



SAND Lab
security, algorithms, networks and data

Understanding Non-Consensual “Undress” applications

Josephine Passananti, Heather Zheng, Ben Zhao
CSCW Digital Intimacy Workshop — Oct 18, 2025



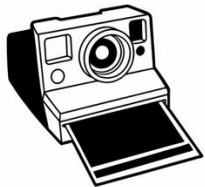
Trigger warning: non-consensual imagery references

Evolution of *Synthetic* Non-Consensual Intimate Imagery (SNCII)

< 2020

SNCII not an “industry” yet

- limited to physical & digital media

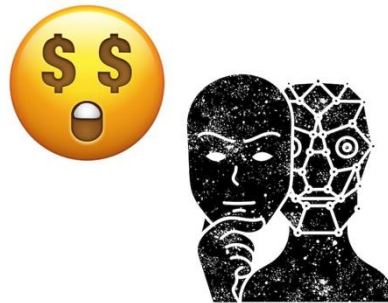


- Leaked photos & videos

2020-2022

GAN based face-swapping tools

- Not widely marketed



- Low-quality, expensive SNCII

2023 >

Generative *Undress AI* tools

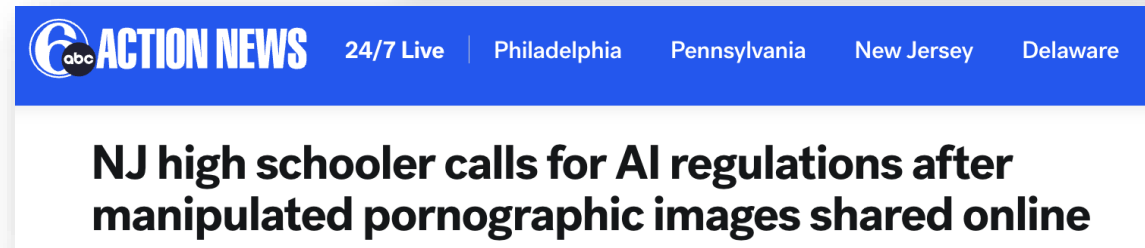
- **2,000% increase** in spam for “undress AI” services
- sites received **over 24 million visitors per month**

- High quality, low cost SNCII

SNCII aka “undress AI” is used maliciously

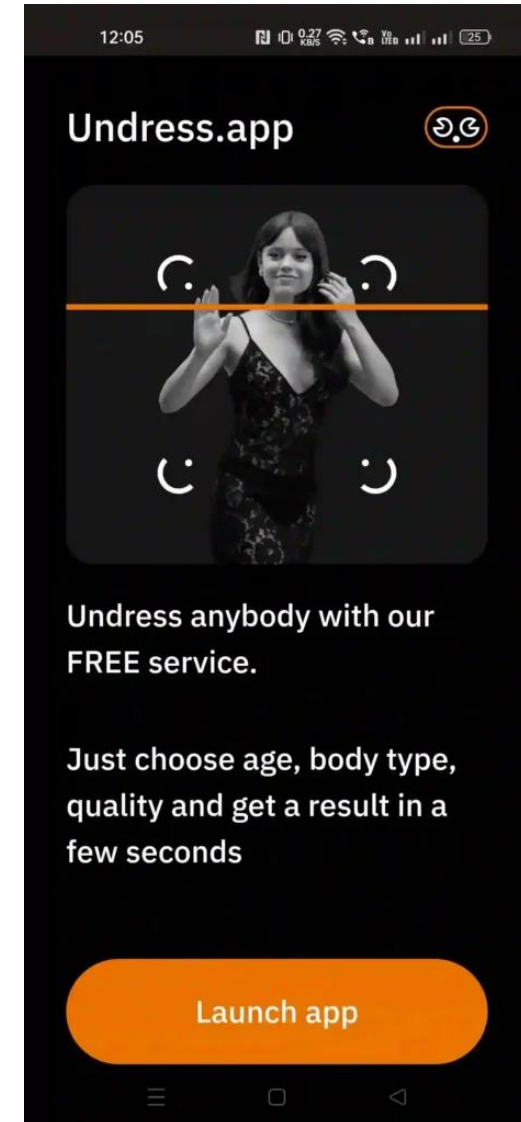
... 2025

Terms and Conditions hold the user responsible for his or her use of the app



Investigating how Undress AI is so pervasive

How easy is it to create an AI nude?



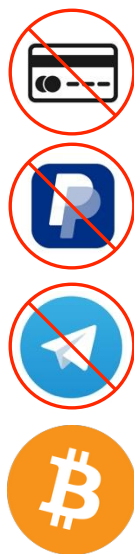
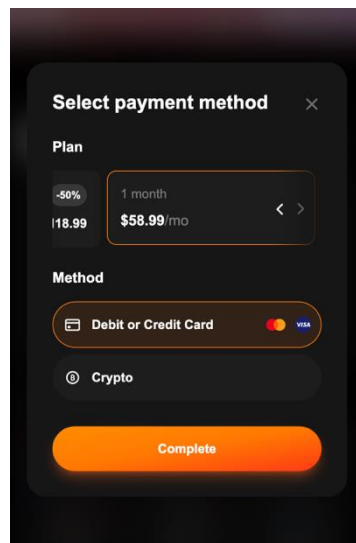
Investigating how Undress AI is so pervasive

Our experience:

1. Email address



2. Payment



How are teenagers buying credits with Crypto?

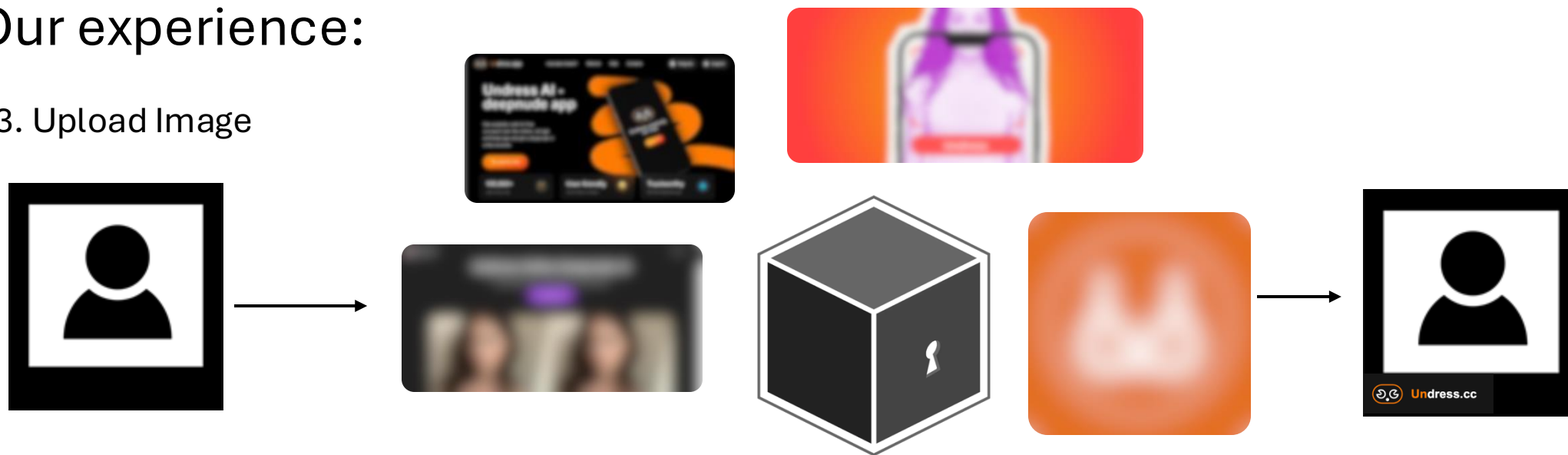


Drug Dealer Model?

Investigating how Undress AI is so pervasive

Our experience:

3. Upload Image



Black Box



How to reverse engineer the pipeline?

- Probe applications **ethically**:



Synthetic Images

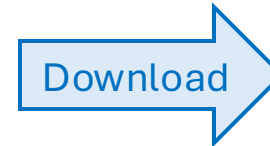
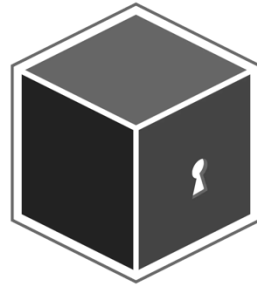


Minimize queries

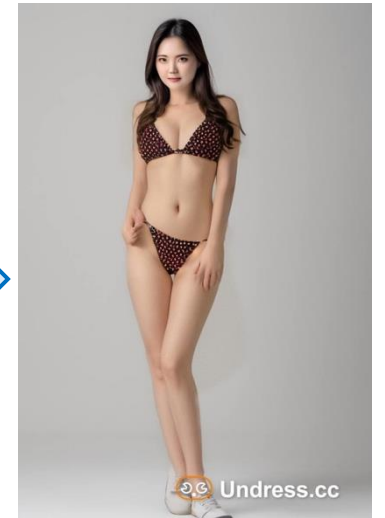
Investigating applications through probing

Undress AI application

Input Image



“Undressed”
Image



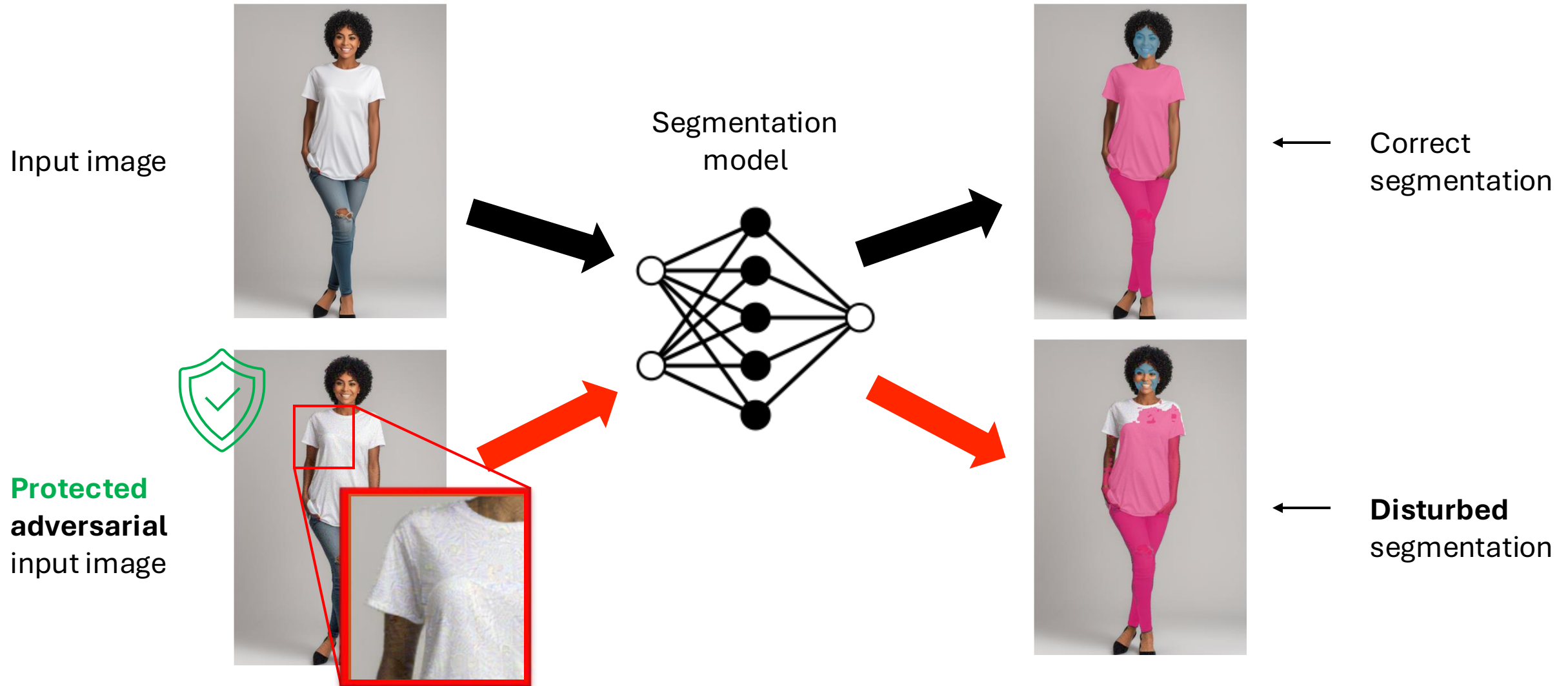
Investigating applications through probing

Undress AI application



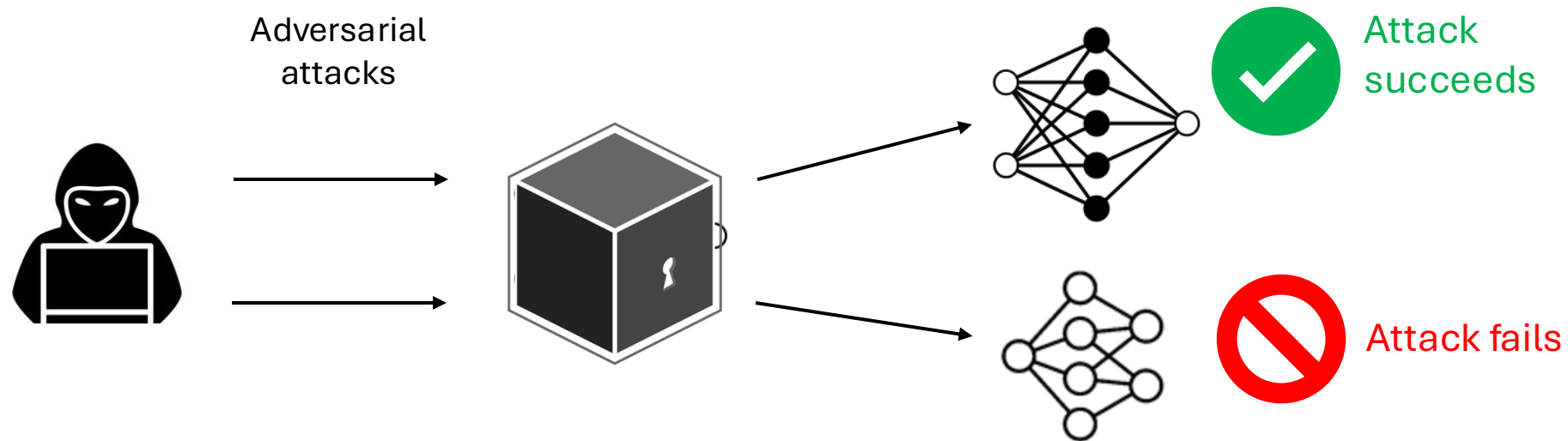
Expensive to train: likely using pre-trained, open-source models

Crafting Adversarial Probes



Running attacks against *black box* models

Successful adversarial attacks require knowledge of the target model



Each successful or failed probe can reveal something about the model

Attacks on AI Undress Applications

Fine line between successful vs unsuccessful attacks



Input



Output



Attack:
Pixel-Based
Segmentation

Attack *fails*



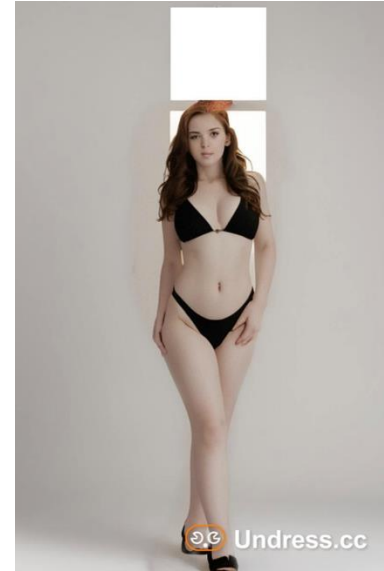
Attack:
Censored Face



Attack *fails*



Attack:
Pixel-Based
Segmentation
+
Censored Face



Attack *succeeds*

Attacks on AI Undress Applications

Fine line between successful vs unsuccessful attacks

Clues: what a successful segmentation attack would look like



Input



Segmentation

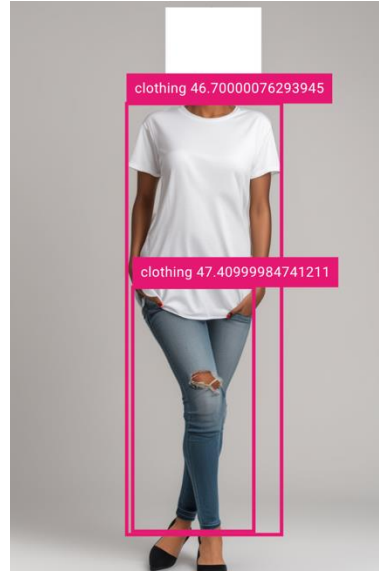


Attack:
Pixel-Based
Segmentation

Attack *fails*



Attack:
Censored Face



Attack *fails*



Attack:
Pixel-Based
Segmentation
+
Censored Face



Attack *succeeds*