



Gradient-Leaks: Understanding Deanonymization in Federated Learning

Tribhuvanesh Orekondy¹, Seong Joon Oh², Yang Zhang³, Bernt Schiele¹, Mario Fritz³

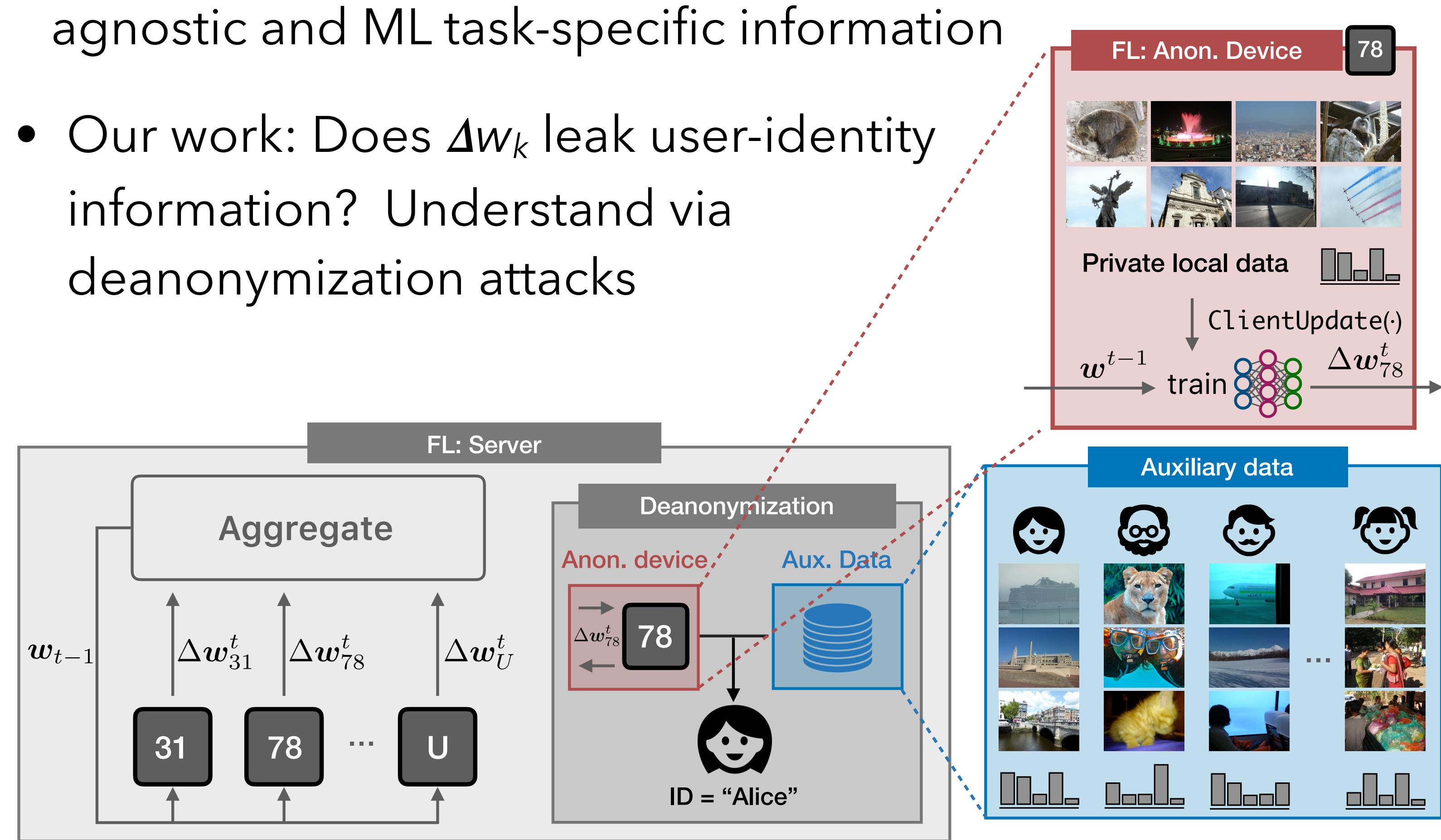
¹ Max Planck Institute for Informatics

² LINE Plus (Naver) Clova ML

³ CISPA Helmholtz Centre for Information Security

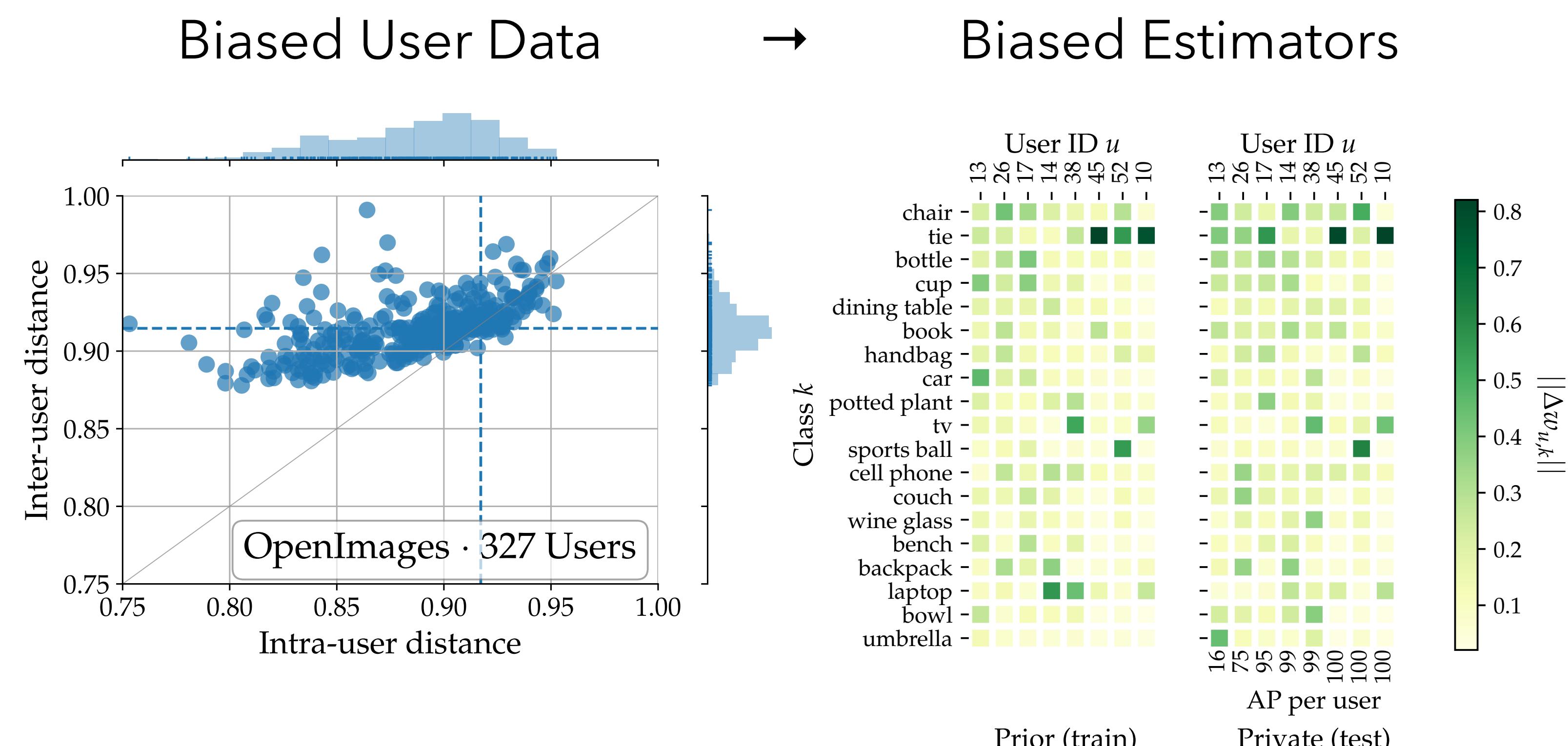
Motivation

- Federated Learning: promising approach to ensure user/data privacy for training ML models on devices e.g., smartphones
- Only model parameter deltas Δw_k are anonymously shared: $\Delta w_k = \text{ClientUpdate}(\text{Private user data})$
- To protect users' privacy, want Δw_k to encode only user-agnostic and ML task-specific information
- Our work: Does Δw_k leak user-identity information? Understand via deanonymization attacks



Deanonymization Attacks in FL

Insight



Threat Model

$$(\Delta w_{\text{anon}}, \Delta w_u^{\text{aux}}) \mapsto \mathbb{P}(\text{anon} = u)$$

=ClientUpdate(Aux. data of user u)

Observed by Attacker

Attack Models

- Re-identification: $f^{\text{re-id}} : \Delta w_{\text{anon}} \mapsto u$ ⇒ Learnt using MLP
- Matching: $f^{\text{mat}} : (\Delta w_{u_i}, \Delta w_{u_j}) \mapsto \mathbb{P}(u_i = u_j)$ ⇒ Siamese Network

Acknowledgment This research was partially supported by the German Research Foundation (DFG CRC 1223)

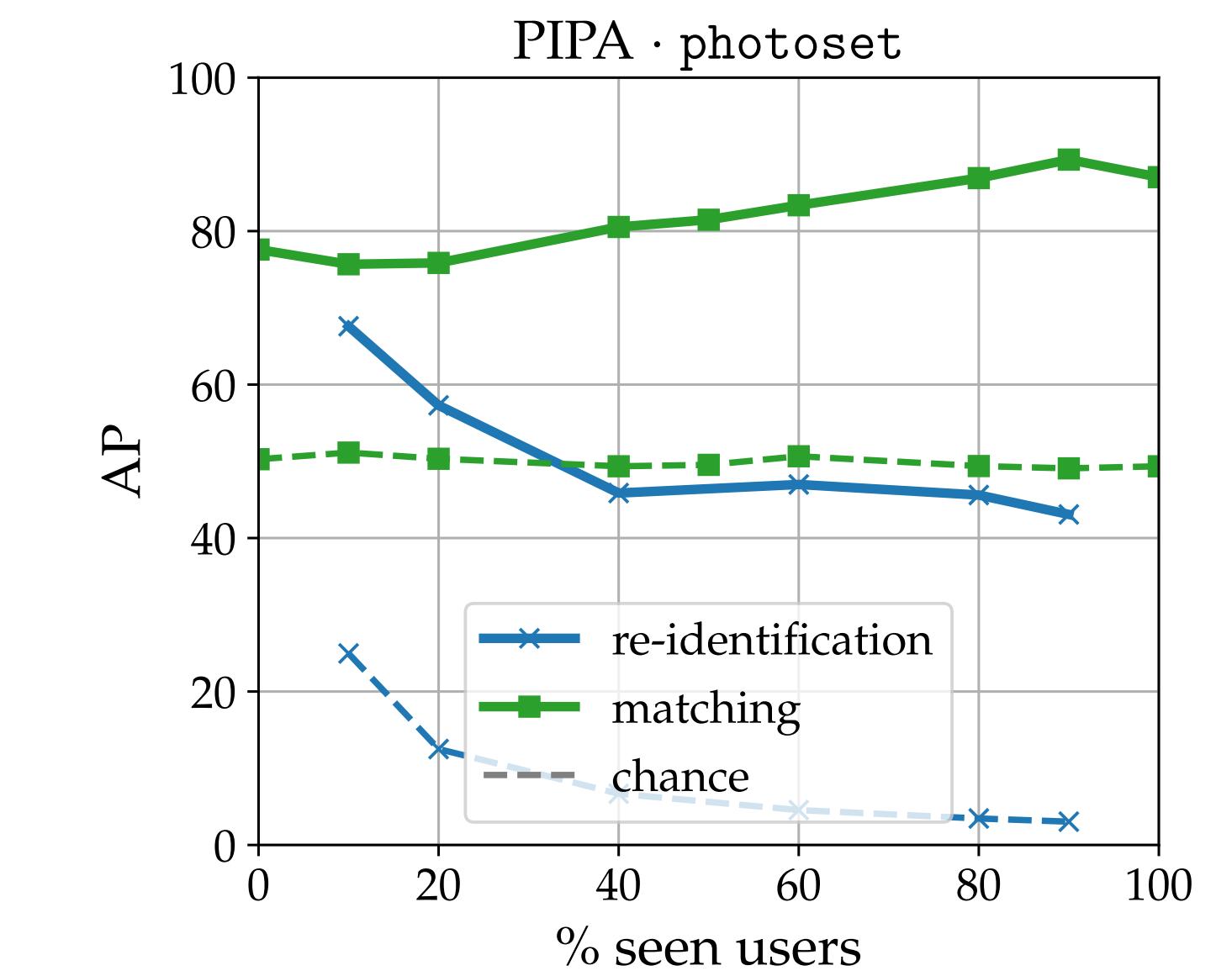
Take-aways

- Task: Deanonymization attack of devices in FL
- Insight: Use user selection bias as a quasi-identification statistical signal to perform deanonymization
- Attacks are effective: 16-91% AP, 19-175x chance-level re-identification performance
- Poses a threat to privacy and anonymity of users participating in FL

Evaluation

How Effective are Deanonymization Attacks?

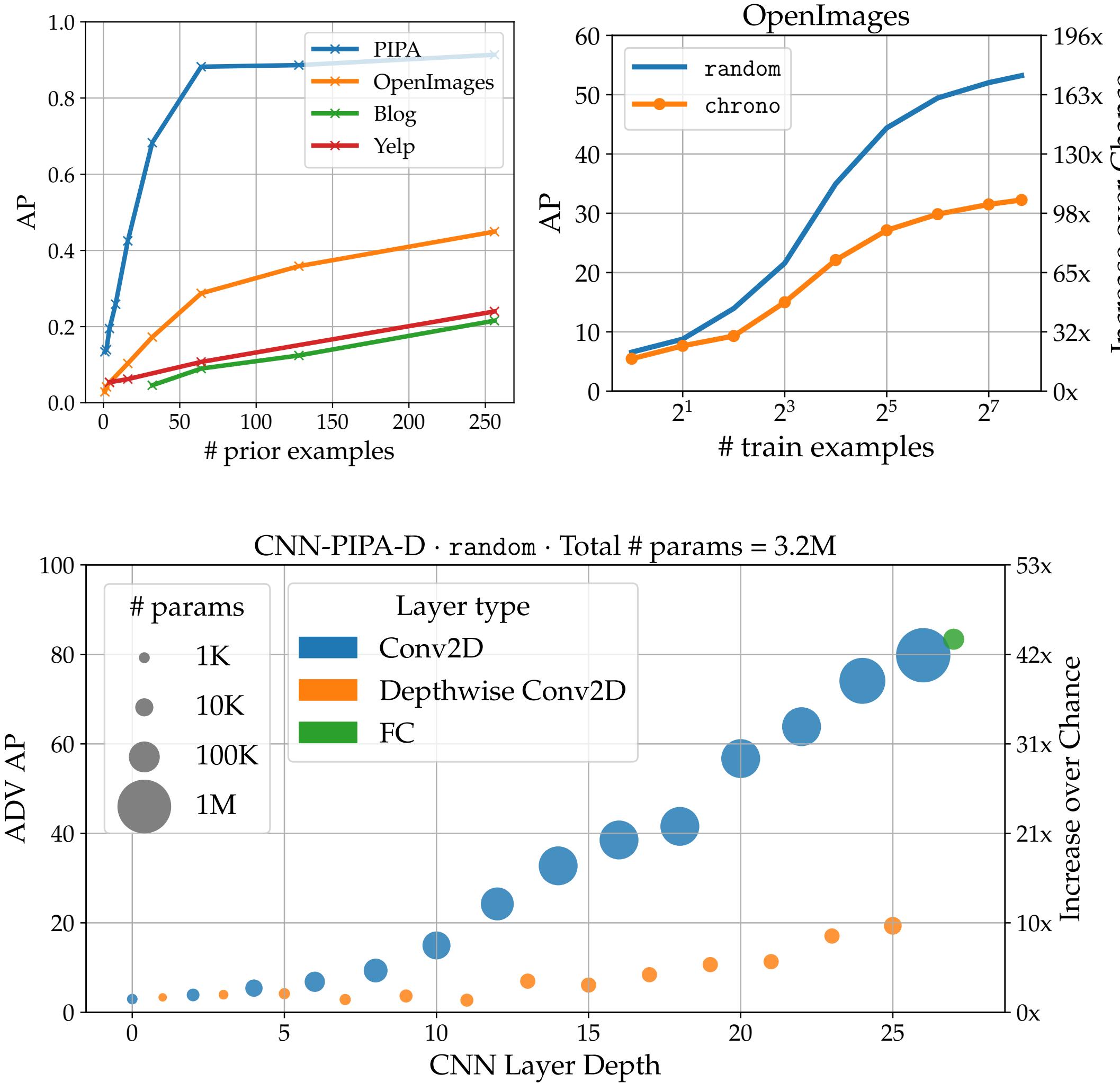
		Re-ID	Matching
PIPA	random	91 (48x)	99.5 (2x)
	photoset	42.2 (22x)	91.2 (1.85x)
Open Images	random	53.7 (175x)	98.2 (1.93x)
	chrono	32.5 (106x)	94.8 (1.93x)
Blog	random	52.9 (29x)	95.3 (1.9x)
	chrono	44.8 (25x)	91.9 (1.89x)
Yelp	random	23.5 (28x)	83.4 (1.7x)
	chrono	16 (19x)	79.3 (1.56x)



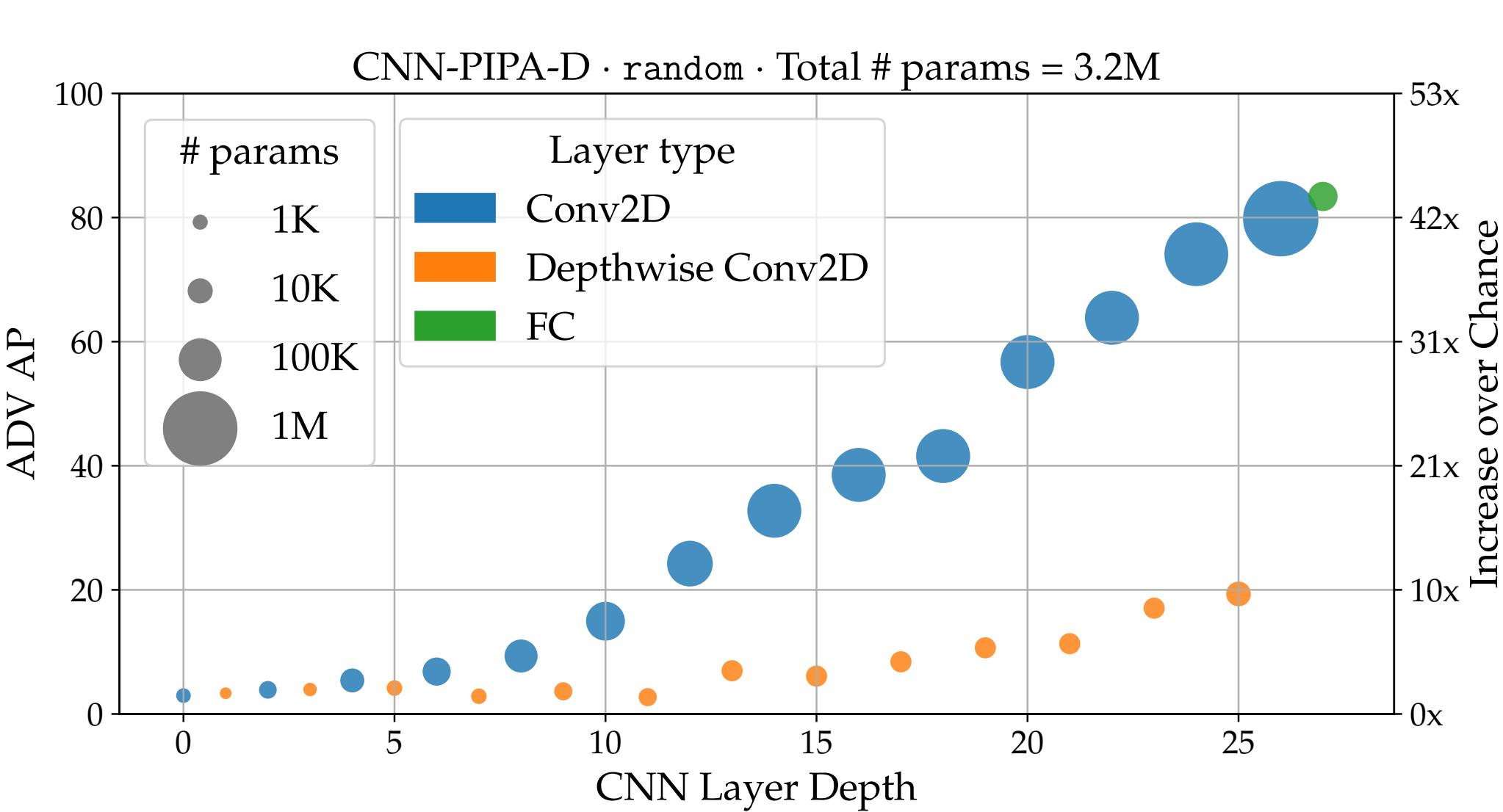
- Closed-world: 19-175x chance-level performance Re-ID performance

- Open-world: Robust to encountering new unseen users at test-time

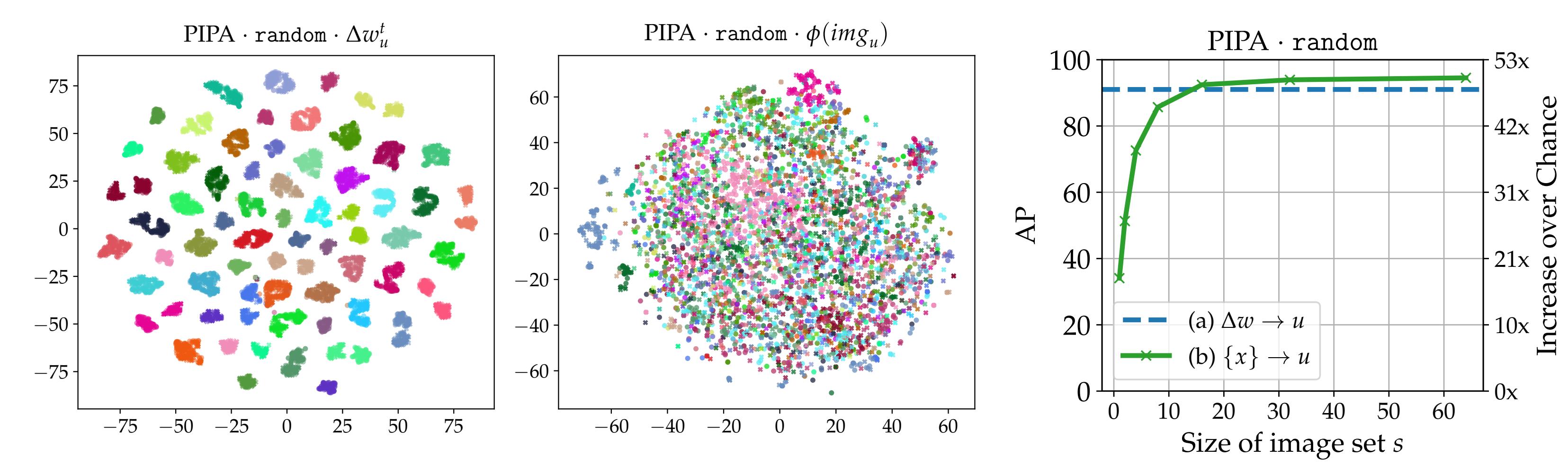
Analysis



- vs. # attacker's training examples: Attacks possible in few-shot training settings



- vs. parameter type: All parameters provide above-chance level information
- FC most informative



- Δw_k encodes aggregated data information → sometimes easier to deanonymize via Δw_k than raw images themselves