

# Data Attribution for Text-to-Image Models by Unlearning Synthesized Images



Sheng-Yu Wang



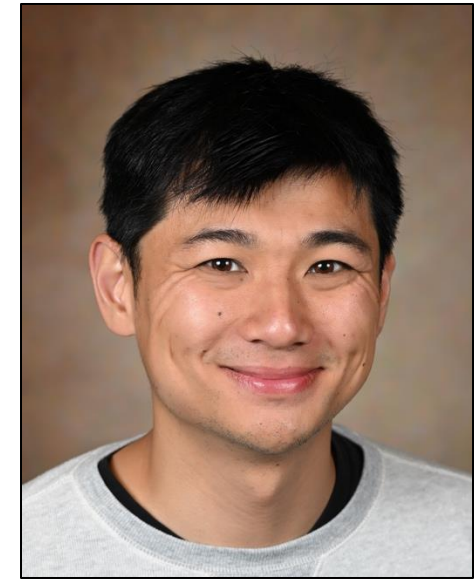
Aaron Hertzmann



Alexei A. Efros



Jun-Yan Zhu

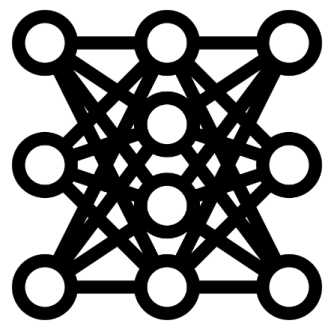


Richard Zhang



In NeurIPS, 2024.





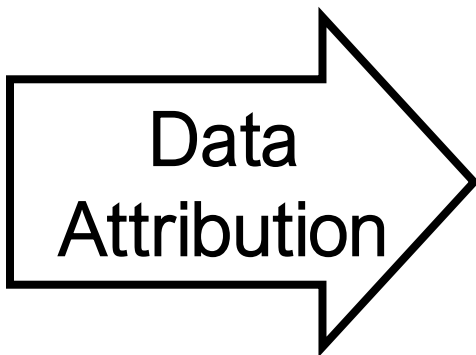
Stable Diffusion



Dataset



GenAI image

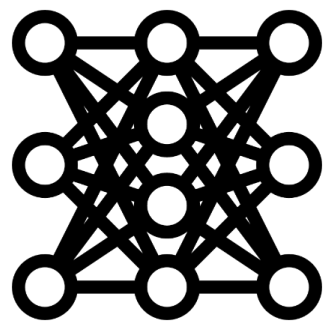


Influence scores

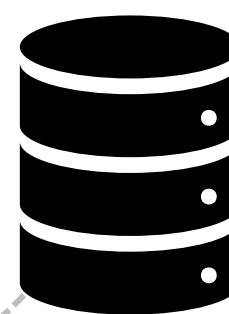
A grid of six small images of spiral staircases, arranged in two rows of three. The top row shows three different views of the staircase. The bottom row shows three more views, with the last one followed by three dots. The entire grid is set against a light grey background.

...

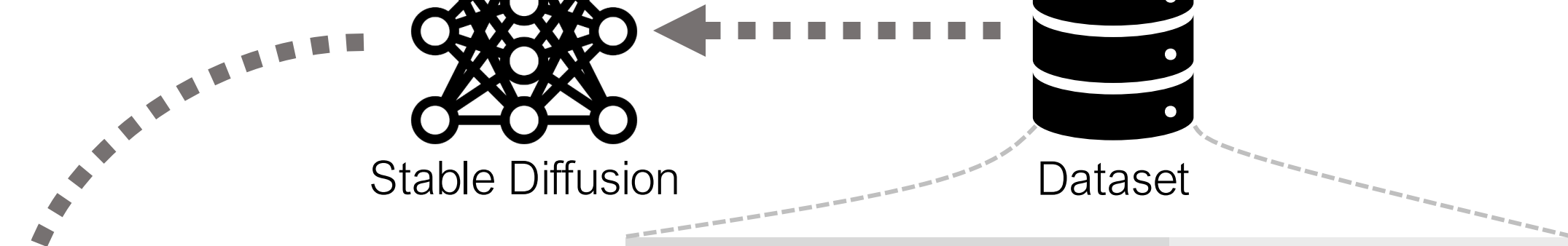




Stable Diffusion



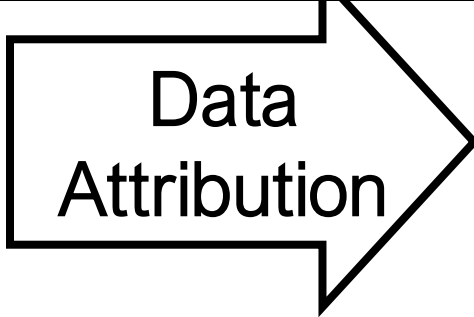
Dataset



Challenge: ground truth influence is unknown...  
Must intervene in the training process



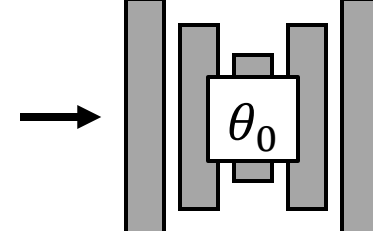
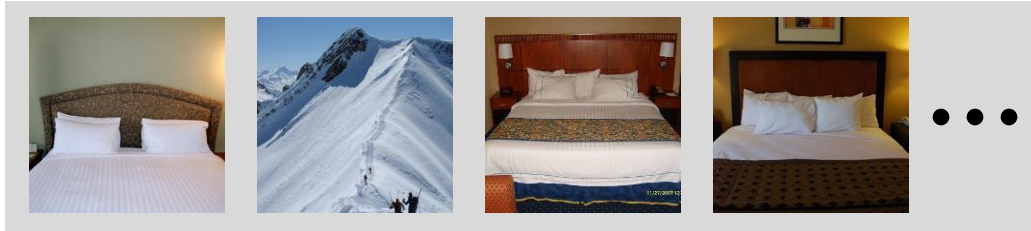
GenAI image



es

# Random subsets

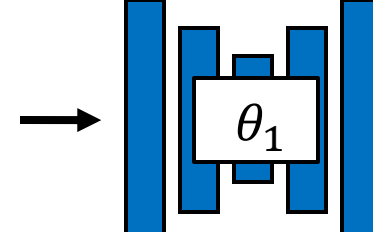
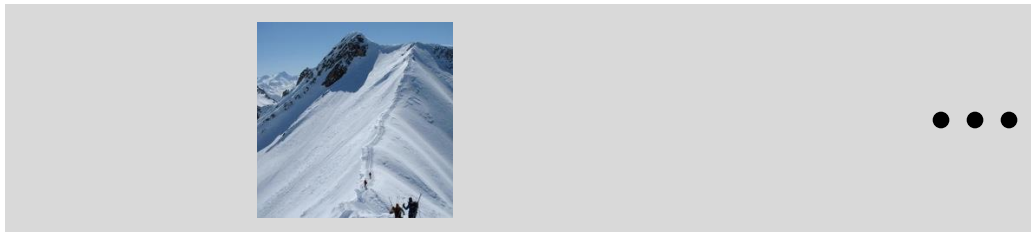
Training dataset



Synthesized

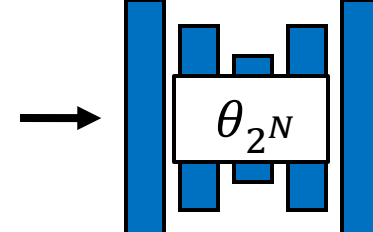
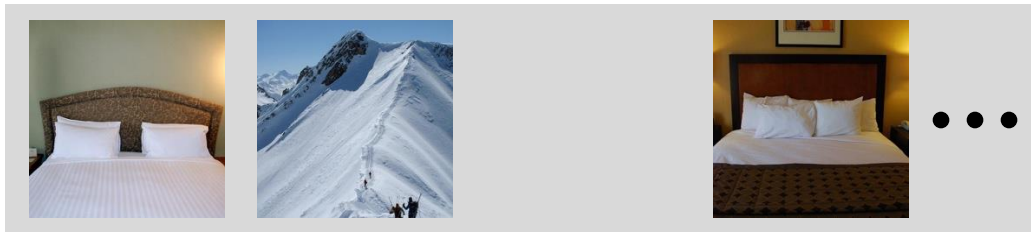


Counterfactual subset 1



...

Counterfactual subset  $2^N$

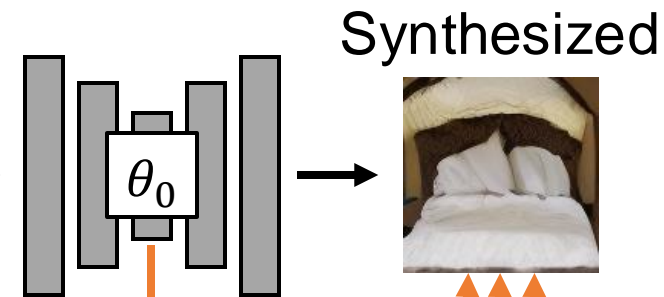
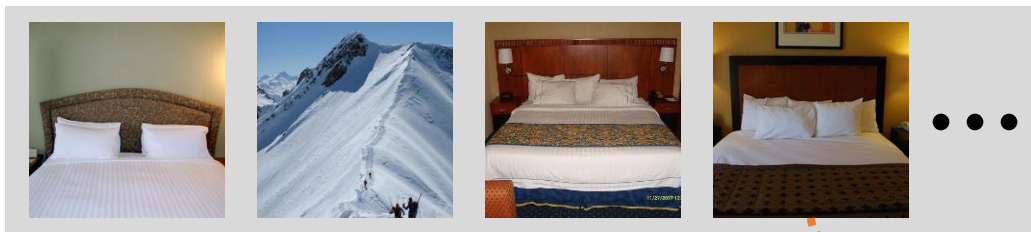


Analyze models

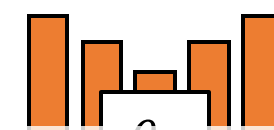
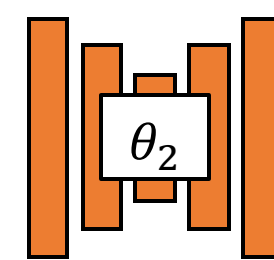
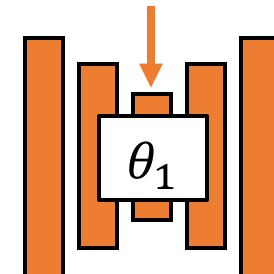
Training  $2^N$  models is too expensive

# Leave-one-out

Training dataset



Unlearning



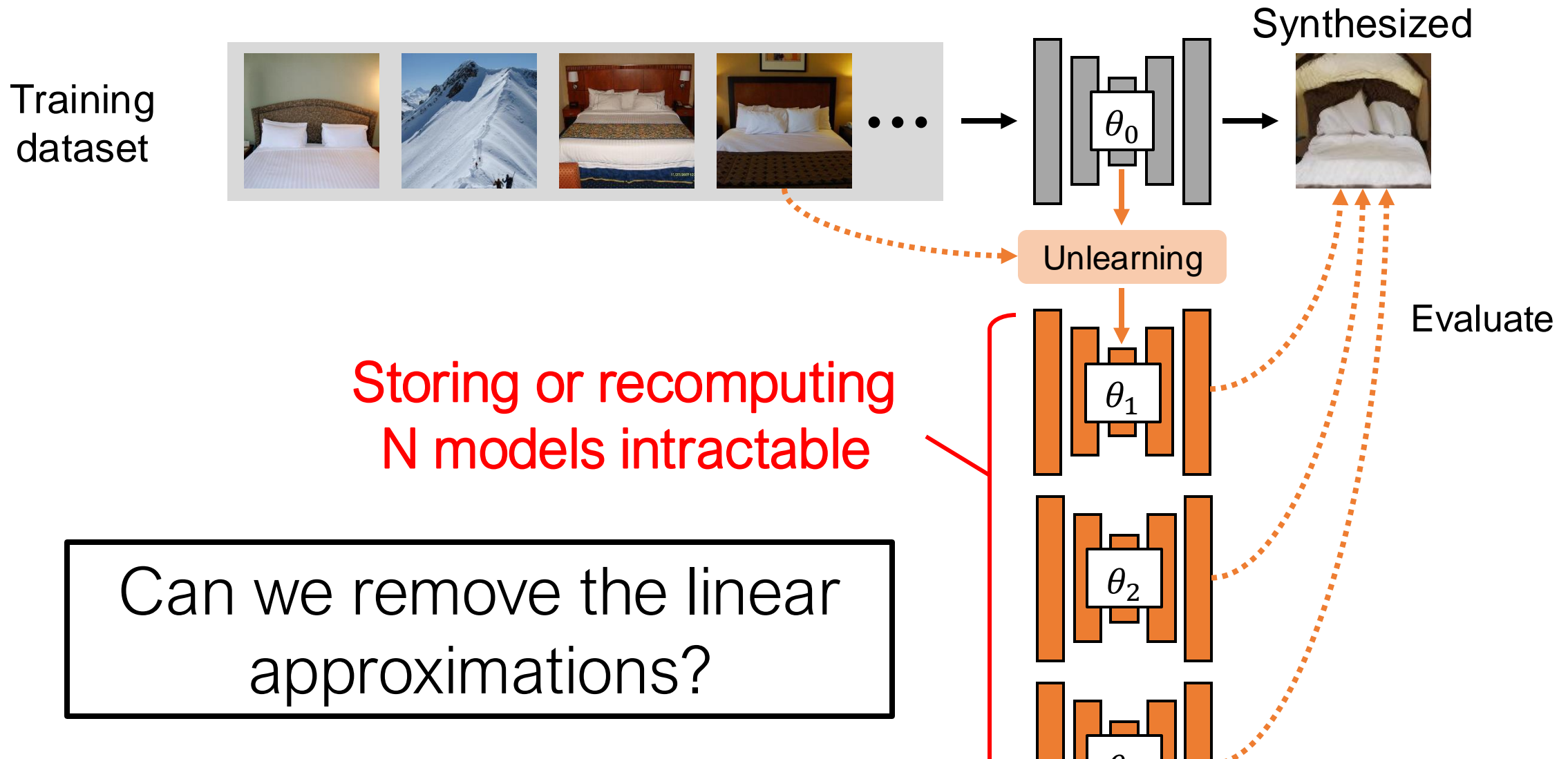
Evaluate

Influence functions: linear approx.  
for unlearning & evaluation

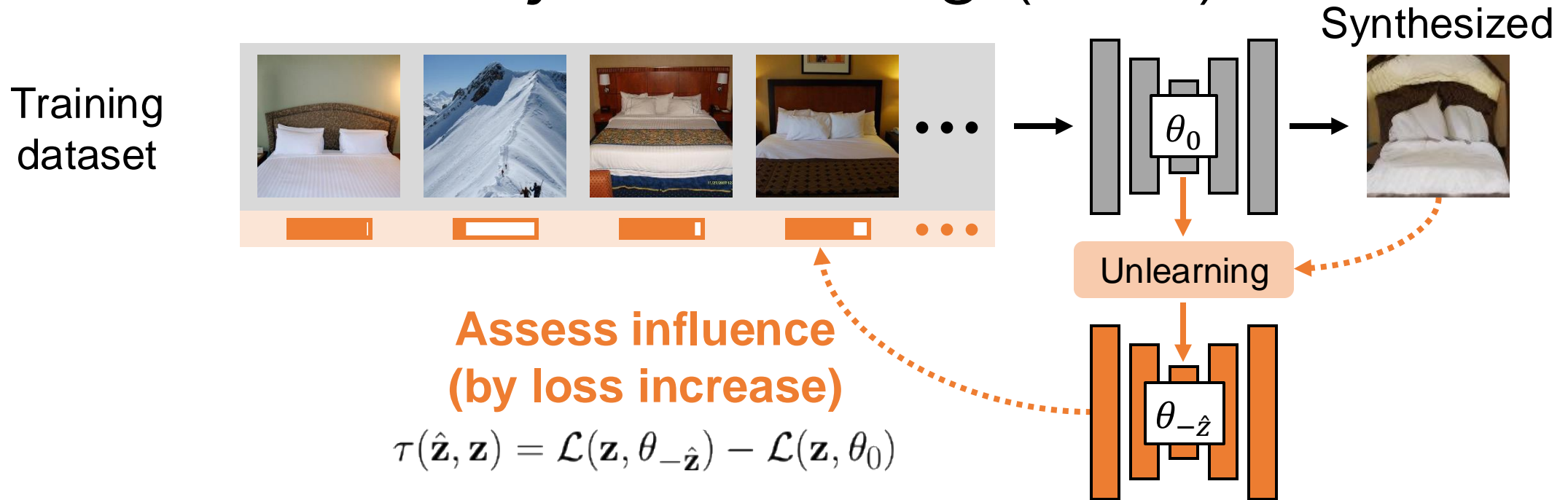
$$\nabla L(\hat{z})^T \underbrace{H_{\theta}^{-1}} \nabla L(x_n)$$

Store a low-dimensional version  
or recompute at test-time

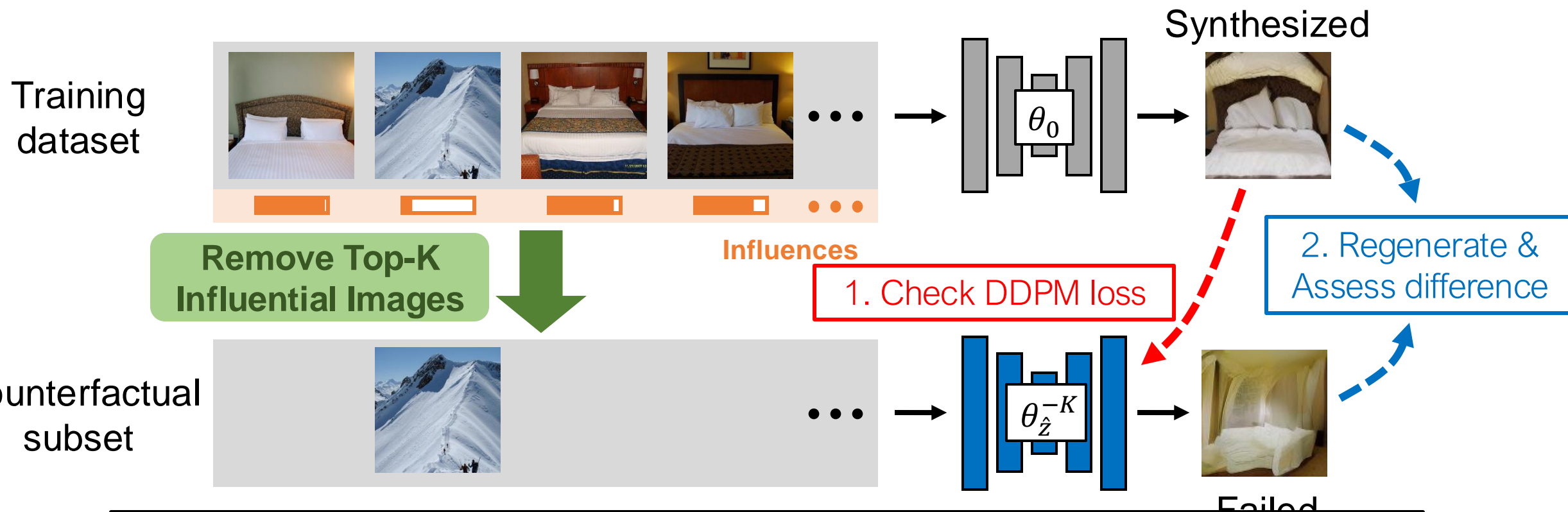
# Leave-one-out



# Attribution by Unlearning (AbU)



# Counterfactual evaluation



If critical training images are identified, removing them should destroy the generation



# MS-COCO results

Remove K=500  
(0.4% of dataset)

Effective removal



“A bus traveling on a freeway next to other traffic.”

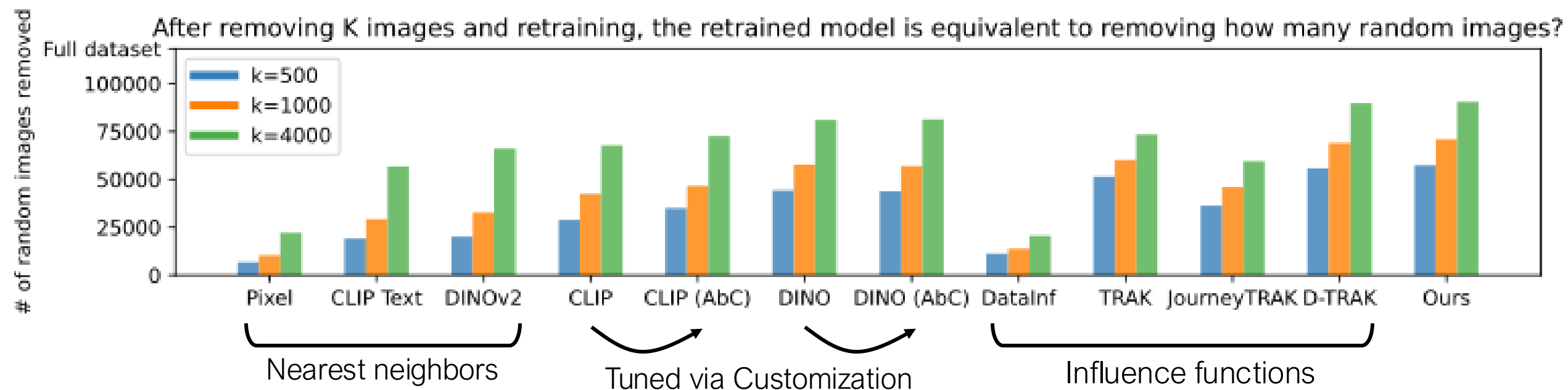
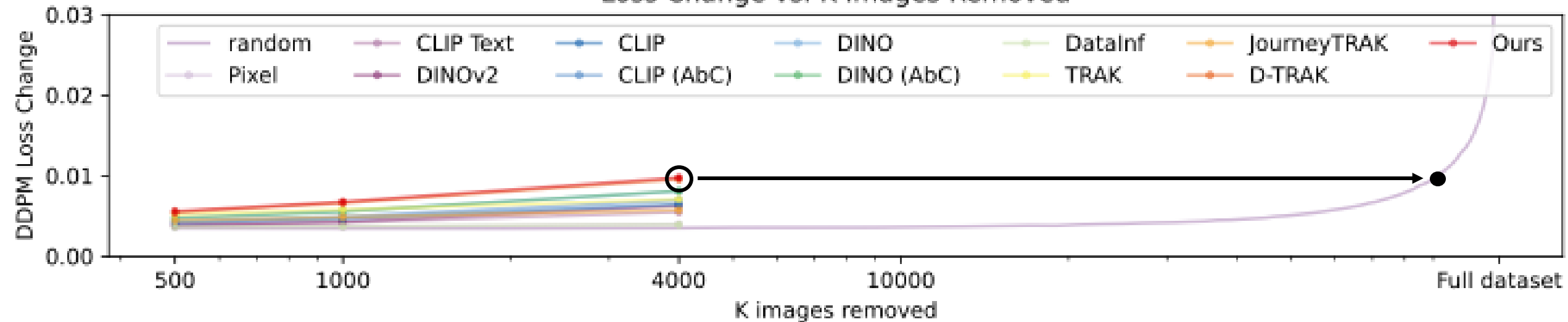


Attribution results

c.f. K. Georgiev, et al. How Training Data Guides Diffusion Models. In ArXiv, 2023.

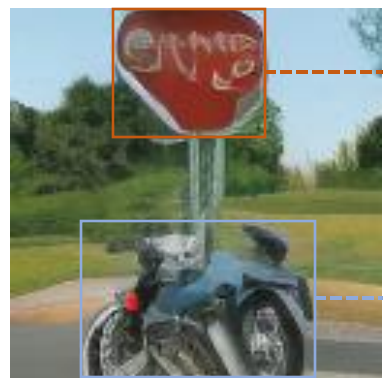
c.f. X. Zheng, et al. Intriguing Properties of Data Attribution on Diffusion Models. In ICLR, 2024.

### Loss Change vs. K Images Removed





# Local attribution



“A motorcycle and a stop sign.”

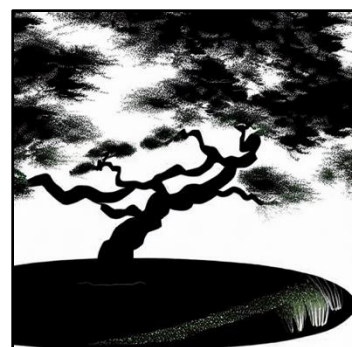
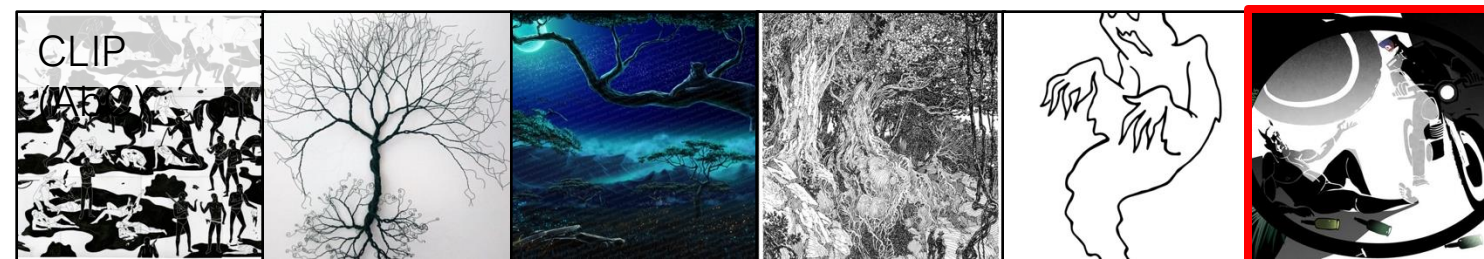
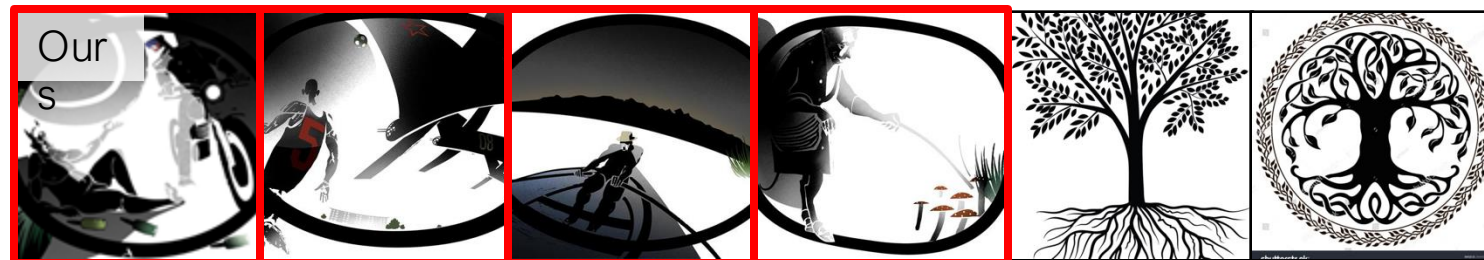


Cropped Queries



Attributed training images

# Customized Model Benchmark



“A picture of tree in the style of V\* art”



Thank you!

<https://peterwang512.github.io/AttributeByUnlearning>