

MLinPL
CONFERENCE 2025

Carnegie
Mellon
University

Auton
Lab

Semantic Label Reconstruction

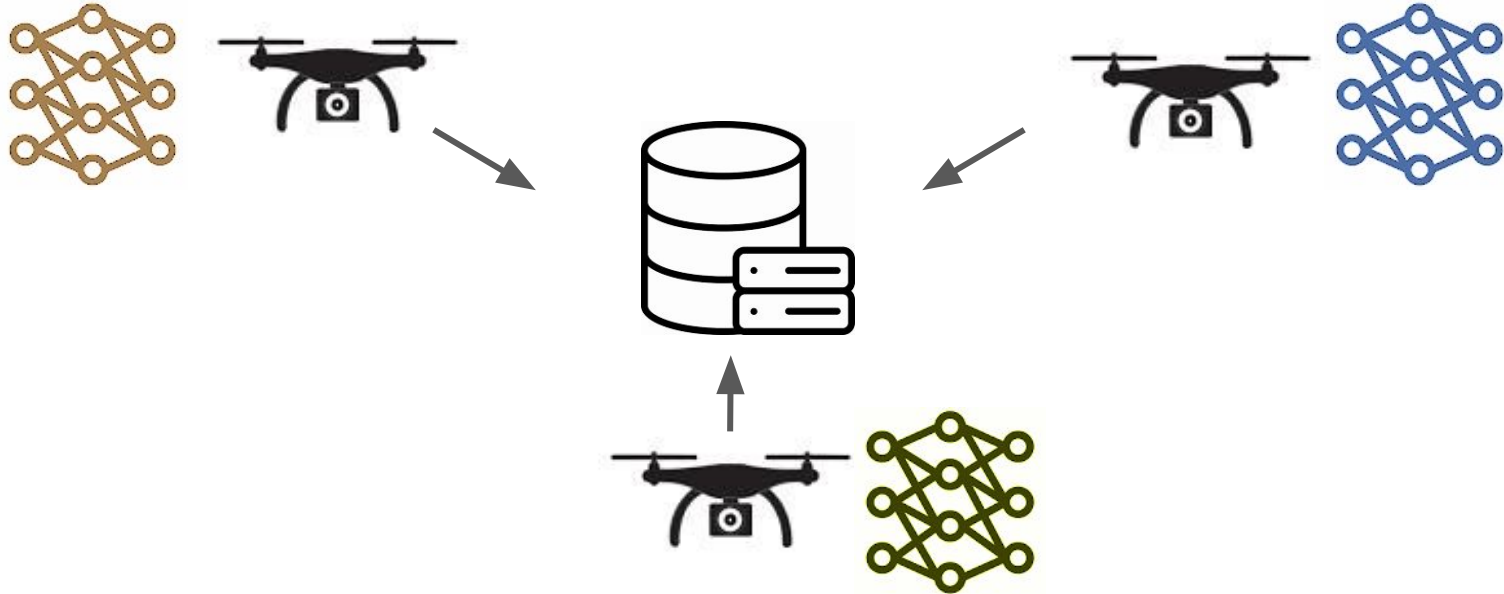
How to Breach Privacy in Federated Learning

Rafał Malcewicz^{1,2}, Ignacy Stępka¹, Abby Turner¹, Artur Dubrawski¹

¹Auton Lab, Carnegie Mellon University

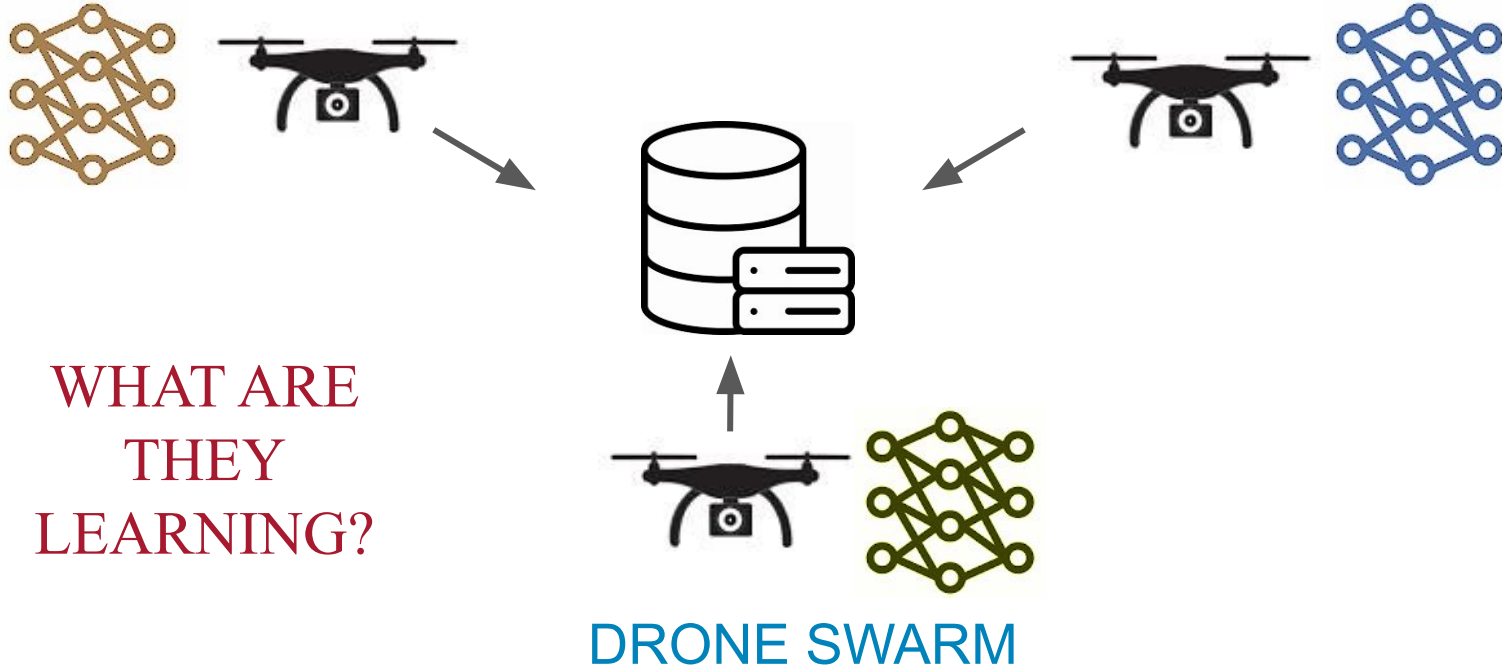
²Aalto University

Why This Matters

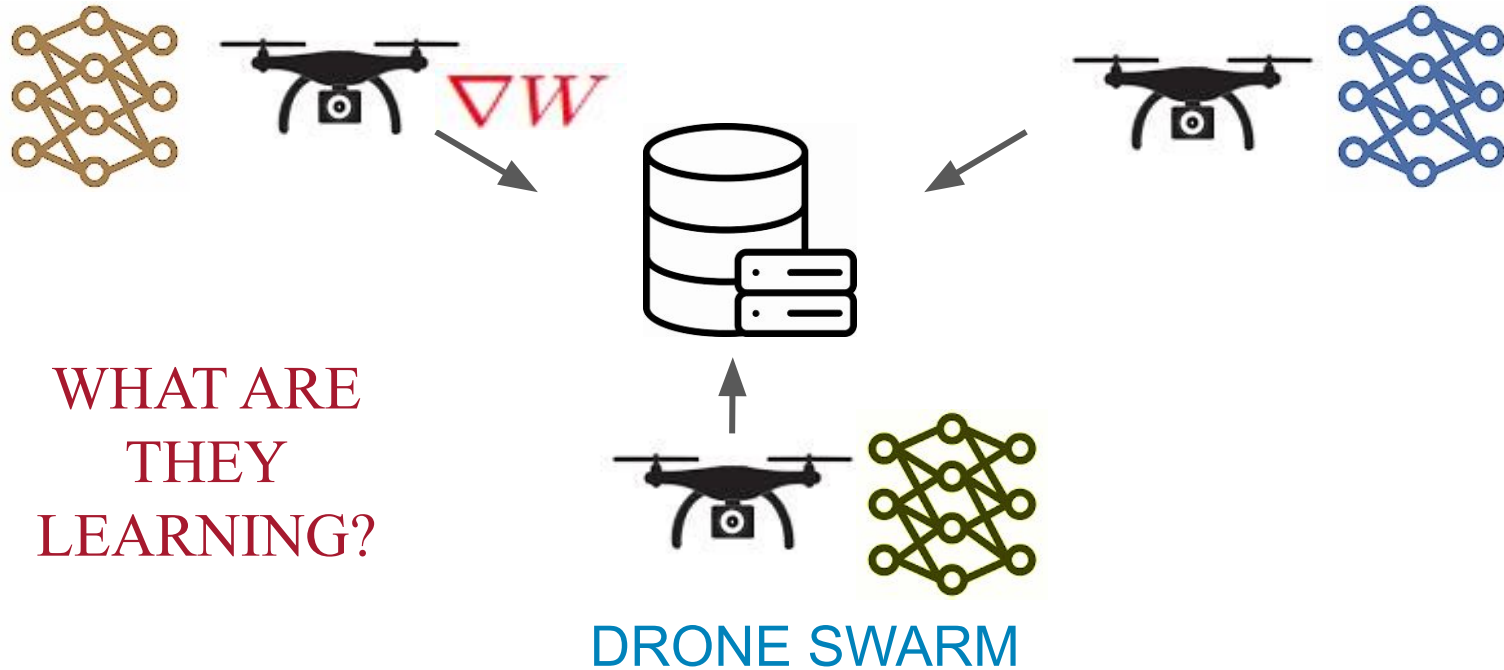


DRONE SWARM

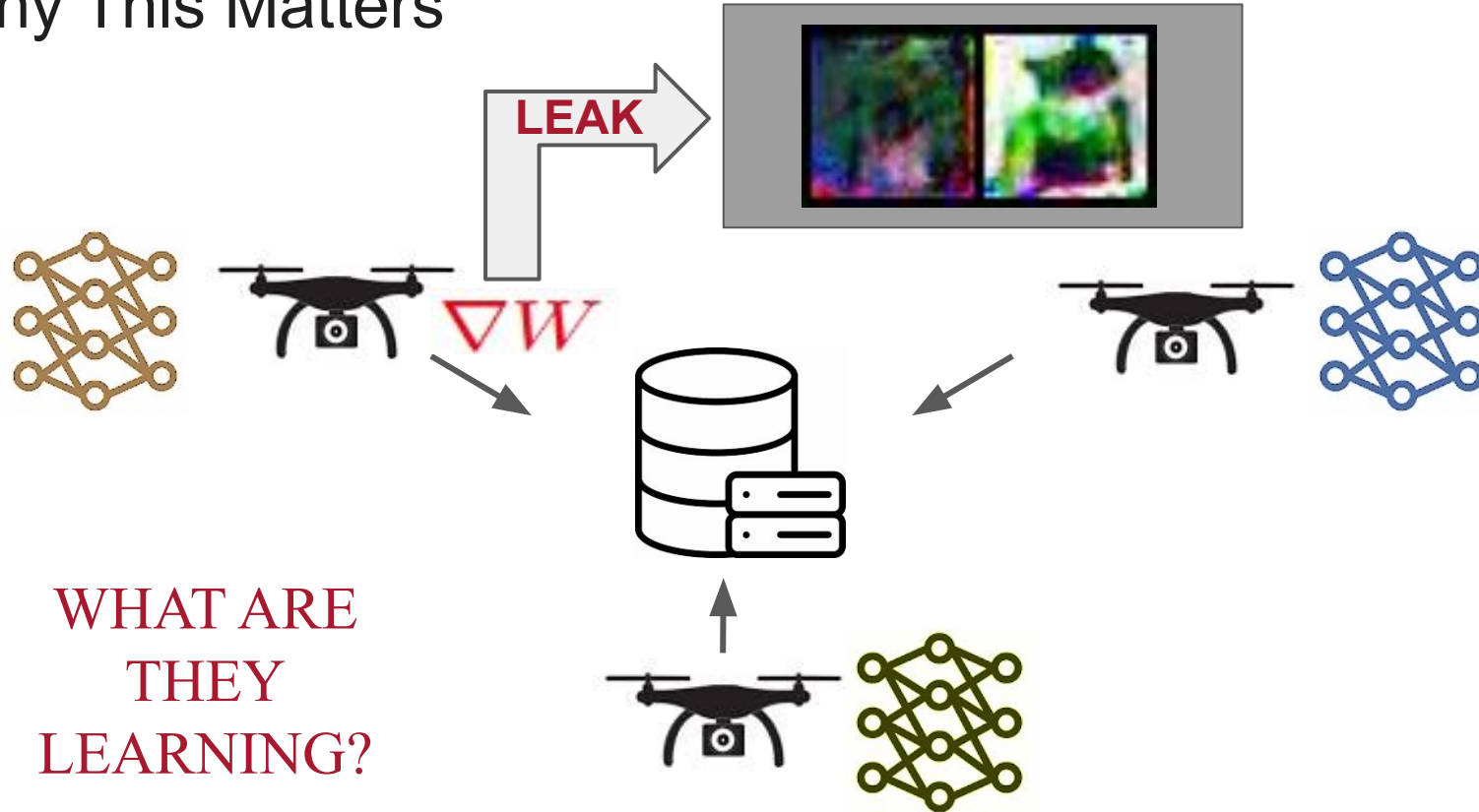
Why This Matters



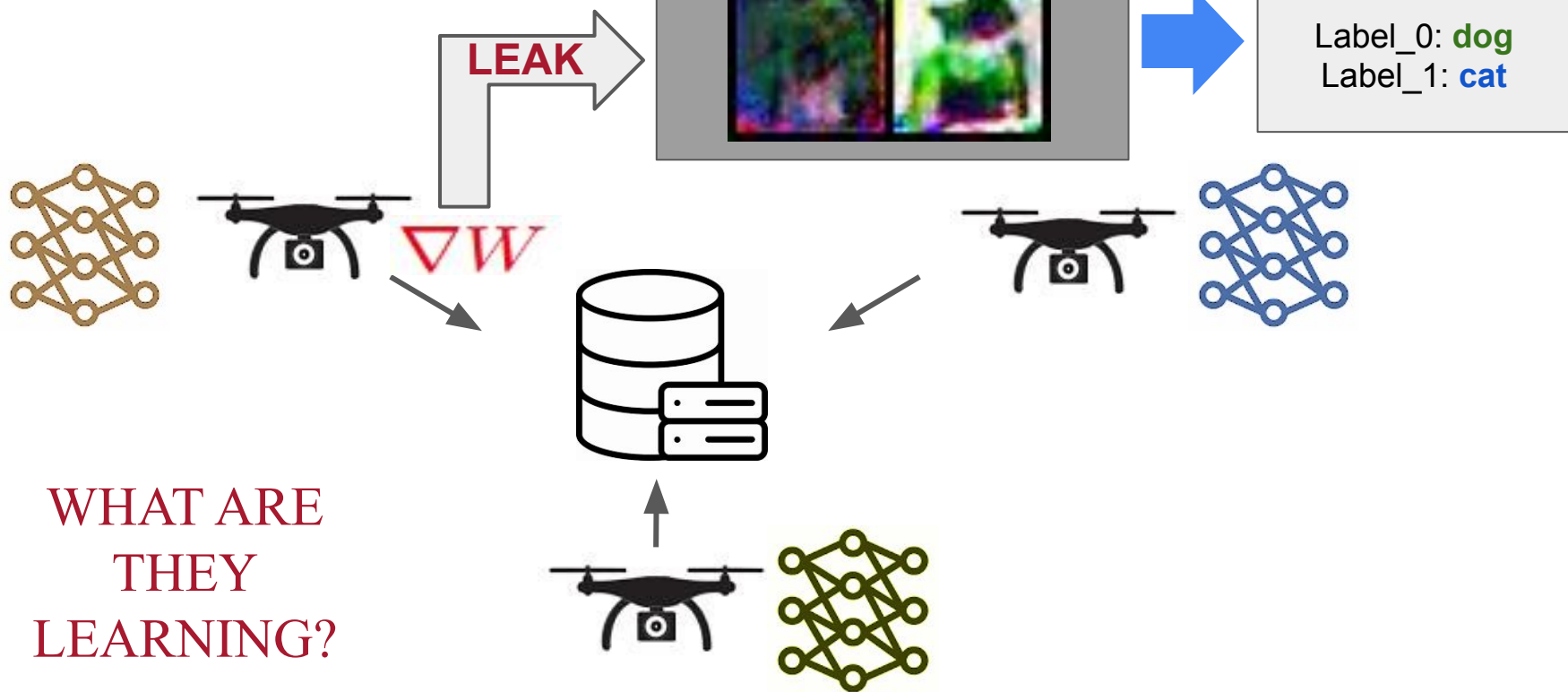
Why This Matters



Why This Matters



Why This Matters



DRONE SWARM

Plan for Today

Background

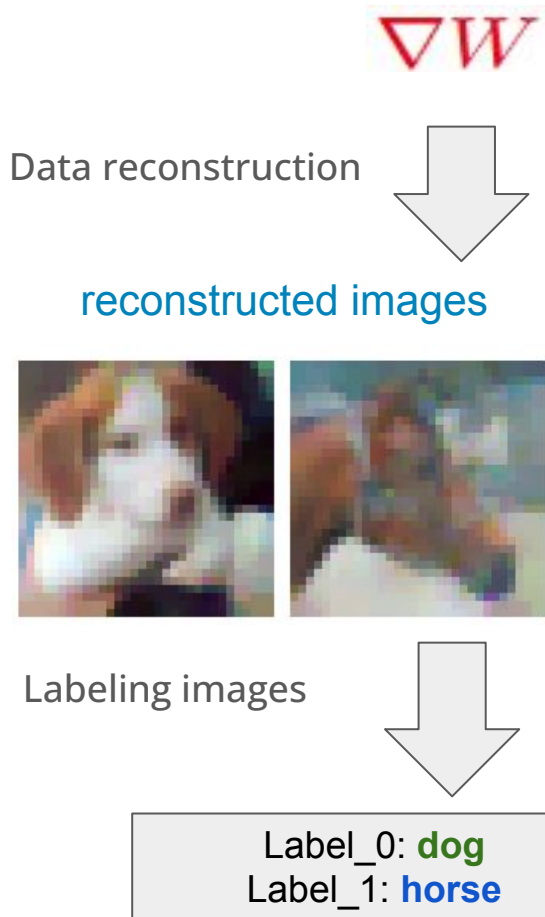
- What are Gradient Inversion Attacks (GIA's)?
- How do we evaluate the success of the reconstruction?

Semantic Label Reconstruction

- Label recovery with CLIP (Contrastive Language-Image Pre-training), Radford et al., 2021 [2]
- CLIP guided reconstruction

From Pixels to Meaning

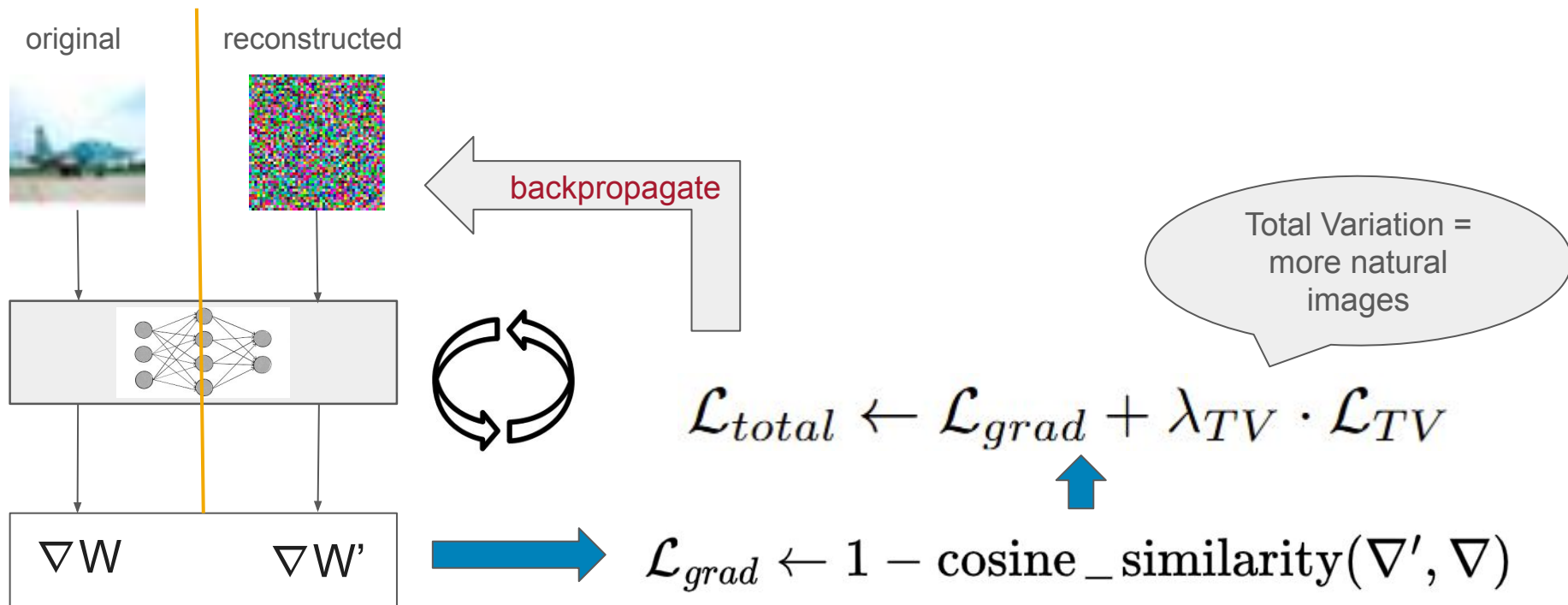
- Start with eavesdropped gradient
- Reconstruct images
- Label the the reconstructed image to retrieve semantic labels of the training data



From Pixels to Meaning



How Gradient Inversion Attack Work



Different Approaches to GIA's

Approach	Paper	Method / Key Idea
Gradient Matching Loss	Zhu et al., 2019 [1]	Used L2 distance
	Geiping et al., 2020 [2]	Used cosine similarity
Additional Loss Terms	Zhu et al., 2019 [1]	None
	Geiping et al., 2020 [2]	Added Total Variation (TV) loss
	Jeon et al., 2021 [3]	Added Batch Normalization statistics
Label Distribution Recovery	Zhao et al., 2020 [4]	Worked only for batch size = 1
	Ma et al., 2023 [5]	Solved system of linear equations
Latent Space Optimization	Fang et al., 2023 [6]	Used GAN to optimize latent space

Different Approaches to GIA's

Approach	Paper	Method / Key Idea
Gradient Matching Loss	Zhu et al., 2019 [1]	Used L2 distance
	Geiping et al., 2020 [2]	Used cosine similarity
Additional Loss Terms	Zhu et al., 2019 [1]	None
	Geiping et al., 2020 [2]	Added Total Variation (TV) loss
	Jeon et al., 2021 [3]	Added Batch Normalization statistics
Label Distribution Recovery	Zhao et al., 2020 [4]	Worked only for batch size = 1
	Ma et al., 2023 [5]	Solved system of linear equations
Latent Space Optimization	Fang et al., 2023 [6]	Used GAN to optimize latent space





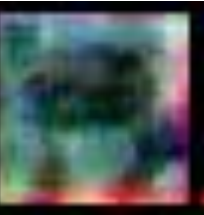


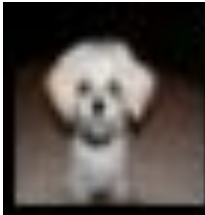


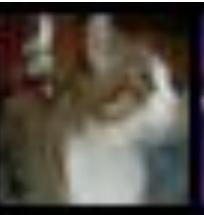

Different Approaches to GIA's

Approach	Paper	Method / Key Idea
Gradient Matching Loss	Zhu et al., 2019 [1]	Used L2 distance
	Geiping et al., 2020 [2]	Used cosine similarity
Additional Loss Terms	Zhu et al., 2019 [1]	None
	Geiping et al., 2020 [2]	Added Total Variation (TV) loss
	Jeon et al., 2021 [3]	Added Batch Normalization statistics
Label Distribution Recovery	Zhao et al., 2020 [4]	Worked only for batch size = 1
	Ma et al., 2023 [5]	Solved system of linear equations
Latent Space Optimization	Fang et al., 2023 [6]	Used GAN to optimize latent space

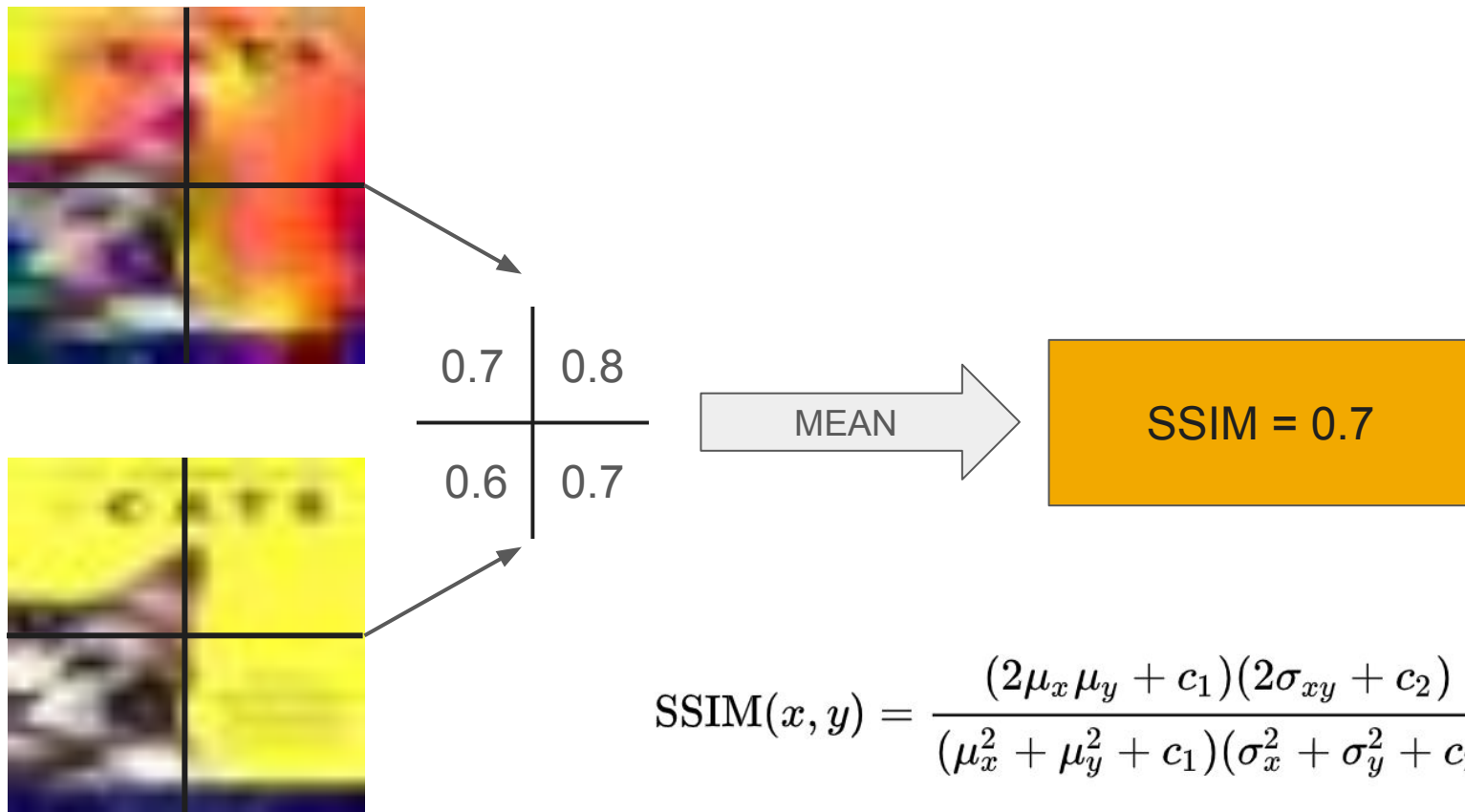
Different Approaches to GIA's

Approach	Paper	Method / Key Idea
Gradient Matching Loss	Zhu et al., 2019 [1]	Used L2 distance
	Geiping et al., 2020 [2]	Used cosine similarity
Additional Loss Terms	Zhu et al., 2019 [1]	None
	Geiping et al., 2020 [2]	Added Total Variation (TV) loss
	Jeon et al., 2021 [3]	Added Batch Normalization statistics
Label Distribution Recovery	Zhao et al., 2020 [4]	Worked only for batch size = 1
	Ma et al., 2023 [5]	Solved system of linear equations
Latent Space Optimization	Fang et al., 2023 [6]	Used GAN to optimize latent space

Results of the Reconstruction

	Batch Size		
	1	2	4
reconstructed		 	   
original		 	   

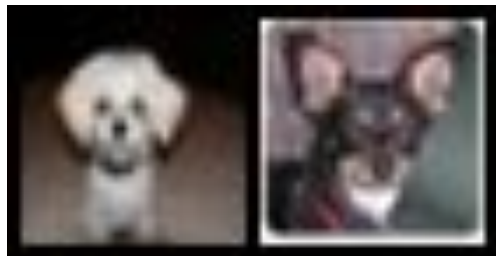
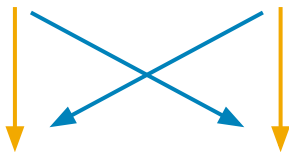
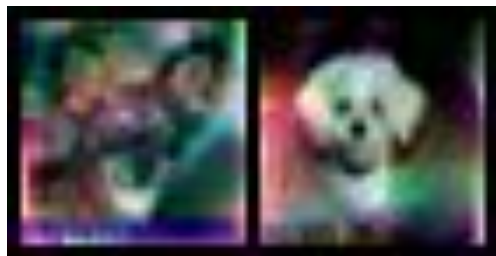
SSIM (Structural Similarity Index Measure)



SSIM (Structural Similarity Index Measure)



Permutation Agnostic SSIM



1st assignment

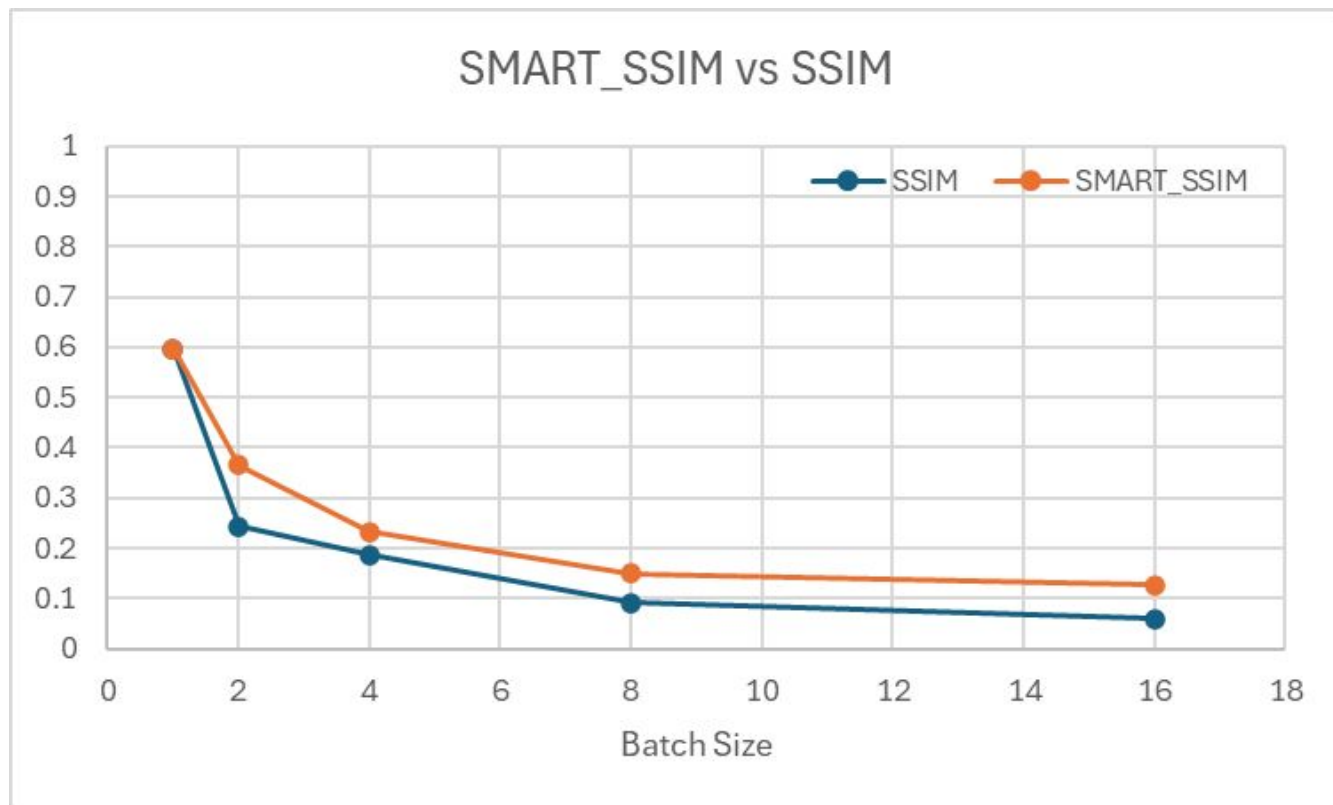
$$SSIM_1 = \frac{SSIM(\text{noisy dog}, \text{clear dog}) + SSIM(\text{clear dog}, \text{noisy dog})}{2}$$

2nd assignment

$$SSIM_2 = \frac{SSIM(\text{noisy dog}, \text{noisy dog}) + SSIM(\text{clear dog}, \text{clear dog})}{2}$$

$$SSIM_{\text{SMART}} = \max(SSIM_1, SSIM_2)$$

Permutation Agnostic SSIM

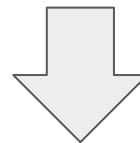


From Pixels to Meaning

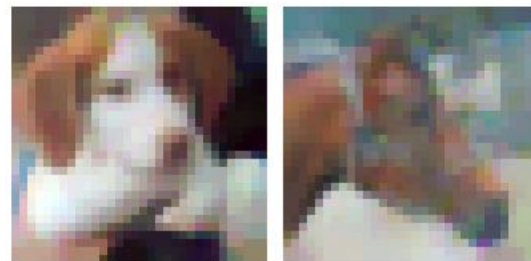


∇W

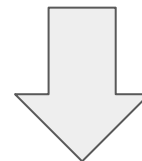
Data reconstruction



reconstructed images

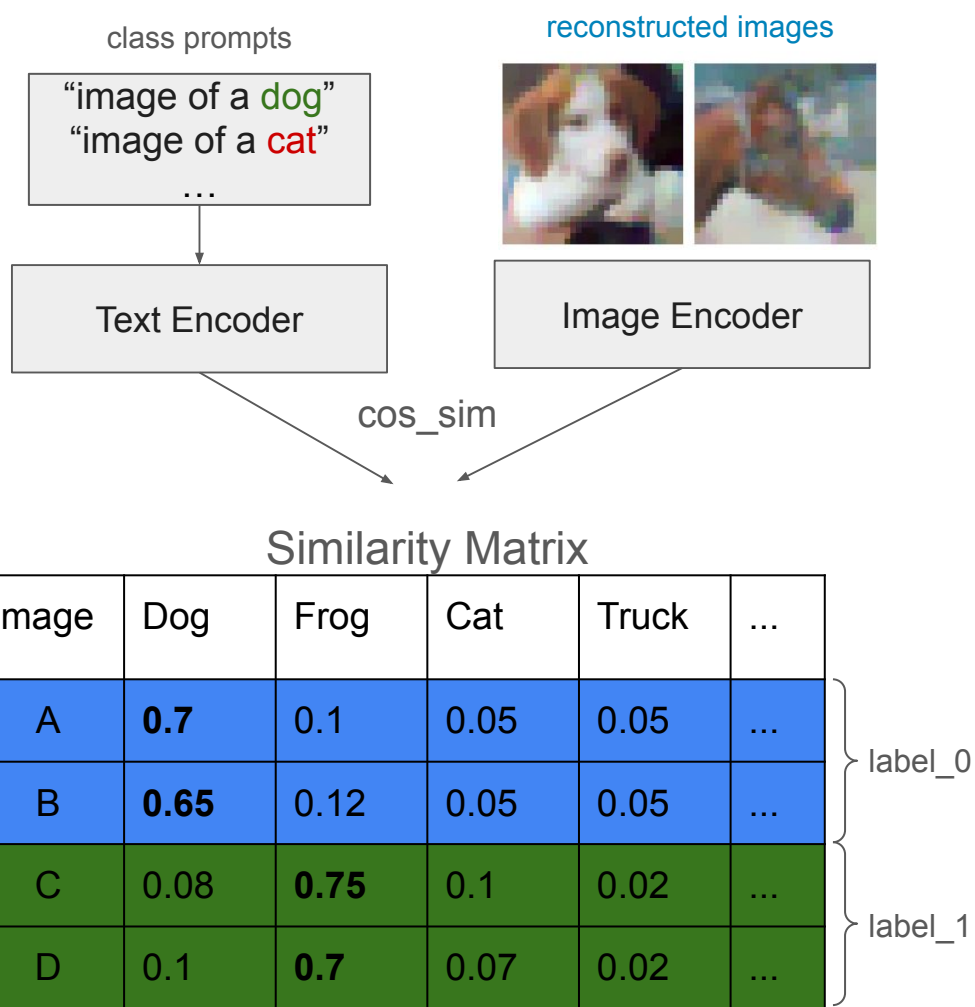


Labeling images



Label_0: **dog**
Label_1: **horse**

CLIP for Label Recovery

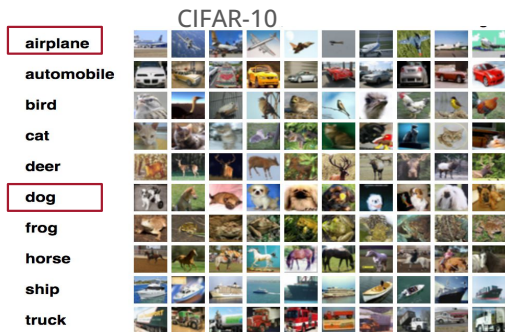


Label	Dog	Frog	Cat	Truck	...
0	0.675	0.11	0.05	0.05	...
1	0.09	0.725	0.085	0.02	...

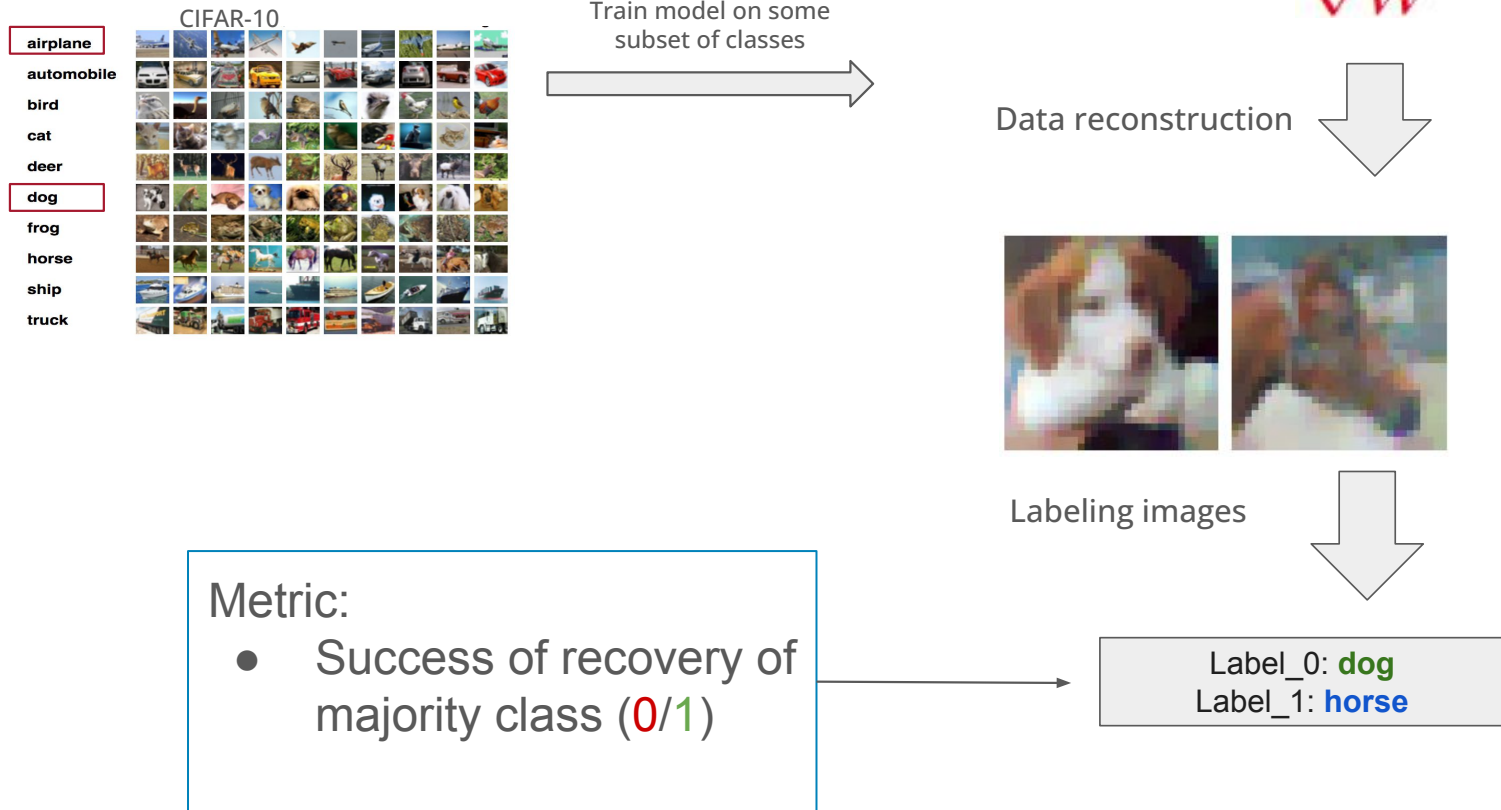
Setup & Evaluation



Setup & Evaluation

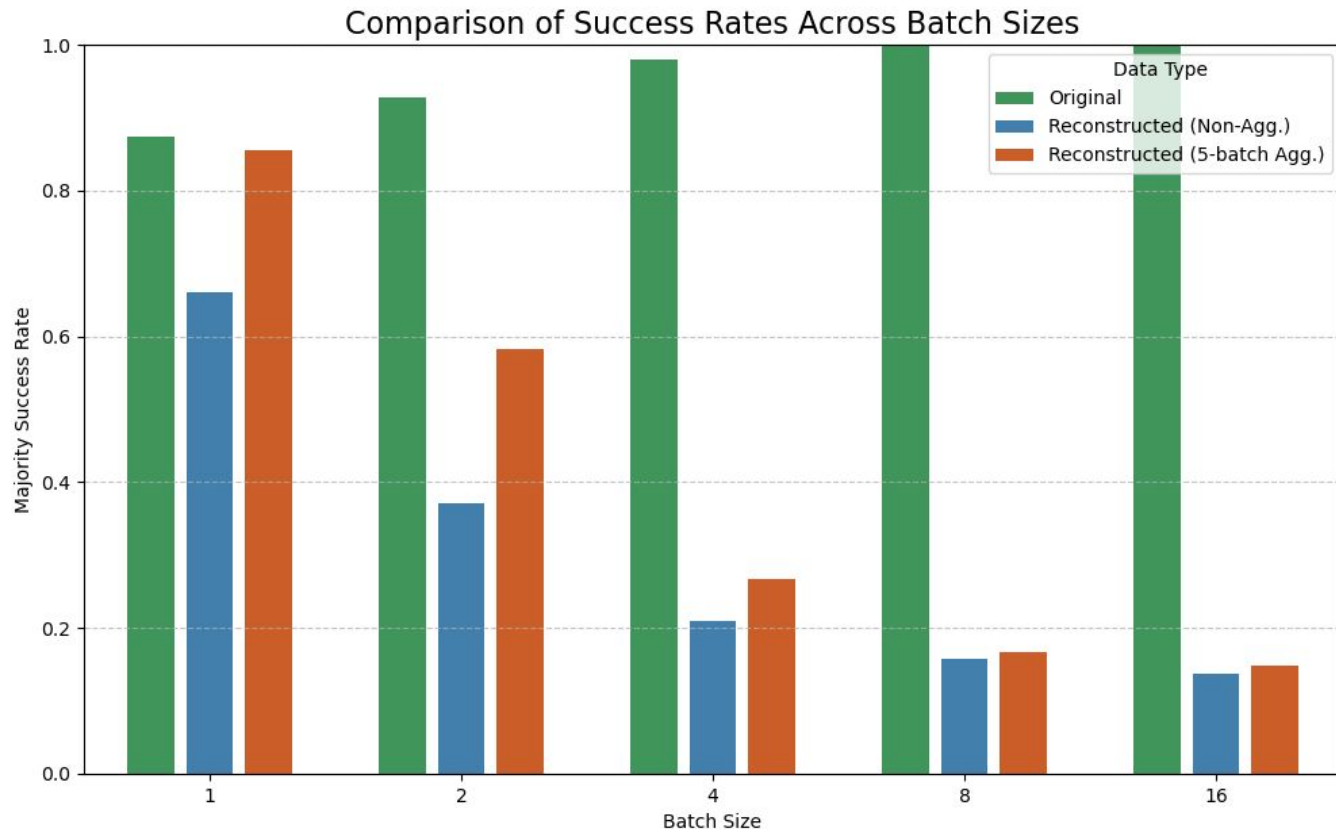


Setup & Evaluation

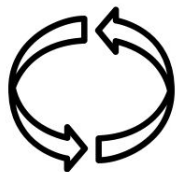
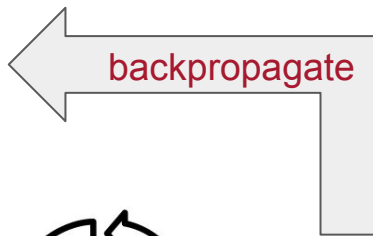
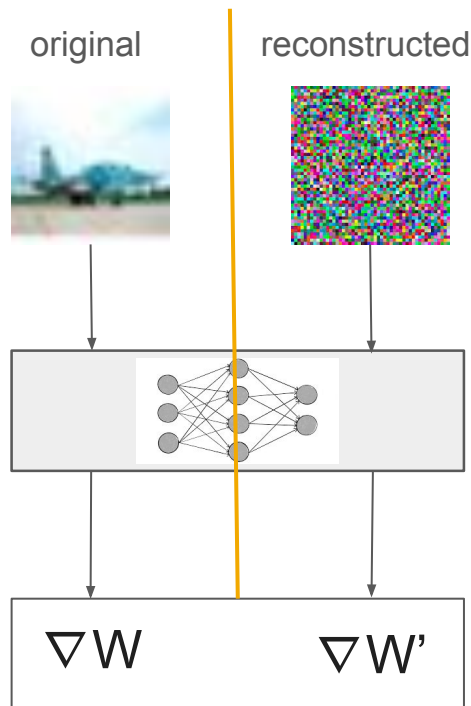


Leakage Drops as Batch Size Increases

- Original - almost perfect
- Aggregation improves quality
- Batch Size 4 still 2x better than random



CLIP Guided Reconstruction



$$\mathcal{L}_{total} \leftarrow \mathcal{L}_{grad} + \lambda_{TV} \cdot \mathcal{L}_{TV}$$



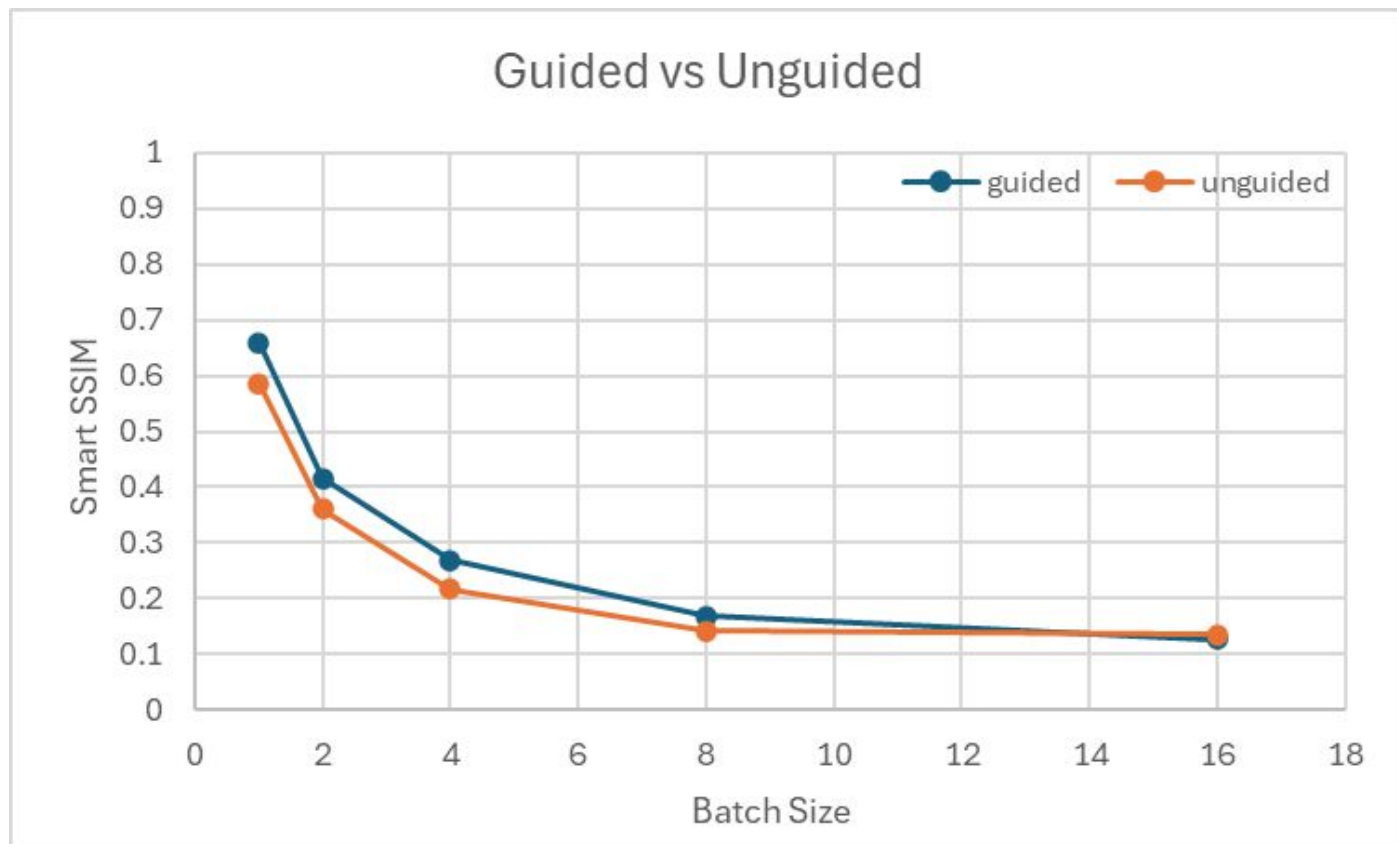
$$\mathcal{L}_{grad} \leftarrow 1 - \text{cosine_similarity}(\nabla', \nabla)$$



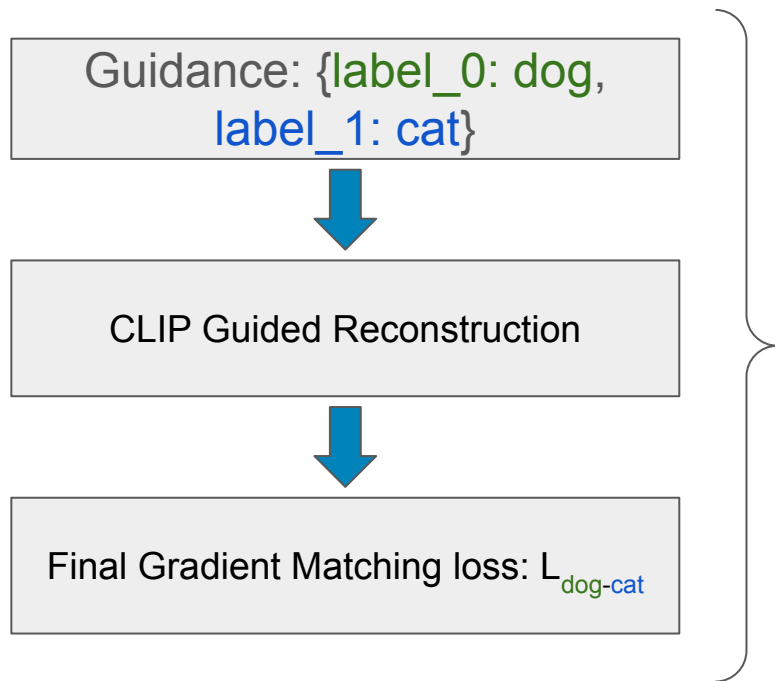
Guidance: "image of an **airplane**"

$$\mathcal{L}_{CLIP} = \frac{1}{N} \sum_{i=1}^N (1 - \text{cosine_similarity}(\text{img}^{(i)}, \text{guidance_prompt}^{(i)}))$$

Oracle Guidance Improves Quality



Recovering Semantic Labels from CLIP Guidance



Repeat for each possible guidance and report one that leads to lowest loss

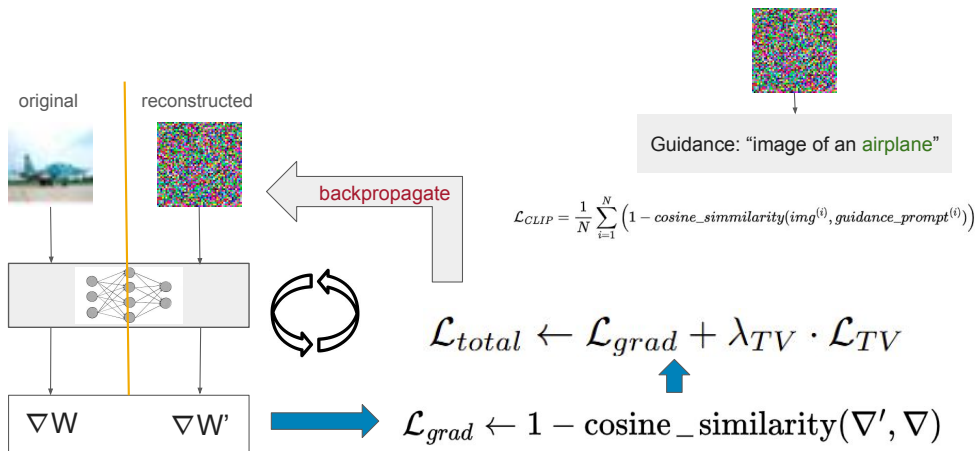
Thank you!

Semantic Label Reconstruction

How to Breach Privacy in Federated Learning

Summary

- Gradients leak semantic information
- Increasing batch size is an effective defensive tool
- Aggregation across multiple batches improves label recovery



References

- [1] Zhu, Ligeng, Zhijian Liu, and Song Han. “Deep Leakage from Gradients.” arXiv, December 19, 2019.
- [2] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Aspell, A., Mishkin, P., Clark, J., Krueger, G., & Sutskever, I. (2021). Learning transferable visual models from natural language supervision. arXiv preprint arXiv:2103.00020
- [3] Geiping, Jonas, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. “Inverting Gradients -- How Easy Is It to Break Privacy in Federated Learning?” arXiv, September 11, 2020.
- [4] Jeon, Jinwoo, jaechang Kim, Kangwook Lee, Sewoong Oh, and Jungseul Ok. “Gradient Inversion with Generative Image Prior.” In Advances in Neural Information Processing Systems, 34:29898–908. Curran Associates, Inc., 2021.
- [5] Zhao, Bo, Konda Reddy Mopuri, and Hakan Bilen. “iDLG: Improved Deep Leakage from Gradients.” arXiv, January 8, 2020.
- [6] Ma, K., Sun, Y., Cui, J., Li, D., Guan, Z., & Liu, J. (2023). Instance-wise batch label restoration via gradients in federated learning. In Proceedings of the International Conference on Learning Representations (ICLR 2023).
- [7] Fang, Hao, Bin Chen, Xuan Wang, Zhi Wang, and Shu-Tao Xia. “GIFD: A Generative Gradient Inversion Method with Feature Domain Optimization.” In Proceedings of the IEEE/CVF International Conference on Computer Vision, 4967–76, 2023.