



The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
UNIVERSITY OF TORONTO

More than Enough is Too Much: Adaptive Defenses against Gradient Leakage in Production Federated Learning

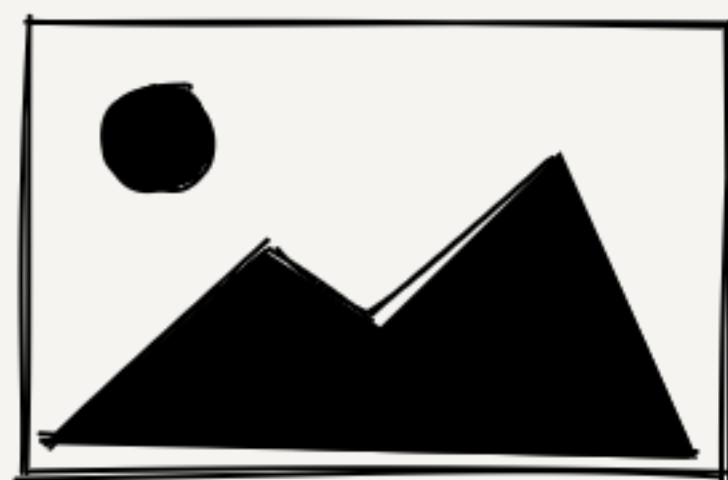
Fei Wang, Ethan Hugh, Baochun Li

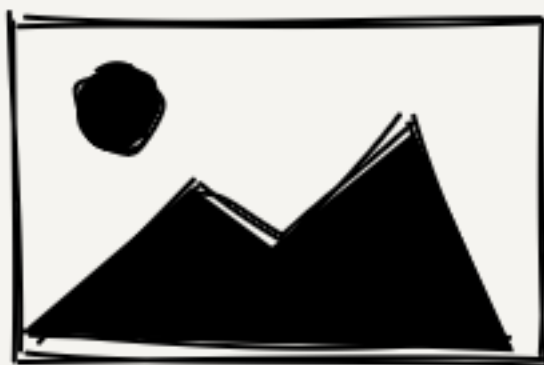
Department of Electrical and Computer Engineering

University of Toronto

INFOCOM'23

A Glance at Gradient Leakage Attacks





Model
 F_w



Prediction



Model
 F_w



Prediction

$$\nabla_w = \frac{\partial}{\partial w} \text{Loss}(\text{Prediction}, \text{Label})$$



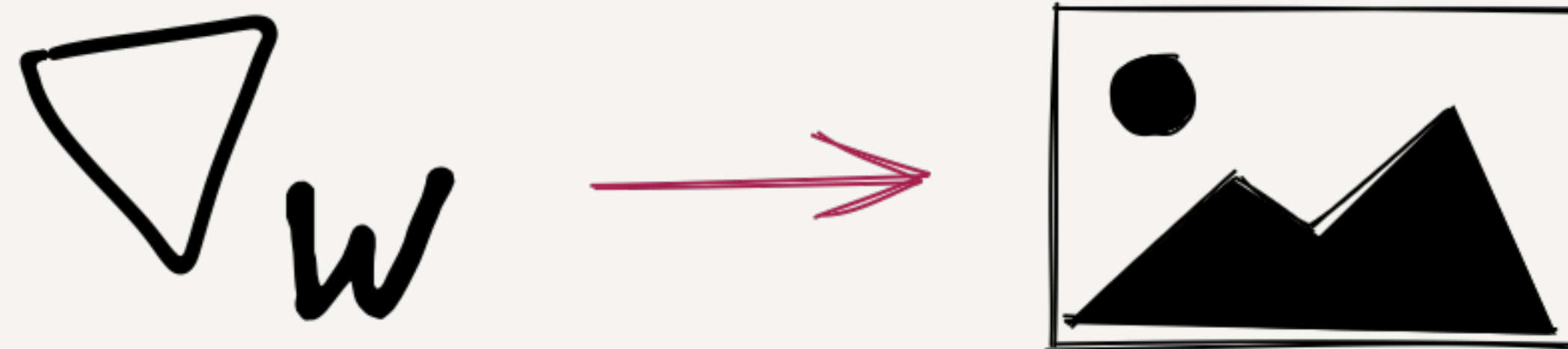
Model
 F_w

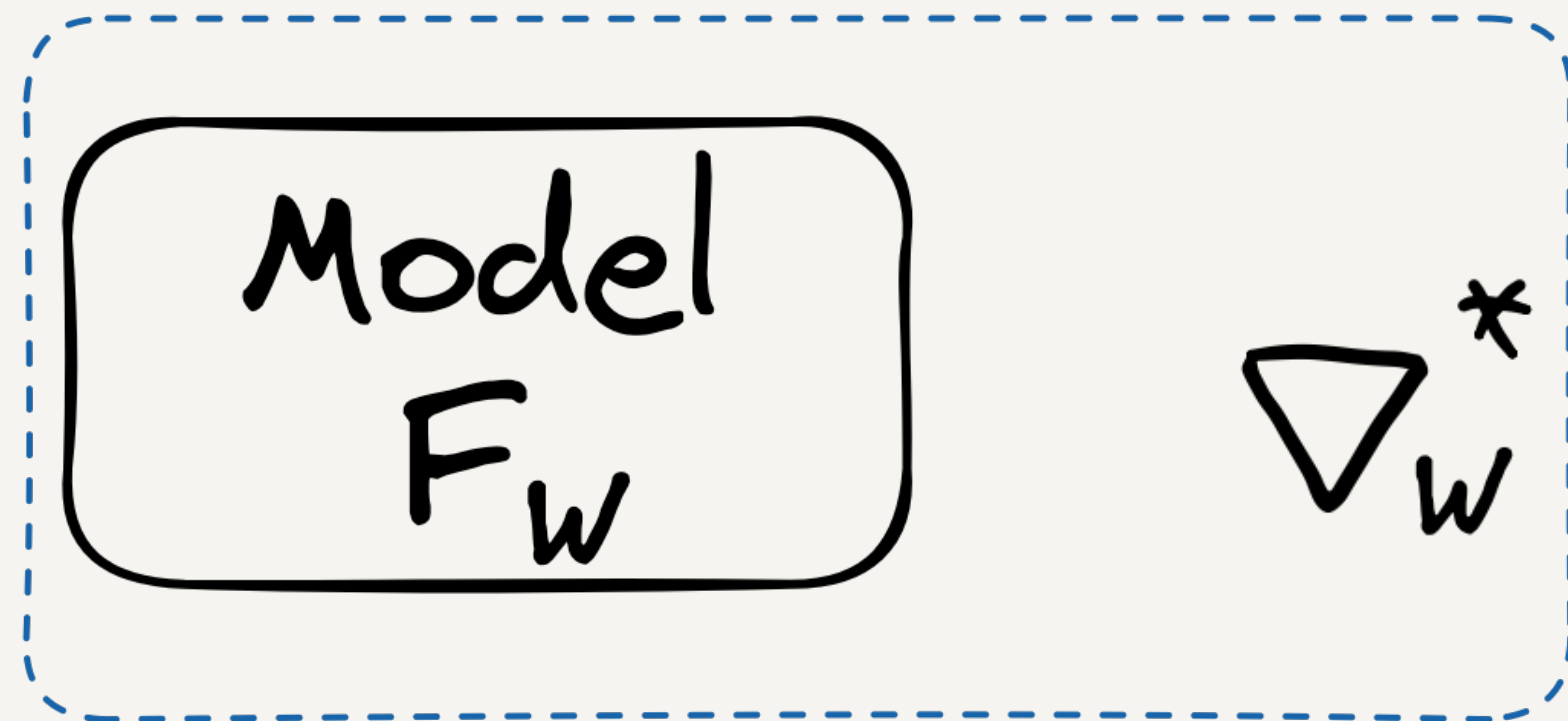


Prediction

$$\nabla_w = \frac{\partial}{\partial w} \text{Loss}(\text{Prediction}, \text{Label})$$

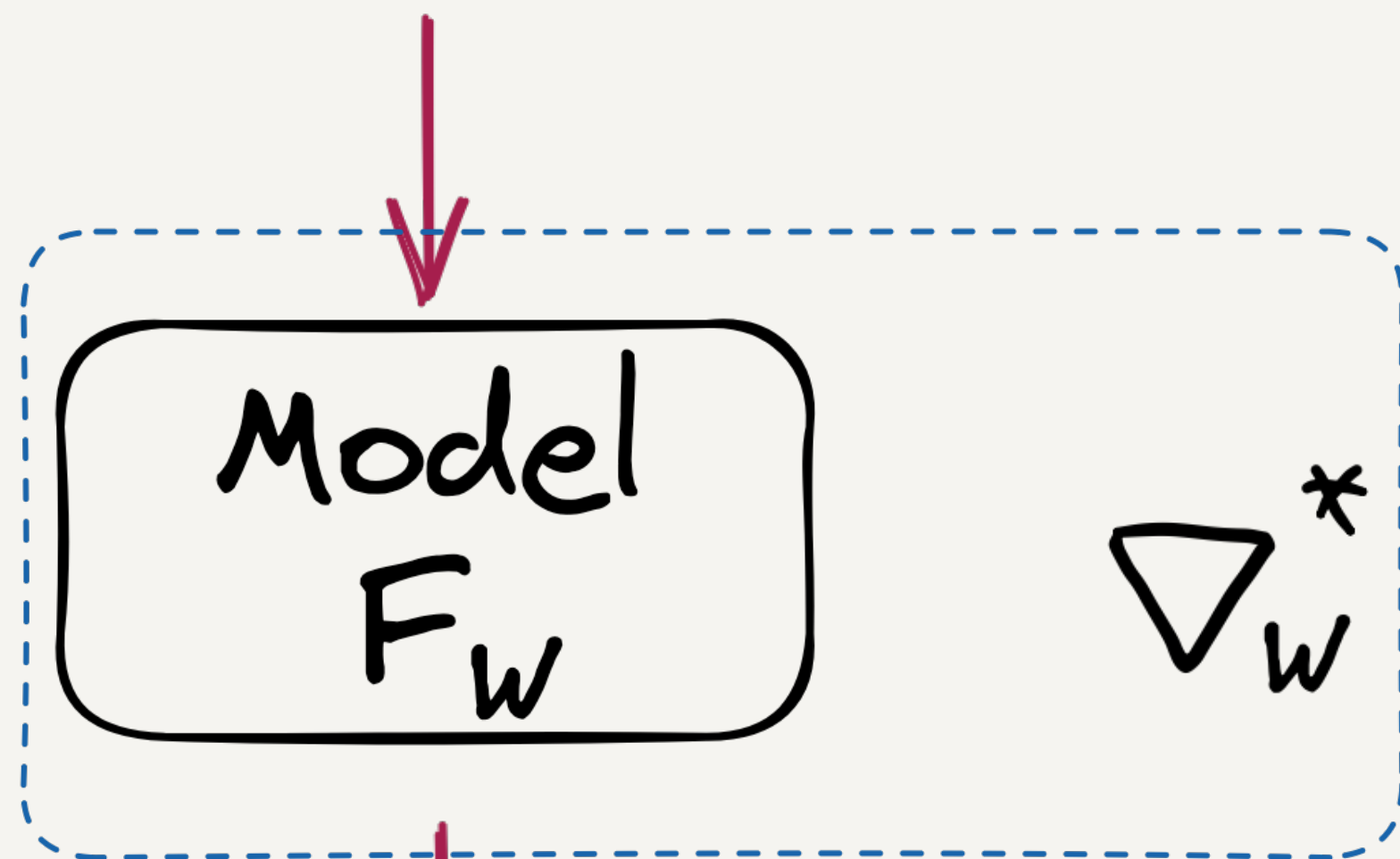
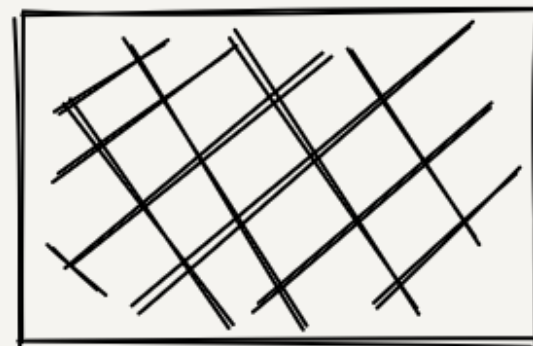
$$w' = w - \eta \nabla_w$$





Dummy Data

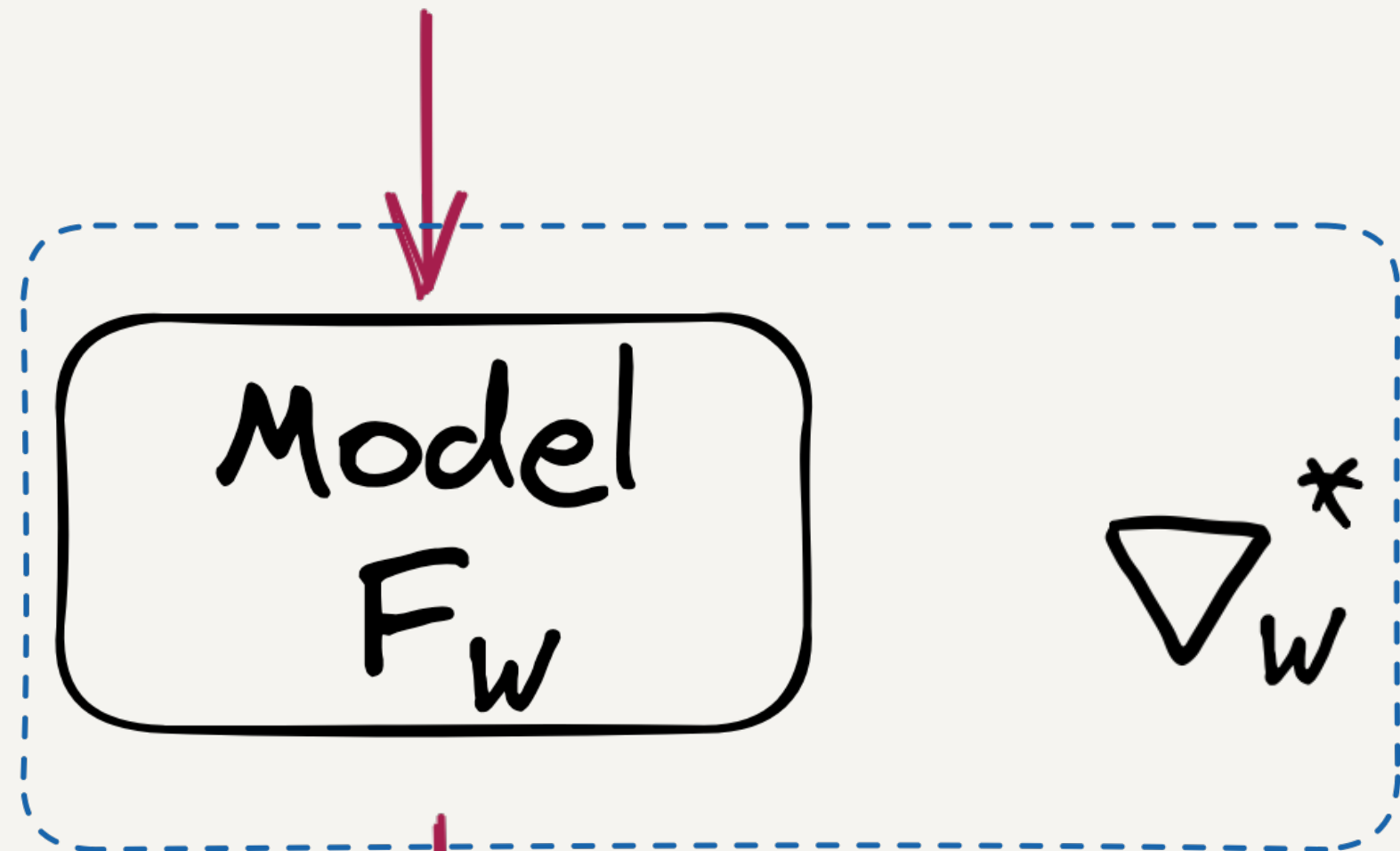
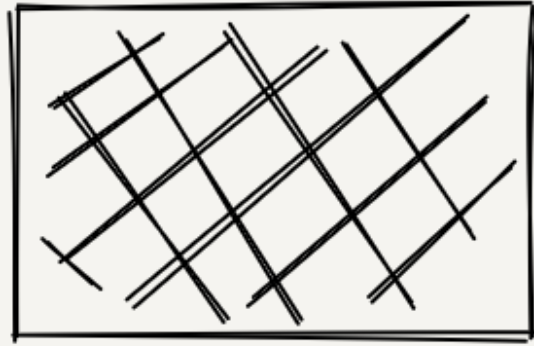
x'



Prediction

Dummy Data

x'



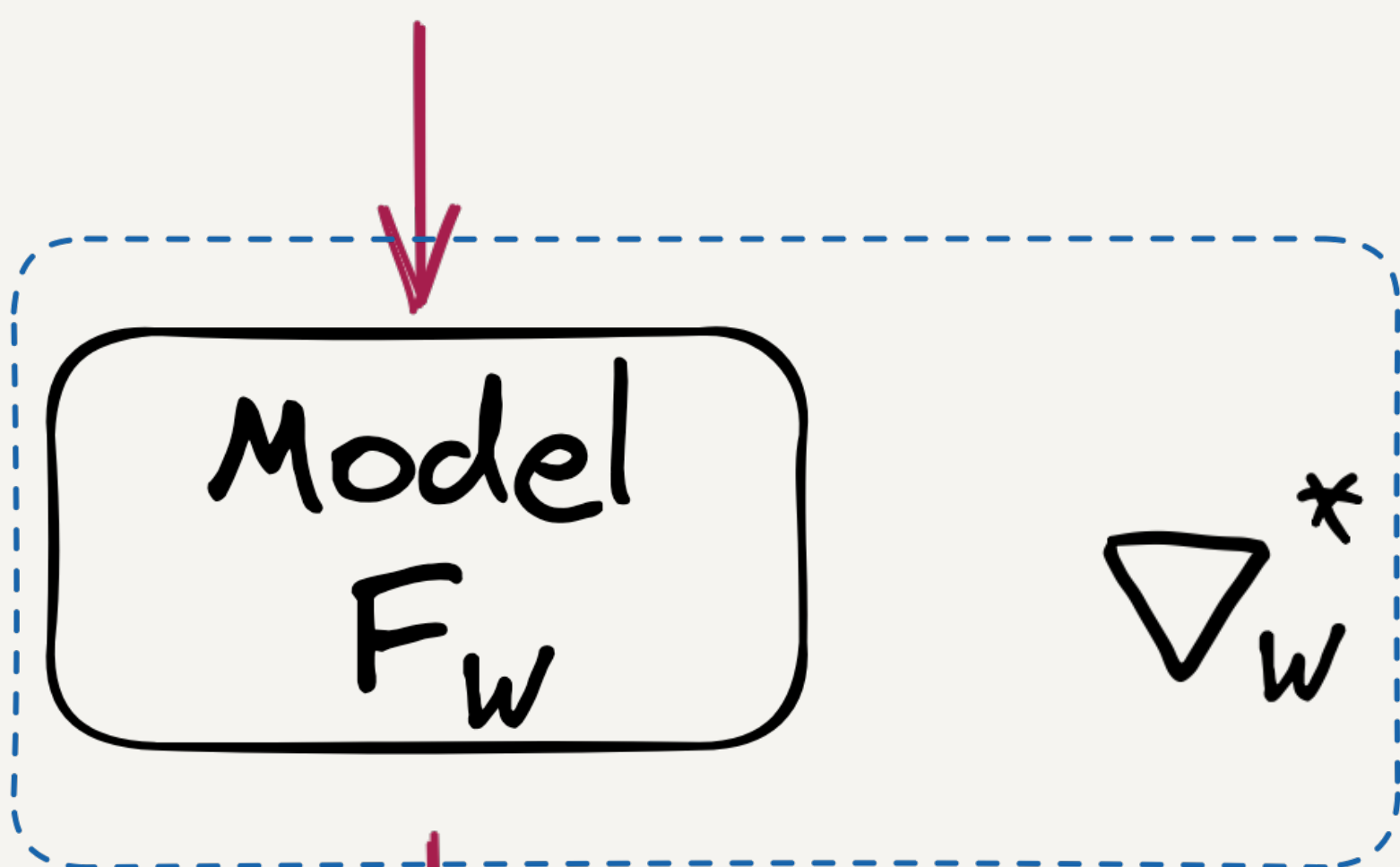
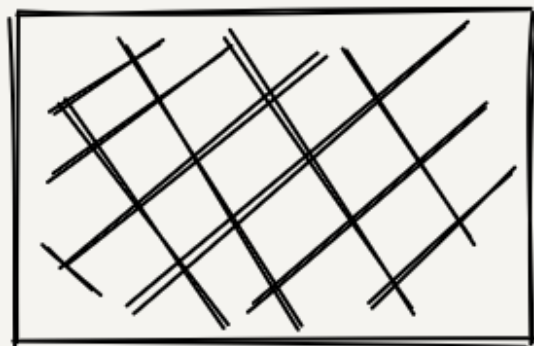
Prediction

$$\nabla_w' = \frac{\partial}{\partial w} \text{Loss}(\text{Prediction}, \text{Dummy Label})$$

$F_w(x')$ y'

Dummy Data

x'



Prediction

Optimization Objective

$$x'^*, y'^* =$$

$$\arg \min_{x', y'} \text{Distance} (\nabla_w^*, \nabla_w')$$

$$\nabla_w^*$$

$$\nabla_w' = \frac{\partial}{\partial w} \text{Loss}(\underbrace{\text{Prediction}}_{F_w(x')}, \underbrace{\text{Dummy Label}}_{y'})$$

Stealing Client's data
in Federated Learning

Stealing Client's data
in Federated Learning

Client's gradients are shared with the server.

Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

Client's gradients are ~~shared~~ with the server.

Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

Client's gradients are ~~shared~~ with the server.

$$\Delta = w' - w$$

Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

Client's gradients are ~~shared~~ with the server.

$$\Delta = w' - w$$

$$\Delta \neq \nabla !$$

Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

$$\nabla = \frac{\Delta}{-\eta}$$

Wei et al., "A Framework for Evaluating Client Privacy Leakages in Federated Learning,"
Proc. European Symposium on Research in Computer Security 2020.

Wu et al., "Fast-Convergent Federated Learning with Adaptive Weighting,"
IEEE Trans. on Cognitive Communications and Networking 2021.

$$\Delta \neq \frac{\Delta}{\eta}$$

Wei et al., "A Framework for Evaluating Client Privacy Leakages in Federated Learning,"

Proc. European Symposium on Research in Computer Security 2020.

Wu et al., "Fast-Convergent Federated Learning with Adaptive Weighting,"

IEEE Trans. on Cognitive Communications and Networking 2021.

$$\nabla \neq \frac{\Delta}{-\eta}$$

There are multiple steps of gradient descent in one round of client's local training.

Wei et al., "A Framework for Evaluating Client Privacy Leakages in Federated Learning,"

Proc. European Symposium on Research in Computer Security 2020.

Wu et al., "Fast-Convergent Federated Learning with Adaptive Weighting,"

IEEE Trans. on Cognitive Communications and Networking 2021.

$$\nabla \neq \frac{\Delta}{-\eta}$$

There are multiple steps of gradient descent in one round of client's local training.

Data samples $\gg 1$

Batch size \ll # Data samples

Epochs > 1

Wei et al., "A Framework for Evaluating Client Privacy Leakages in Federated Learning,"

Proc. European Symposium on Research in Computer Security 2020.

Wu et al., "Fast-Convergent Federated Learning with Adaptive Weighting,"

IEEE Trans. on Cognitive Communications and Networking 2021.

$$\nabla \neq \frac{\Delta}{-\eta}$$

There are multiple steps of gradient descent in one round of client's local training.

More sophisticated gradient descent algorithms are routinely used.

Data samples $\gg 1$

Batch size \ll # Data samples

Epochs > 1

Wei et al., "A Framework for Evaluating Client Privacy Leakages in Federated Learning,"

Proc. European Symposium on Research in Computer Security 2020.

Wu et al., "Fast-Convergent Federated Learning with Adaptive Weighting,"

IEEE Trans. on Cognitive Communications and Networking 2021.

$$\nabla \neq \frac{\Delta}{-\eta}$$

There are multiple steps of gradient descent in one round of client's local training.

Data samples $\gg 1$

Batch size \ll # Data samples

Epochs > 1

More sophisticated gradient descent algorithms are routinely used.

Momentum

Weight decay

Learning rate scheduler

...

Wei et al., "A Framework for Evaluating Client Privacy Leakages in Federated Learning,"

Proc. European Symposium on Research in Computer Security 2020.

Wu et al., "Fast-Convergent Federated Learning with Adaptive Weighting,"

IEEE Trans. on Cognitive Communications and Networking 2021.

$$x'^*, y'^* = \arg \min_{x', y'} \text{Distance} (\nabla_w^*, \nabla_w')$$

$$x'^*, y'^* =$$

$$\arg \min_{x', y'} \text{Distance} \left(\cancel{\Delta_w^*}, \Delta_w' \right)$$
$$\Delta_w^*, \Delta_w' ?$$

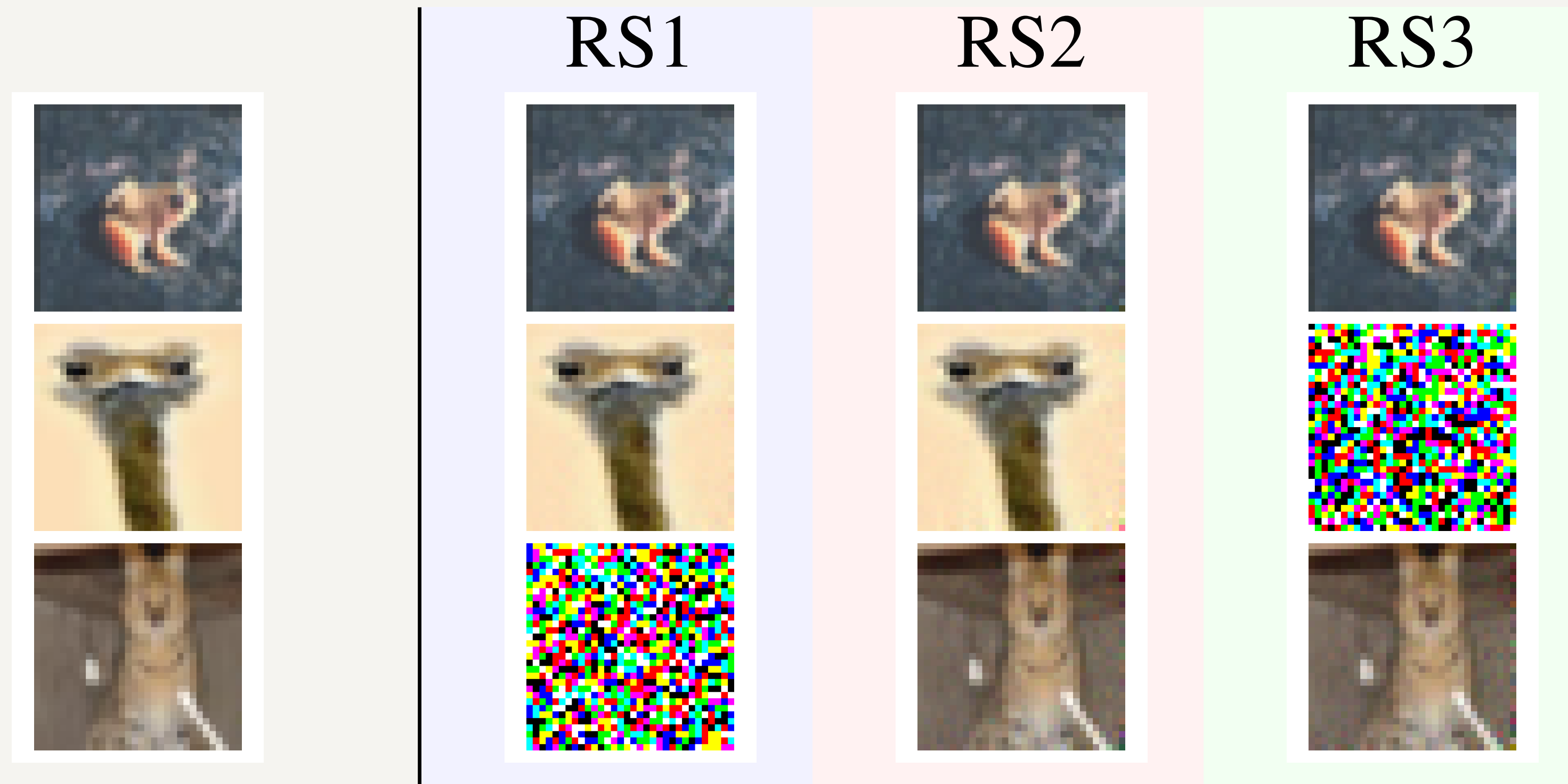
$$x'^*, y'^* = \arg \min_{x', y'} \text{Distance} \left(\cancel{\Delta_w^*}, \cancel{\Delta_w'} \right)$$

$$\Delta_w^*, \Delta_w' ?$$

To realize the same gradient descent process using the dummy data instead, the server requires a series of prior knowledge.

Ground truth

Untrained network with explicit initialization



Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

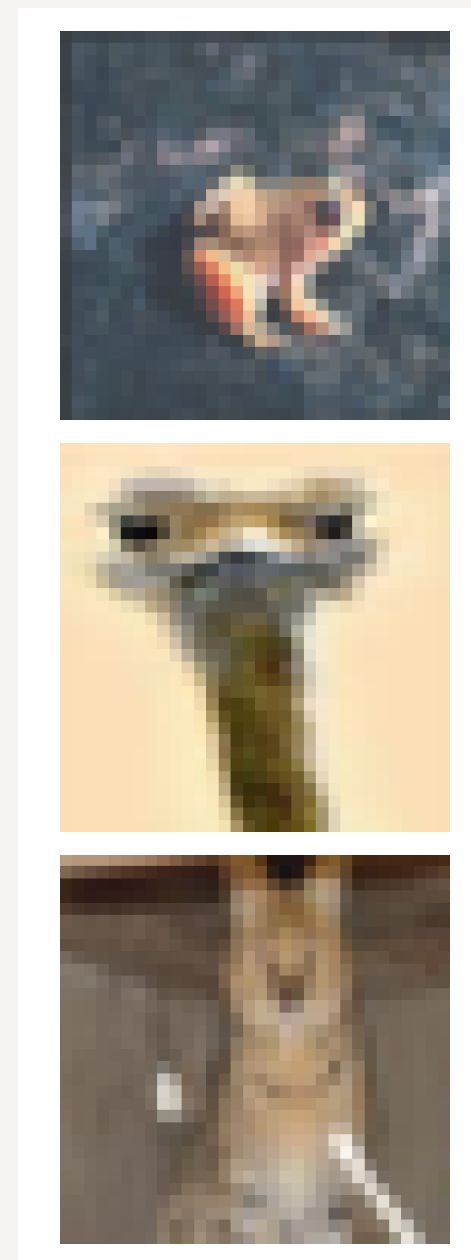
Geiping et al., "Inverting Gradients — How Easy Is It to Break Privacy in Federated Learning," NeurIPS 2020.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

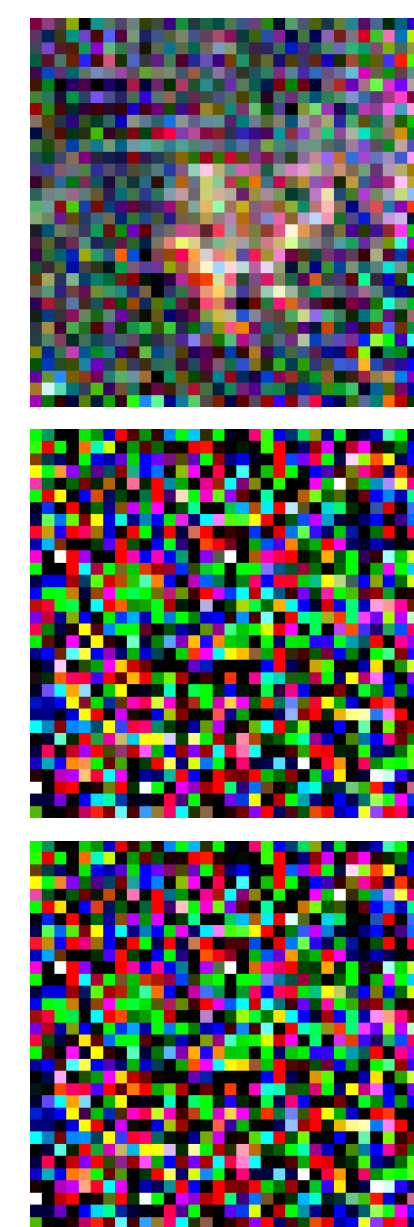
Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

Ground truth

Untrained network with default PyTorch initialization



RS1



RS2



RS3



Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

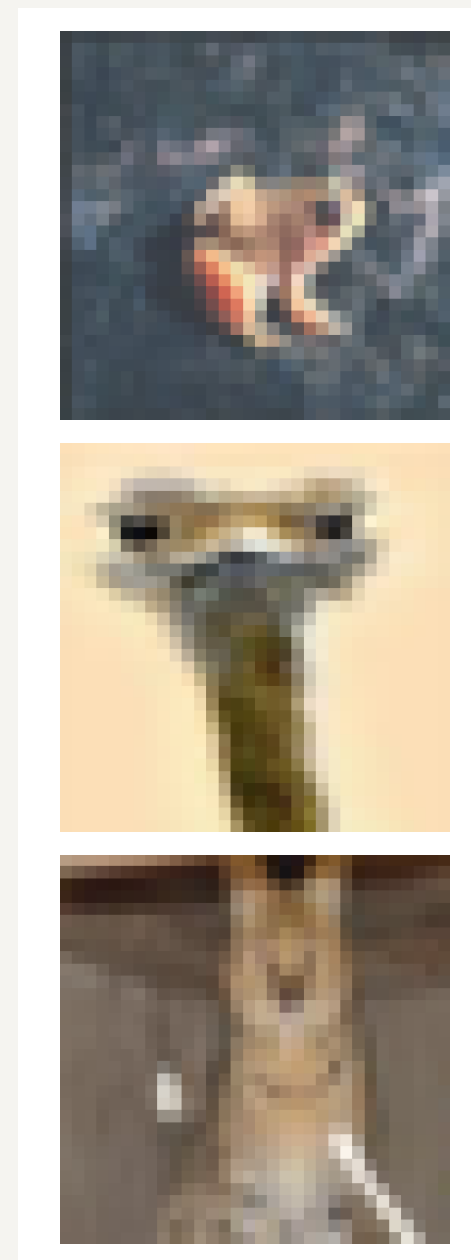
Geiping et al., "Inverting Gradients — How Easy Is It to Break Privacy in Federated Learning," NeurIPS 2020.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

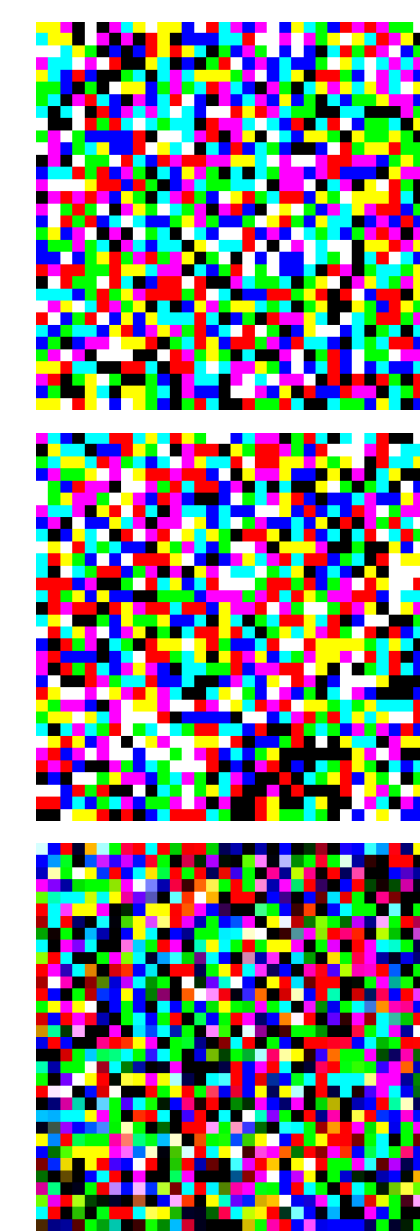
Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

Ground truth

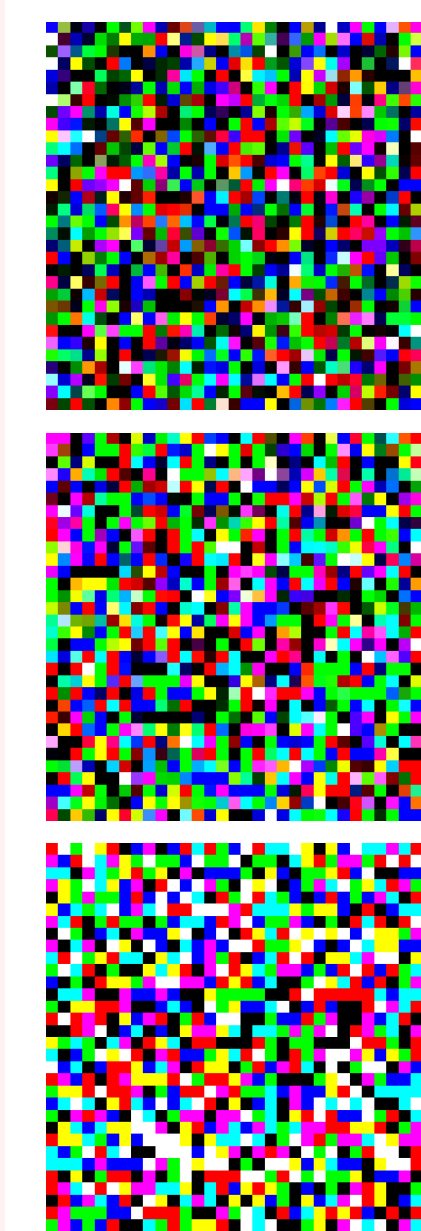
**Trained network pre-trained
with the same data**



RS1



RS2



RS3



Zhu et al., "Deep Leakage from Gradients," NeurIPS 2019.

Geiping et al., "Inverting Gradients — How Easy Is It to Break Privacy in Federated Learning," NeurIPS 2020.

Zhao et al., "iDLG: Improved Deep Leakage from Gradients," arXiv 2020.

Jeon et al., "Gradient Inversion with Generative Image Prior," NeurIPS 2021.

Outpost: Our Lightweight Defense

Sufficient and self-adaptive protection

Sufficient and self-adaptive protection

Selective perturbation

Sufficient and self-adaptive protection

The diagonal of the
Fisher information matrix

Selective perturbation

Sufficient and self-adaptive protection

The diagonal of the
Fisher information matrix

Selective perturbation

Sufficient and self-adaptive protection

Different intensity

The diagonal of the
Fisher information matrix

Selective perturbation

Sufficient and self-adaptive protection

Different intensity

The range of weight values

The diagonal of the
Fisher information matrix

Selective perturbation

Sufficient and self-adaptive protection

Different intensity

The range of weight values

The diagonal of the
Fisher information matrix

Selective perturbation

Gradually decreased likelihood

Sufficient and **self-adaptive** protection

Different intensity

The range of weight values

The diagonal of the Fisher information matrix

With each local step of gradient descent

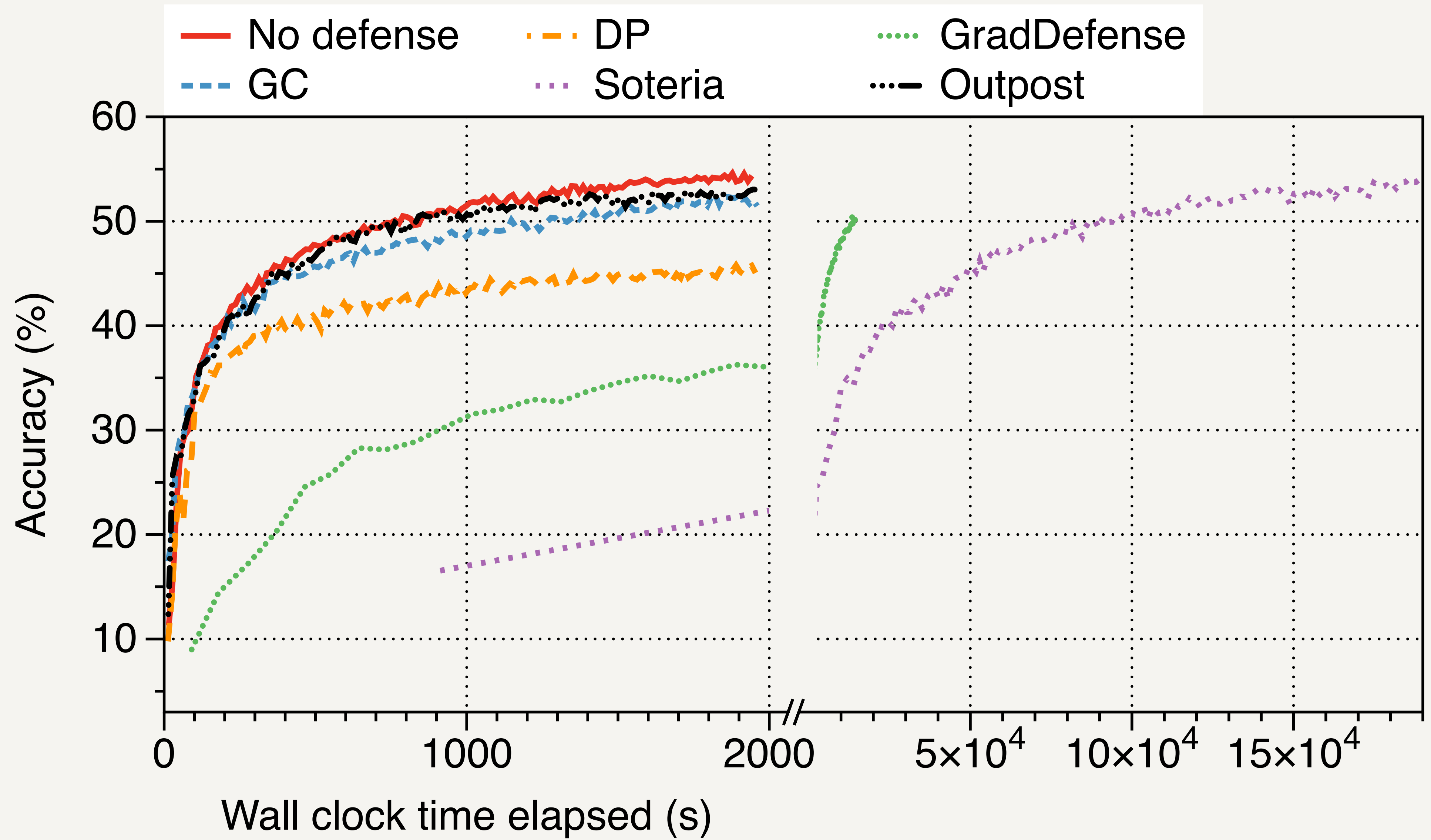
Gradually decreased likelihood

Selective perturbation

Sufficient and **self-adaptive** protection

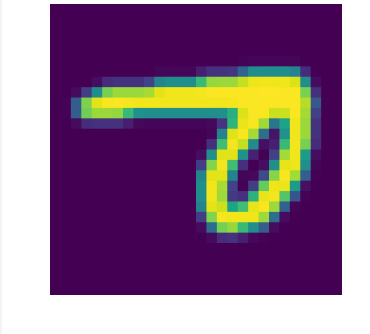
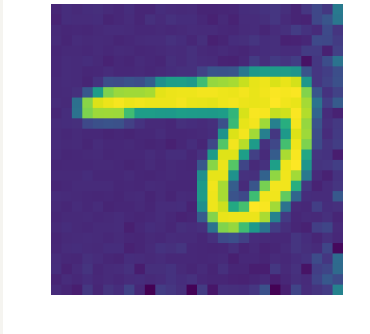
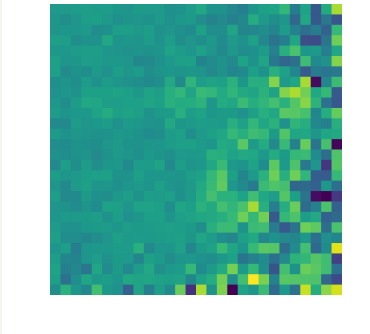
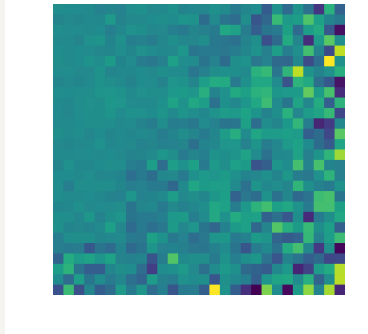
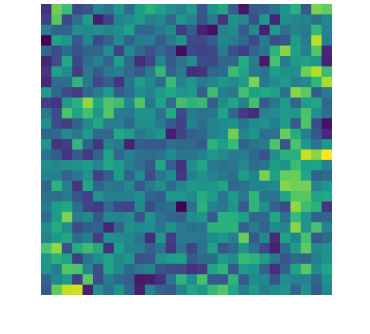
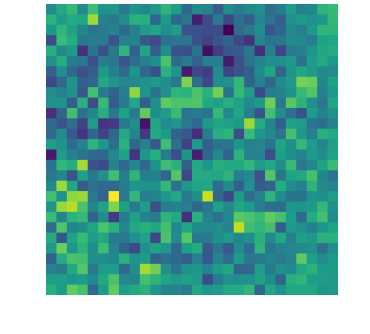
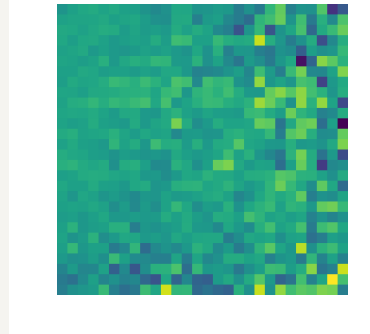
Different intensity

The range of weight values

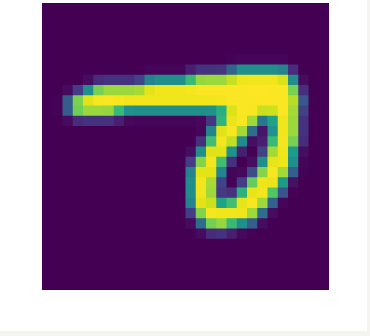
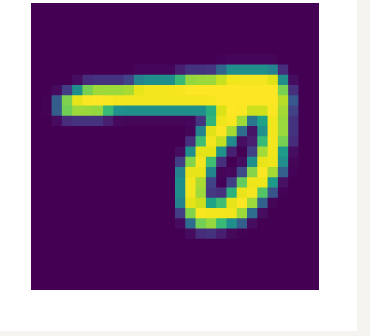
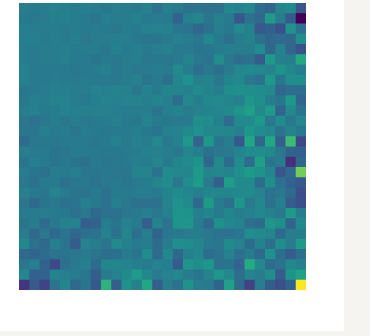
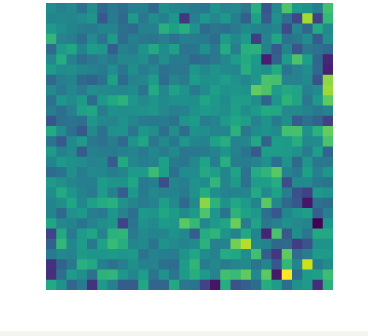
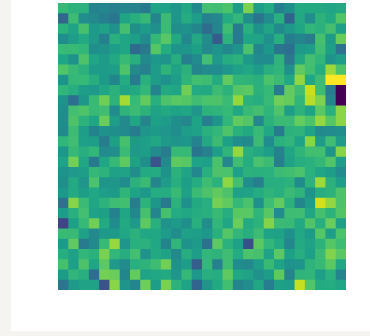
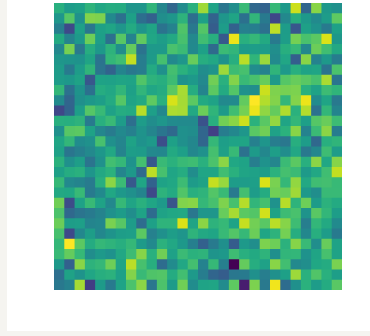
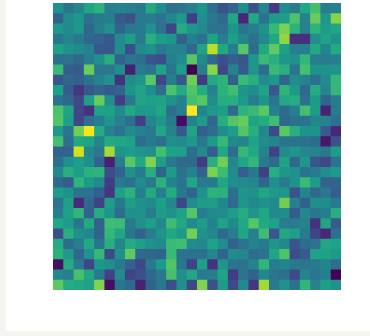



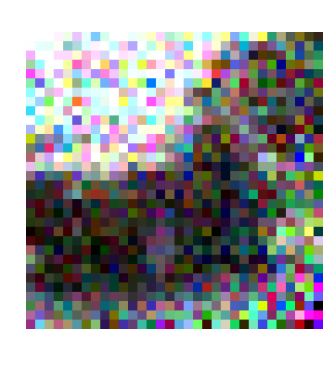
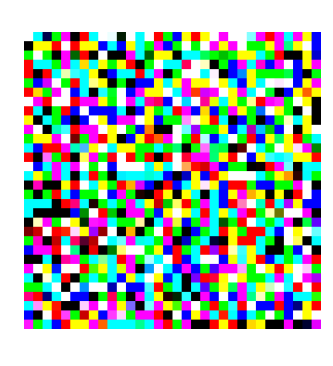
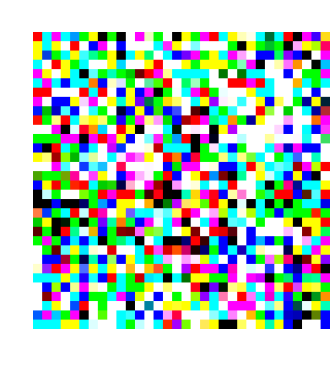
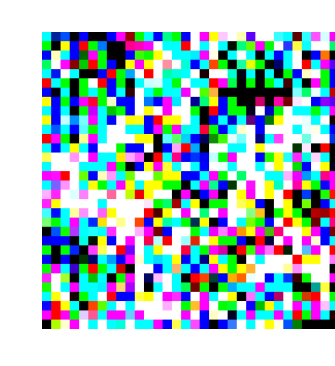
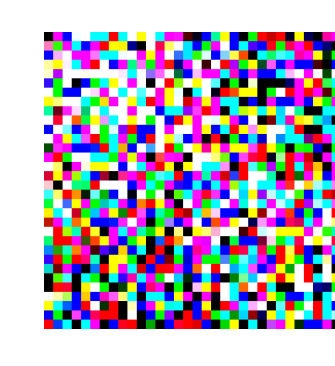
Sun et al., "Soteria: Provable Defense against Privacy Leakage in Federated Learning from Representation Perspective," CVPR 2021.

Wang et al., "Protect Privacy from Gradient Leakage Attack in Federated Learning," INFOCOM 2022.



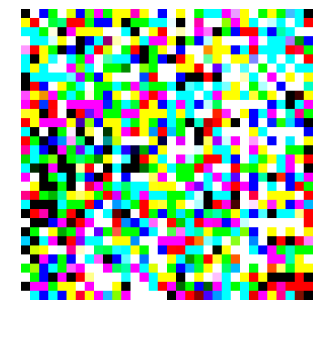


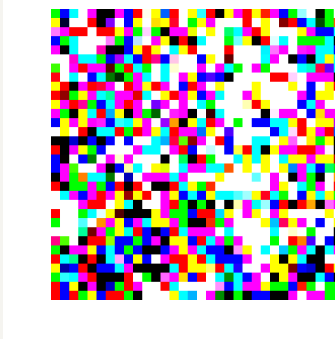
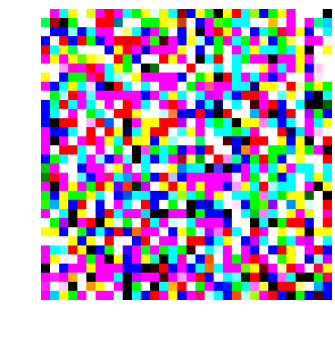
	No defense	GC	DP	Soteria	GD	OUTPOST
MSE \uparrow	$6.6e-3$	13.96	113.63	95.19	32.57	77.05
LPIPS \uparrow	$7.1e-2$	0.55	0.60	0.63	0.64	0.58
SSIM \downarrow	0.99	0.30	0.19	$3.3e-2$	$1.0e-2$	0.13
						

[Scenario 1] EMNIST: $E = 1, n = 1, B = 1$

	No defense	GC	DP	Soteria	GD	OUTPOST
MSE \uparrow	$2.6e-7$	199.08	297.84	296.76	360.98	294.678
LPIPS \uparrow	$5.8e-7$	0.60	0.66	0.63	0.64	0.68
SSIM \downarrow	1.00	0.32	$6.5e-2$	$4.1e-2$	$1.7e-2$	$1.6e-2$
						

	No defense	GC	DP	Soteria	GD	OUTPOST
MSE \uparrow	$5.1e-2$	7.83	34.08	25.91	11.46	13.10
LPIPS \uparrow	0.53	0.77	0.77	0.76	0.74	0.77
SSIM \downarrow	0.57	$4.4e-2$	$3.6e-2$	$5.5e-2$	$2.2e-2$	$2.1e-2$
						

[Scenario 3] CIFAR-10: $E = 1, n = 1, B = 1$

	No defense	GC	DP	Soteria	GD	OUTPOST
MSE \uparrow	$5.9e-5$	27.50	34.51	25.91	56.66	35.24
LPIPS \uparrow	$1.8e-3$	0.76	0.78	0.76	0.77	0.77
SSIM \downarrow	0.99	$2.6e-2$	$2.9e-2$	$5.5e-2$	$1.7e-2$	$3.4e-2$
						

Plato: A New Framework for
Scalable Federated Learning Research

<https://github.com/TL-System/plato>

A thorough investigation of
gradient leakage attacks in
production federated learning

A thorough investigation of
gradient leakage attacks in
production federated learning

Significantly weakened!

A thorough investigation of
gradient leakage attacks in
production federated learning

Significantly weakened!

Outpost: a defense mechanism

A thorough investigation of
gradient leakage attacks in
production federated learning

Significantly weakened!

Outpost: a defense mechanism
sufficient with minimal sacrifice!