

Evaluating Text-to-Image Models

Shobhita Sundaram

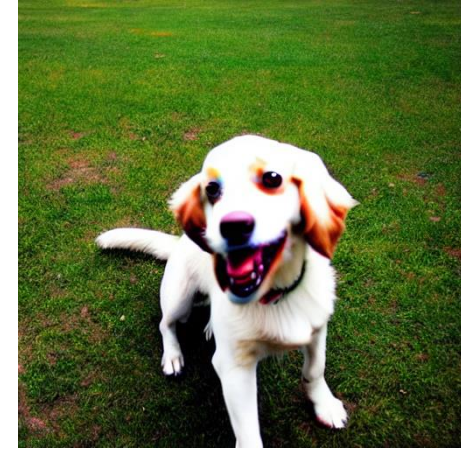
"Generate a photo of a dog playing outside"



- ✓ Shows a dog
- ✗ Not a photo



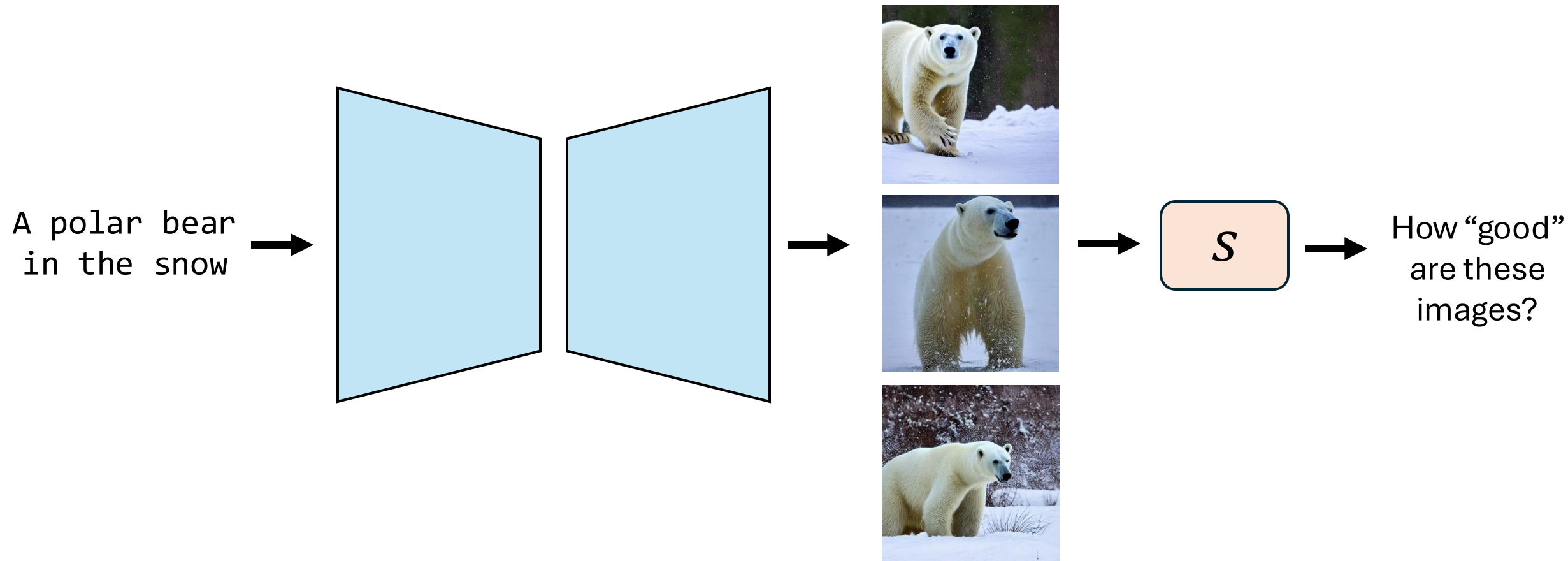
- ✓ Dog is playing
- ✓ Aesthetically pleasing



- ✓ Is a photo
- ✗ Strange lighting/artifacts

How do we evaluate generative models and their outputs?

Evaluating T2I models



Agenda

- What are the current image evaluation metrics?
- What are the best/most popular metrics for T2I models?
- How do you design a good evaluation metric that reflects human preferences?

Agenda

- **What are the current image evaluation metrics?**
- What are the best/most popular metrics for T2I models?
- How do you design a good evaluation metric that reflects human preferences?

What are the tools for image evaluation?

	Low-Level	High-Level
Unary/Holistic $s(x)$	Blurriness, No-Reference IQA	PickScore, ImageReward
Image Similarity $s(x, x_{ref})$	PSNR, SSIM, LPIPS, DISTS	DreamSim
Distribution $s(p(x)); s(p(x), p_{ref})$	InceptionScore, FID, CMMD	
Cross-Modal Similarity $s(x, y_{ref})$	SOA, CLIPScore	

	Low-Level	High-Level
Unary/Holistic $s(x)$	Blurriness, No-Reference IQA	PickScore, ImageReward

“A cat on a propaganda poster”



“A demon exiting through a portal...”









Agenda

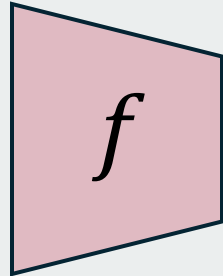
- What are the current image evaluation metrics?
- **What are the best/most popular metrics for T2I models?**
- How do you design a good evaluation metric that reflects human preferences?

	Low-Level	High-Level
Unary/Holistic $s(x)$	Blurriness, No-Reference IQA	PickScore, ImageReward
Similarity $s(x, x_{ref})$	PSNR, SSIM, LPIPS, DISTS	DreamSim
Distribution $s(p(x), p_{ref})$	FID, InceptionScore, CMMD	
Text-Alignment $s(x, y_{ref})$	SOA, CLIPScore	

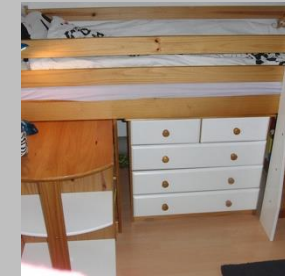
Why compare image distributions?

Caption	Generated Image	Real Image
A shoe rack with some shoes and a dog sleeping on them.		
Bunk bed with a narrow shelf sitting underneath it		
A table full of food such as peas and carrots, bread, salad and gravy		

How do we compare image distributions?

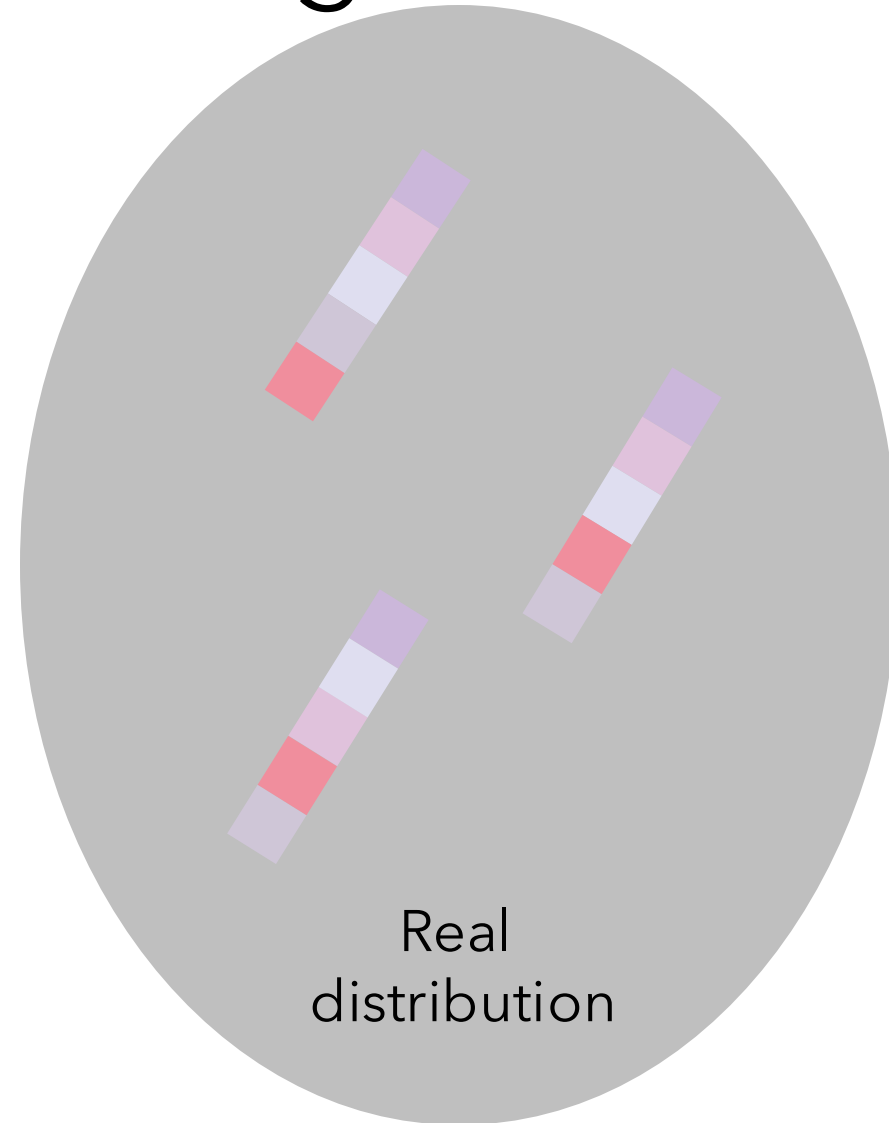


Generated
distribution



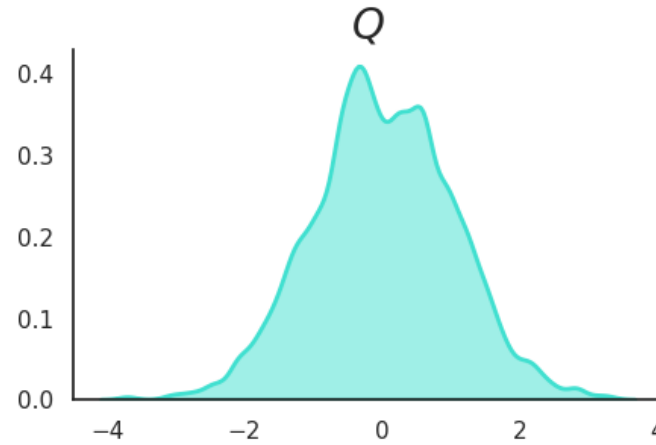
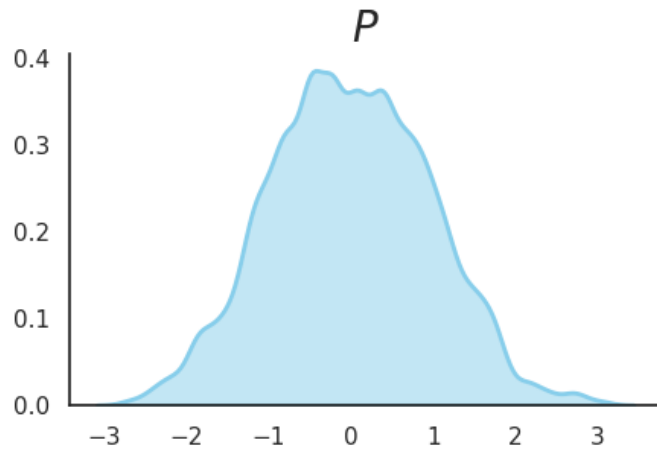
Real
distribution

How do we compare image distributions?



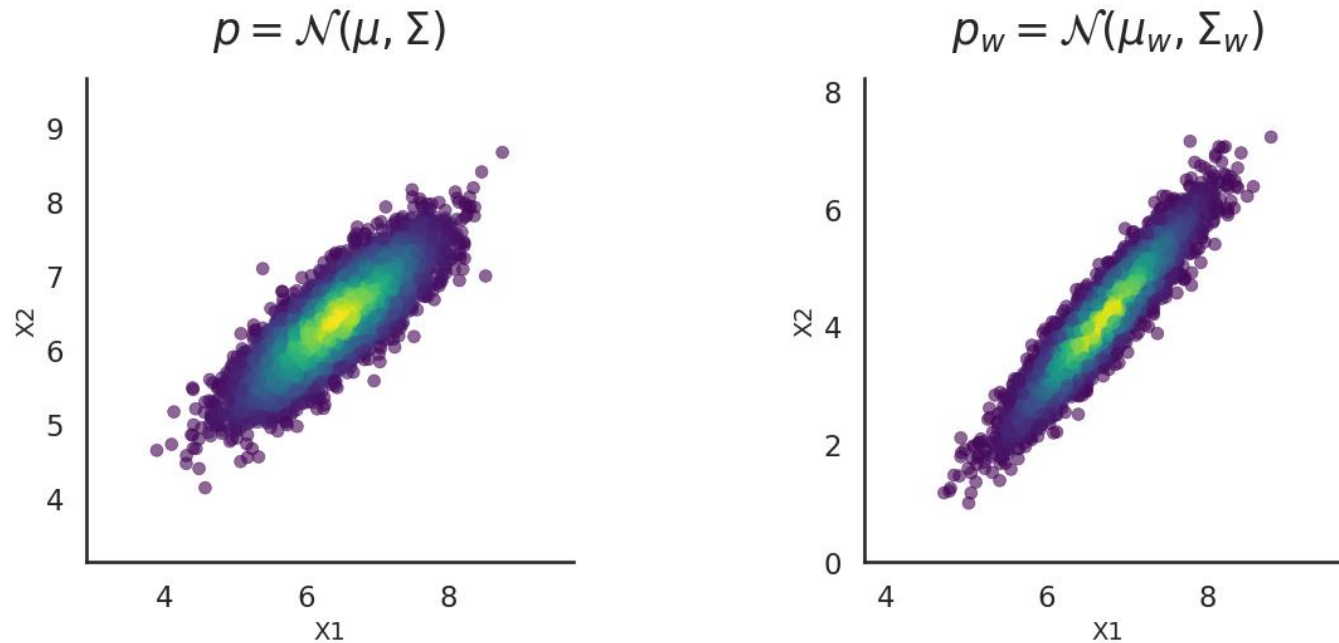
How do we compare image distributions?

Fréchet Distance (= Wasserstein-2 Distance)



$$W_2^2(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \mathbb{E}_{(X, Y) \sim \gamma} [\|X - Y\|_2^2]$$

How do we compare image distributions?



$$d_F(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\mu_w, \Sigma_w)) = \|\mu - \mu_w\|_2^2 + \text{Tr} \left(\Sigma + \Sigma_w - 2(\Sigma^{1/2} \Sigma_w \Sigma^{1/2})^{1/2} \right)$$

Fréchet Distance between Multivariate Gaussians

Fréchet Inception Distance (FID)

- Fréchet distance between Inception V3 embeddings of our real and generated images.
- **Advantages:**
 - Comparing images embedded in a meaningful representation space
 - Sensitive to both quality and diversity
 - Some GAN studies have shown correlation with human judgements^{1,2,3}

1. Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Proc. NIPS*, 2017.

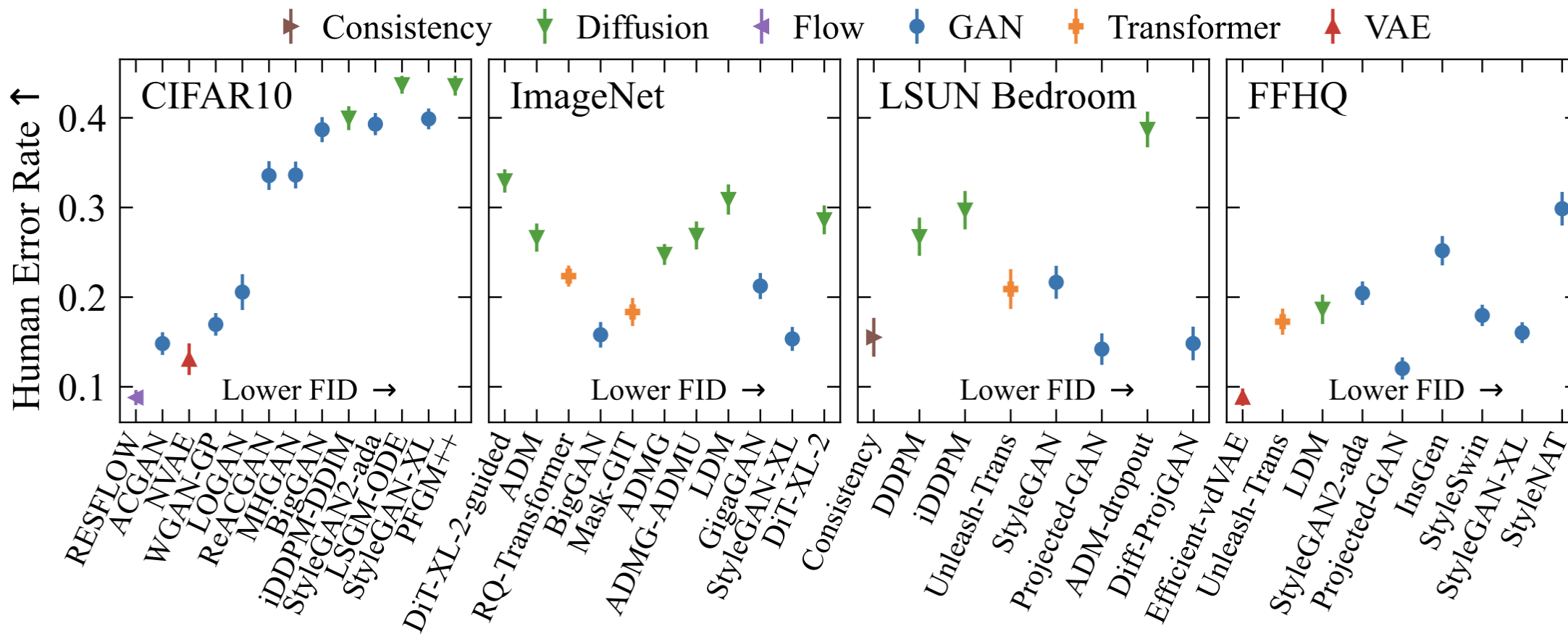
2. Weinberger. An Empirical Study on Evaluation Metrics of Generative Adversarial Networks

3. Mario Lucic, Karol Kurach, Marcin Michalski, S. Gelly, and O. Bousquet. Are GANs Created Equal? A Large-Scale Study. In *Proc. NeurIPS*, 2018.

Slide adapted from "Rethinking FID: Towards a Better Evaluation Metric for Image Generation", Sadeep Jayasuma

Fréchet Inception Distance (FID)

- Fréchet distance between Inception V3 embeddings of our real and generated images.
- **Disadvantages**
 - InceptionV3 only trained on ImageNet (~1M images)
 - Gaussian assumption (often untrue)
 - Need to estimate a large (2048x2048) covariance matrix
 - Biased estimator¹



CMMD

CLIP + Maximum Mean Discrepancy

- CLIP Embeddings
 - Trained on ~400M training images & complex scenes

CMMD

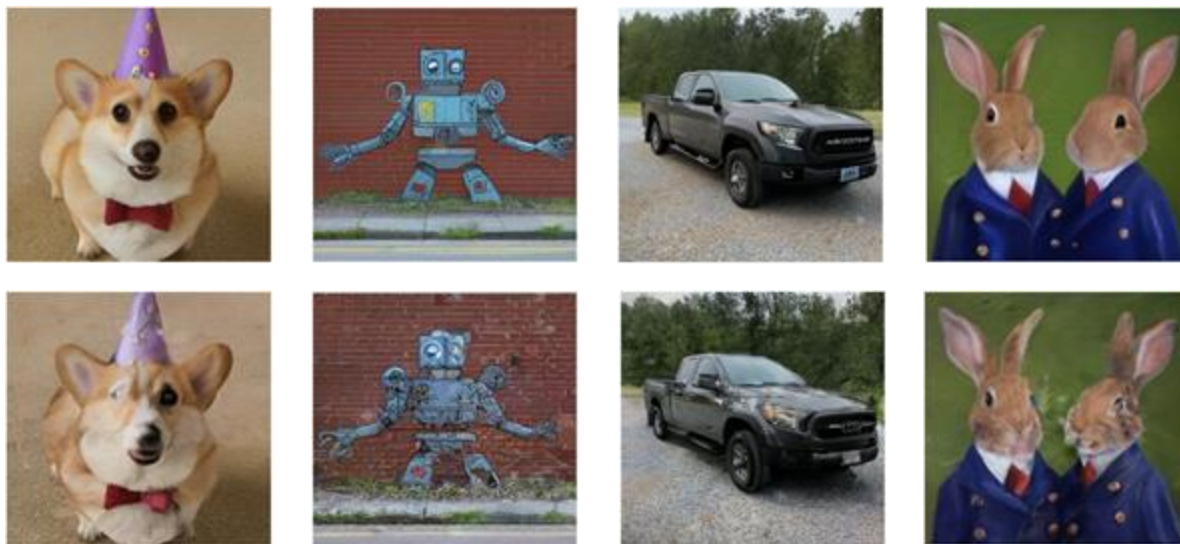
CLIP + Maximum Mean Discrepancy

- CLIP Embeddings
 - Trained on ~400M training images & complex scenes
- MMD Distance

$$\hat{\text{dist}}_{\text{MMD}}^2(X, Y) = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(\mathbf{x}_i, \mathbf{y}_j)$$

- No distributional assumptions
- Sample efficient
- Unbiased estimator

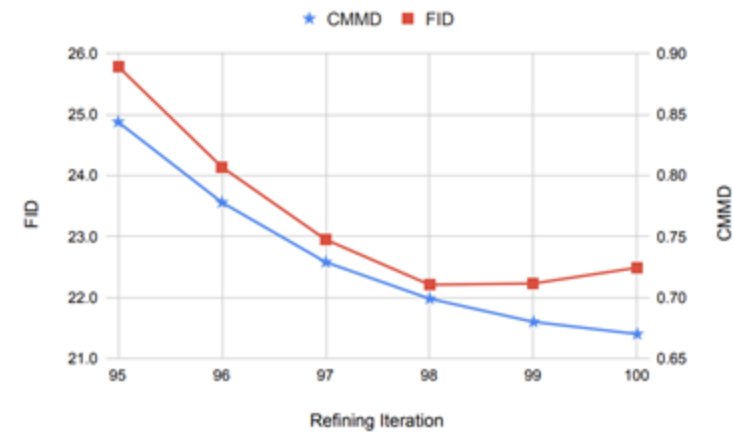
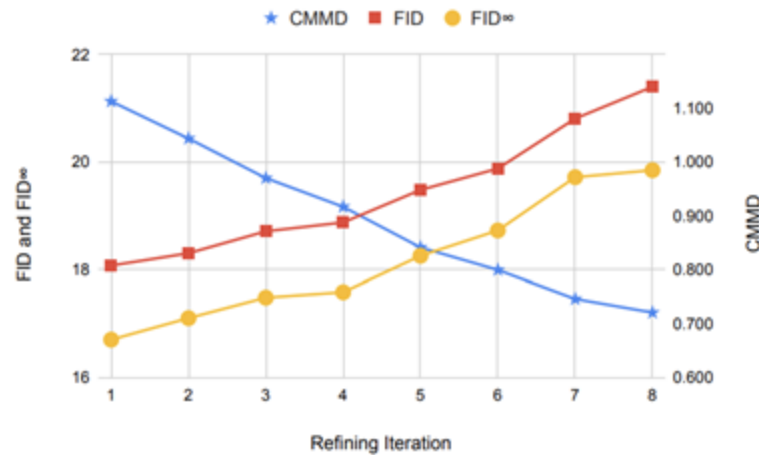
CMMD: Human Evaluation



Model	Model-A	Model-B
FID	21.40	18.42
FID _∞	20.16	17.19
KID	0.0105	0.0080
CMMD	0.721	0.951
Human rater preference	92.5%	6.9%

Table 3. *Human evaluation. FID and KID contradict human evaluation while CMMD agrees. Lower is better for all metrics.*

Measuring Model Improvements



Agenda

- What are the current image evaluation metrics?
- What are the best/most popular evaluation metrics for T2I models?
- **How do you design a good evaluation metric that reflects human preferences?**

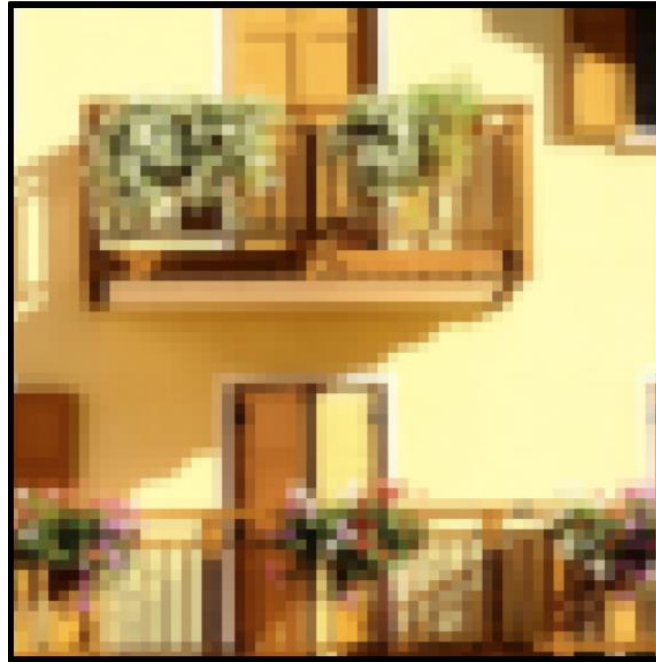
	Low-Level	High-Level
Holistic $s(x)$	Blurriness, No-Reference IQA	PickScore, ImageReward
Similarity $s(x, x_{ref})$	PSNR, SSIM, LPIPS, DISTS	DreamSim
Distribution $s(p(x)); s(p(x), p_{ref})$	FID, InceptionScore, CMMD	
Text-Alignment $s(x, y_{ref})$	SOA, CLIPScore	



Which patch is more similar to the middle?



< Clap >



Humans

L2/PSNR

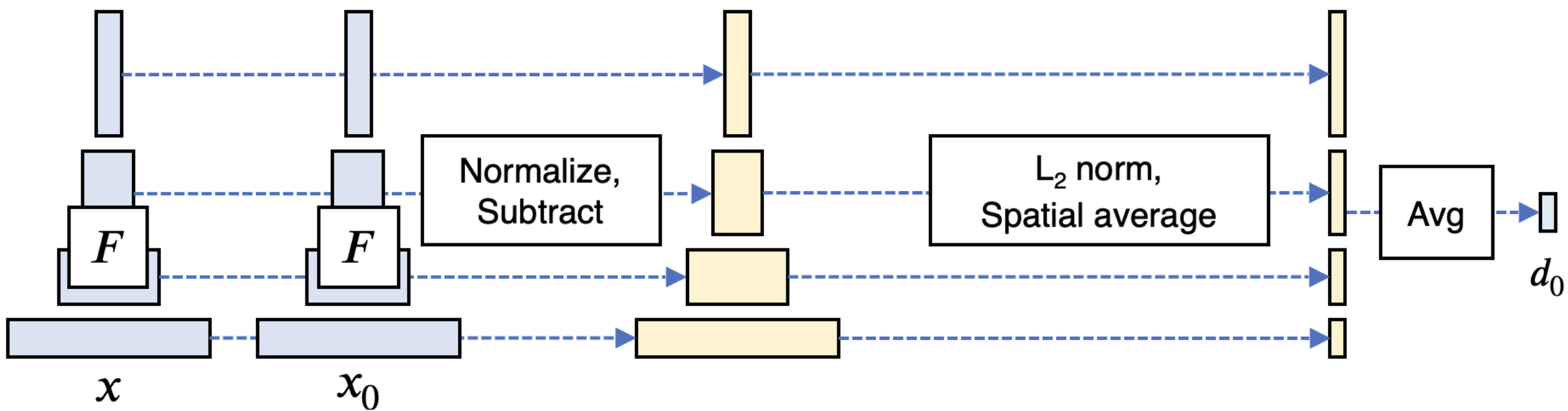
SSIM/FSIMc

Deep Networks?

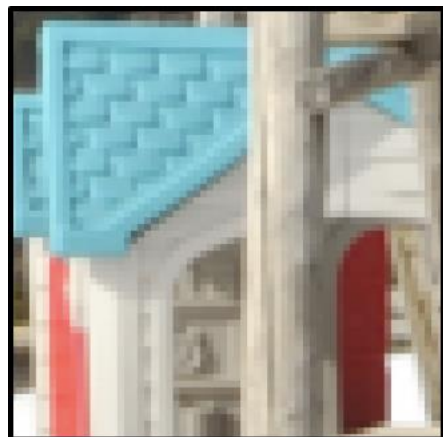


< C ✓ >

Deep Networks as a Perceptual Metric

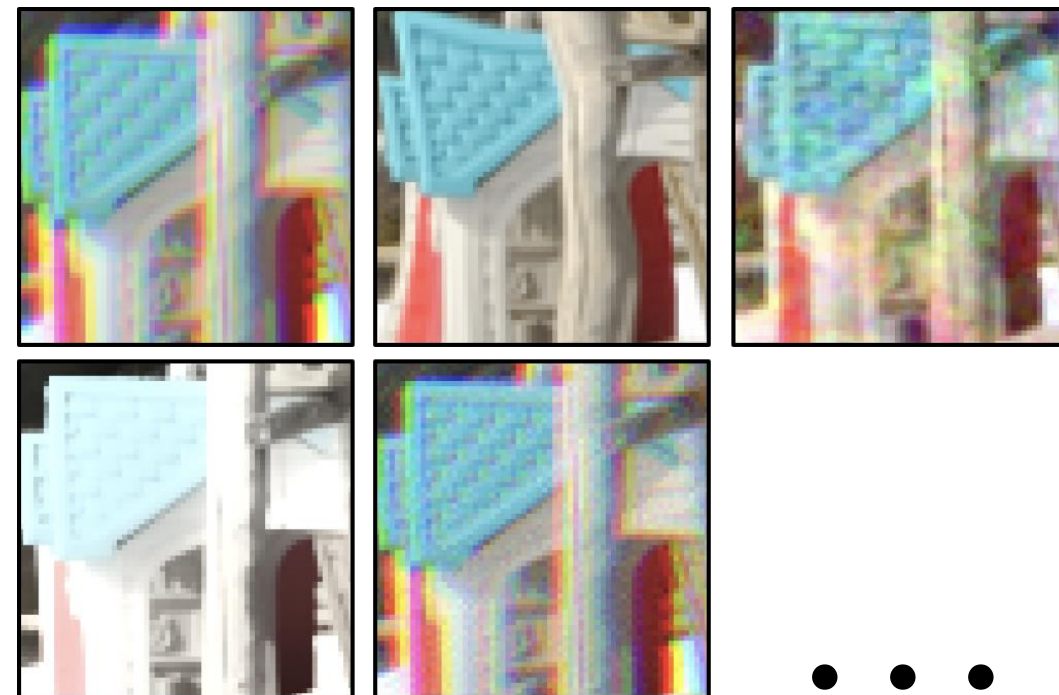


Distortions



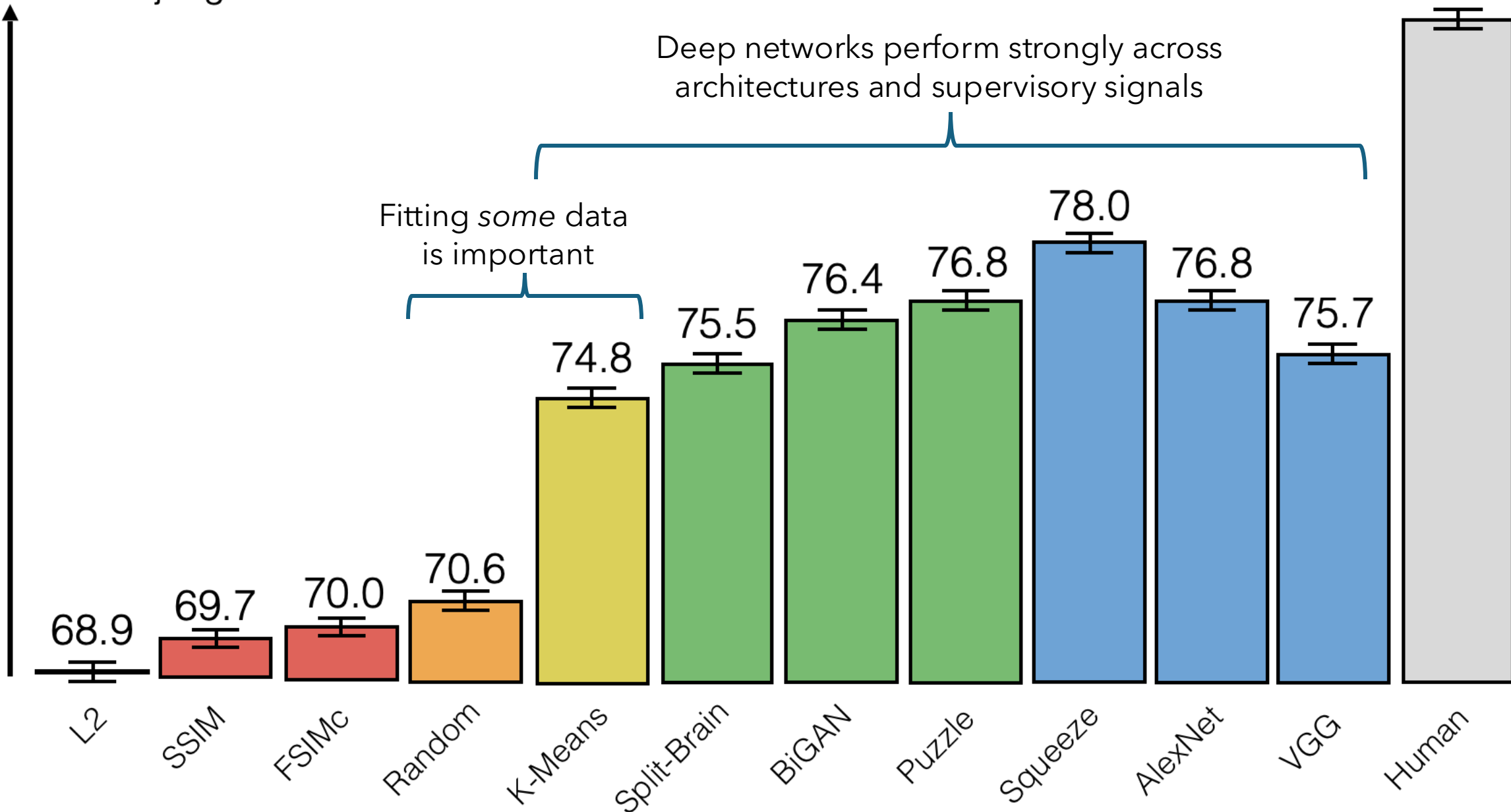
Original Patch

Noise
Photometric
Spatial warps
Compression
Blur



Distorted Patches

% agreement with human judges



Deep networks perform strongly across architectures and supervisory signals

Fitting some data is important

How different are *these* images?



DreamSim: Learning New Dimensions of Human Visual Similarity using Synthetic Data

 <https://dreamsim-nights.github.io/> 



Stephanie Fu^{*1}



Netanel Y. Tamir^{*2}



Shobhita Sundaram^{*1}



Lucy Chai¹



Richard Zhang³



Tali Dekel²



Phillip Isola¹

^{*}Equal contribution, order decided by random seed



Which image, A or B, is more similar to the reference?

A



Reference



B



< LPIPS >

Humans'

Which image, A or B, is more similar to the reference?

A



Reference



B



 Humans >

<  LPIPS

 DINO

 CLIP

Which image, A or B, is more similar to the reference?

A



Reference



B



< LPIPS >

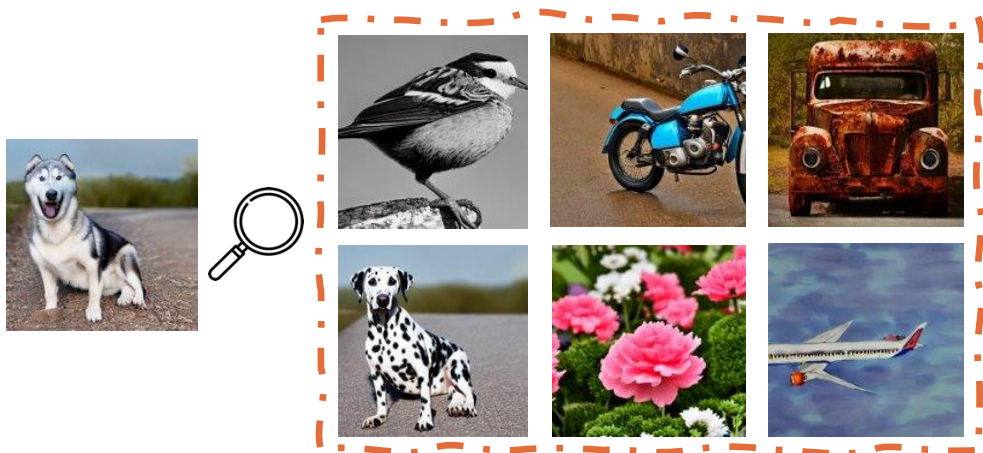
DINO CLIP

Humans'

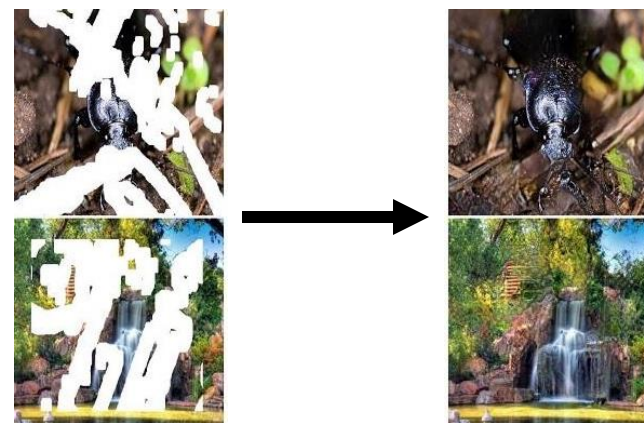
DreamSim

$$f\left(\text{img}_1, \text{img}_2\right) = d$$

Image retrieval



Loss function

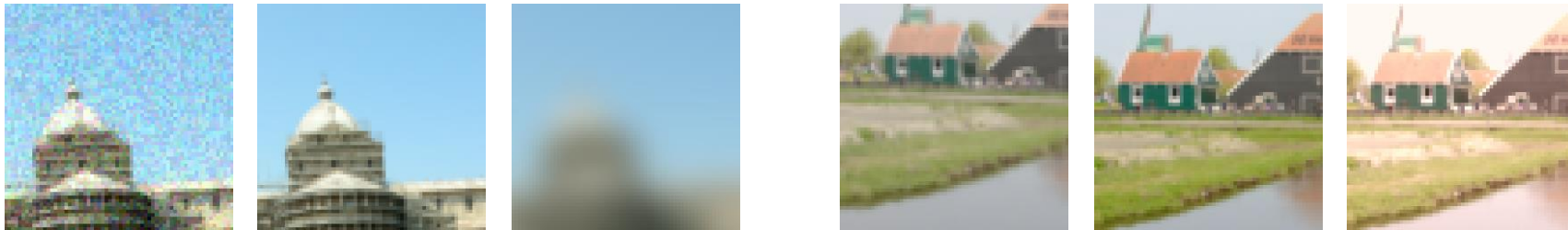


Liu et al, Image Inpainting for Irregular Holes Using Partial Convolutions, *ECCV 2018*

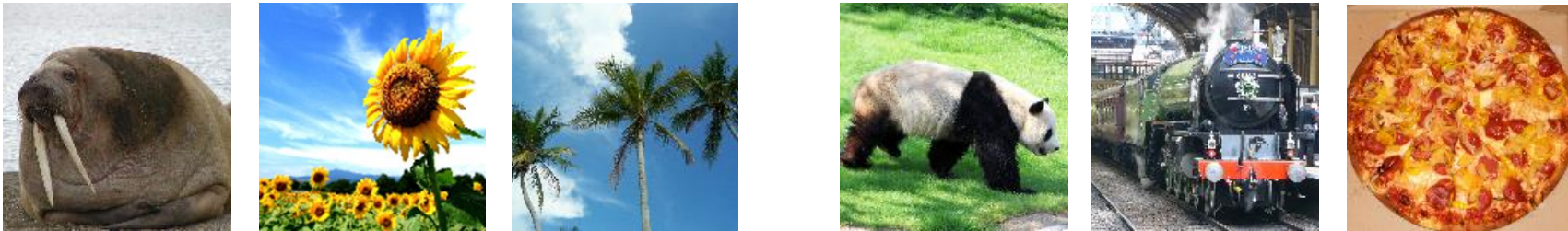
Perceptual similarity datasets

We can improve f by finetuning on perceptual similarity datasets

- BAPPS - images & low-level variations (blurring, saturation, shifting, etc..)

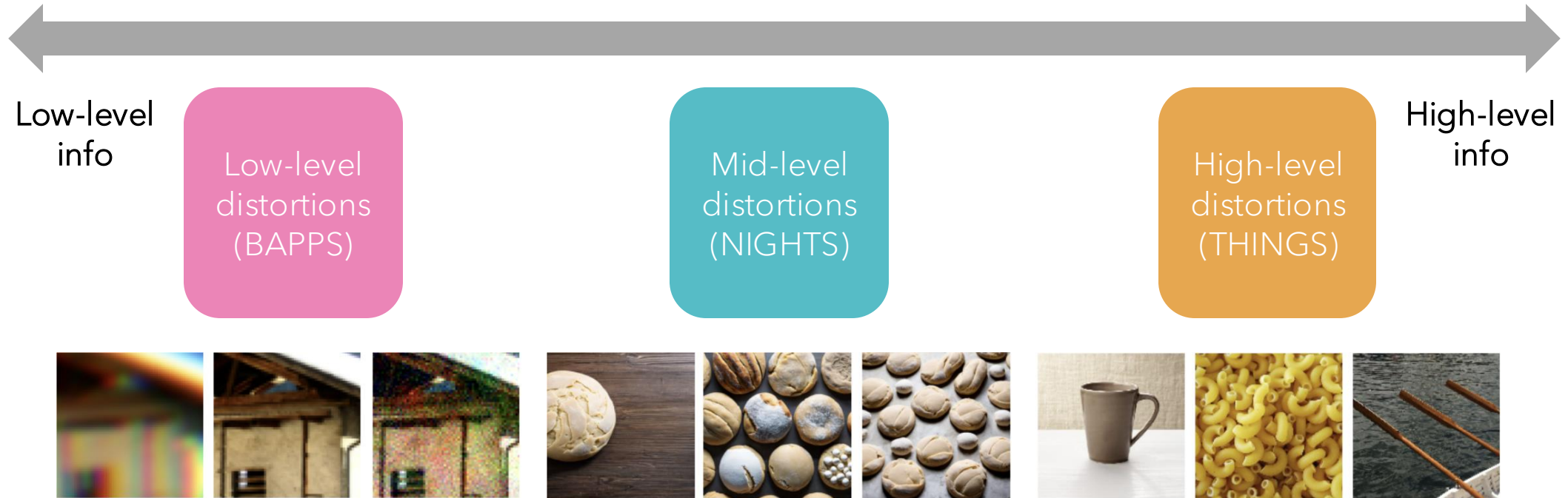


- THINGS - images depicting classes (more conceptual)



These datasets don't capture the variations we saw in our experiment!

Perceptual similarity datasets





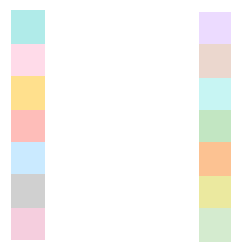
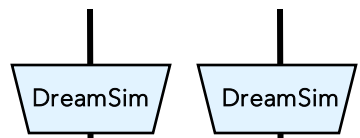
Low-level
info

Low-level
distortions
(BAPPS)

Mid-level
distortions
(NIGHTS)

High-level
distortions
(THINGS)

High-level
info



$D (\text{bread image} , \text{bread image})$

NIGHTS – Novel Image Generations with Human-Tested Similarity

Goal: create a dataset of triplets which exhibit changes in **mid-level** information

“An image of
a **ski lodge**”



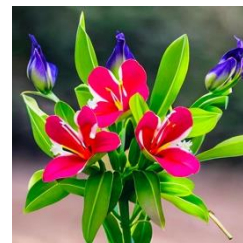
Stable
Diffusion



3 seeds

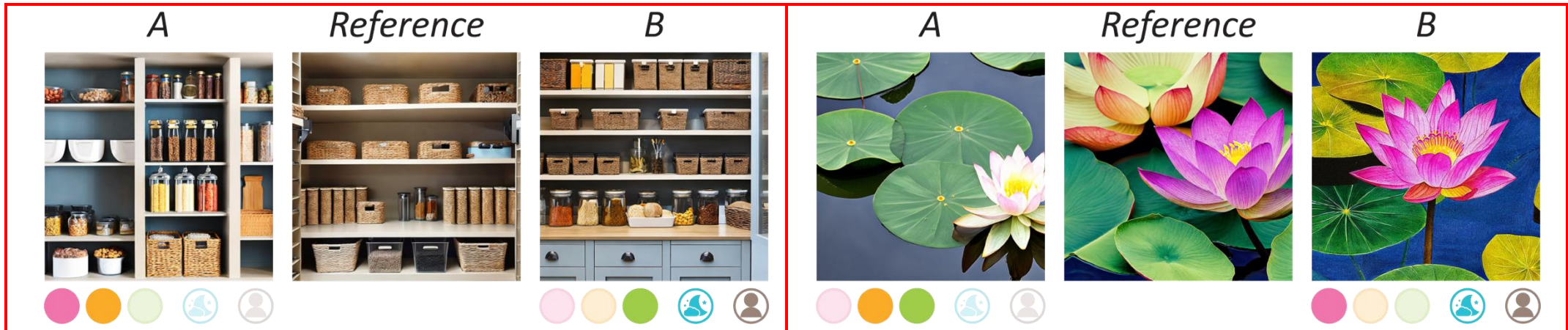
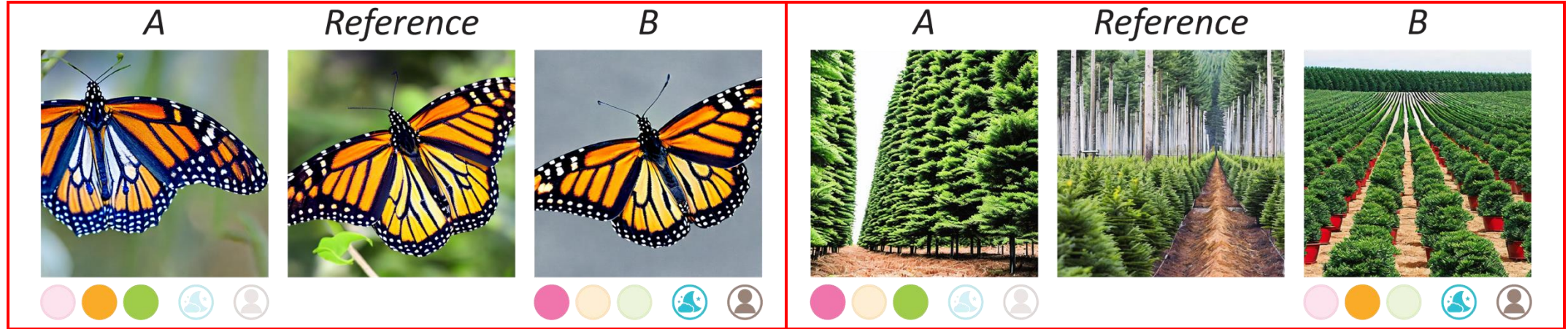
Two-alternative forced choice (2AFC) test

Which image, A or B, is more similar to the
reference?



- ~20k **synthetic** image triplets with unanimous human votes
- Average of 7 votes per triplet
- Classes taken from ImageNet, Food-101, SUN397, etc.

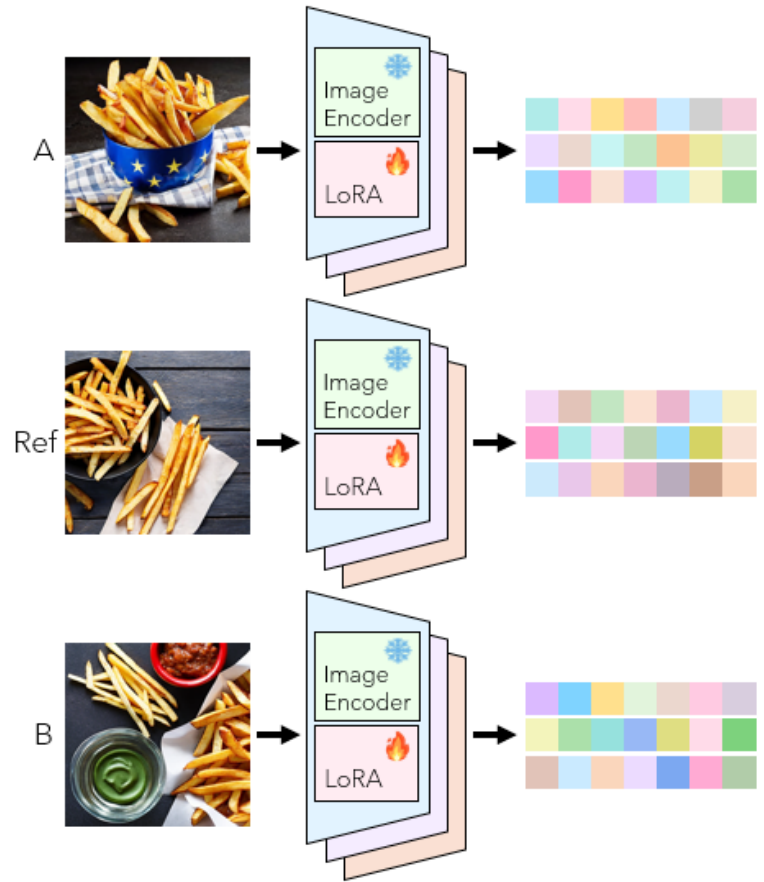
Examples of NIGHTS triplets



 LPIPS  DINO  CLIP  DreamSim  Humans

Training & Inference

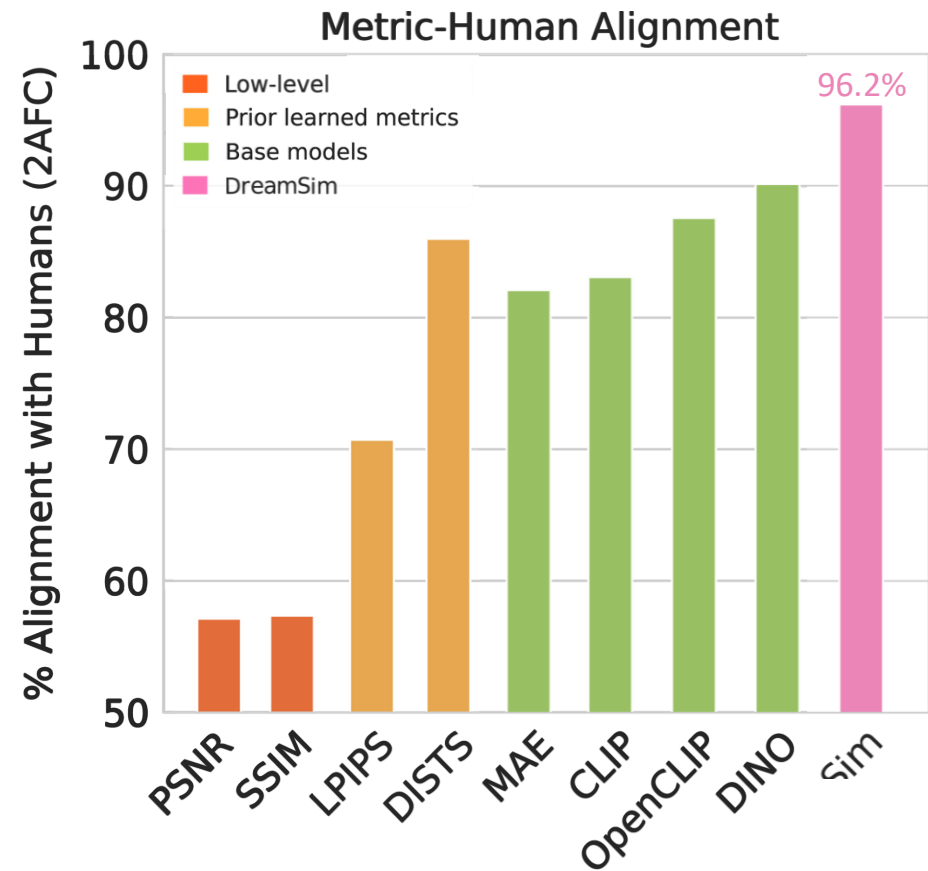
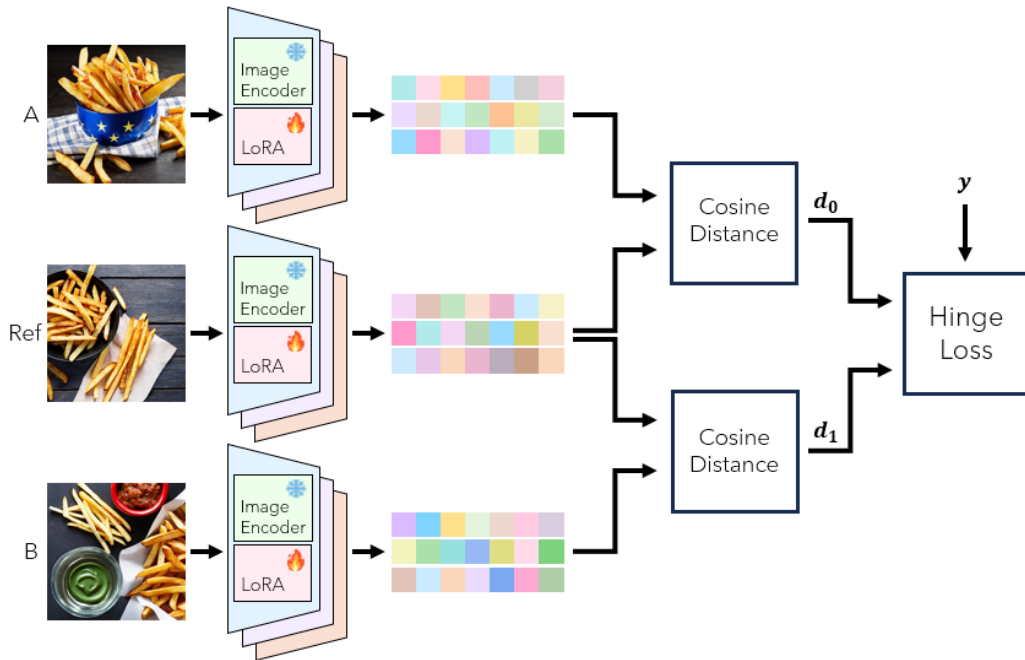
Training: use hinge loss between distances (= triplet loss between embeddings)



Use Low-Rank Adaptation (LoRA)
Tunes 0.5% of ViT parameters

Inference: cosine distance between embeddings of two images

Training & Inference



Nearest Neighbors

Input



Nearest Neighbors

LPIPS



DISTS



OpenCLIP



DINO



Ours



Nearest Neighbors (COCO + ImageNet-R)

Input



Nearest Neighbors

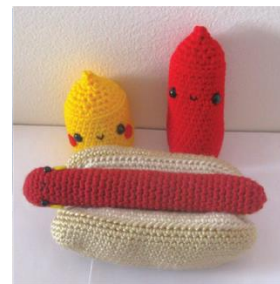
LPIPS



DISTS



OpenCLIP



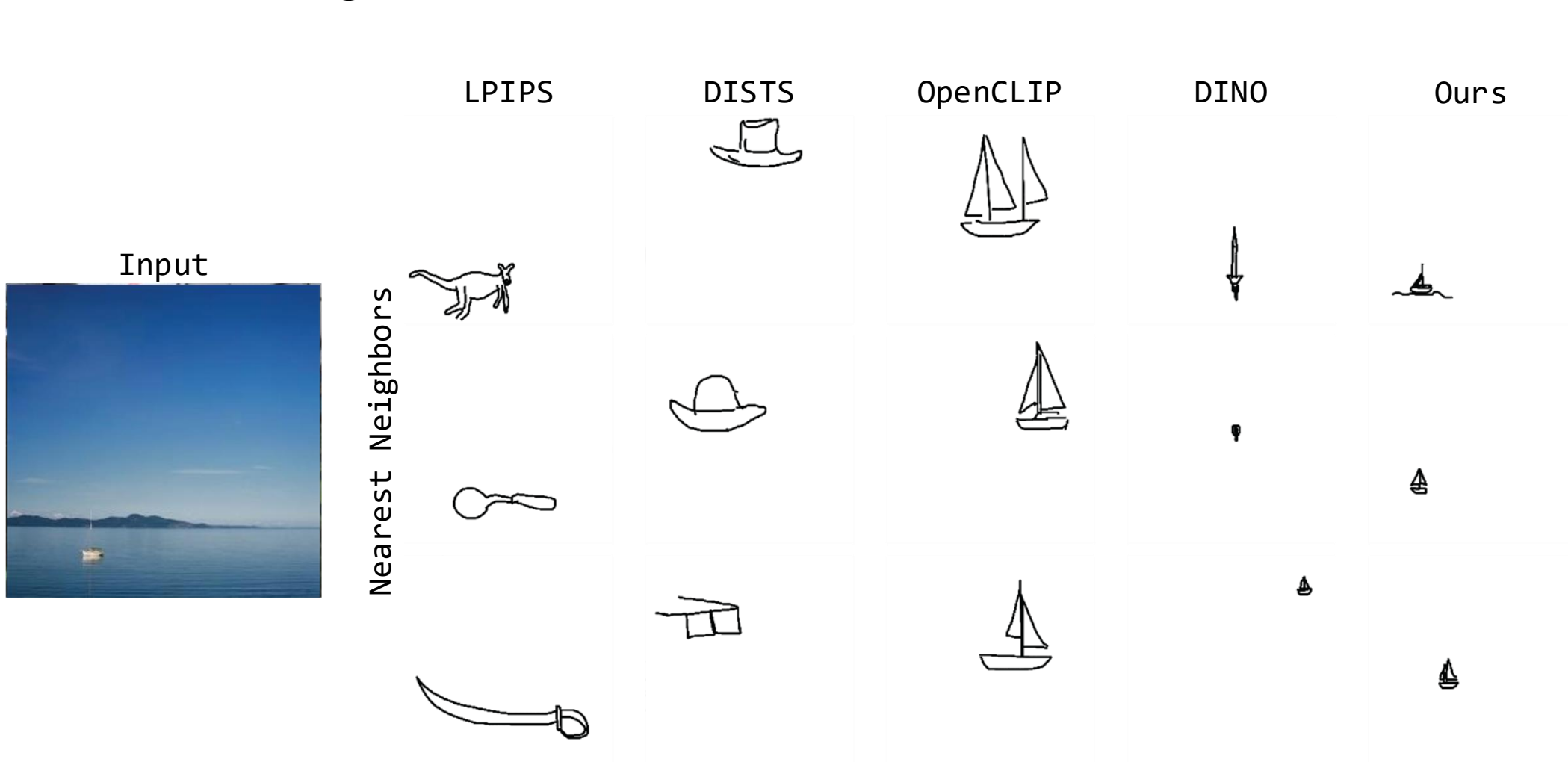
DINO



Ours



Nearest neighbors (Photos → Sketches)



Generation

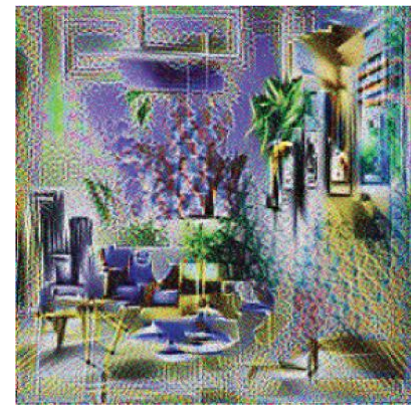
Target



OpenCLIP



DINO



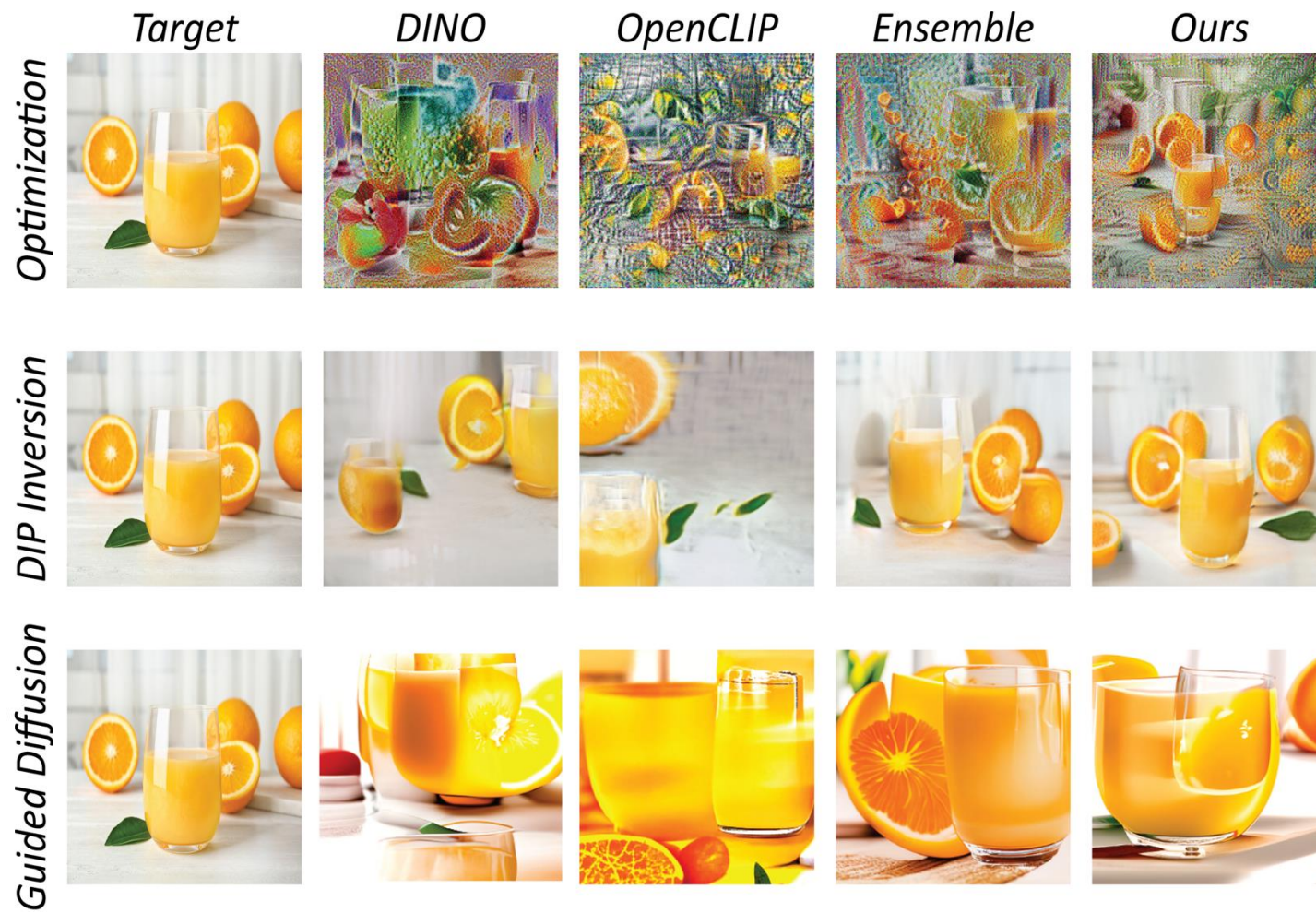
Ours



Guided Diffusion Optimization



Inversion



Evaluating Generated Images

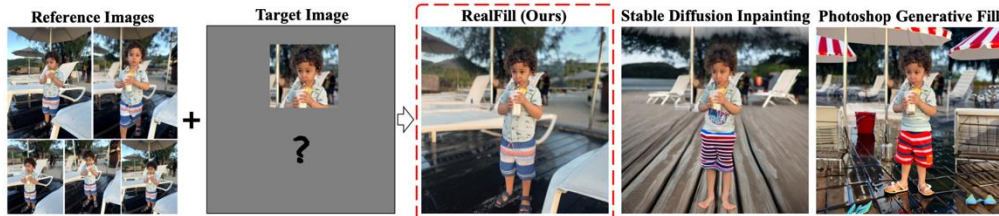
RealFill: Reference-Driven Generation for Authentic Image Completion

LUMING TANG, Cornell University, US
 NATANIEL RUIZ, Google Research, US
 QINGHAO CHU, Google Research, US
 YUANZHEN LI, Google Research, US
 ALEKSANDER HOŁYŃSKI, Google Research, US
 DAVID E. JACOBS, Google Research, US
 BHARATH HARIHARAN, Cornell University, US
 Yael Pritch, Google Research, Israel
 NEAL WADHWA, Google Research, US
 KFIR ABERMAN, Snap Research, US
 MICHAEL RUBINSTEIN, Google Research, US

SC

Zhixu

¹Carr



American Culture

Nigerian Culture

Korean Culture

Adi Haviv¹

Shahar Sarfaty¹

Uri Hachohen²

Niva Elkin-Koren²

Roï Livni³

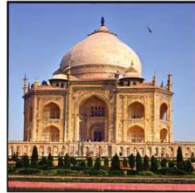
Amit H Bermano¹

Original
Stable Diffusion



Stereot

SCoFT (Ours)



(a) "Photo of a traditional building, in [Culture]"

Customizing Text-to-Image Models with a Single Image Pair

Image Generation

Every Image is Worth a Thousand Words: Quantifying Originality in Stable Diffusion

Sheng-Yu Wang¹ Nupur Kumari¹
 David Bau² Jun-Yan Zhu¹

ig¹, Zhifei Zhang², Zhe Lin², Scott Cohen², Brian Price²,
 ang², Soo Ye Kim², He Zhang², Wei Xiong², Daniel Aliaga¹
 Purdue University¹, Adobe Research²



Conclusion

	Low-Level	High-Level
Unary/Holistic $s(x)$	Blurriness, No-Reference IQA	PickScore, ImageReward
Image Similarity $s(x, x_{ref})$	PSNR, SSIM, LPIPS, DIST _S	DreamSim
Distribution $s(p(x)); s(p(x), p_{ref})$	InceptionScore, FID, CMMD	
Cross-Modal Similarity $s(x, y_{ref})$	SOA, CLIPScore	

What's Next?

- How can evaluation metrics be incorporated more directly into generation pipelines?
 - RLHF
 - Reward functions

What's Next?

- How can evaluation metrics be incorporated more directly into generation pipelines?
 - RLHF
 - Reward functions
- Multiple different eval metrics v. one holistic eval metric?

What's Next?

- How can evaluation metrics be incorporated more directly into generation pipelines?
 - RLHF
 - Reward functions
- Multiple different eval metrics v. one holistic eval metric?
- Cross-model alignment