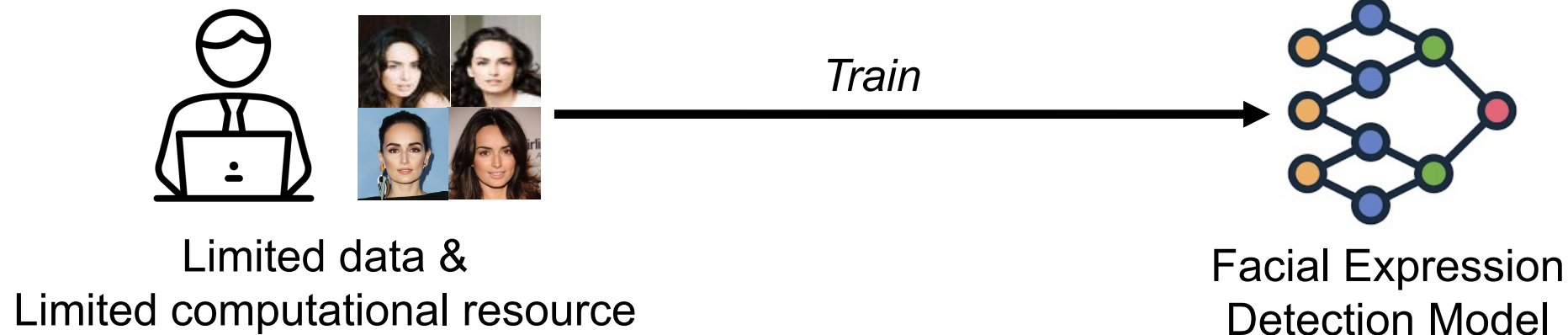# Crafter: Facial Feature Crafting against Inversion-based Identity Theft on Deep Models

Shiming Wang, Zhe Ji, Liyao Xiang, Hao Zhang, Xinbing Wang, Chenghu Zhou, Bo Li
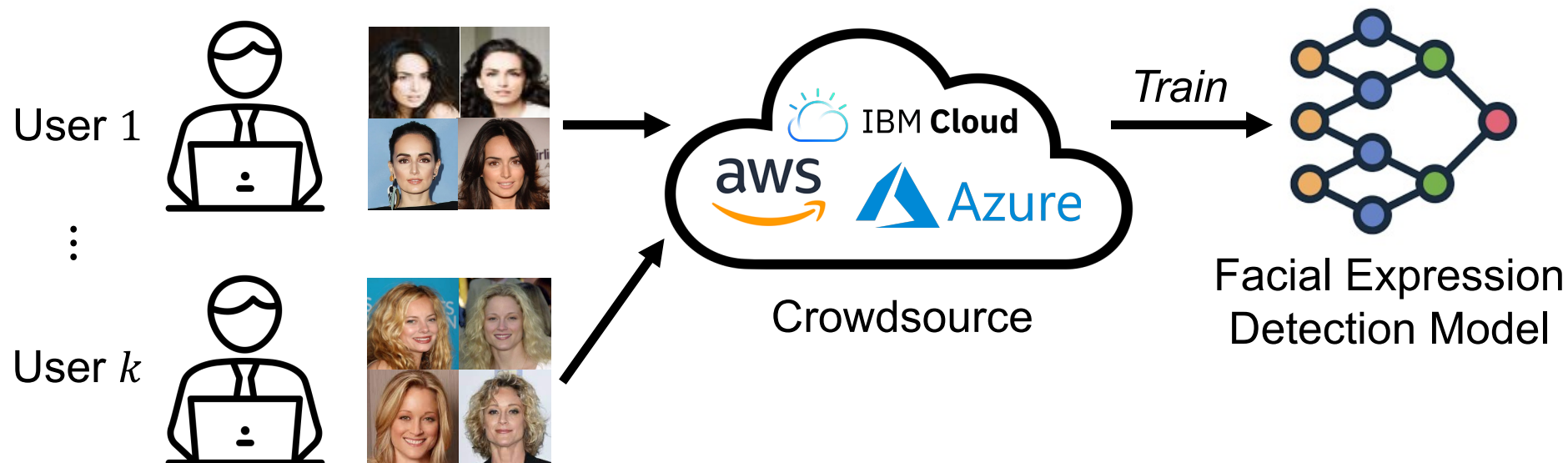
NDSS SYMPOSIUM/2024

Presented by
Internet Society

#NDSSSymposium2024

# Machine Learning as a Service

**Example 1:** Teaching **deep learning task.**



Limited data &
Limited computational resource

*Train*

Facial Expression
Detection Model

# Machine Learning as a Service

**Example 1: Training deep learning task.**



User 1

User $k$

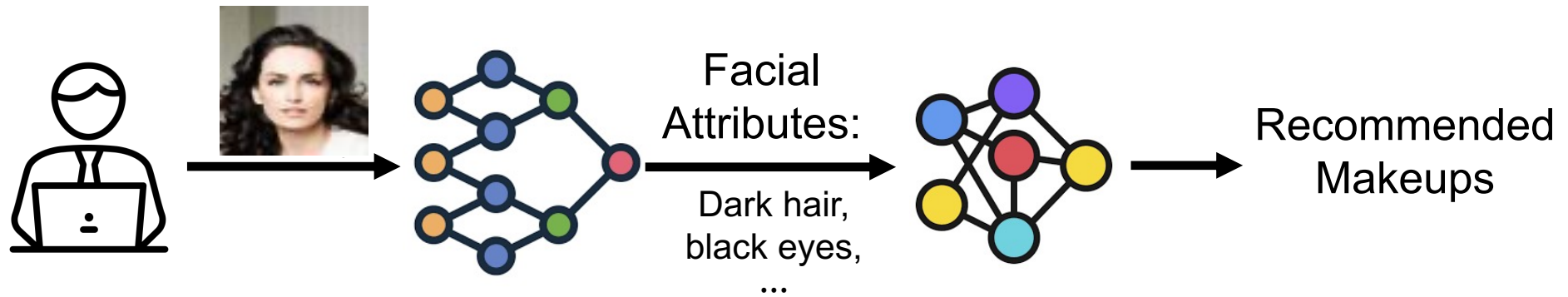Crowdsource
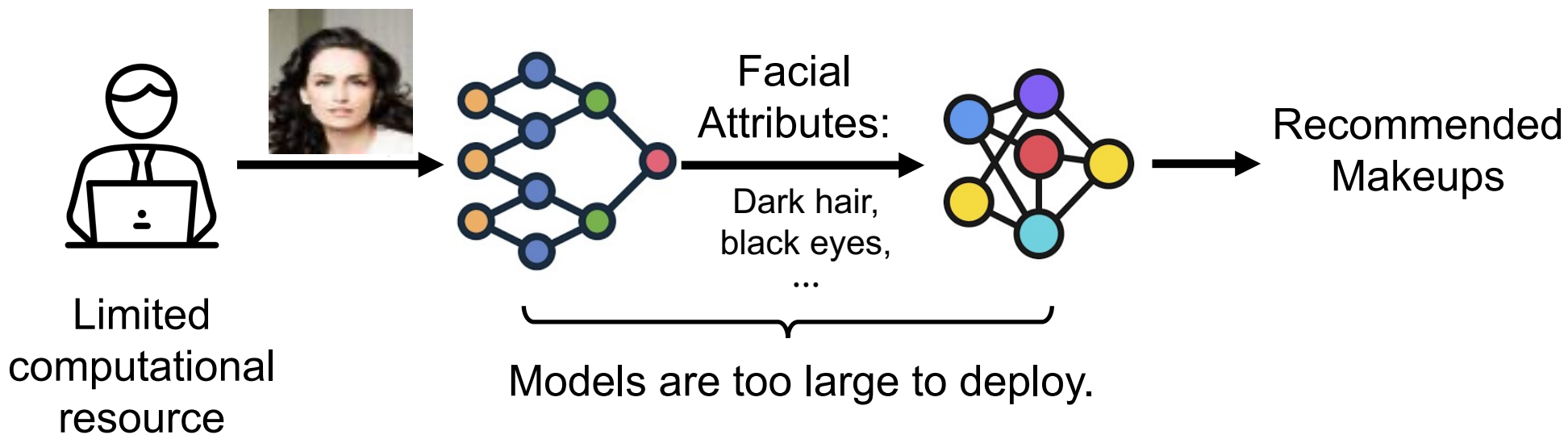
*Train*

Facial Expression
Detection Model

# Machine Learning as a Service

**Example 2: Inference deep learning task.**

# Machine Learning as a Service

**Example 2: Inference deep learning task.**



Facial Attributes:

Dark hair, black eyes, ...

Recommended Makeups

Limited computational resource

Models are too large to deploy.

# Machine Learning as a Service

**User**                    **Cloud**

"Alice"



Private image $X$.
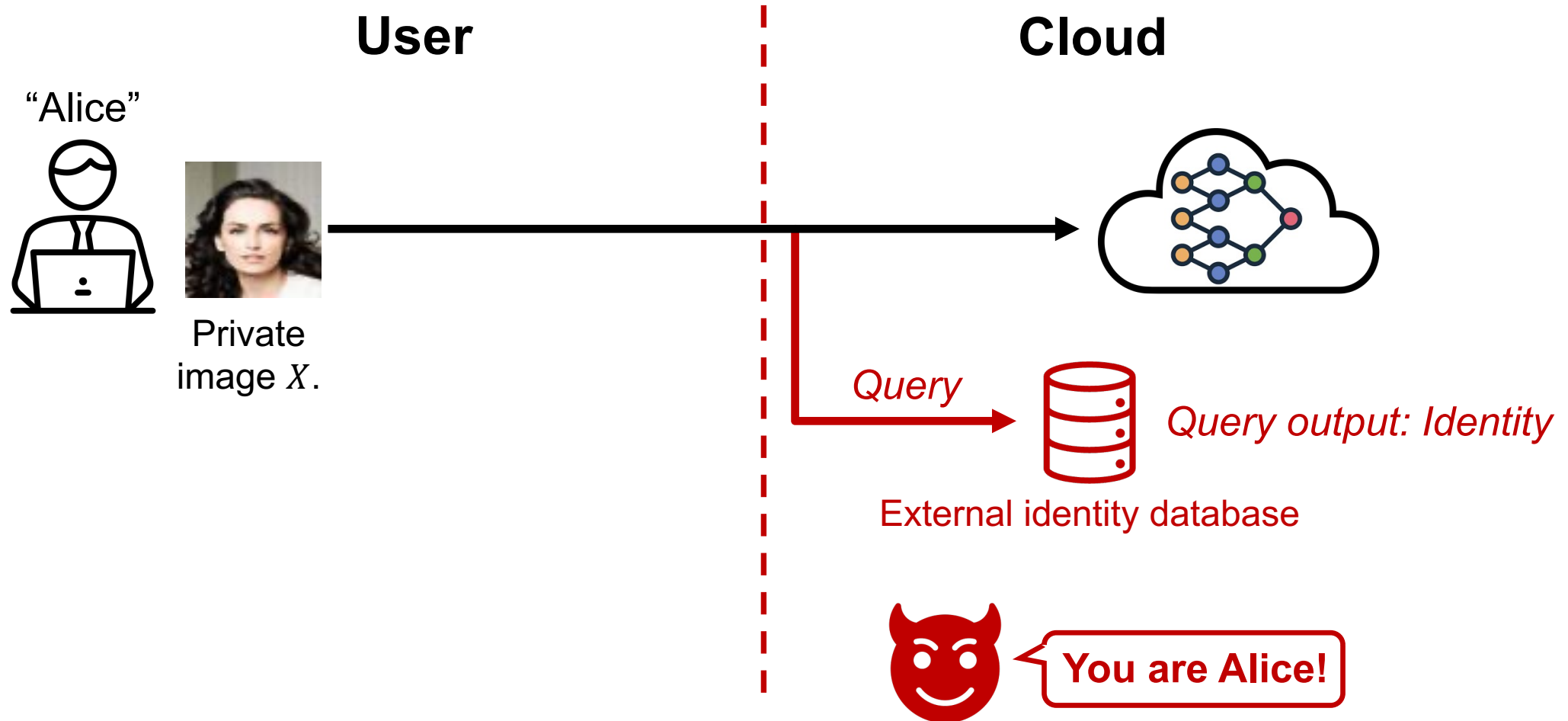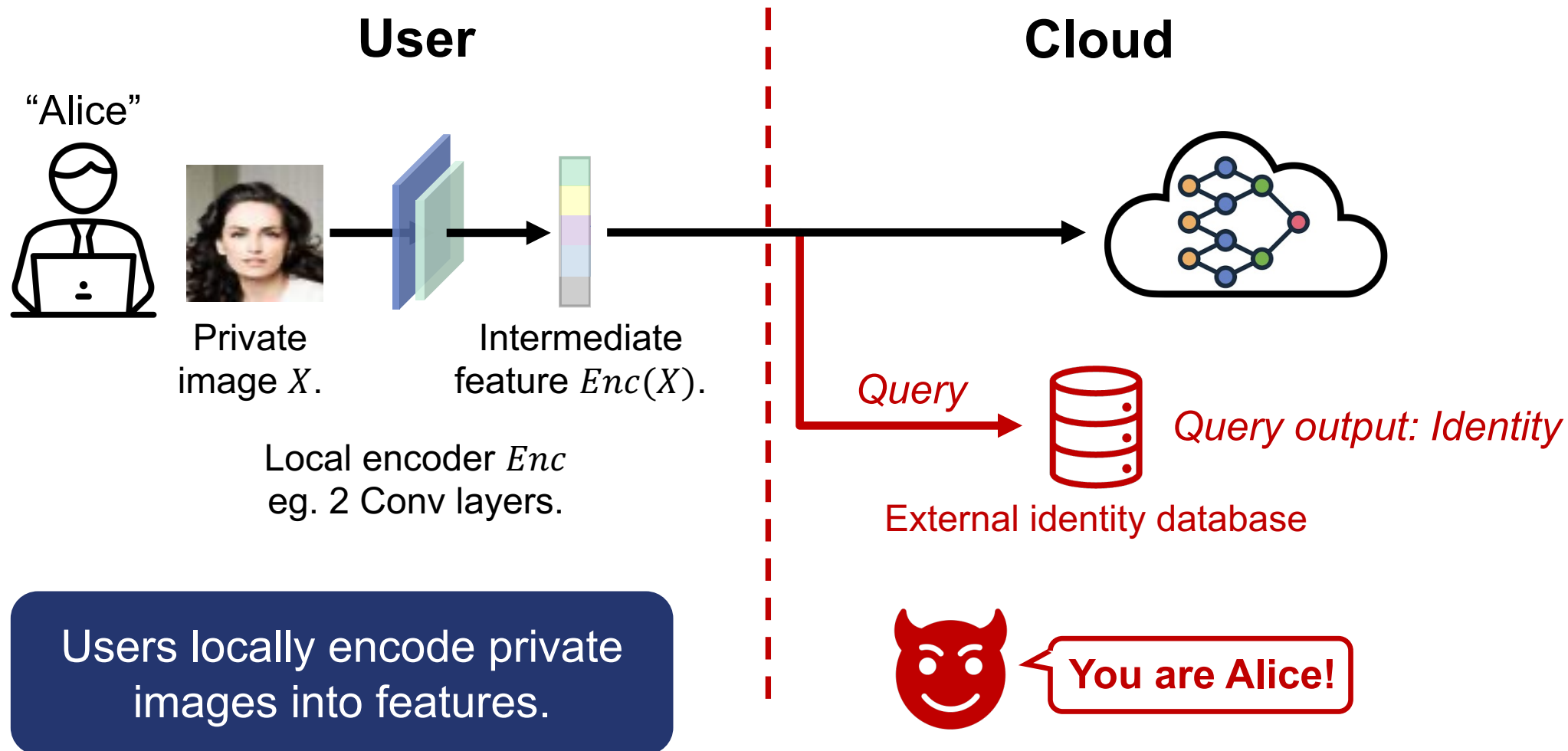
Users are motivated to share their facial images with the cloud.

# Privacy Concern: Identity Theft

# Privacy Concern: Identity Theft

**User**

**Cloud**

"Alice"

Private
image $X$.

Intermediate
feature $Enc(X)$.

*Inversion
attack*

Inverted
image $\hat{X}$.

# Privacy Concern: Identity Theft

# Privacy Concern: Identity Theft



**User**

"Alice"

Private image $X$.

Intermediate feature $Enc(X)$.

Local encoder $Enc$ 2 Conv layers.

**Cloud**

*Inversion attack*

*Query* → *ID*

Inverted image $\hat{X}$.

# Defending Inversion-based Identity Theft

**Previous Defense**:

*AdvLearn*[1], *Disco*[2], *TIPRDC*[3]

- Vulnerable against adaptive attacks;

- Fail to balance privacy & utility;

- Limited application scenarios.

[1] Xiao et al. "Adversarial learning of privacy-preserving and task-oriented representations ", 2020
[2] Singh et al. "Disco: Dynamic and invariant sensitive channel obfuscation for deep neural networks ", 2021
[3] Li et al. "Tiprdc: task-independent privacy-respecting data crowdsourcing framework for deep learning with anonymized intermediate representations ", 2020

# In Our Work

**Crafter Defense**:

User-end feature crafting that protects identity info against various inversion attacks, while preserving data utility.

Threat Model

Intuitions & Design

Evaluation

**Black-box inversion attack**:

- Access to public images; query access to the local $Enc.$



User

"Alice"

Private image $X$.

Local encoder $Enc$

Intermediate feature $Enc(X)$.

# Threat Model

**Black-box inversion attack**:

- Access to public images; query access to the local $Enc$.
- Train a decoder network $Dec$.

# Threat Model

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.



User

Private image $X$.

Local encoder $Enc$.

Intermediate feature $Enc(X)$.

# Threat Model

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.
- Optimize over the inverted image.



User

Private image $X$.

Local encoder $Enc$

Intermediate feature $Enc(X)$.

Attacker

Inverted image $\hat{X}$.

$Enc(\hat{X})$

# Threat Model

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.
- Optimize over the inverted image.



User

Private image $X$.

Local encoder $Enc$

Intermediate feature $\underline{Enc(X)}$.

Attacker

Inverted image $\hat{X}$.

$\underline{Enc(\hat{X})}$

Compute the feature-level difference.

Presented

Inte
Soci

# Threat Model

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.
- Optimize over the inverted image.



User

Private image $X$.

Local encoder $Enc$

Intermediate feature $\underline{Enc(X)}$.

Attacker

Inverted image $\hat{X}$.

$\underline{Enc(\hat{X})}$

Compute the feature-level difference.

# Threat Model

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.
- Optimize over the inverted image.

$+$ Pretrained public generator $G$.



User

Private image $X$.

Local encoder $Enc$

Intermediate feature $\underline{Enc(X)}$.

*In practice:*

**Attacker** 😈

$Enc(\hat{X})$

Compute the feature-level difference.

Inverted image $\hat{X}$.

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.
- Optimize over the inverted image.

$+$ Pretrained public generator $G$.



*In practice:*

**User**

Private image $X$.

Local encoder $Enc$

Intermediate feature $\underline{Enc(X)}$.

Compute the feature-level difference.

$\underline{Enc(\hat{X})}$

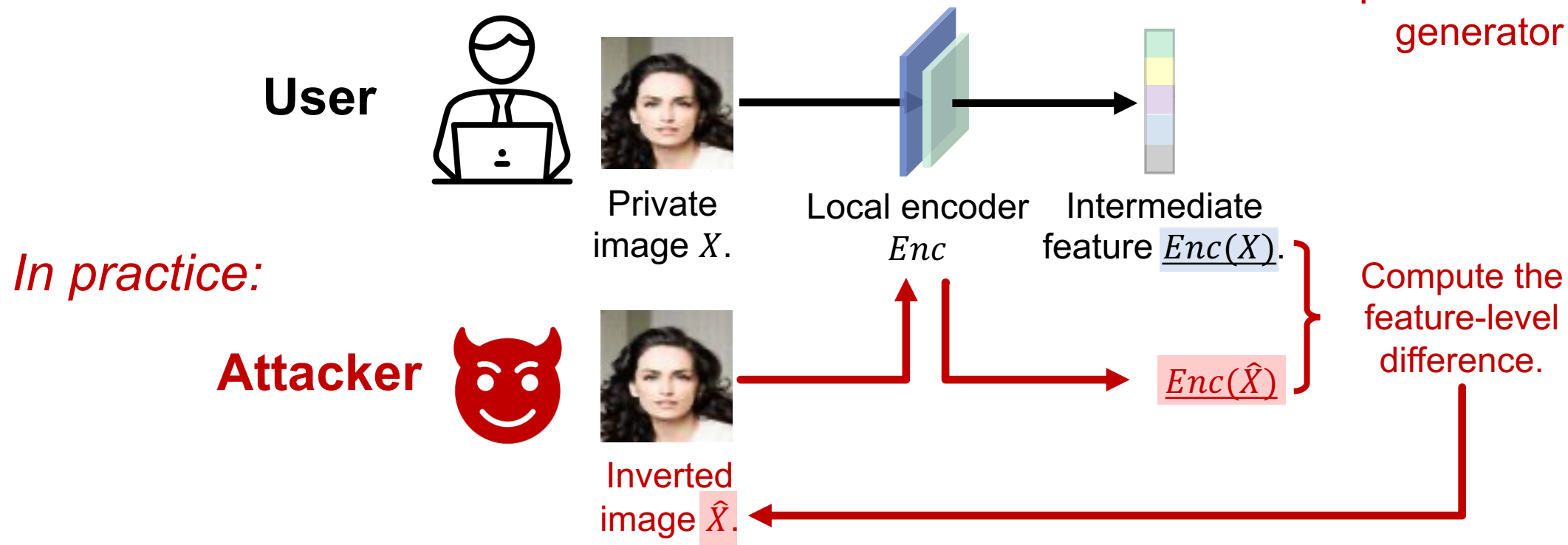Latent vector $z$. Generator $G$. Inverted $\hat{X} = G(z)$.
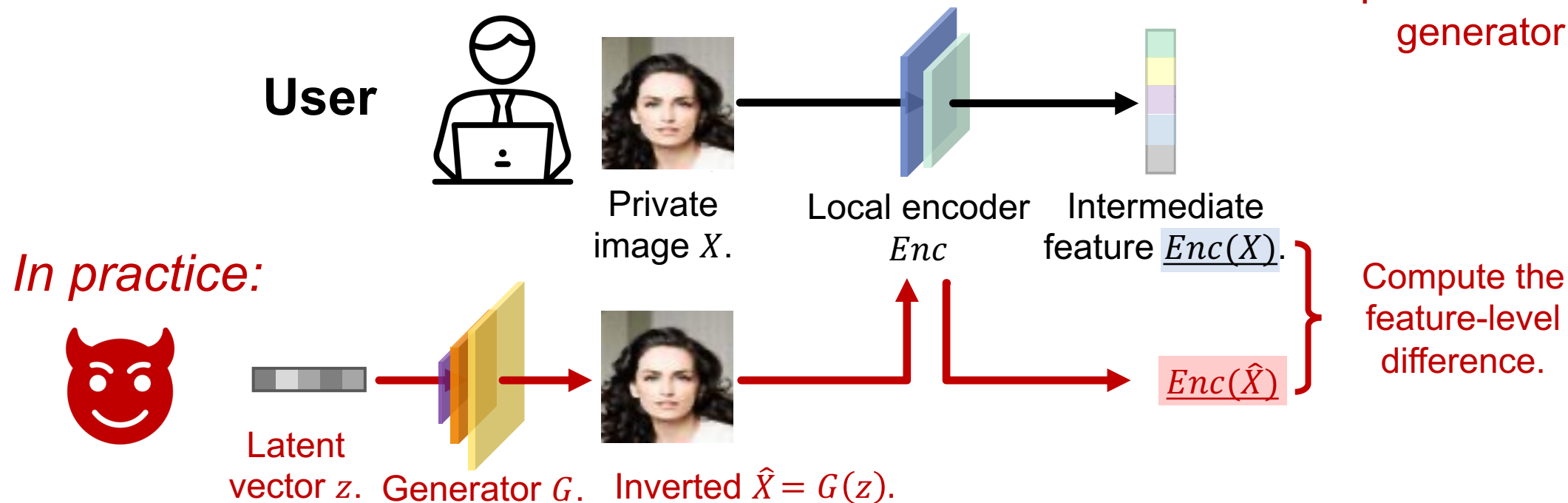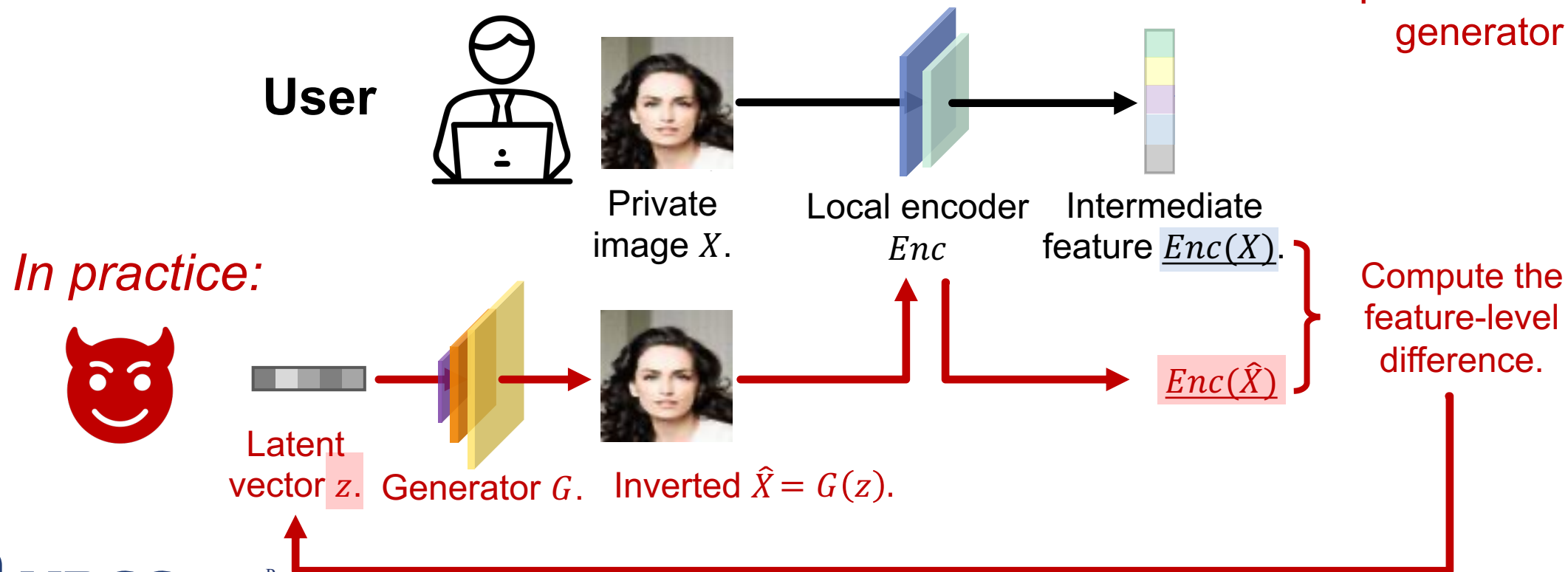
Presented

Inte Soci

**White-box inversion attack**:

- Access to public images; access to the local $Enc$ and its parameters.
- Optimize over the ~~inverted image~~ latent vector.

$+$ Pretrained public generator $G$.



**User**

Private image $X$.

Local encoder $Enc$

Intermediate feature $\underline{Enc(X)}$.

Compute the feature-level difference.

*In practice:*

Latent vector $z$.  Generator $G$.  Inverted $\hat{X} = G(z)$.

$\underline{Enc(\hat{X})}$

# Defense Intuitions

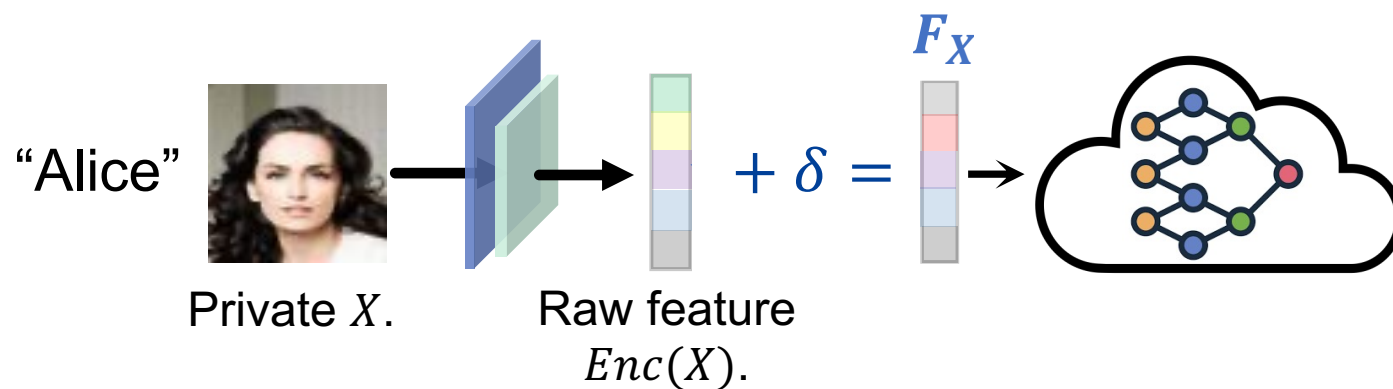**Privacy goal:** Inverted image does not look like Alice.

**Utility goal:** Feature completes cloud tasks well.

# Defense Intuitions

**Privacy goal:** Inverted image does not look like Alice.

**Utility goal:** Feature completes cloud tasks well.

**General intuition:** Perturb the feature.



"Alice"

Private $X$.

Raw feature $Enc(X)$.

$F_X$

$+ \delta =$

# Defense Intuitions

**Privacy goal:** Inverted image does not look like Alice.

**Utility goal:** Feature completes cloud tasks well.

**General intuition:** Perturb the feature.

- (*Privacy*) Mislead a simulated inversion attacker.

"Alice"

$F_X$

$+ \delta =$

Private $X$.

Raw feature $Enc(X)$.

*Inversion attack*

Inverted $\hat{X}$.

Presented by
Internet Society

**Privacy goal:** Inverted image does not look like Alice.

**Utility goal:** Feature completes cloud tasks well.
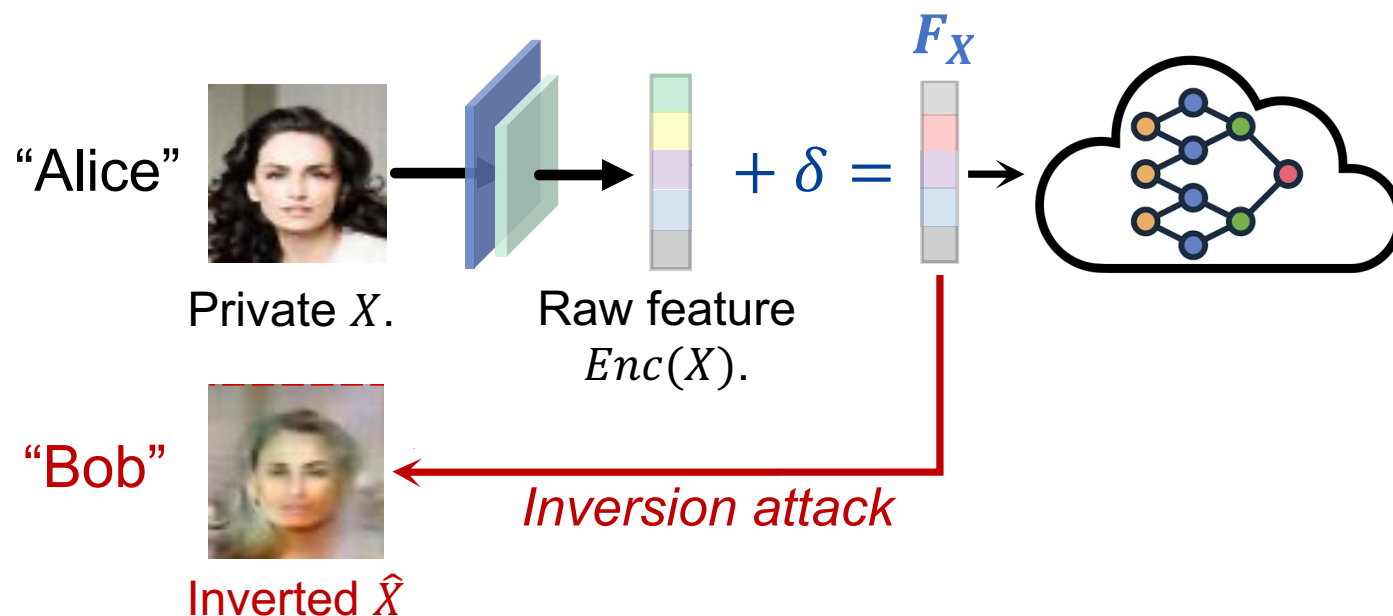
**General intuition:** Perturb the feature.
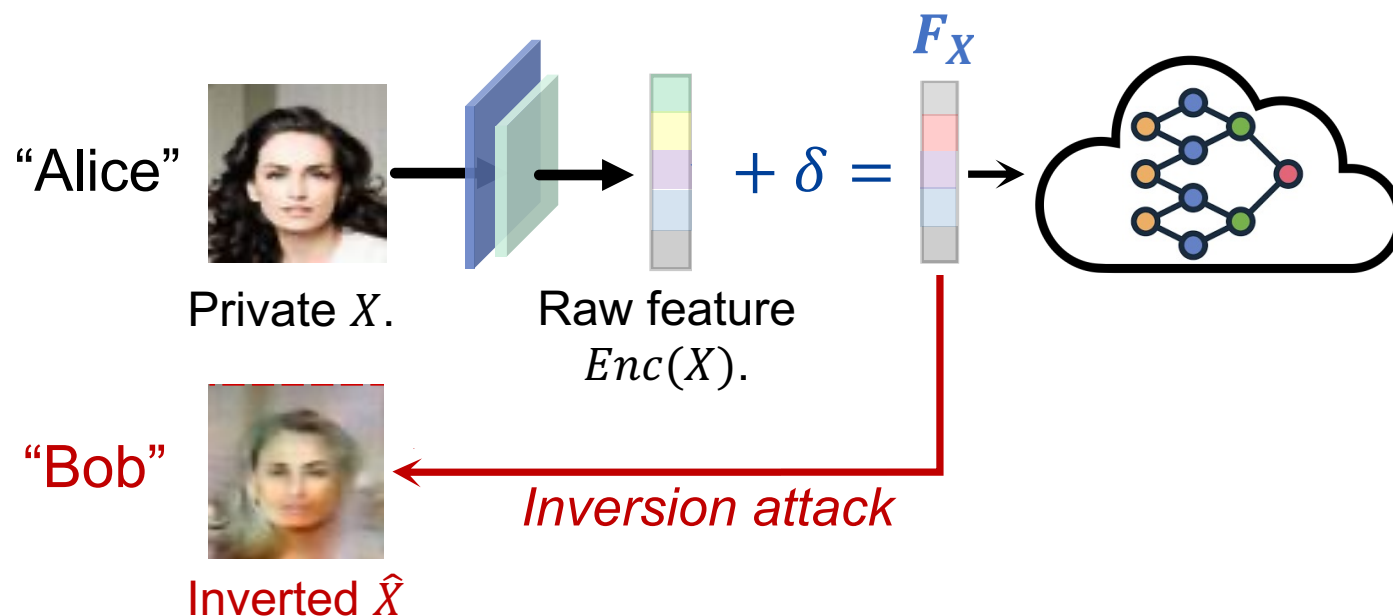
- (*Privacy*) Mislead a simulated inversion attacker.

- (*Utility*) Keep the perturbation small.

"Alice"

$F_X$

$+ \delta =$

Private $X$.

Raw feature $Enc(X)$.

*Inversion attack*

Inverted $\hat{X}$.

# Defense Intuitions (Utility)

**Utility loss:** $L_{utility}$ = perturbation magnitude.

**Preserves utility:** Cloud model is robust against minor perturbation.

**Utility task agnostic:** $L_{utility}$ independent from cloud model

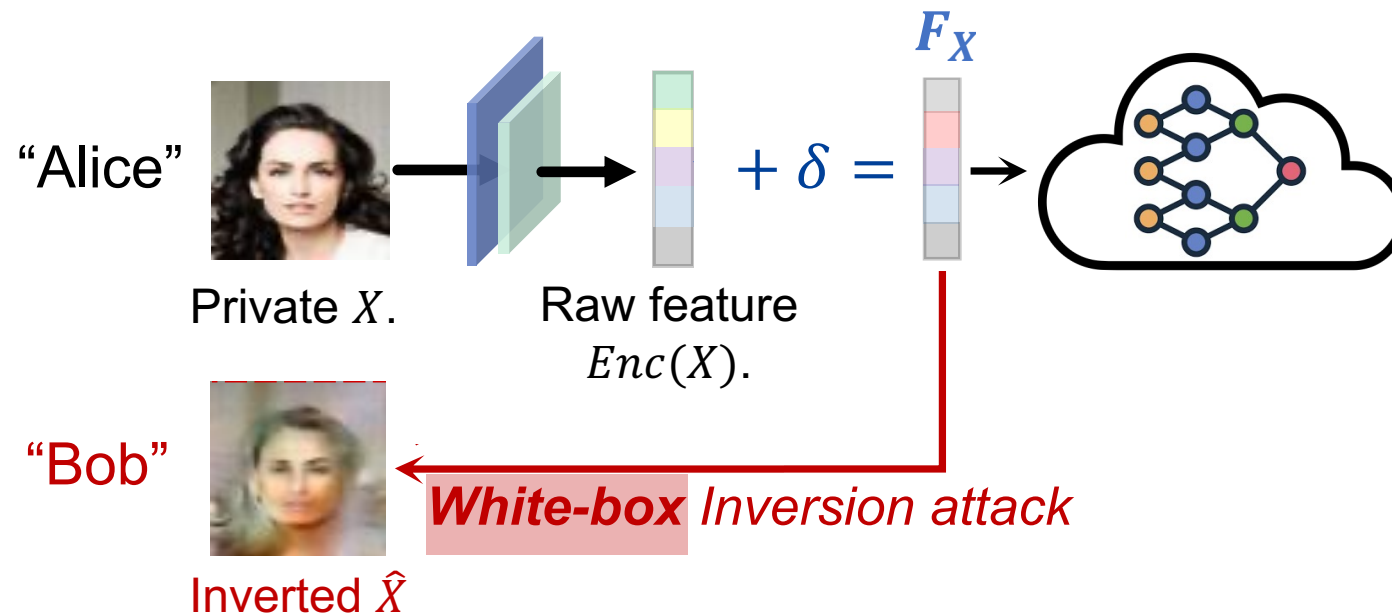$\rightarrow$ deployable as a plug-in.

# Defense Intuitions (Privacy)

**Challenge 1:** *Robust against both black- & white-box inversion.*

# Defense Intuitions (Privacy)

**Challenge 1:** *Robust against both black- & white-box inversion.*

**Intuition:** White-box attack is stronger; simulate a **white-box attacker**.



"Alice"

Private $X$.    Raw feature
$Enc(X)$.

$F_X$

$+ \delta =$

*White-box Inversion attack*

Inverted $\hat{X}$.

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against <u>adaptive attacks</u>.*

*Attacker tries to bypass a fixed defense.*

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against <u>adaptive attacks</u>.*

*Attacker tries to bypass a fixed defense.*

*If a defense is not robust:* <span style="color:red">*false security, meaningless!*</span>

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against <u>adaptive attacks</u>.*

*Attacker tries to bypass a fixed defense.*

*If a defense is not robust: <span style="color:red">false security, meaningless!</span>*

**Previous defense:** Push the attacker away from the private image.

"Stay Away"    Tit for tat between attacker & defense.

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against <u>adaptive attacks</u>.*
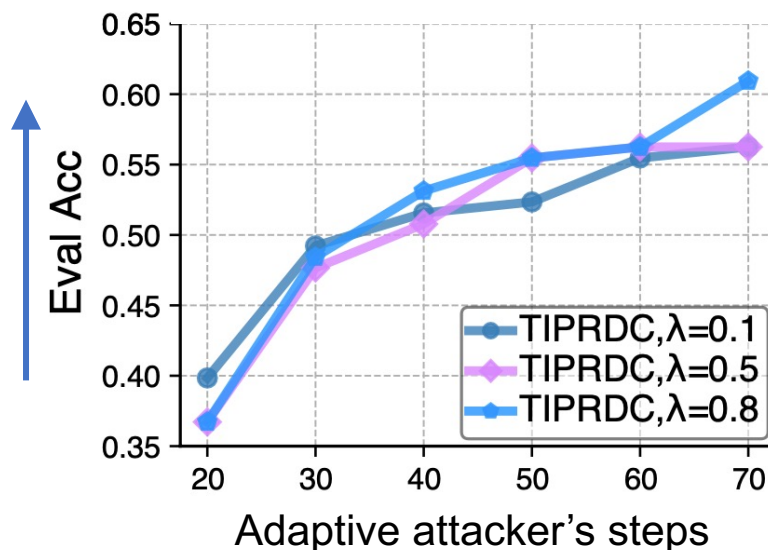
*Attacker tries to bypass a fixed defense.*

*If a defense is not robust: false security, meaningless!*

**Previous defense:** Push the attacker away from the private image.

"Stay Away"

Tit for tat between attacker & defense.



*Weaker privacy*

Eval Acc vs Adaptive attacker's steps

Legend: TIPRDC,$\lambda=0.1$; TIPRDC,$\lambda=0.5$; TIPRDC,$\lambda=0.8$

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against <u>adaptive attacks</u>.*

*Attacker tries to bypass a fixed defense.*

*If a defense is not robust: <span style="color:red">false security, meaningless!</span>*

**Previous defense:** Push the attacker away from the private image.

"Stay Away"

Tit for tat between attacker & defense.

Why is "Stay Away" vulnerable against adaptive attacks?

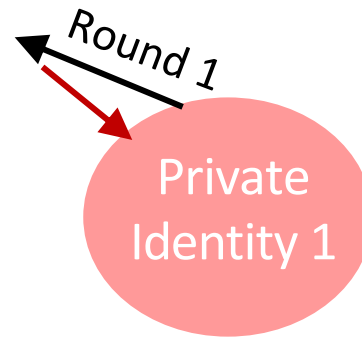# Defense Intuitions (Privacy)

A game view:

Attack  ⟶

Defense  ⟶

Private Identity 1

# Defense Intuitions (Privacy)

A game view:

Attack ──→

Defense ──→

**Conventional:**
stay away from private identity

Round 1

Private Identity 1

# Defense Intuitions (Privacy)

A game view:

Attack →

Defense →

**Conventional**: stay away from private identity
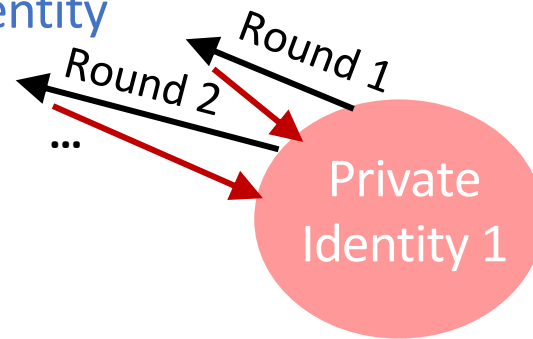
Round 1

Round 2

...

Private Identity 1

# Defense Intuitions (Privacy)

A game view:

**Attack** →
Defense →

**Conventional**: stay away from private identity
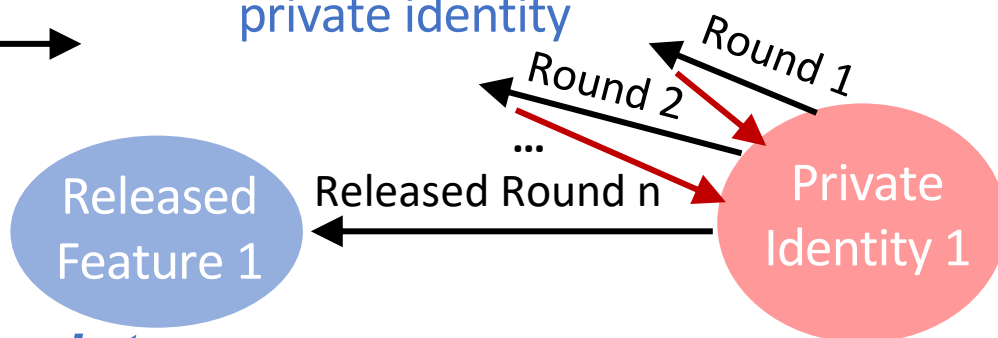
Round 1

Round 2

...

Released Round n

Released Feature 1

Private Identity 1
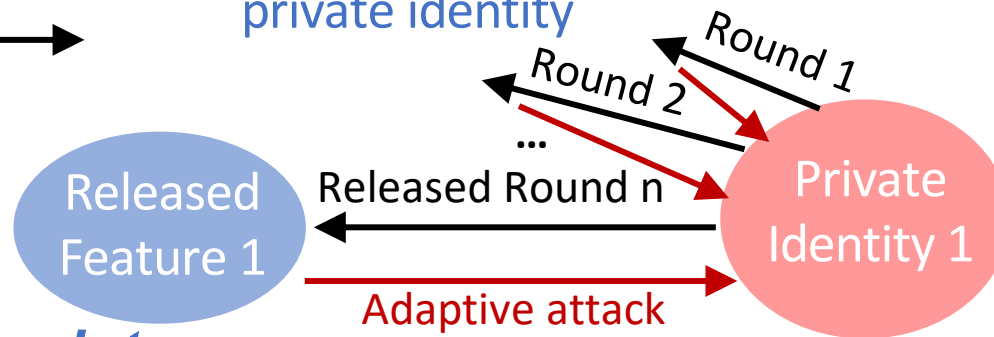
*A stationary point.*
*But not equilibrium (there is NO equilibrium in reality).*

# Defense Intuitions (Privacy)

A game view:

Attack →

Defense →

**Conventional**:
stay away from
private identity

Round 1

Round 2

...

Released Round n

Released Feature 1

Private Identity 1

Adaptive attack

*A stationary point.*
*But not equilibrium (there is NO equilibrium in reality).*

# Defense Intuitions (Privacy)

A game view:

Attack →

Defense →
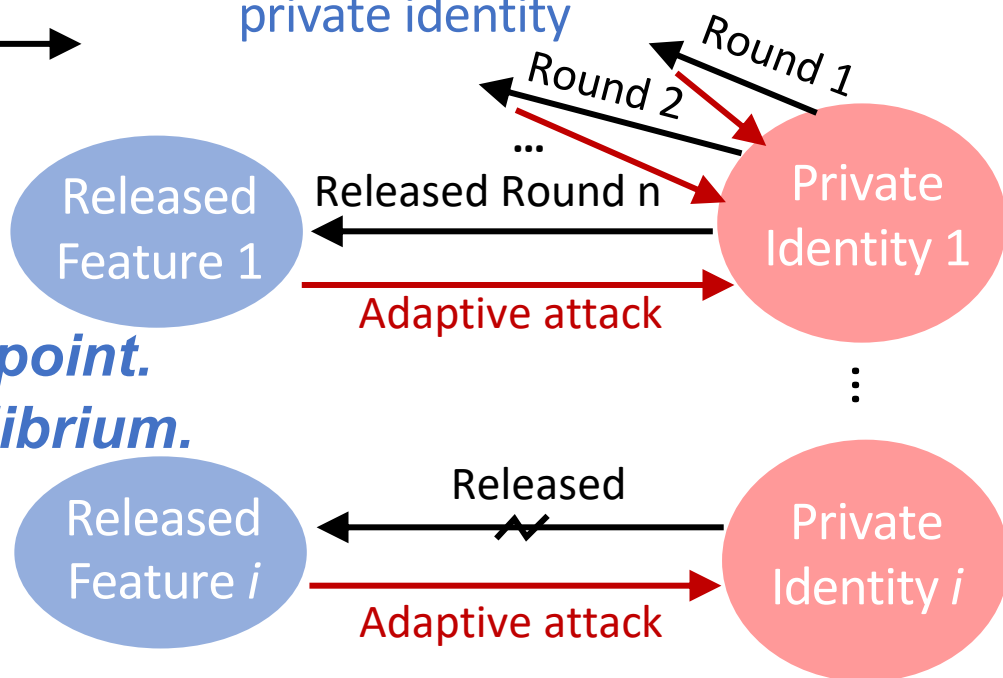
**Conventional**: stay away from private identity

Round 1

Round 2

...

Released Round n

Released Feature 1 ← Private Identity 1

Adaptive attack →

*A stationary point.*
*But not equilibrium.*

Released Feature *i* ← Released → Private Identity *i*

Adaptive attack →

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against adaptive attacks.*

*Attacker tries to bypass a fixed defense.*

*If a defense is not robust: false security, meaningless!*

**Our Intuition:** Limit attacker's **knowledge gain** from the exposed feature.

"Get Close"

# Defense Intuitions (Privacy)

**Challenge 2:** *Robust against <u>adaptive attacks</u>.*

*Attacker tries to bypass a fixed defense.*

*If a defense is not robust:* *false security, meaningless!*

**Our Intuition:** Limit attacker's **knowledge gain** from the exposed feature.

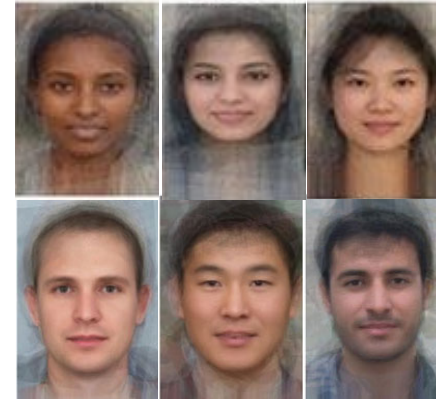"Get Close"         = **Prior** vs. **Posterior**

# Defense Intuitions (Privacy)

**Prior**: "*Average face*", public face distribution.

# Defense Intuitions (Privacy)

**Prior**: "*Average face*", public face distribution.

Contains no private ID info.

# Defense Intuitions (Privacy)

**Prior**: "*Average face*", public face distribution.
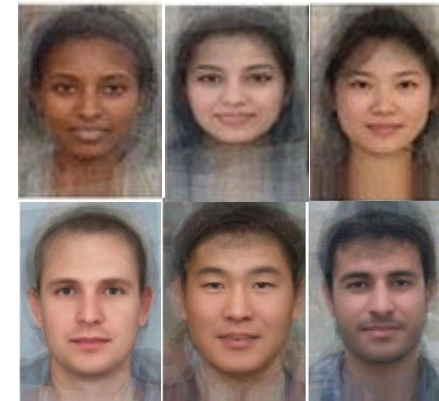
Contains no private ID info.

We use: $G(z_{random})$

Public generator    Random latent vectors
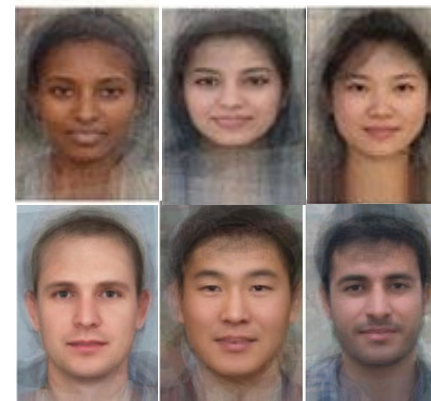
# Defense Intuitions (Privacy)

**Prior**: "*Average face*", public face distribution.

Contains no private ID info.

We use: $G(z_{random})$

Public generator    Random latent vectors

**Posterior**: Image $\hat{X}$ inverted from feature $F_X$.

**"Get Close"** : Minimize distance between prior & posterior.

# Defense Intuitions (Privacy)

**"Get Close"** : Minimize distance between prior & posterior.

We use: Earth-Mover distance $EMD$.

**Privacy loss:** $L_{privacy} = EMD$ between inverted image $\hat{X}$ & average face.

# Defense Intuitions (Privacy+Utility)

"Get Close" : Minimize distance between prior & posterior.

We use: Earth-Mover distance $EMD$.

**Privacy loss:** $L_{privacy} = EMD$ between inverted image $\hat{X}$ & average face.

**Combine utility & privacy:** $L_{combined} = \beta \cdot L_{privacy} + L_{utility}$

| "Get Close" | **:** Minimize distance between prior & posterior. |

We use: Earth-Mover distance $EMD$.

**Privacy loss:** $L_{privacy} = EMD$ between inverted image $\hat{X}$ & average face.

**Combine utility & privacy:** $L_{combined} = \beta \cdot L_{privacy} + \boxed{L_{utility}}$

Find a feature perturbation that 1) is small;

# Defense Intuitions (Privacy+Utility)

"Get Close" **:** Minimize distance between prior & posterior.

We use: Earth-Mover distance $EMD$.

**Privacy loss:** $L_{privacy} = EMD$ between inverted image $\hat{X}$ & average face.

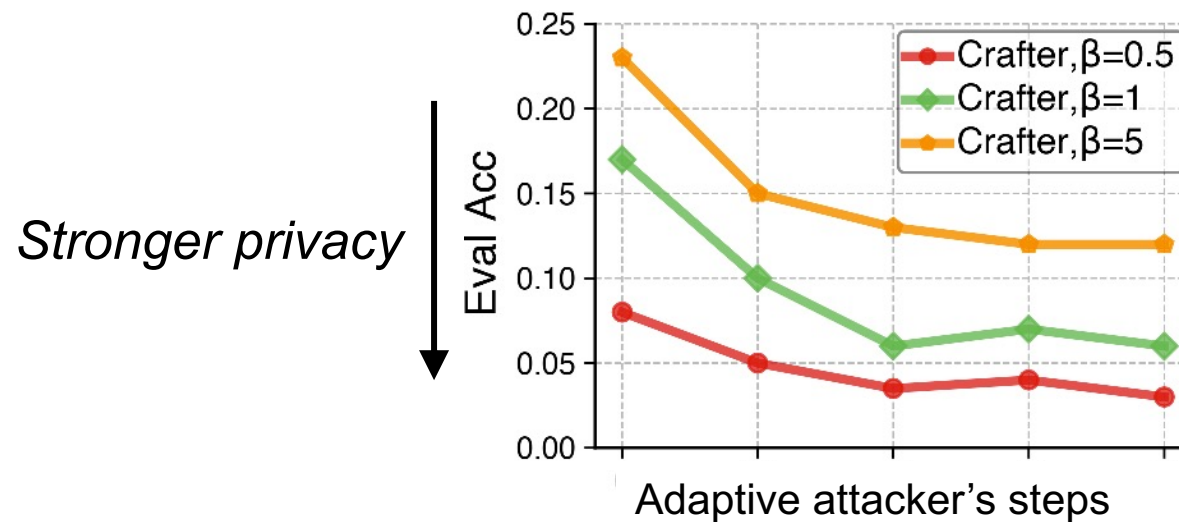**Combine utility & privacy:** $L_{combined} = β \cdot L_{privacy} + L_{utility}$

Find a feature perturbation that 1) is small;

2) draws inverted image close to public average faces.

Presented by
Internet Society

# Defense Intuitions (Privacy+Utility)

**"Get Close"** ✓ : Minimize distance between prior & posterior.

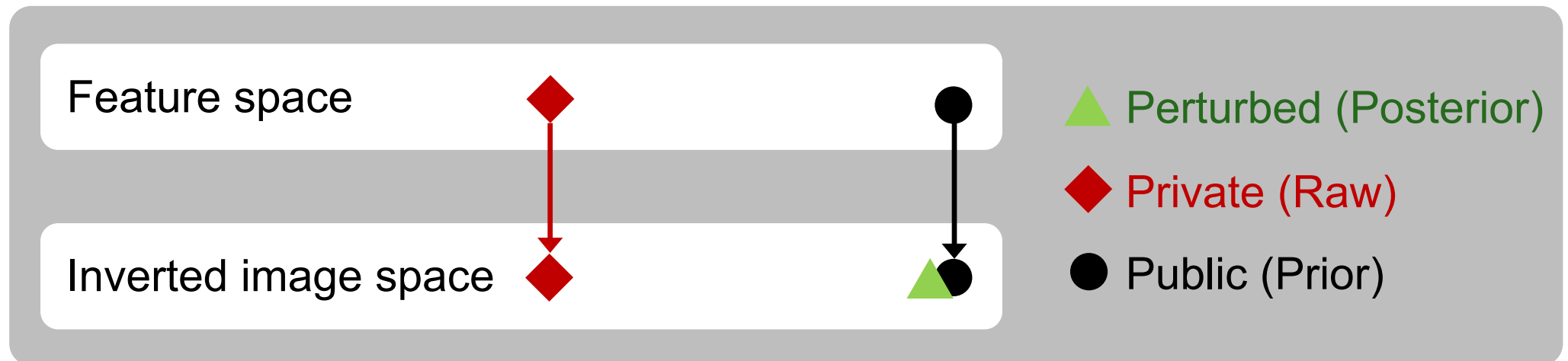

*Stronger privacy* ↓

Why is "Get Close" robust against adaptive attacks?

# Defense Intuitions (Privacy+Utility)

"Get Close" ✓

Find a feature perturbation that

1) is small;
2) draws inverted image close to public average faces.

Feature space     ◆           ●

Inverted image space     ◆           ▲●

▲ Perturbed (Posterior)
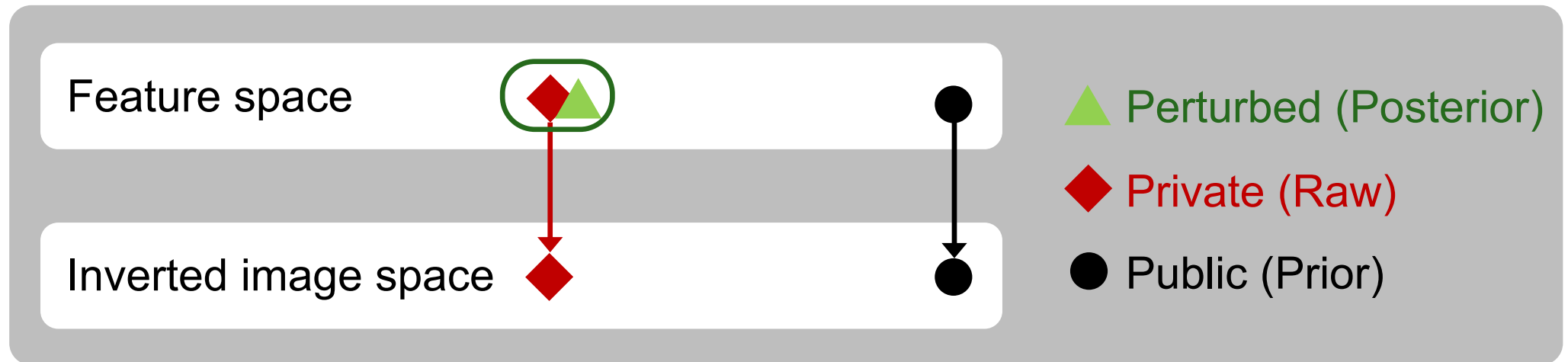
◆ Private (Raw)

● Public (Prior)

# Defense Intuitions (Privacy+Utility)

"Get Close" ✓   Find a feature perturbation that
1) is small;
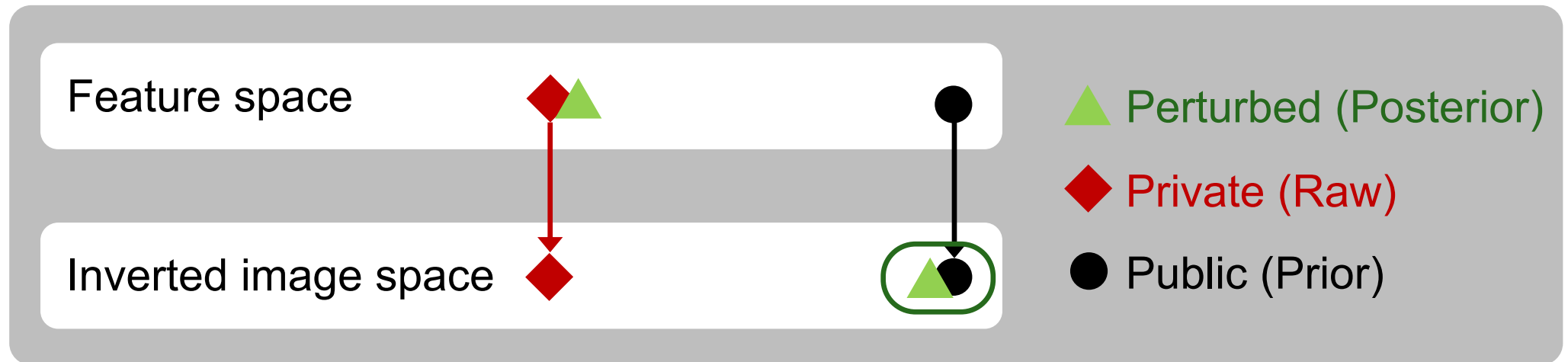2) draws inverted image close to public average faces.

Feature space   🔺 Perturbed (Posterior)

🔶 Private (Raw)

Inverted image space   ⚫ Public (Prior)

# Crafter

# Crafter



A game view:

Attack →
Defense →

**Conventional:** stay away from private identity

**Our Crafter:** get close to non-private prior

*The perturbed feature poisons adaptive attackers!*

*A stationary point. But not equilibrium.*

Round 1
Round 2
...

Released Round n

Private Identity 1

Released Feature 1

Adaptive attack

Released

Non-private Prior

Adaptive attempt

Private Identity i

Released Feature i

Adaptive attack

Released

# Crafter



A game view:

Attack ➡ (red arrow)
Defense ➡ (black arrow)

**Conventional**: stay away from private identity

**Our Crafter**: get close to non-private prior

*The perturbed feature poisons adaptive attackers!*

*A stationary point. But not equilibrium.*

Round 1
Round 2
...
Released Round n

Released Feature 1

Adaptive attack

Private Identity 1

Released

Adaptive result

Adaptive attempt

Non-private Prior

Released

Adaptive attack

Released Feature *i*

Private Identity *i*

# Crafter



A game view:

Attack →
Defense →

**Conventional**: stay away from private identity

**Our Crafter**: get close to non-private prior

*The perturbed feature poisons adaptive attackers!*

Round 1
Round 2
...
Released Round n

Private Identity 1

Released

Adaptive result

Released Feature 1

Adaptive attack

Adaptive attempt

Non-private Prior

*A stationary point. But not equilibrium.*

Released Feature *i*

Released
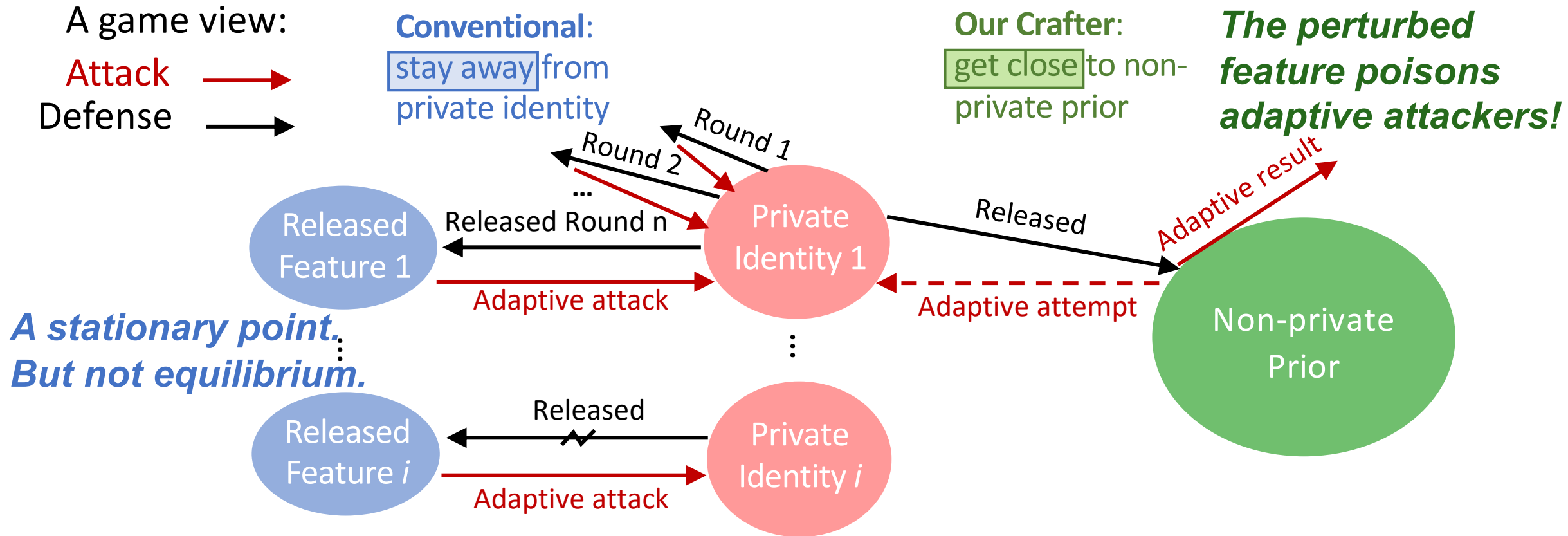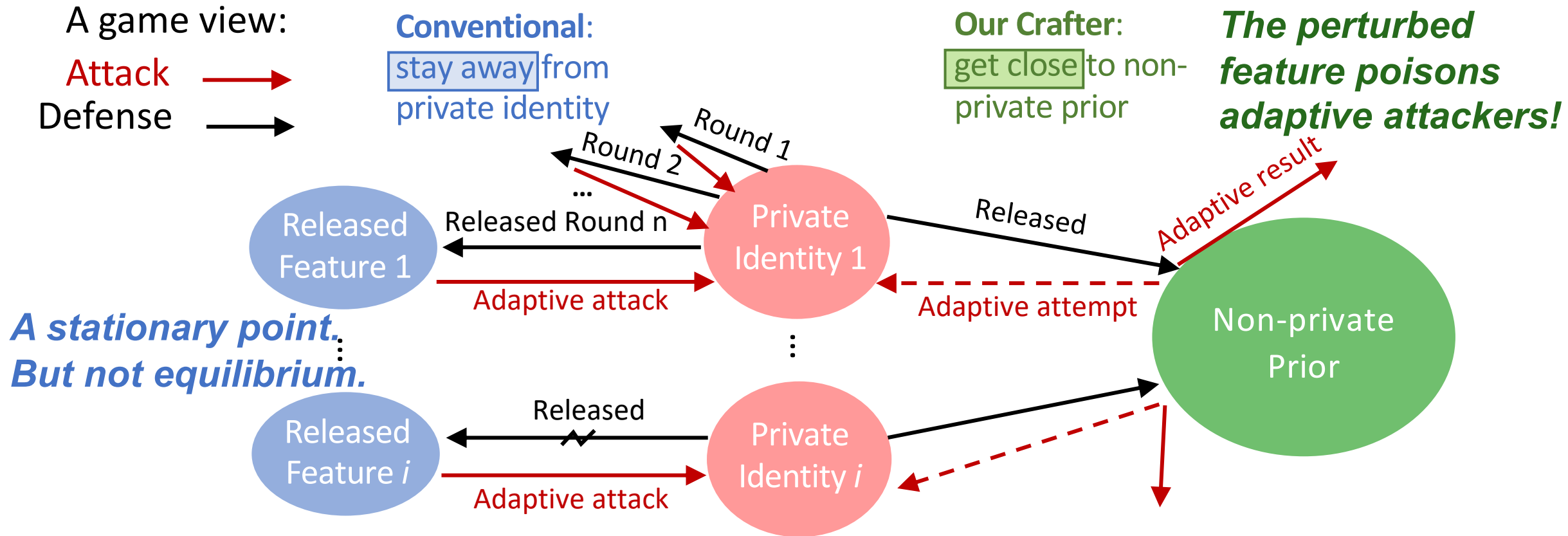
Private Identity *i*

Adaptive attack

# Evaluation

## Datasets

### CelebA (64*64)

- 40 binary utility attributes

### LFW (128*128)

- 10 binary utility attributes

### VGGFace2 (112*112)

- 5-class hair color utility attribute

## Baselines

### AdvLearn

- Deployment scenario

### Disco

- Deployment scenario
- Improves upon AdvLearn with a pruner

### TIPRDC

- Development scenario

Xiao et al. "Adversarial learning of privacy-preserving and task-oriented representations ", 2020

Singh et al. "Disco: Dynamic and invariant sensitive channel obfuscation for deep neural networks ", 2021

Li et al. "Tiprdc: task-independent privacy-respecting data crowdsourcing framework for deep learning with anonymized intermediate representations ", 2020

# Evaluation

**Tradeoff parameter.**

- **AdvLearn:** {0.1, 0.5, 0.8}

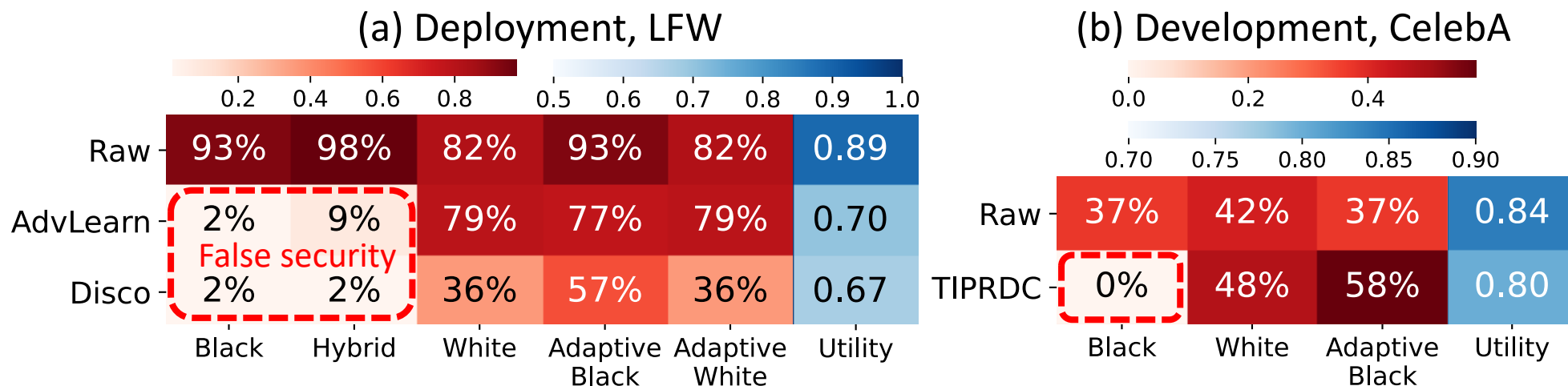- **Disco:** {0.2, 0.6, 0.8}

- **TIPRDC:** {0.1, 0.5, 0.8}

**Privacy Metrics.**

- **Eval Acc:** identification accuracy of the inverted images.

- **Feature Similarity:** cosine similarity between of the raw & inverted images.

- **SSIM:** pixel-level resemblance between the raw & inverted images.

- **Human study:** 35 human feedbacks, Macro-F1 score of reidentification.

|         | Black | Hybrid |
|---------|-------|--------|
| Raw     | 0.93  | 0.9    |
| AdvLearn| 0.02  | 0.0    |
| Disco   | 0.02  | 0.0    |
| Crafter | 0.32  | 0.3    |

|         | Black | Hybrid | White |
|---------|-------|--------|-------|
| Raw     | 0.50  | 0.50   | 1.00  |
| AdvLearn| 0.50  | 0.50   | 1.00  |
| Disco   | 0.50  | 0.50   | 1.00  |
| Crafter | 0.50  | 0.50   | 1.00  |

|        | Black | White | Adaptive Black |    |    |
|--------|-------|-------|----------------|----|----|
| Raw    | 0%    | 0%    | 0%             | 0% | 1% |
| TIPRDC | 0%    | 0%    | 1%             | 0% | 1% |
| Crafter| 0%    | 0%    | 0%             | 0% | 1% |

# Evaluation



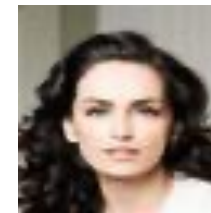(a) Deployment — Original, False security, Adv Learning, Crafter; Black, White, Adaptive Black, Adaptive White
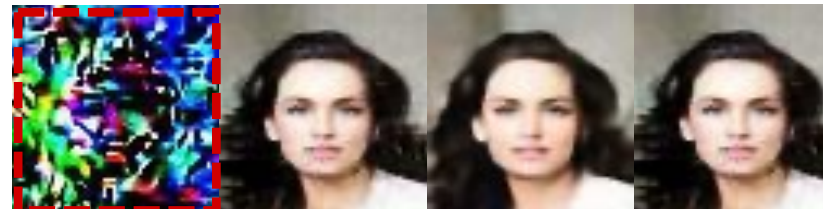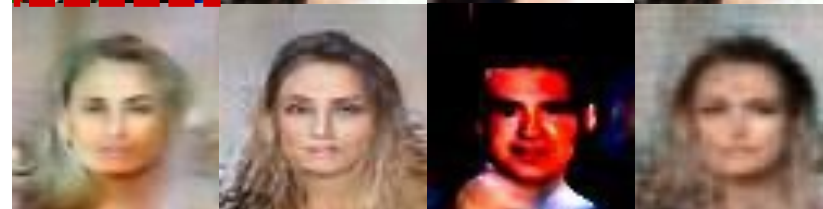
(b) Development — Original, False security, TIPRDC, Crafter; Black, White, Adaptive Black, Adaptive White

# Crafter: Facial Feature Crafting against Inversion-based Identity Theft on Deep Models

Shiming Wang, Zhe Ji, Liyao Xiang, Hao Zhang,
Xinbing Wang, Chenghu Zhou, Bo Li

Code Available @GitHub 👉