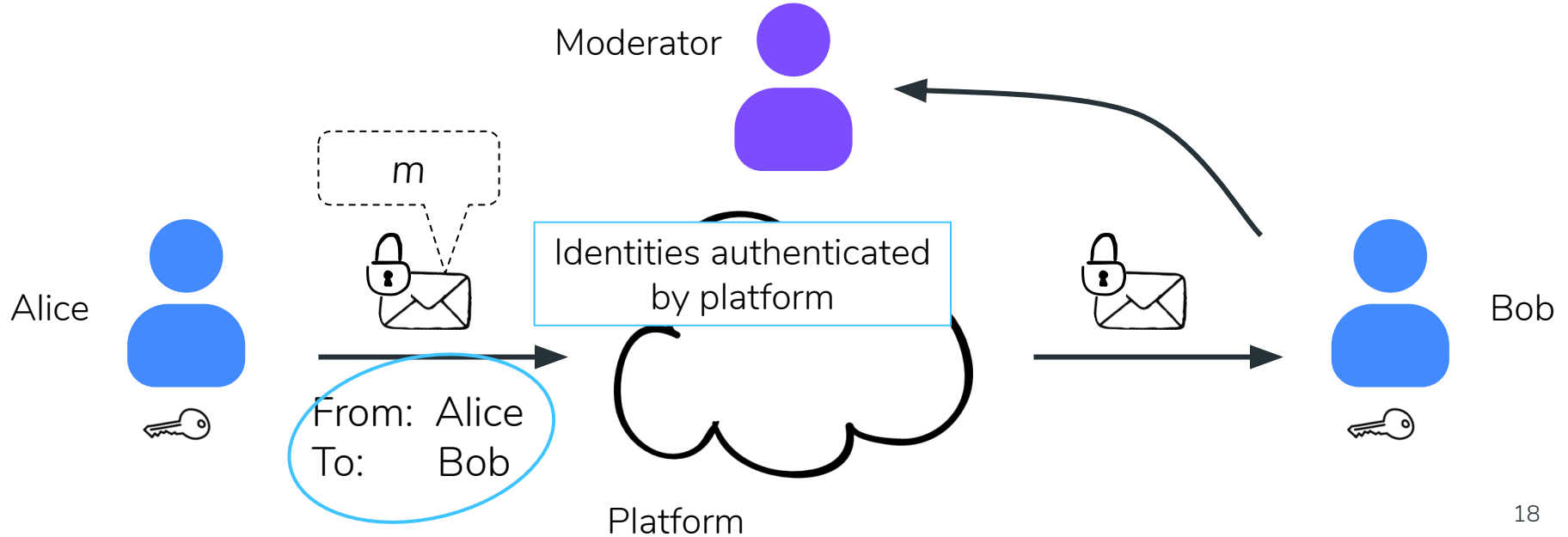
- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption



# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

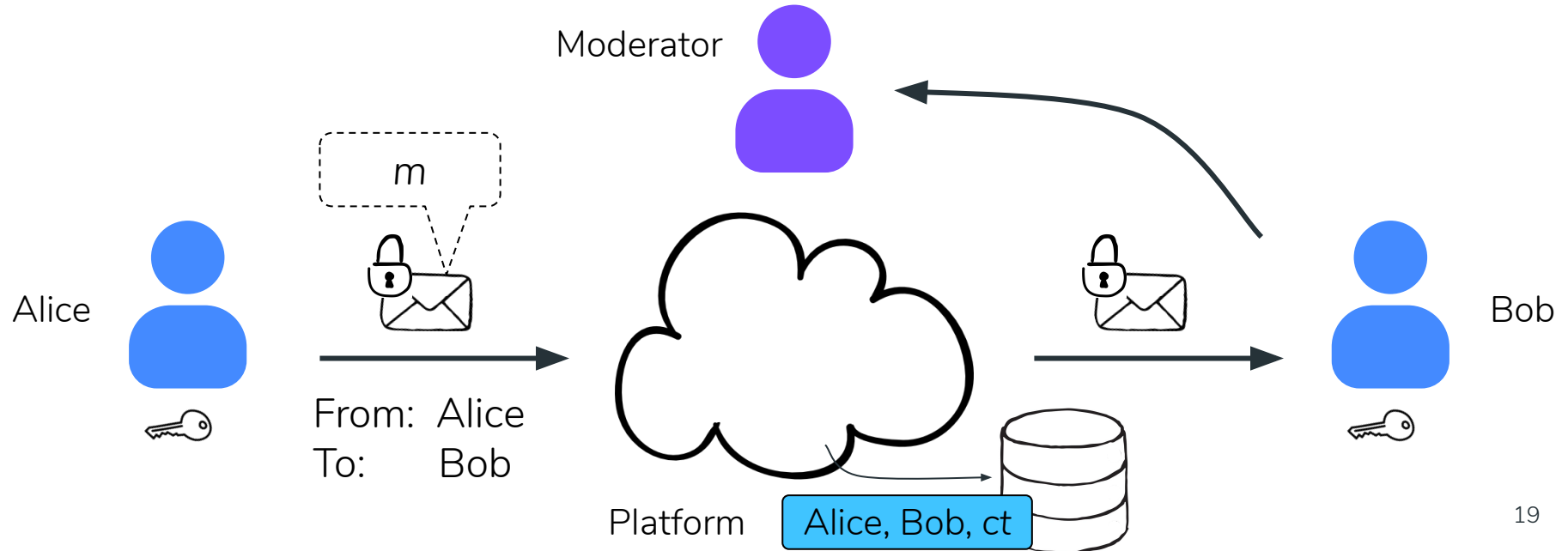
- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption



# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

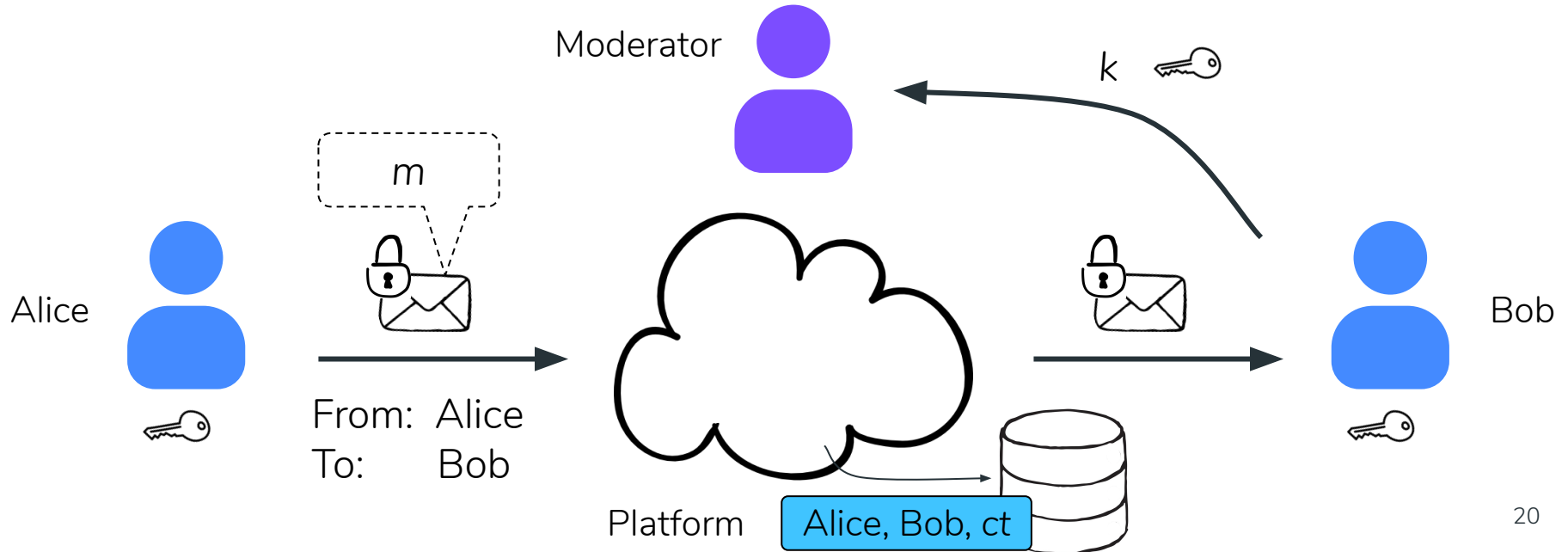
- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption



# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

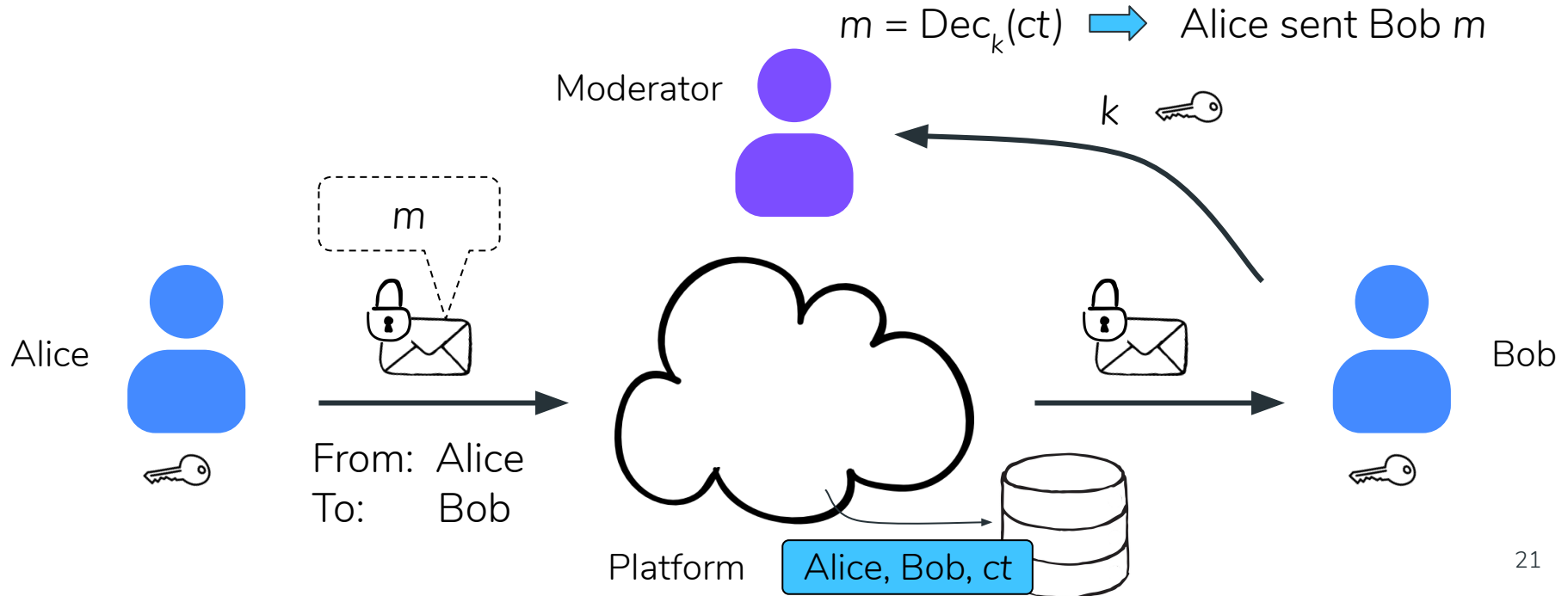
- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption



# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

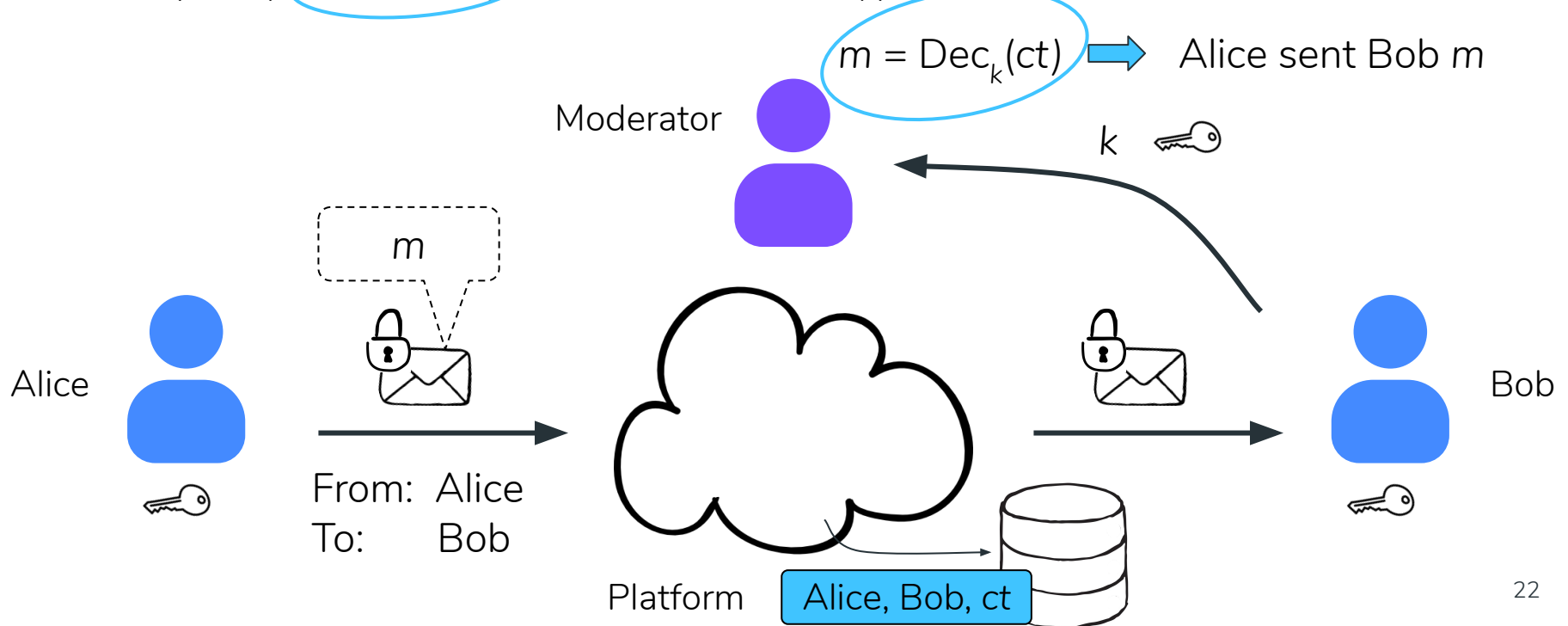
- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption



# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

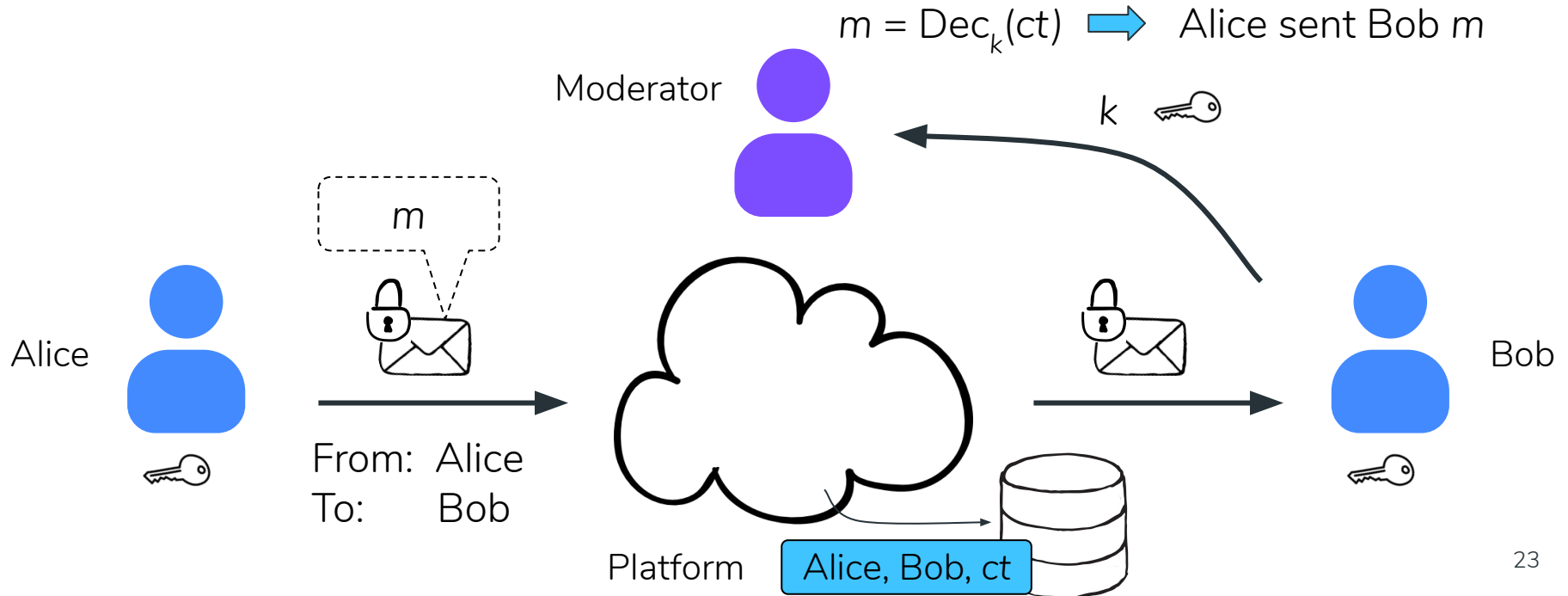
- Content-based moderation of encryption that is NOT metadata-private
- Compactly-committing authenticated encryption



# Prior work on moderation in E2E encryption

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

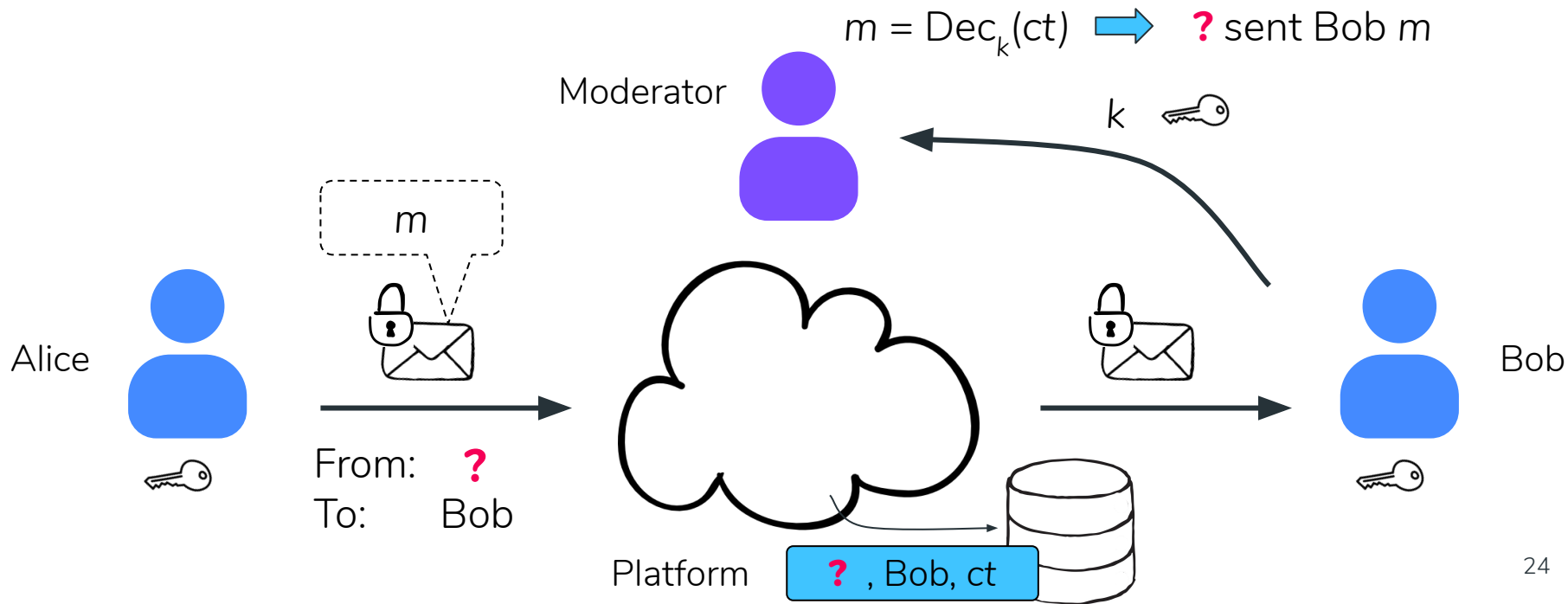
- Content-based moderation of encryption that is **NOT metadata-private**
- Compactly-committing authenticated encryption



# Message franking for metadata-private setting?

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is **NOT metadata-private**
- Compactly-committing authenticated encryption

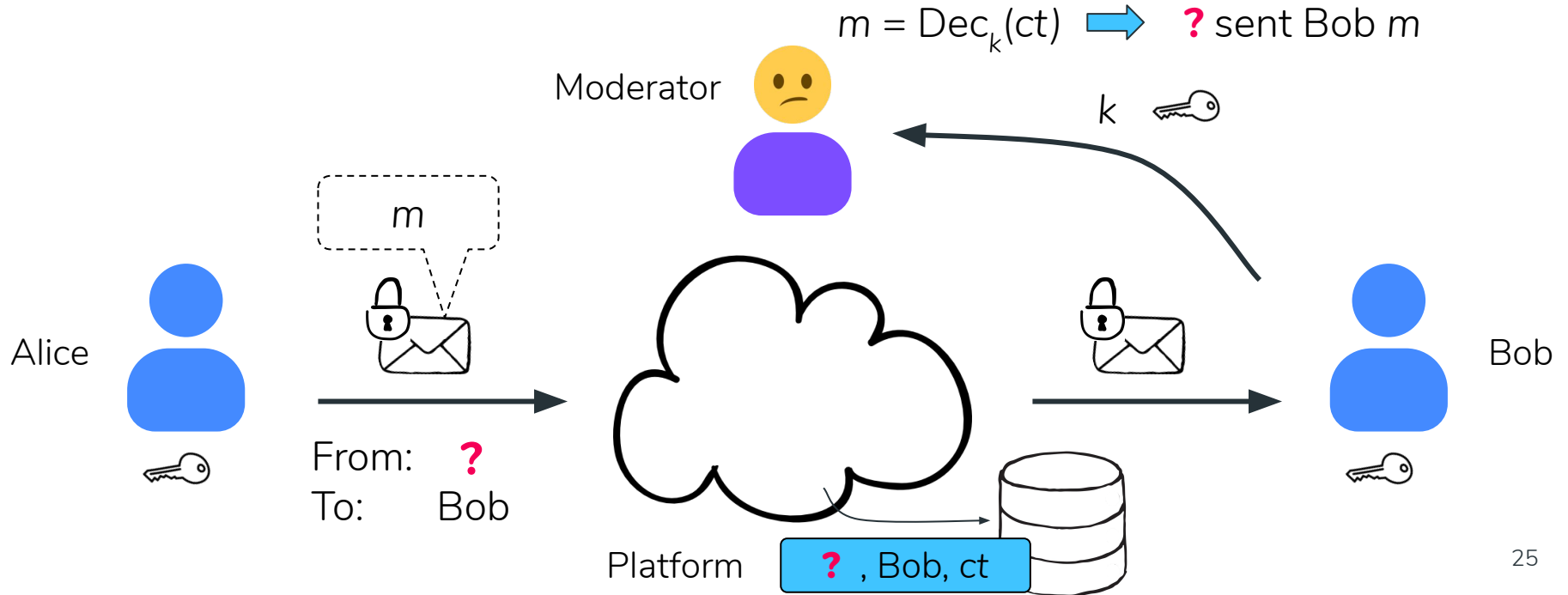




# Message franking for metadata-private setting?

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is **NOT metadata-private**
- Compactly-committing authenticated encryption

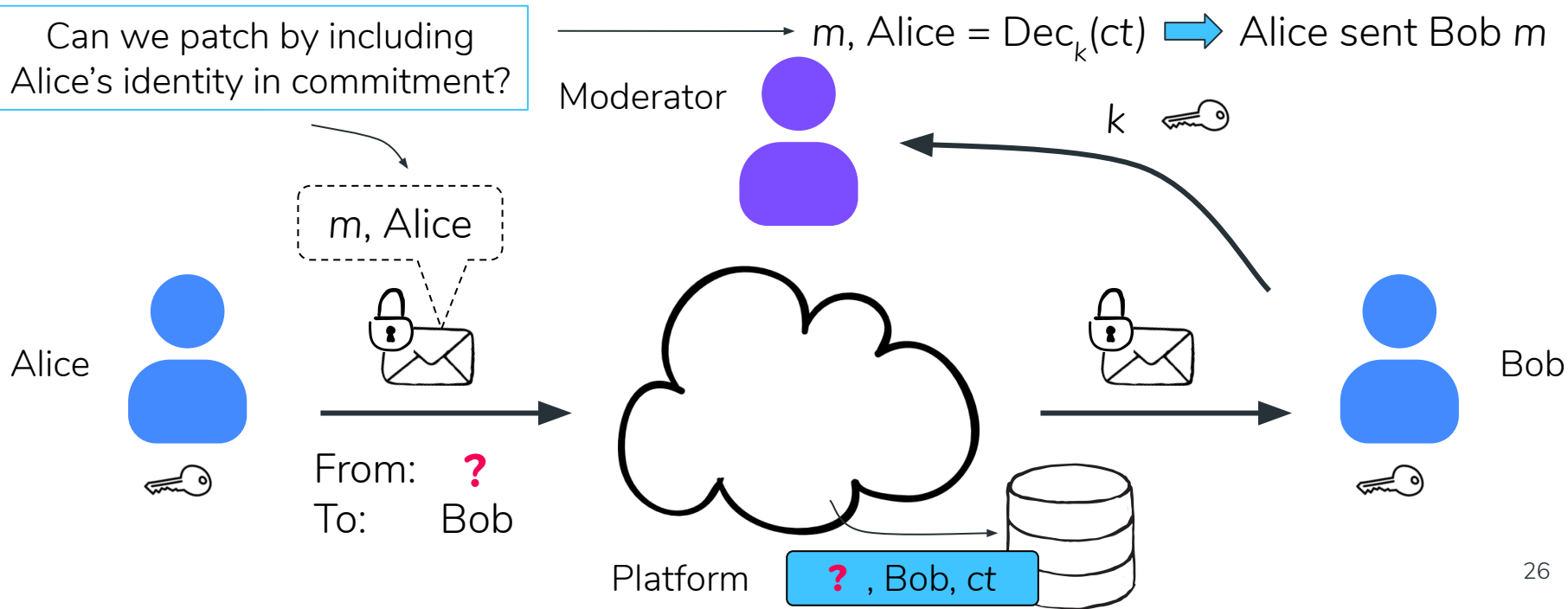


# Message franking for metadata-private setting?

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is **NOT metadata-private**
- Compactly-committing authenticated encryption

Can we patch by including Alice's identity in commitment?

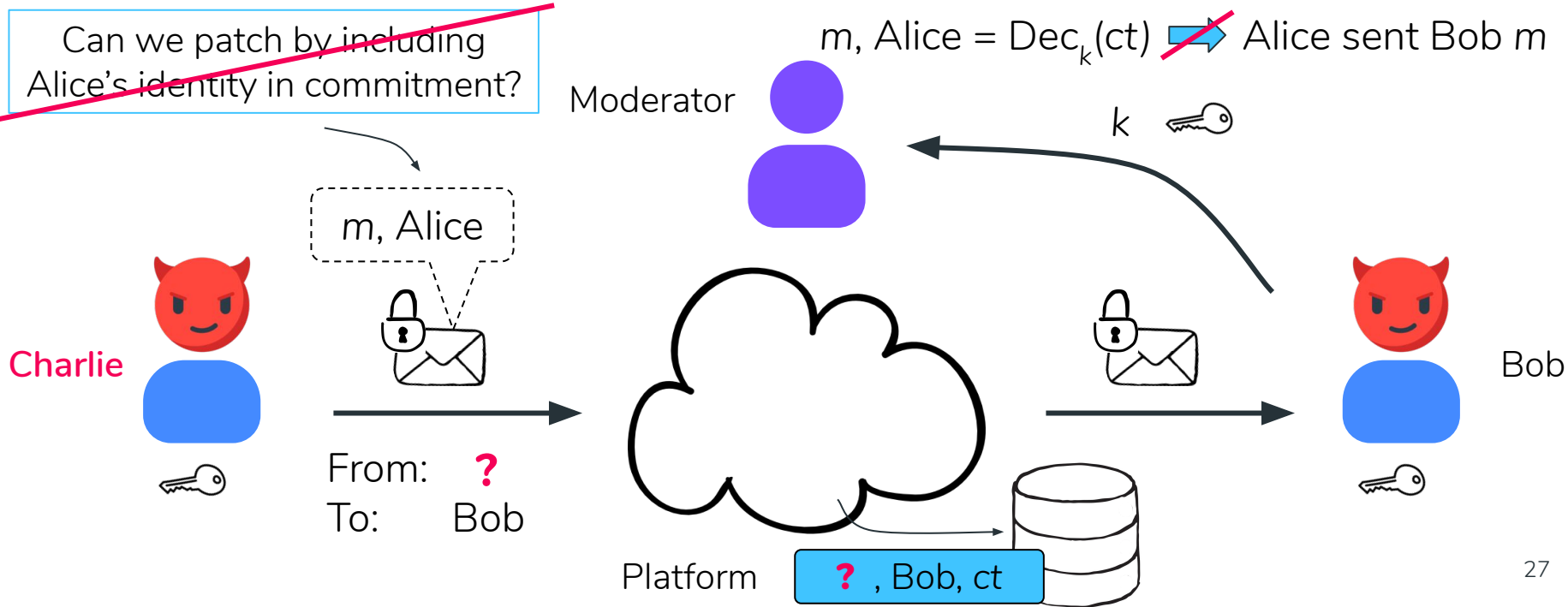


# Message franking for metadata-private setting?

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is **NOT metadata-private**
- Compactly-committing authenticated encryption

Can we patch by including Alice's identity in commitment?

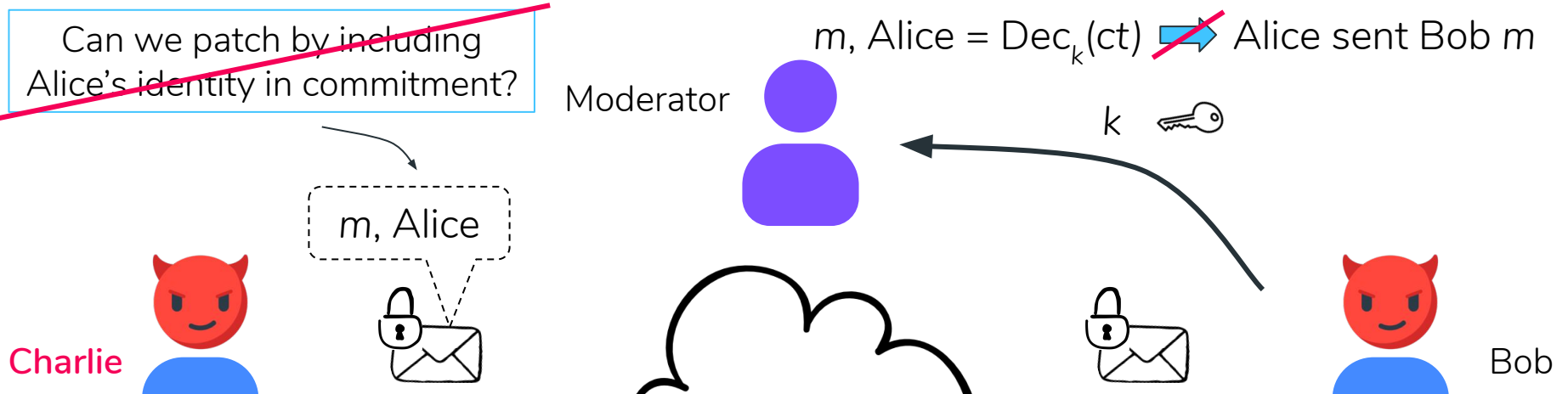


# Message franking for metadata-private setting?

Message franking [FB white paper '17], [GLR CRYPTO'17], [DGRW CRYPTO'18]

- Content-based moderation of encryption that is **NOT metadata-private**
- Compactly-committing authenticated encryption

~~Can we patch by including Alice's identity in commitment?~~



Core problem: Alice's identity not cryptographically bound to message content

# AMFs: High level idea

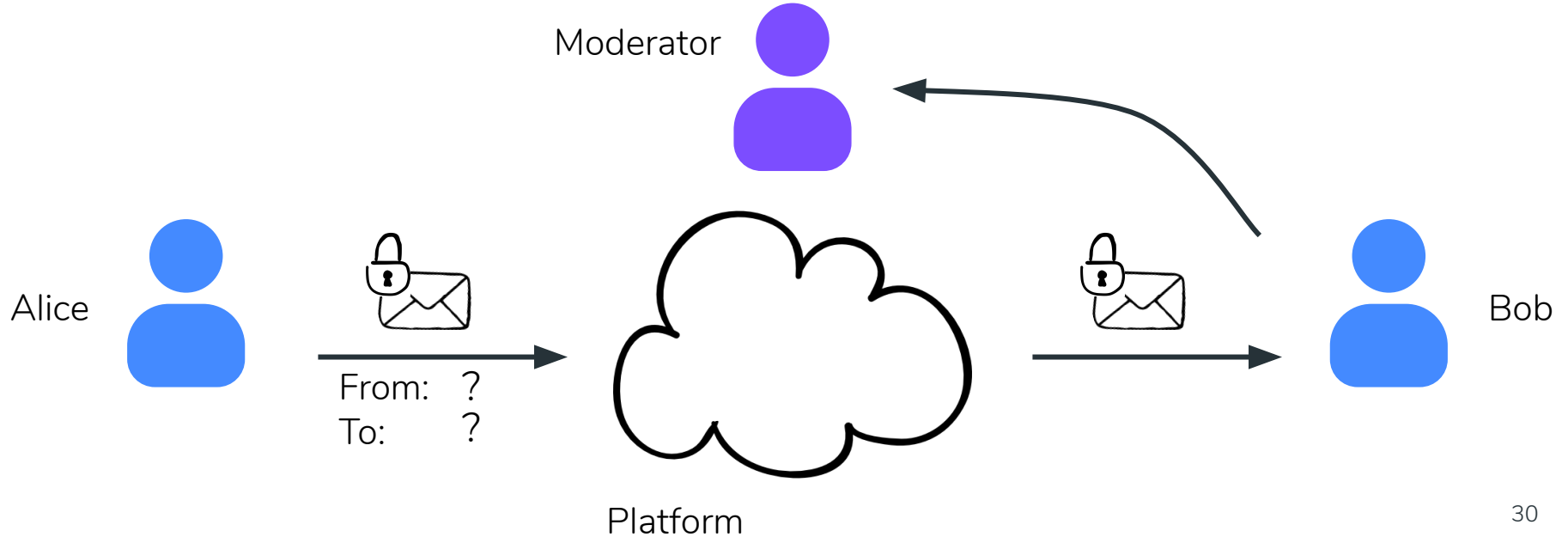
Specialized digital signature scheme that provides:

- Accountability
- Deniability

# AMFs: High level idea

Specialized digital signature scheme that provides:

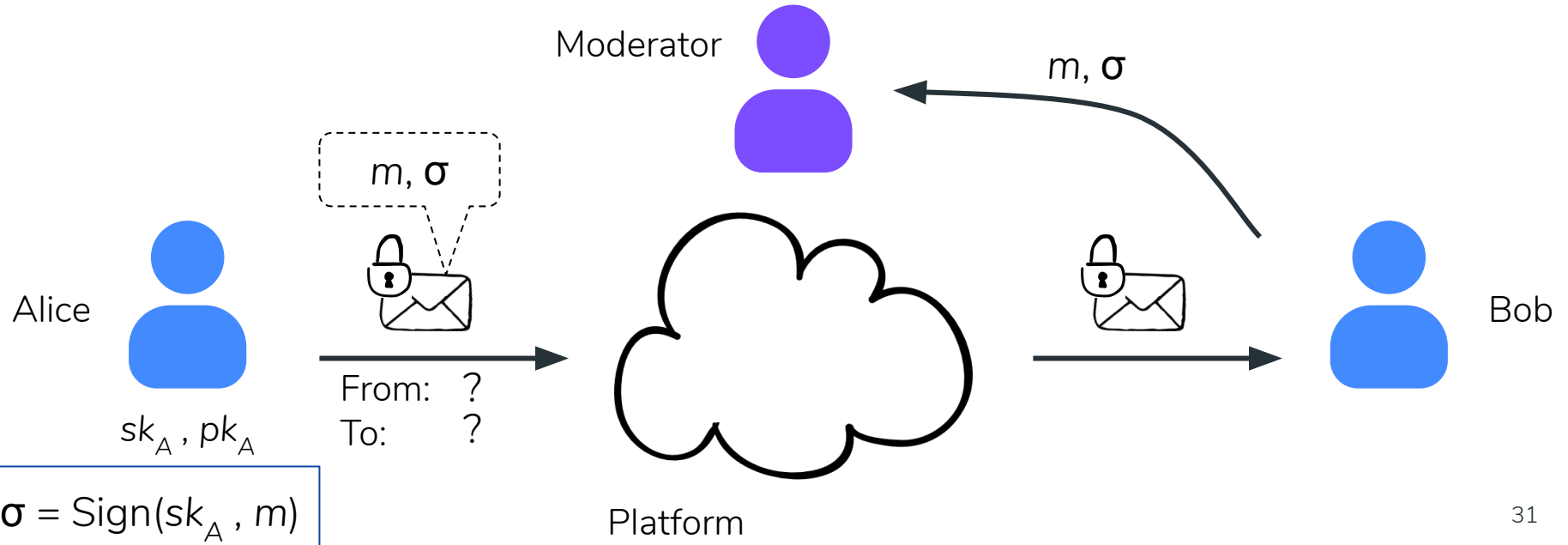
- Accountability
- Deniability



# AMFs: High level idea

Specialized digital signature scheme that provides:

- Accountability
- Deniability

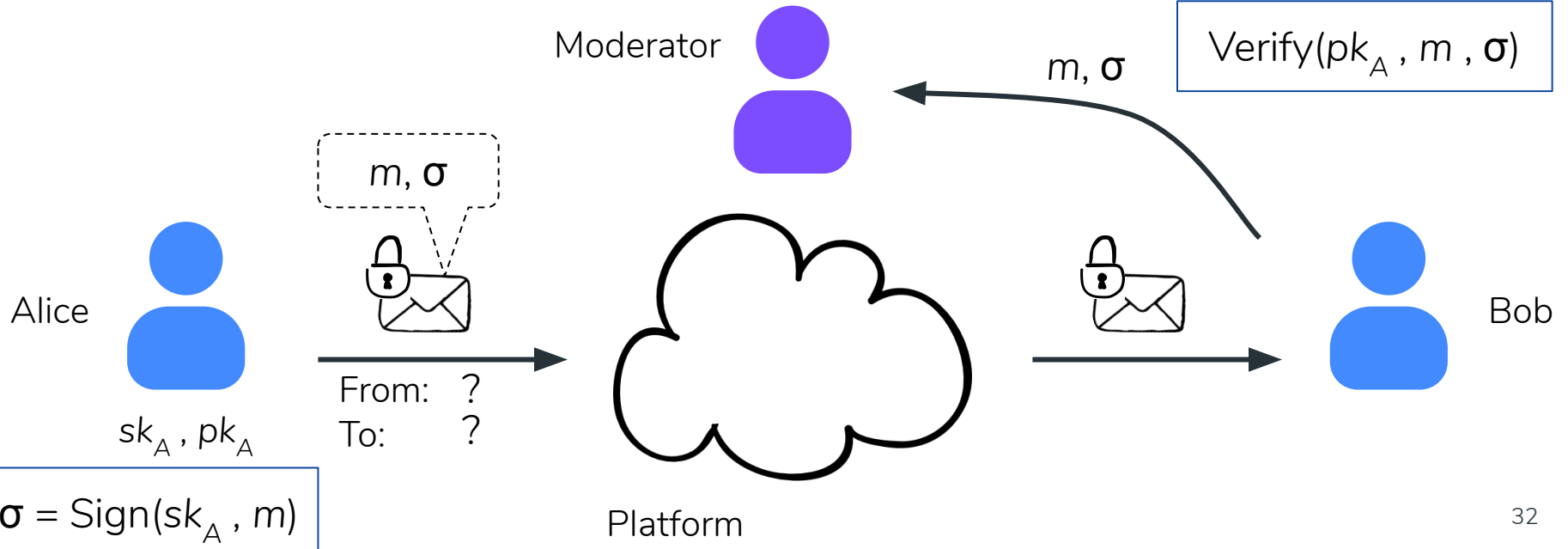


# AMFs: High level idea

Specialized digital signature scheme that provides:

- Accountability
- Deniability

Standard digital signatures provide accountability





# AMFs: High level idea

Specialized digital signature scheme that provides:

- Accountability
- Deniability

Standard digital signatures provide accountability ...but not deniability

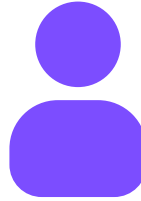
“Public”



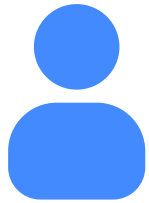
Verify( $pk_A, m, \sigma$ )

$m, \sigma$

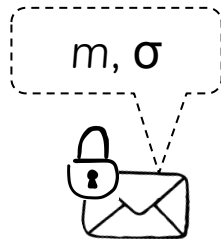
Moderator



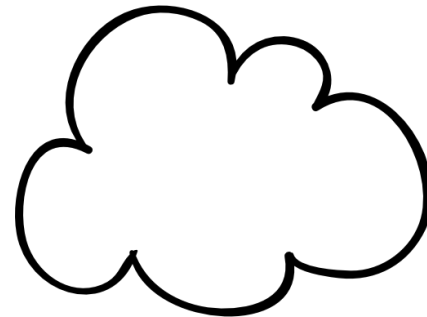
Alice



$sk_A, pk_A$



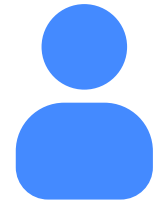
From: ?  
To: ?



Platform



Bob



$\sigma = \text{Sign}(sk_A, m)$

# Starting point: Designated-verifier signatures

Digital signatures where only one party can verify [JSI EUROCRYPT '96]

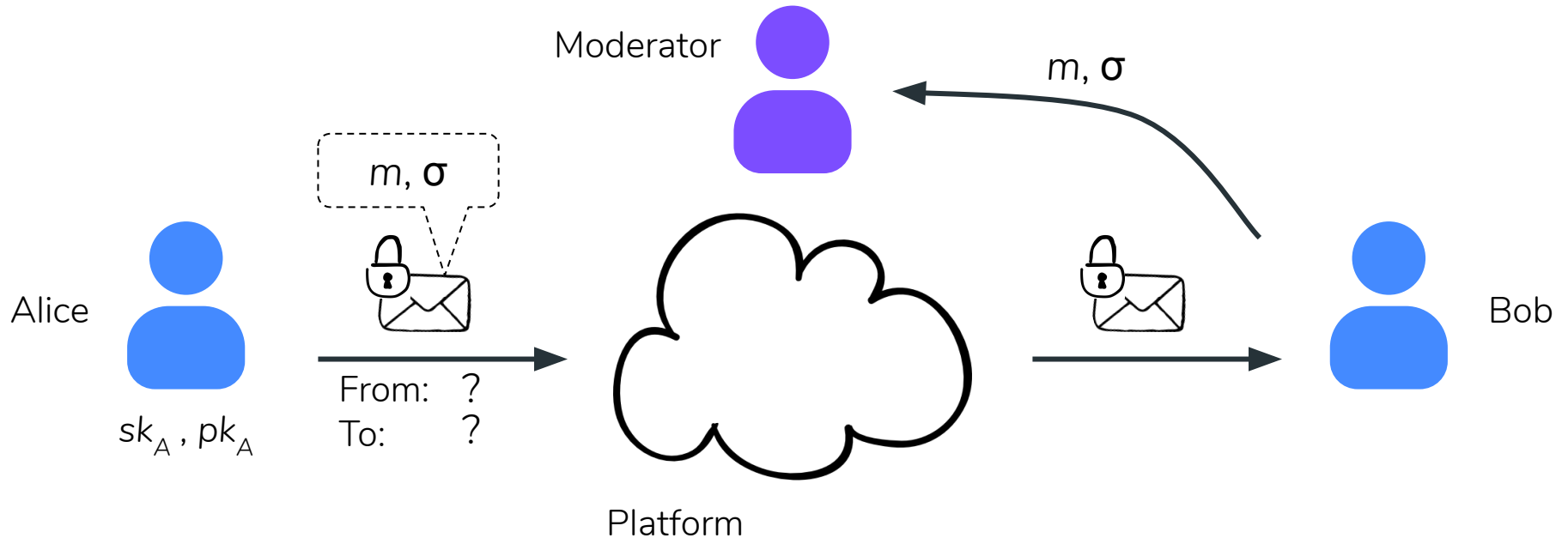
# Starting point: Designated-verifier signatures

Digital signatures where only one party can verify [JSI EUROCRYPT '96]

- Accountability  
Designated verifier can't be fooled by forgery
- Deniability  
There exists forgery algorithm that fools everyone else

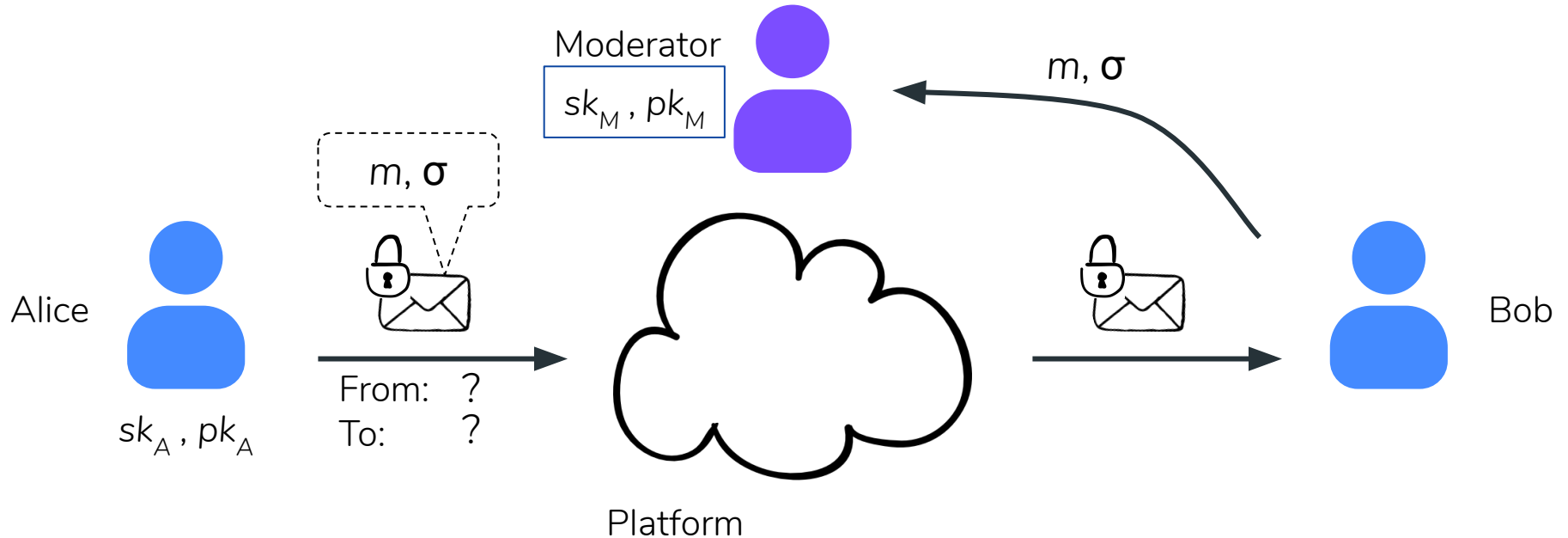
# Starting point: Designated-verifier signatures

Idea: Designating the moderator as a verifier?



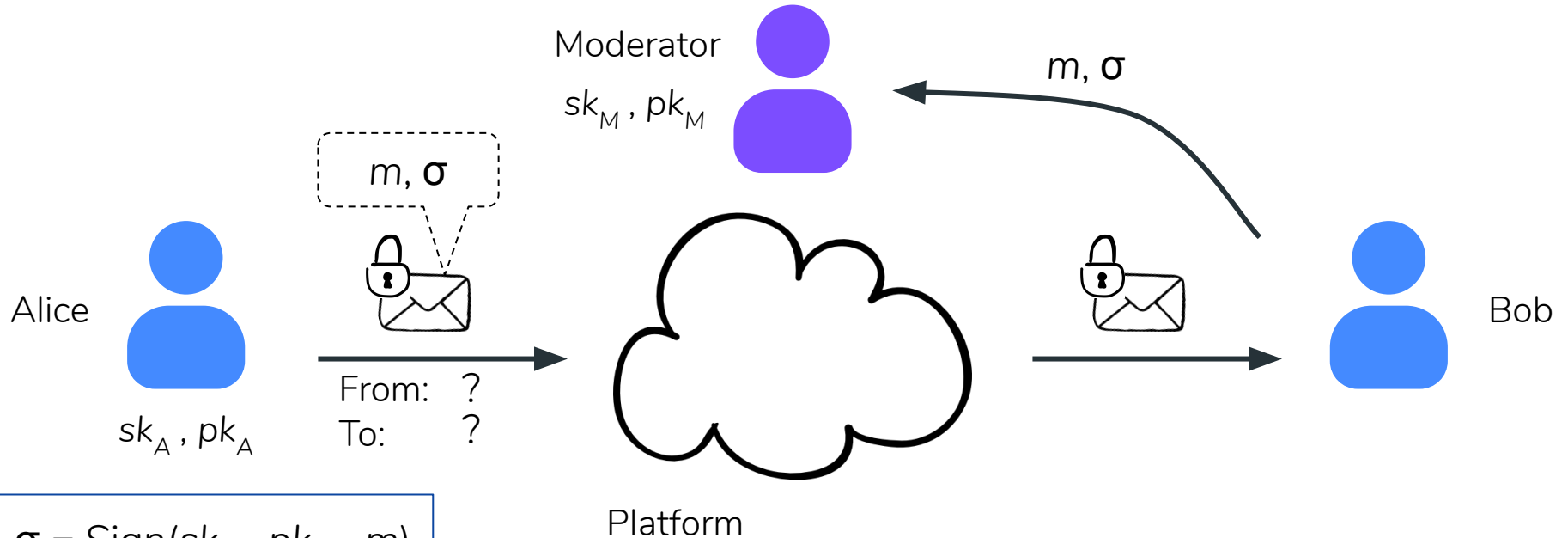
# Starting point: Designated-verifier signatures

Idea: Designating the moderator as a verifier?



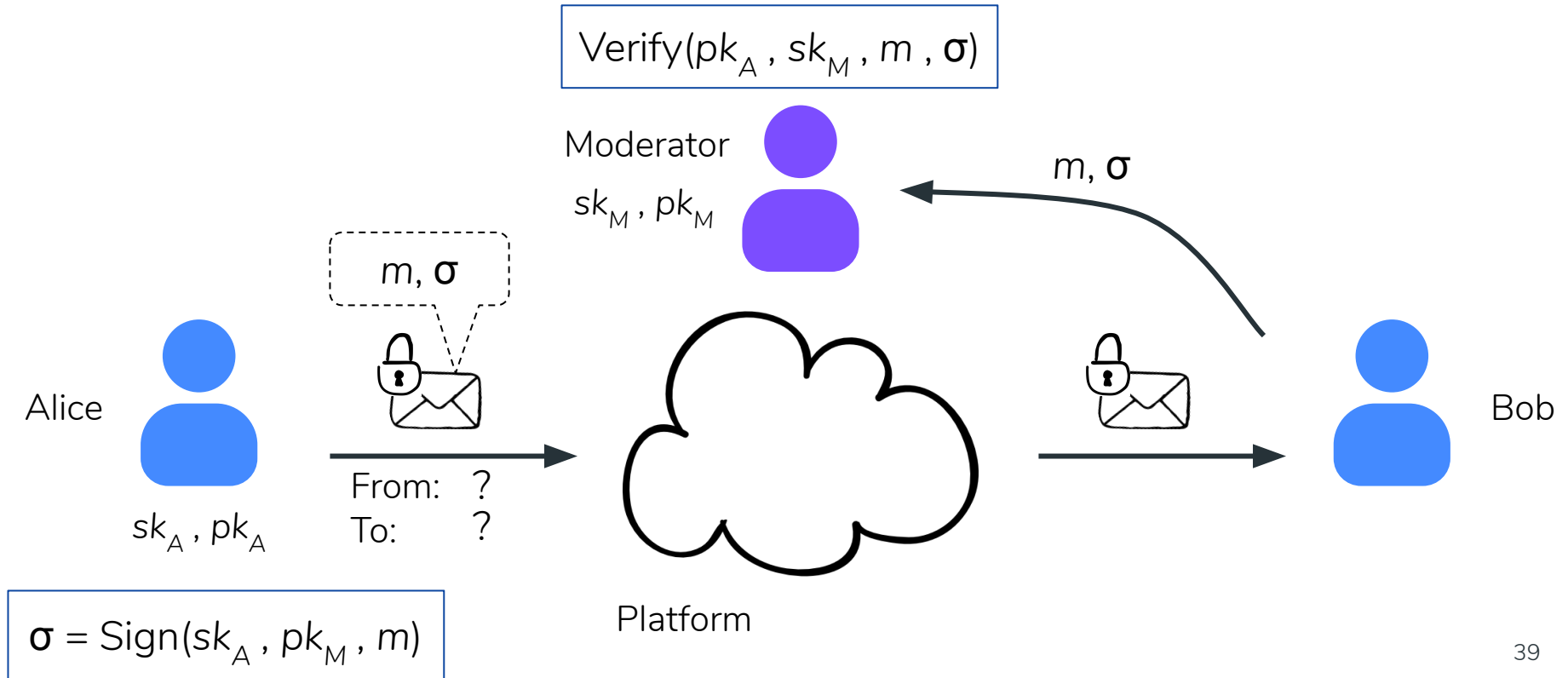
# Starting point: Designated-verifier signatures

Idea: Designating the moderator as a verifier?



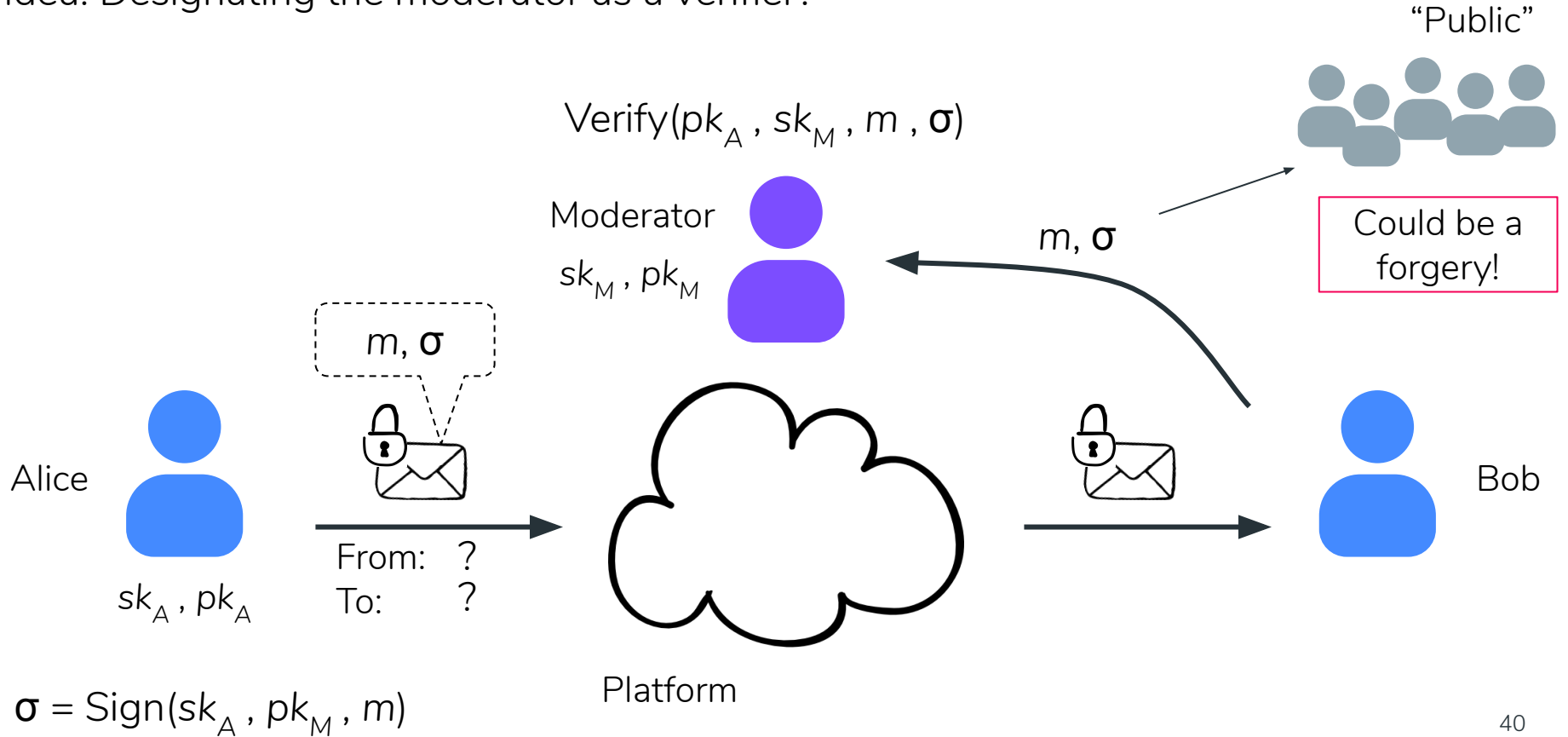
# Starting point: Designated-verifier signatures

Idea: Designating the moderator as a verifier?



# Starting point: Designated-verifier signatures

Idea: Designating the moderator as a verifier?

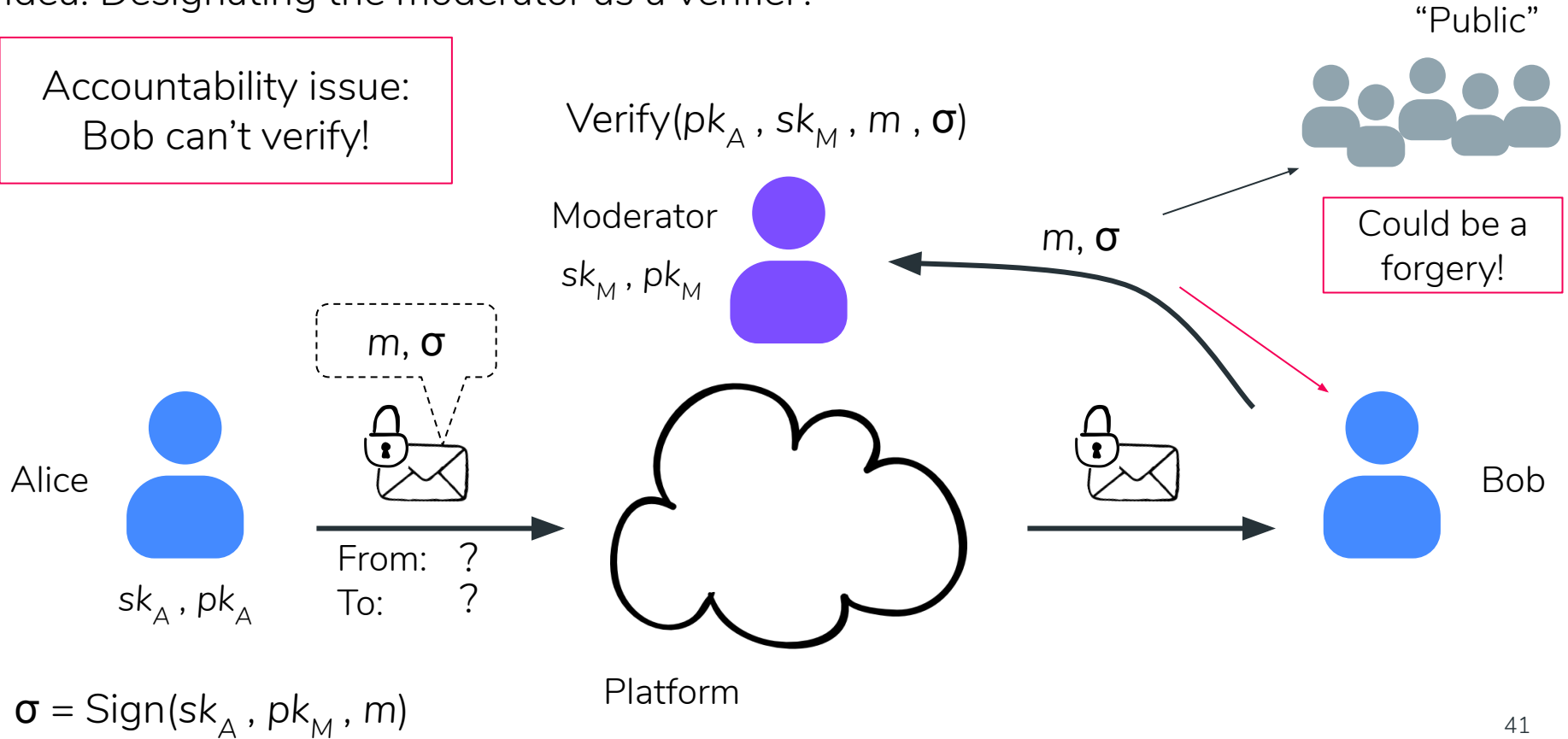




# Starting point: Designated-verifier signatures

Idea: Designating the moderator as a verifier?

Accountability issue:  
Bob can't verify!

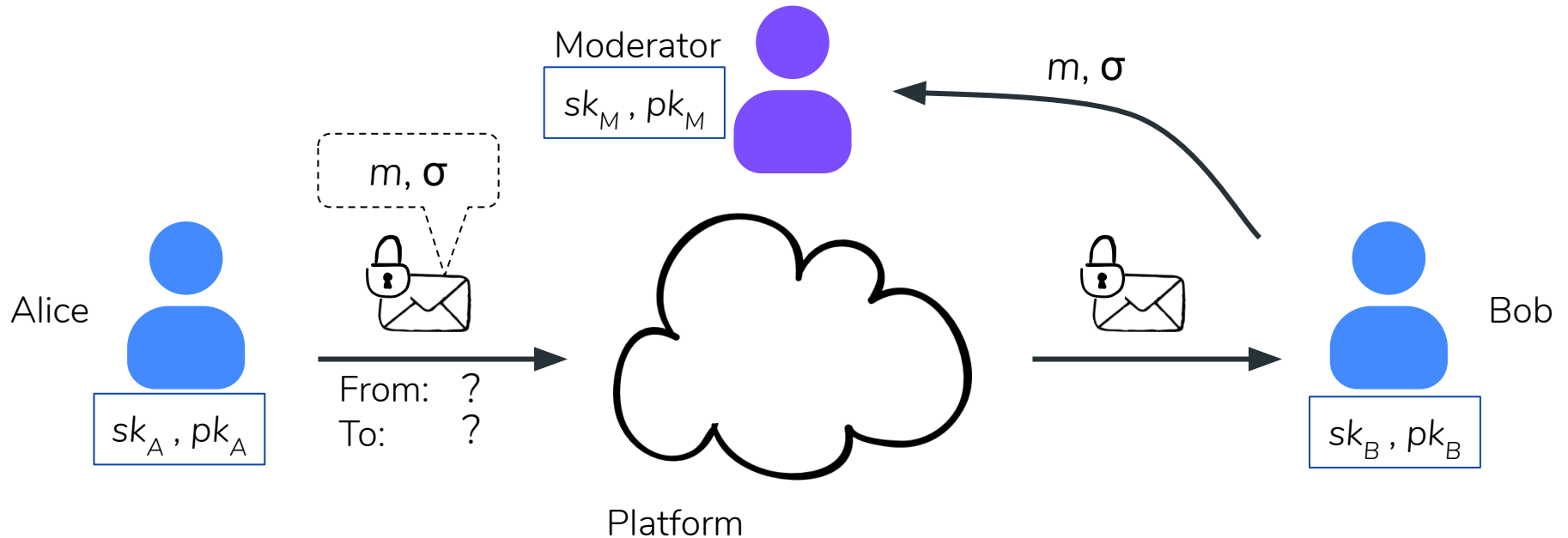


# AMFs: Include recipient as verifying party

Solution: Designate Bob as verifier of proof that signature to moderator will succeed

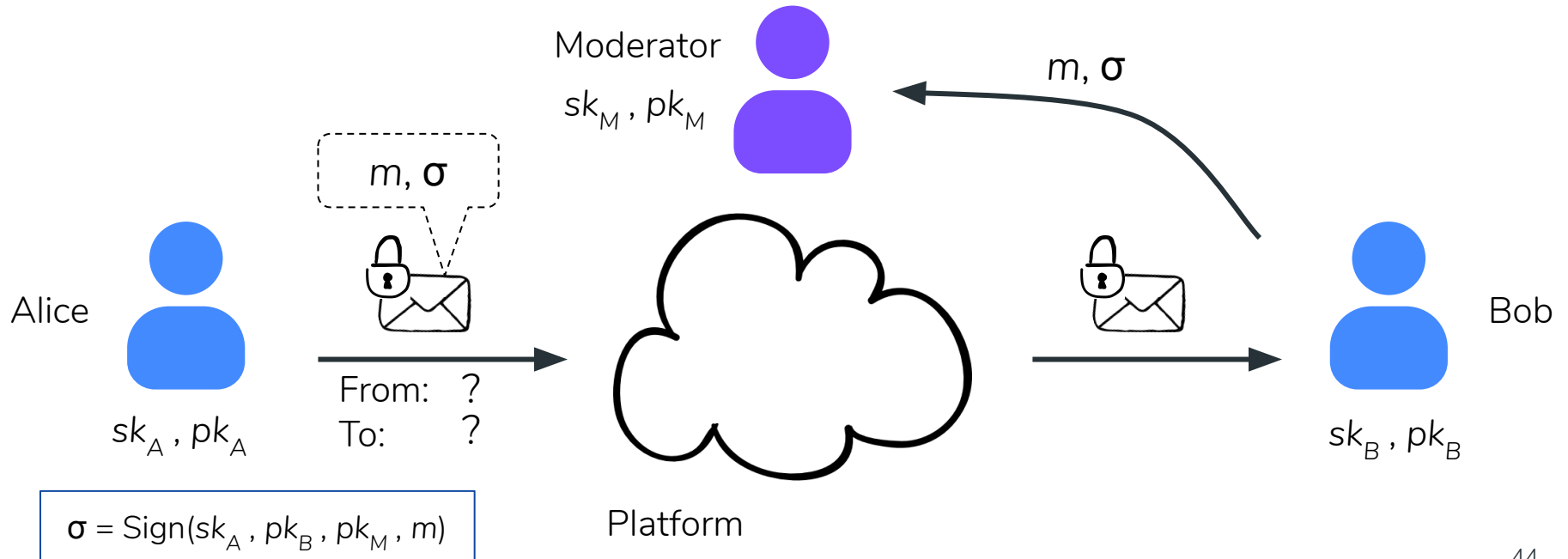
# AMFs: Include recipient as verifying party

Solution: Designate Bob as verifier of proof that signature to moderator will succeed



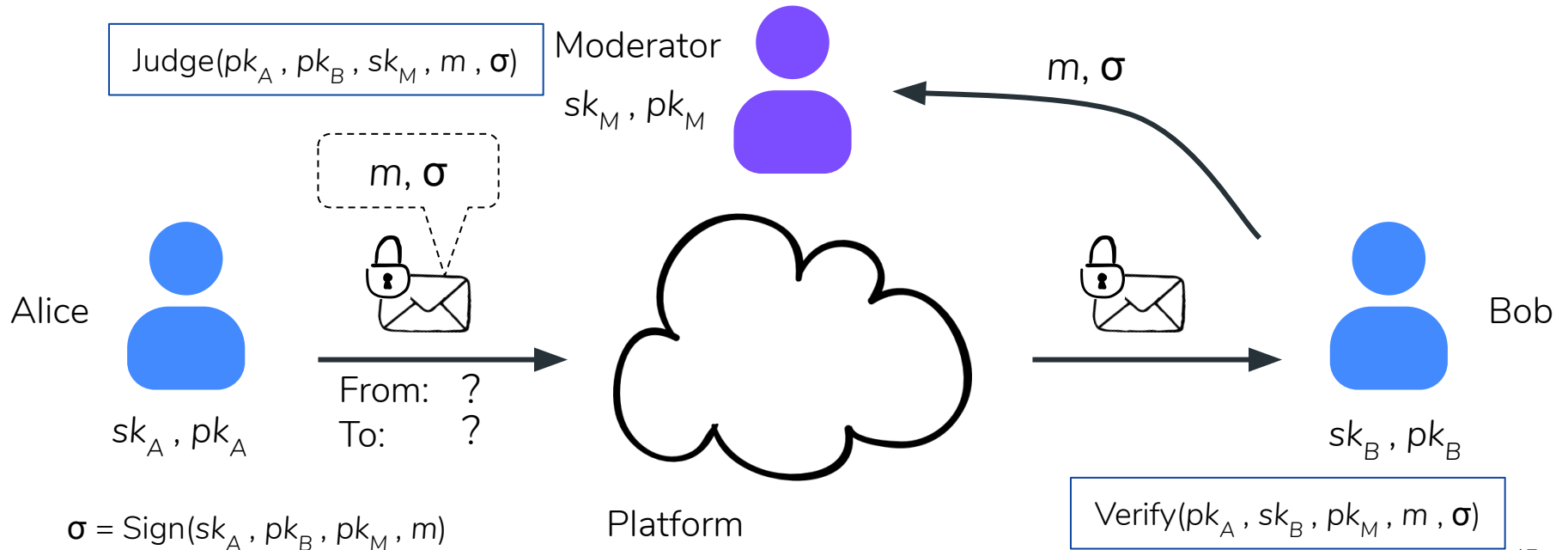
# AMFs: Include recipient as verifying party

Solution: Designate Bob as verifier of proof that signature to moderator will succeed



# AMFs: Include recipient as verifying party

Solution: Designate Bob as verifier of proof that signature to moderator will succeed

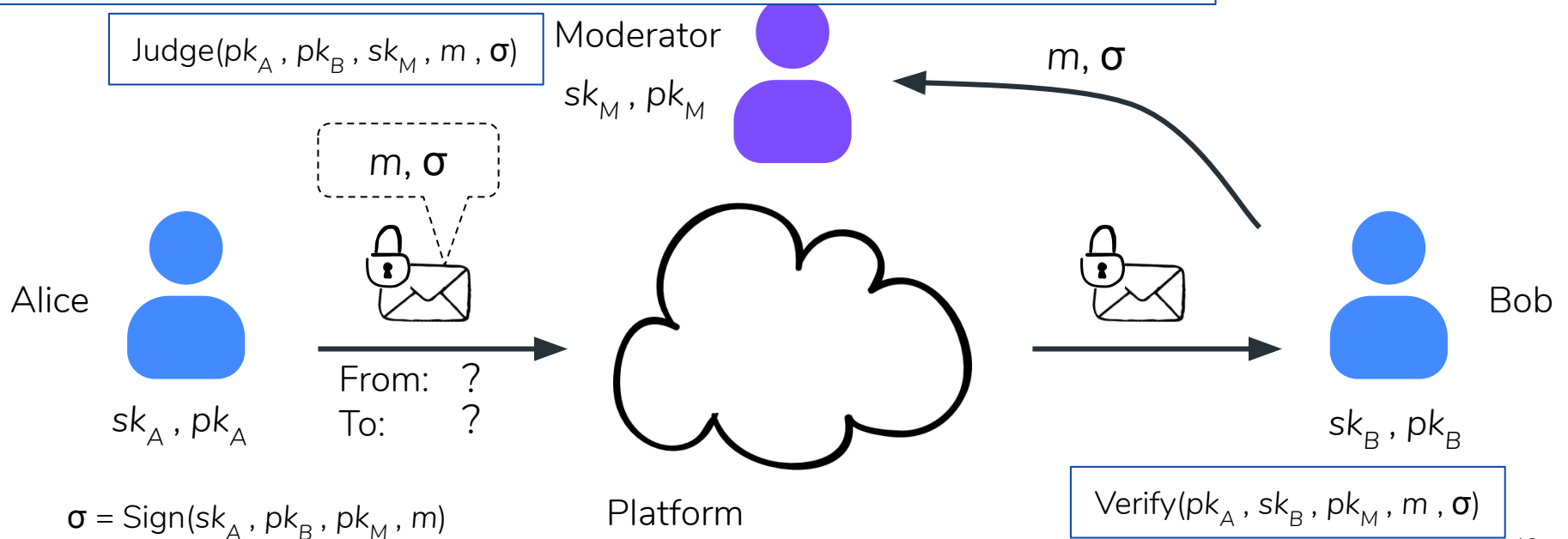


# AMFs: Include recipient as verifying party

Solution: Designate Bob as verifier of proof that signature to moderator will succeed

Accountability notions

- **Receiver binding:** Bob can't frame Alice for a message she did not send
- **Sender binding:** Alice can't send Bob a message that evades moderation



# Deniability landscape: “Who can trick whom?”

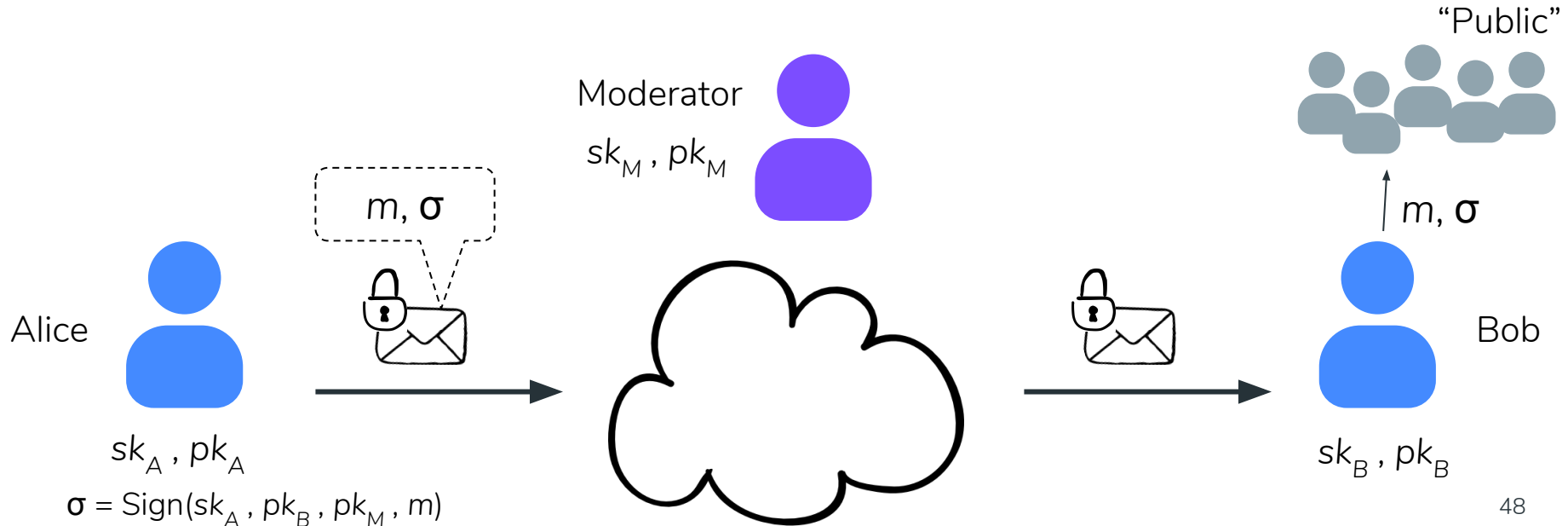


# Deniability landscape: “Who can trick whom?”

Forger

Distinguisher  $D$

$$\sigma' = \text{Forge}(pk_A, sk_B, pk_M, m) \xrightarrow{\sigma \approx_D \sigma'} pk_A, pk_B, pk_M$$



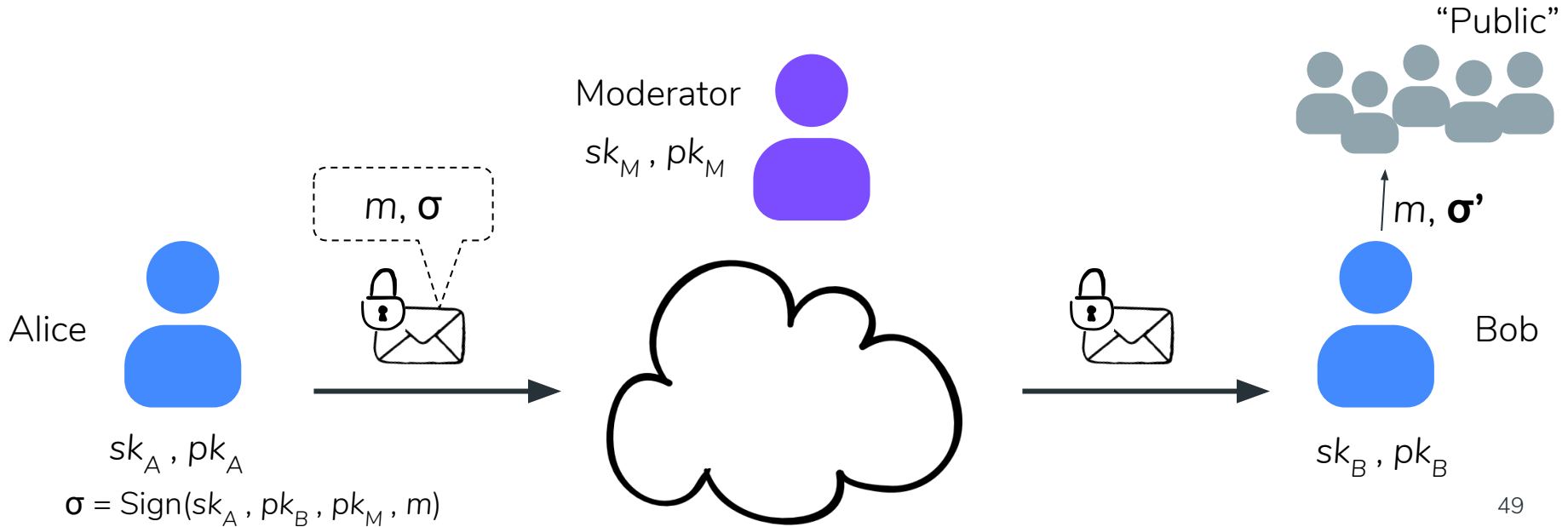


# Deniability landscape: “Who can trick whom?”

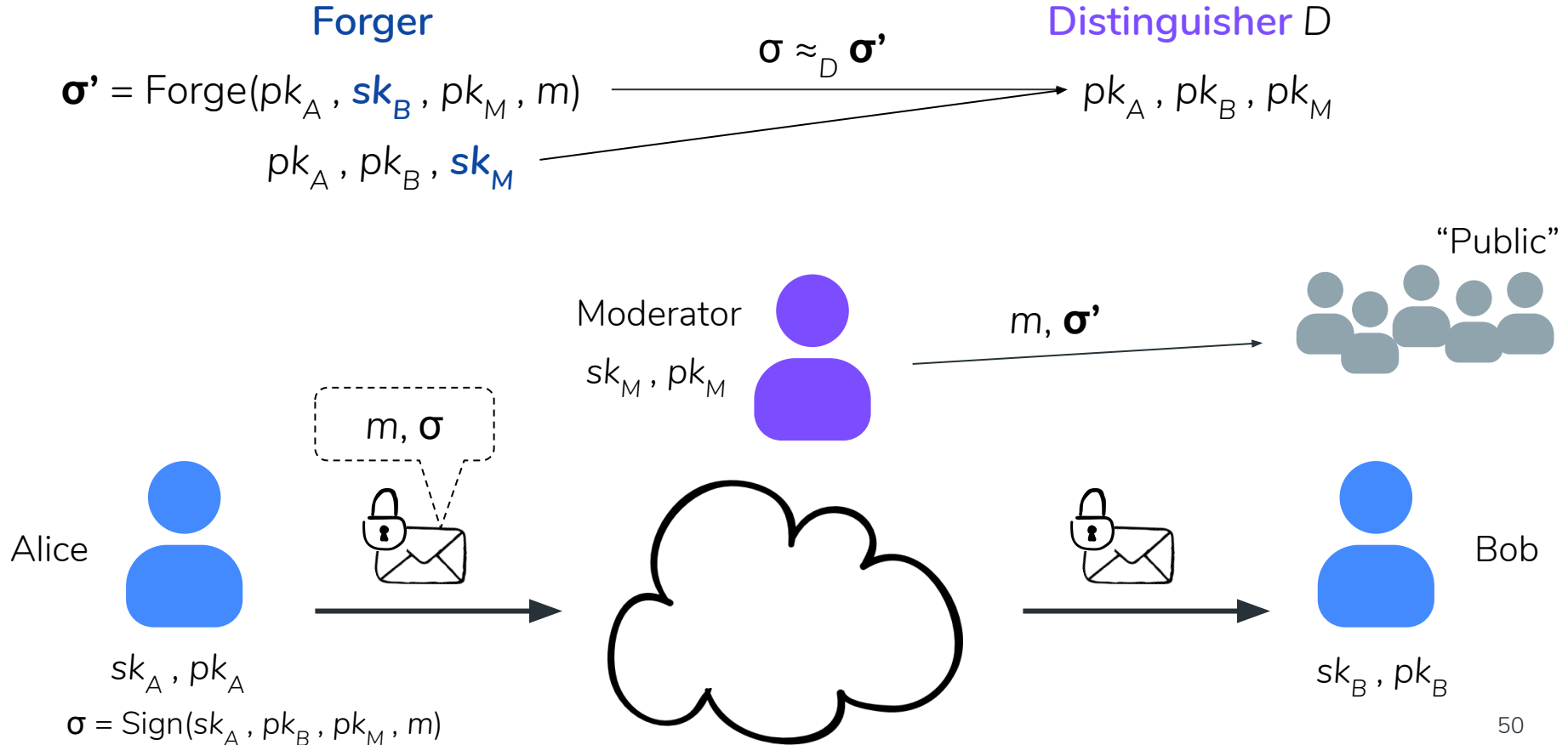
Forger

Distinguisher  $D$

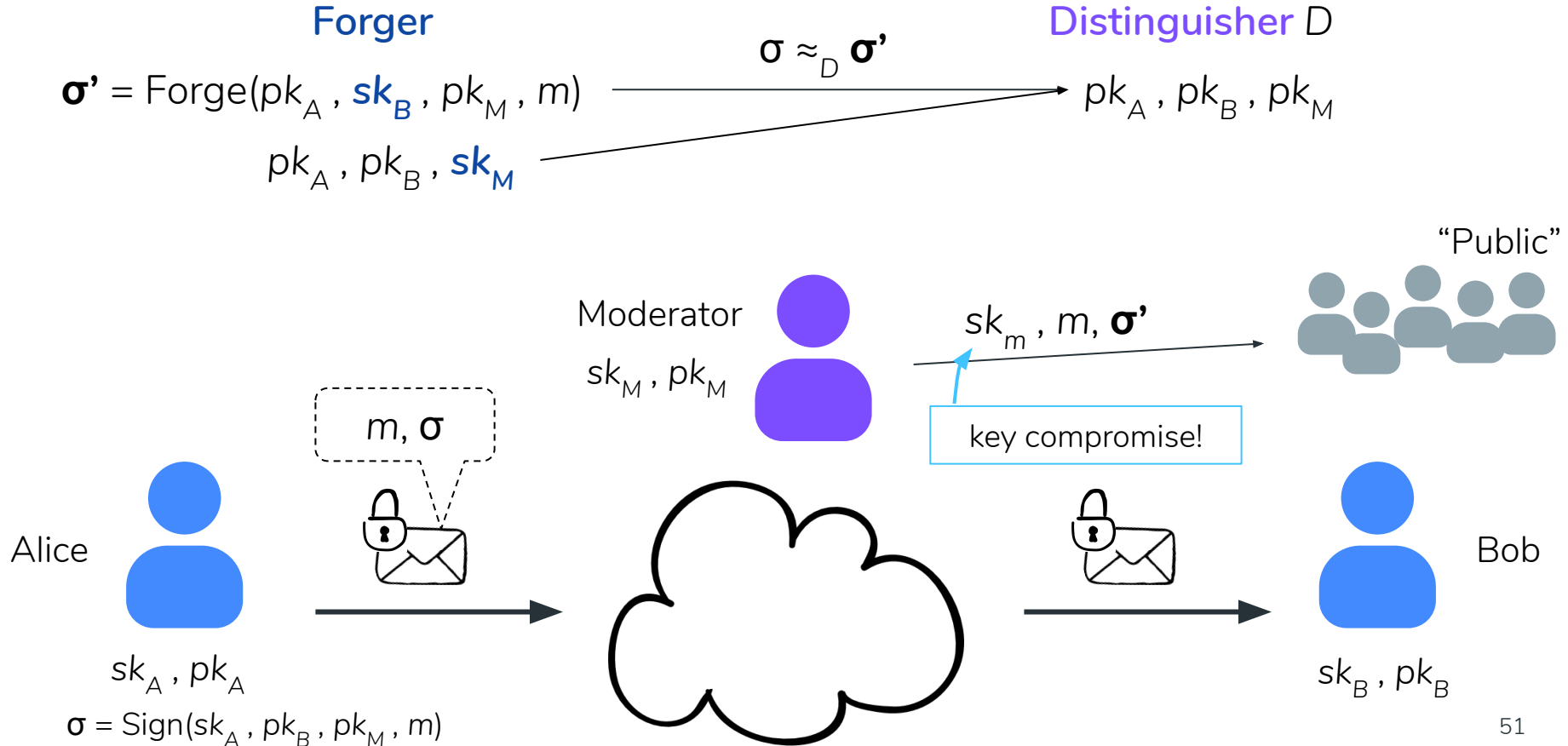
$$\sigma' = \text{Forge}(pk_A, sk_B, pk_M, m) \xrightarrow{\sigma \approx_D \sigma'} pk_A, pk_B, pk_M$$



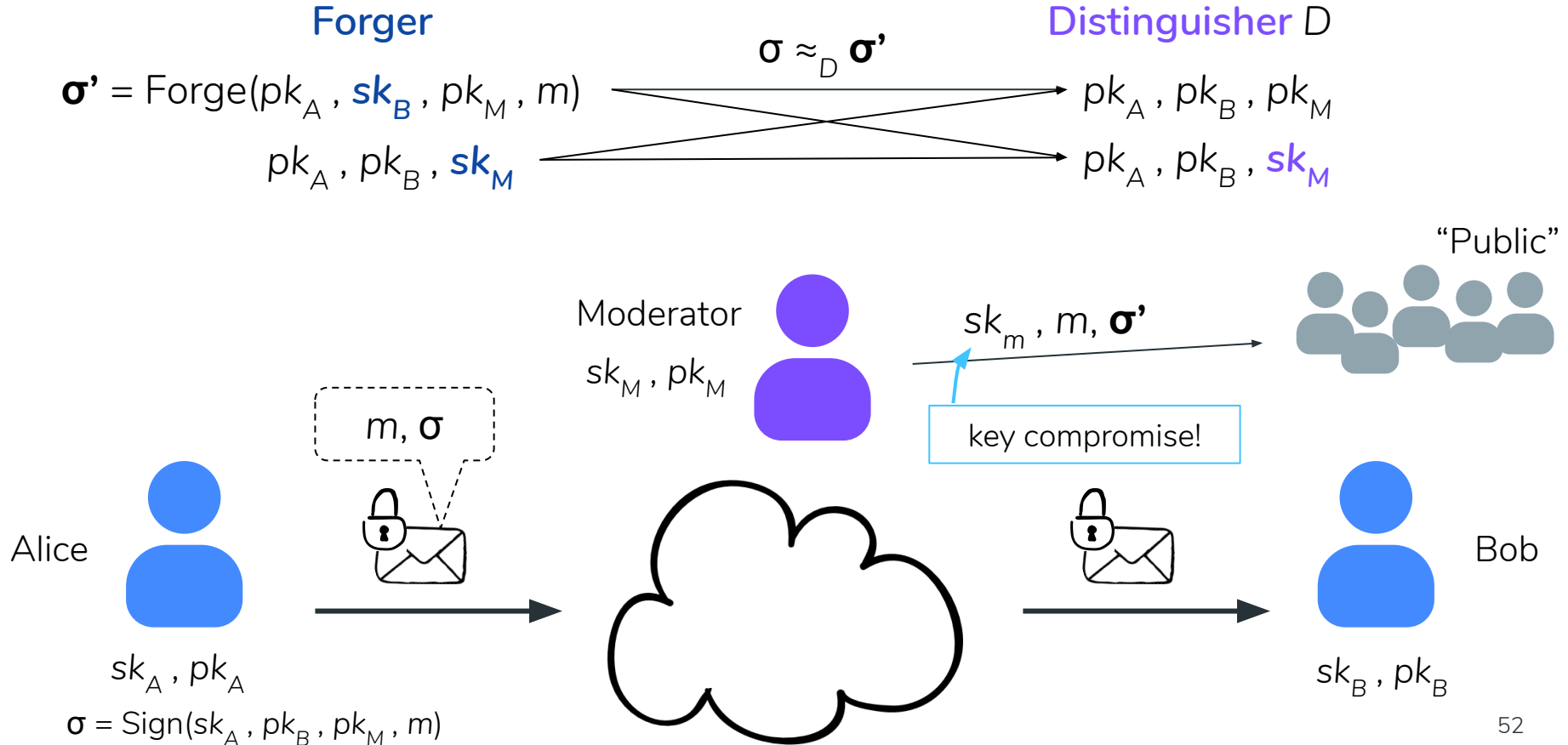
# Deniability landscape: “Who can trick whom?”



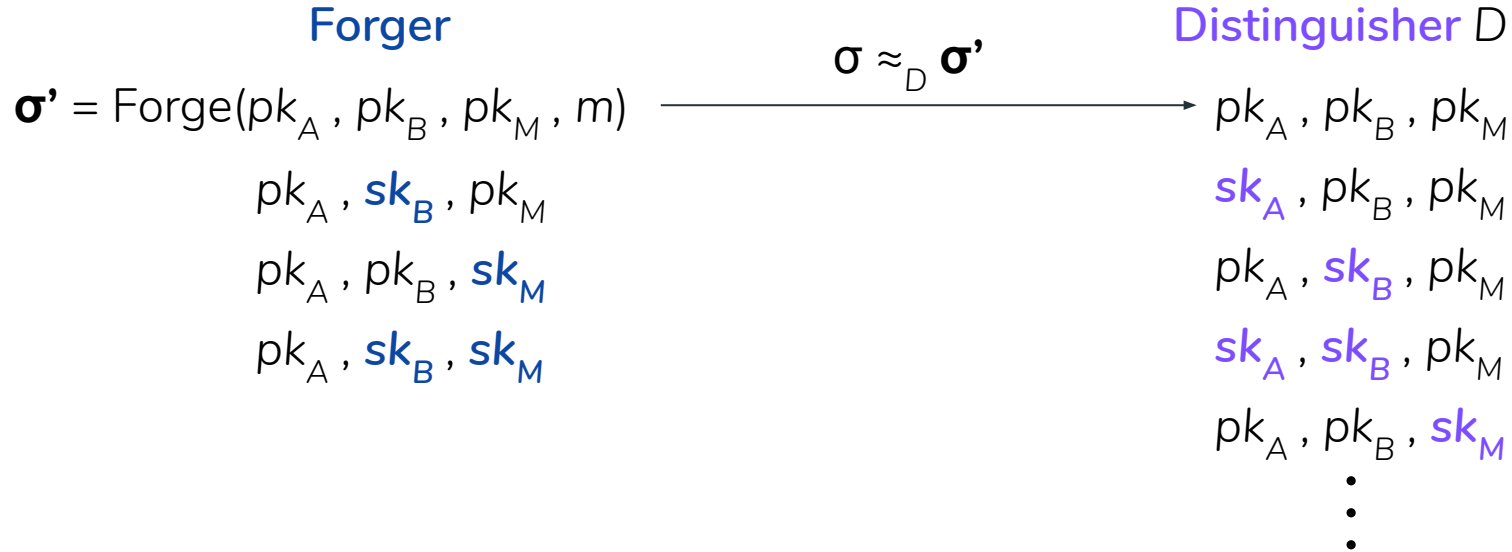
# Deniability landscape: “Who can trick whom?”



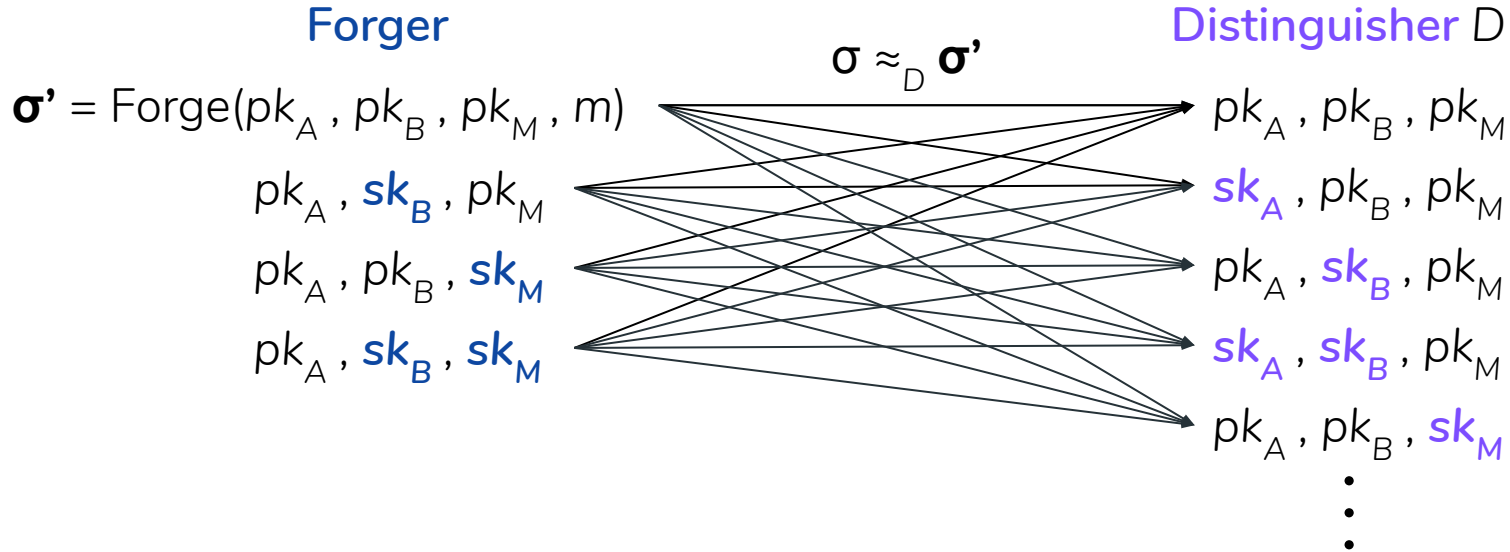
# Deniability landscape: “Who can trick whom?”



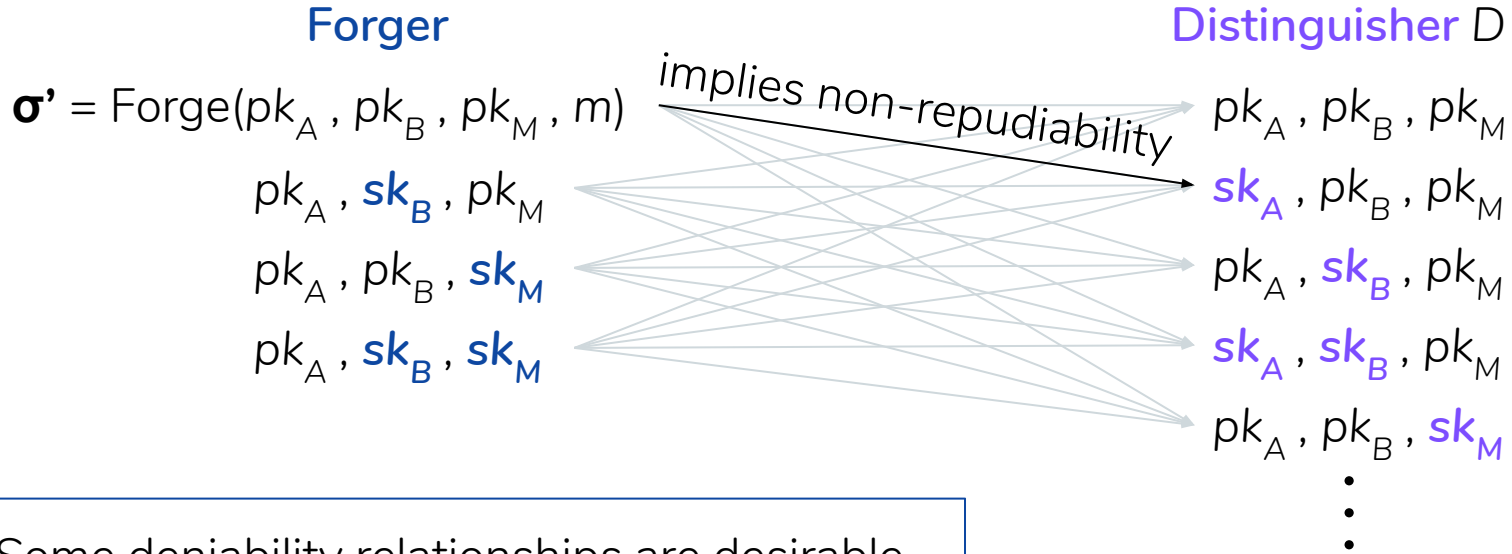
# Deniability landscape: “Who can trick whom?”



# Deniability landscape: “Who can trick whom?”

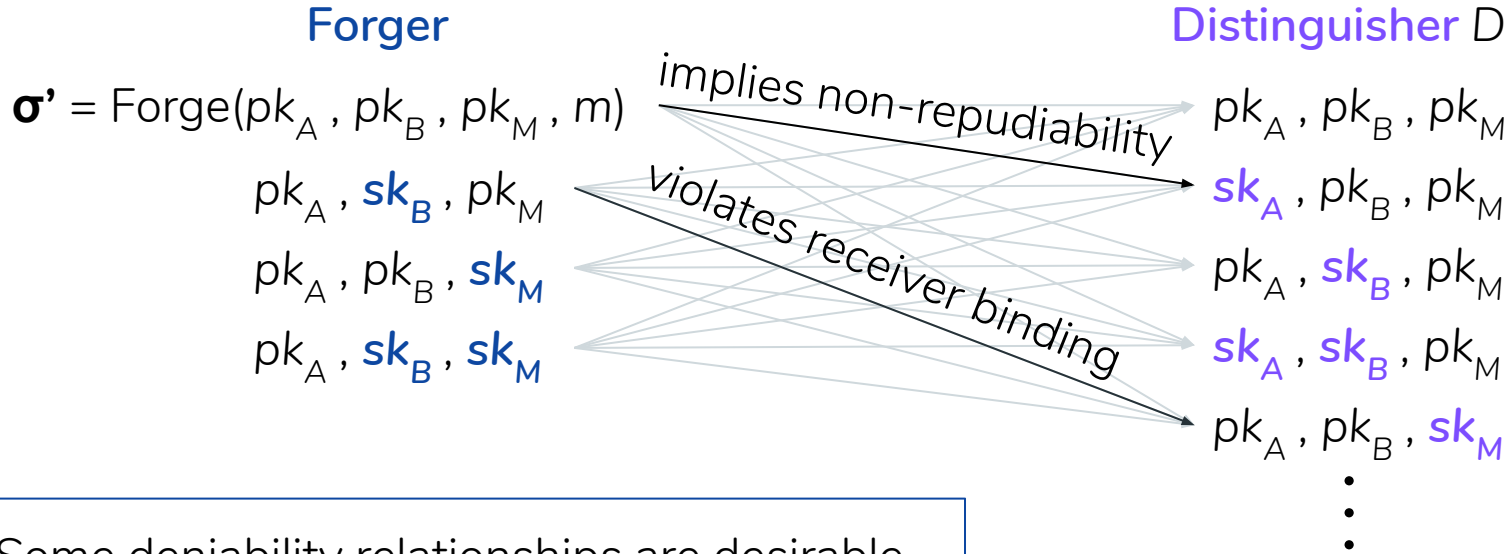


# Deniability landscape: “Who can trick whom?”



Some deniability relationships are desirable

# Deniability landscape: “Who can trick whom?”



Some deniability relationships are desirable

Others contradict directly with accountability



# Deniability landscape: “Who can trick whom?”

Distinguisher

			Distinguisher								
			0	1	0	1	0	1	0	1	
Forger	$sk_A$		0	1	0	1	0	1	0	1	
		$sk_B$	0	0	1	1	0	0	1	1	
			$sk_M$	0	0	0	0	1	1	1	1
	0	0	0			●	●	◆	◆	●◆	●◆
	0	1	0					◆	◆	◆	◆
0	0	1									
0	1	1									

● : Incompatible with unforgeability  
 ◆ : Incompatible with receiver binding

# Deniability landscape: “Who can trick whom?”

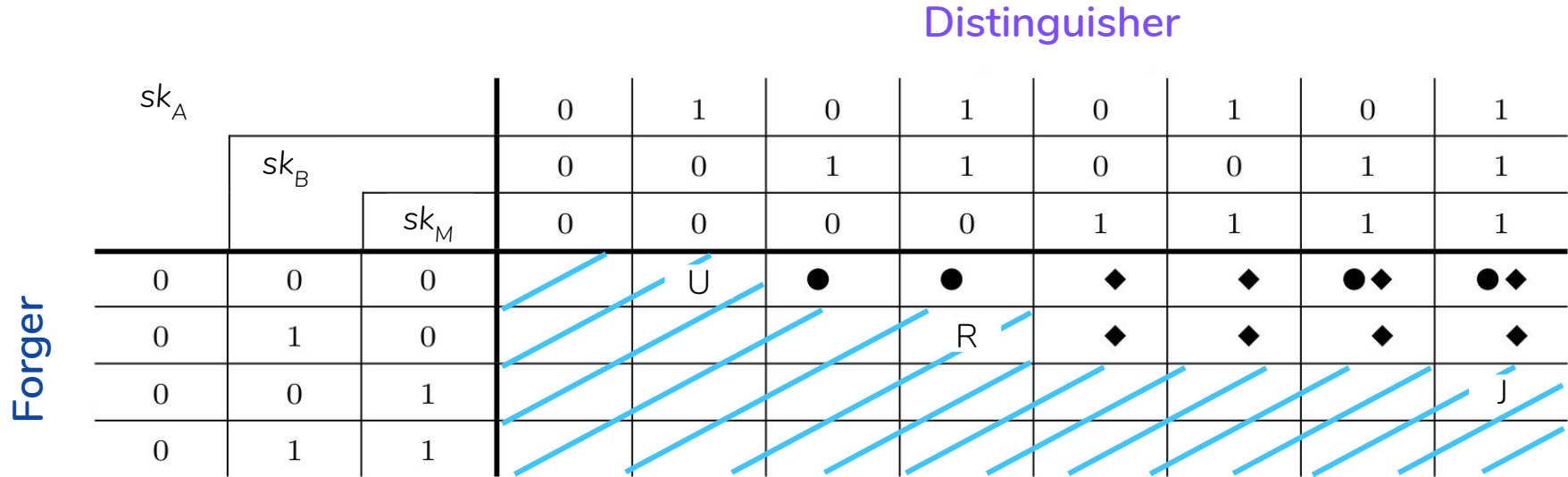
Distinguisher

			Distinguisher								
			0	1	0	1	0	1	0	1	
Forger	$sk_A$										
	$sk_B$	$sk_M$									
		0	1	0	1	0	1	0	1	0	1
	0	0	0	U	●	●	◆	◆	●◆	●◆	
0	1	0			R	◆	◆	◆	◆		
0	0	1							J		
0	1	1									

U : Universal deniability  
 R : Receiver compromise deniability  
 J : Judge compromise deniability

● : Incompatible with unforgeability  
 ◆ : Incompatible with receiver binding

# Deniability landscape: “Who can trick whom?”



U : Universal deniability  
 R : Receiver compromise deniability  
 J : Judge compromise deniability

● : Incompatible with unforgeability  
 ◆ : Incompatible with receiver binding

This represents only one possible set of tradeoffs!

# Summary of AMF goals

Specialized digital signature scheme that provides:

- Accountability
  - Receiver binding
  - Sender binding
- Deniability
  - Universal deniability
  - Receiver compromise deniability
  - Judge compromise deniability

# Our Construction

- Proof of knowledge of carefully-crafted expression of discrete log relationships
- Create signature by adding message via Fiat-Shamir transform

# Our Construction

- Proof of knowledge of carefully-crafted expression of discrete log relationships
- Create signature by adding message via Fiat-Shamir transform

Example of signature proof of knowledge (SPK) notation:

Standard digital signature (Schnorr)

$$\sigma \leftarrow_{\$} SPK \{t : pk_A = g^t\}$$

# Our Construction

- Proof of knowledge of carefully-crafted expression of discrete log relationships
- Create signature by adding message via Fiat-Shamir transform

Example of signature proof of knowledge (SPK) notation:

Standard digital signature (Schnorr)

$$\sigma \leftarrow_{\$} SPK \{t : pk_A = g^t\}$$

---

$\Sigma$ -Protocol Proof  
of Knowledge

**Prover**

$\xrightarrow{\text{com}}$

**Verifier**

$\xleftarrow{\text{chal}}$

$\xrightarrow{\text{resp}}$

SPK via Fiat-Shamir

**Prover**

$\xrightarrow{\text{com}}$

**Verifier**

$\text{chal} = H(\text{com}, m)$

$\xrightarrow{\text{resp}}$

# Our Construction

$$\begin{aligned} &\alpha \leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow \text{pk}_M^\alpha \\ &\beta \leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow \text{pk}_B^\beta \\ &\sigma \leftarrow_{\$} \text{SPK} \left\{ t, u, v, w : \underbrace{(\text{pk}_A = g^t \vee J = g^u)}_{\text{DV signature to moderator}} \wedge \underbrace{((J = (\text{pk}_M)^v \wedge E_J = g^v) \vee R = g^w)}_{\text{DV proof to Bob}} \right\}, E_J, E_R \end{aligned}$$



# Our Construction

$$\begin{aligned} \alpha &\leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha \\ \beta &\leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta \\ \sigma &\leftarrow_{\$} SPK \left\{ t, u, v, w : \underbrace{(pk_A = g^t \vee J = g^u)}_{\text{DV signature to moderator}} \wedge \underbrace{((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)}_{\text{DV proof to Bob}} \right\}, E_J, E_R \end{aligned}$$

# Our Construction

$$\begin{aligned} \alpha &\leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha \\ \beta &\leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta \\ \sigma &\leftarrow_{\$} SPK \{t, u, v, w : (pk_A = g^t \vee J = g^u) \wedge ((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)\}, E_J, E_R \end{aligned}$$

“What Alice is proving  
to the moderator”

DV proof to Bob

DV signature to moderator

# Our Construction

$$\alpha \leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha$$
$$\beta \leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta$$
$$\sigma \leftarrow_{\$} SPK \{t, u, v, w : (pk_A = g^t \vee J = g^u) \wedge ((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)\}, E_J, E_R$$

“What Alice is proving  
to the moderator”

“What allows other  
parties to forge”

DV proof to Bob

DV signature to moderator

# Our Construction

$$\alpha \leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha$$

$$\beta \leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta$$

$$\sigma \leftarrow_{\$} SPK \{t, u, v, w : (pk_A = g^t \vee J = g^u) \wedge ((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)\}, E_J, E_R$$

Moderator accepts if  $pk_M, E_J, J$  form a Diffie-Hellman triple

“What Alice is proving to the moderator”

“What allows other parties to forge”

DV proof to Bob

DV signature to moderator

# Our Construction

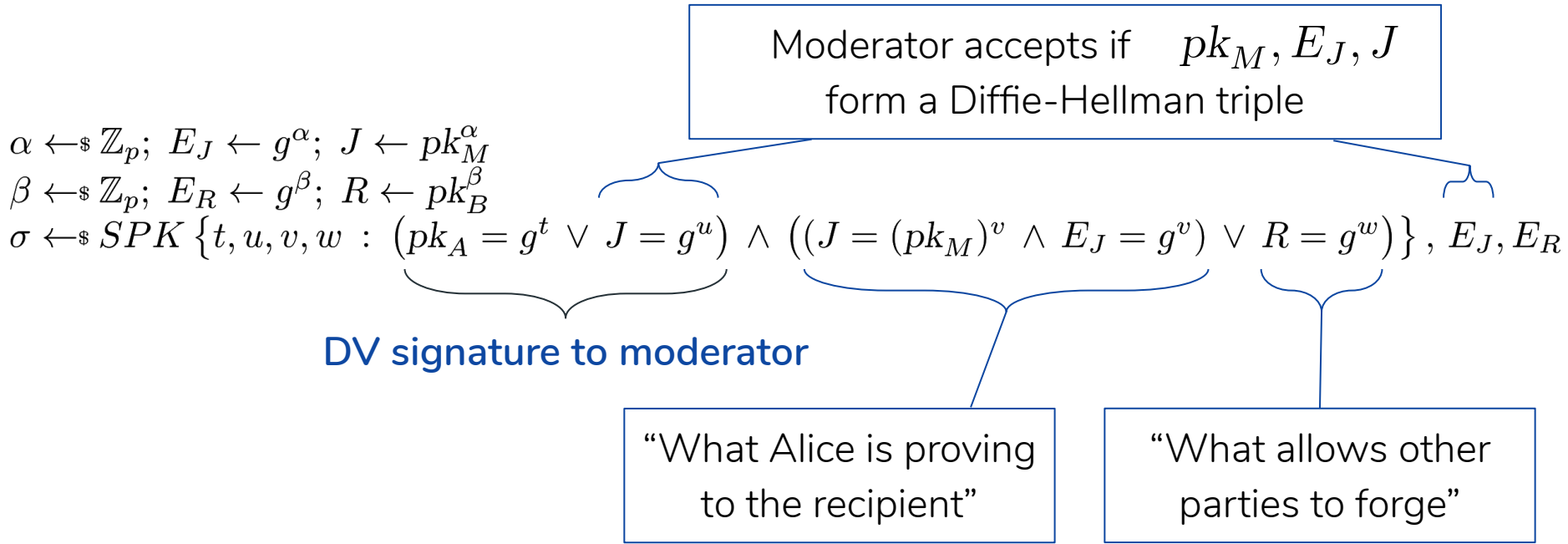
Moderator accepts if  $pk_M, E_J, J$   
form a Diffie-Hellman triple

$$\alpha \leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha$$
$$\beta \leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta$$
$$\sigma \leftarrow_{\$} SPK \{t, u, v, w : (pk_A = g^t \vee J = g^u) \wedge ((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)\}, E_J, E_R$$

DV signature to moderator

DV proof to Bob

# Our Construction



# Our Construction

$$\alpha \leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha$$
$$\beta \leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta$$
$$\sigma \leftarrow_{\$} SPK \{t, u, v, w : (pk_A = g^t \vee J = g^u) \wedge ((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)\}, E_J, E_R$$

DV signature to moderator

Moderator accepts if  $pk_M, E_J, J$   
form a Diffie-Hellman triple

Alice is proving Diffie-Hellman  
relationship to Bob!

DV proof to Bob

# Our Construction

Moderator accepts if  $pk_M, E_J, J$   
form a Diffie-Hellman triple

$$\alpha \leftarrow_{\$} \mathbb{Z}_p; E_J \leftarrow g^\alpha; J \leftarrow pk_M^\alpha$$
$$\beta \leftarrow_{\$} \mathbb{Z}_p; E_R \leftarrow g^\beta; R \leftarrow pk_B^\beta$$
$$\sigma \leftarrow_{\$} SPK \{t, u, v, w : (pk_A = g^t \vee J = g^u) \wedge ((J = (pk_M)^v \wedge E_J = g^v) \vee R = g^w)\}, E_J, E_R$$

**DV signature to moderator**

Accountability

- Moderator can attribute signature to sender
- Recipient can verify moderator will accept signature

Deniability

- Signature supports multiple forgery algorithms for various key compromise scenarios

Alice is proving Diffie-Hellman  
relationship to Bob!

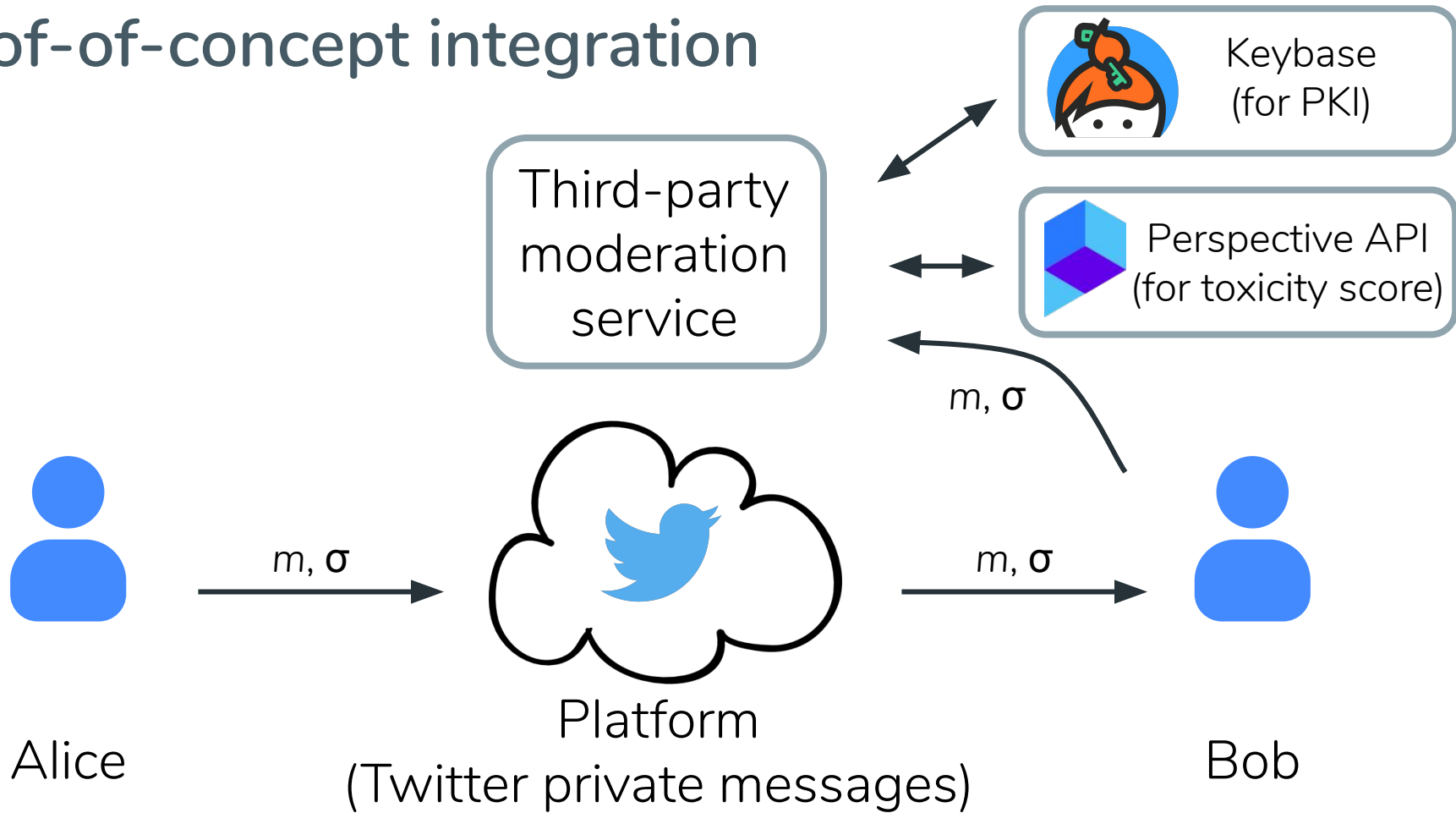
**DV proof to Bob**



# Implementation

- Implemented in Python 3 using petlib (OpenSSL bindings)
- Fast and efficient
  - < 500 bytes for P-256 (9 group elements + 6 scalars)
  - < 10 ms for P-256
- Available at [github.com/julialen/asymmetric-message-franking](https://github.com/julialen/asymmetric-message-franking)

# Proof-of-concept integration



# Our contributions

- Asymmetric Message Franking (AMF)
  - new cryptographic primitive for content moderation of metadata-private messaging
  - formal accountability and deniability security notions for content moderation
- Construction based on “designated-verifier” signatures
- Implementation and proof-of-concept integration
  - Available at [github.com/juliaalen/asymmetric-message-franking](https://github.com/juliaalen/asymmetric-message-franking)