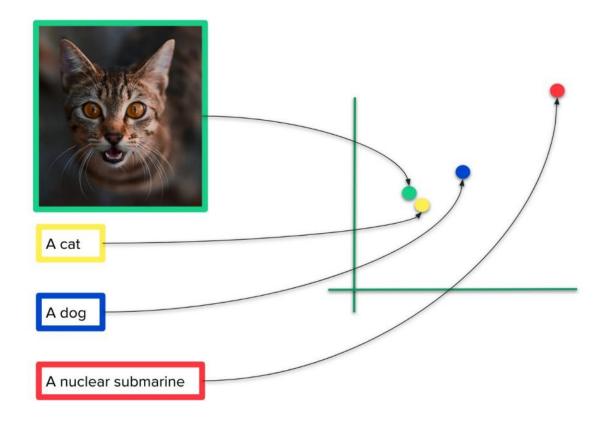
Session #18: Self-Supervised Vision-Lang Models

Thursday, October 27 CSCI 601.771: Self-supervised Statistical Models







The Toronto skyline with Google brain logo written in fireworks.



Teddy bears swimming at the Olympics 400m Butterfly event.

Stability AI, the startup behind Stable Diffusion, raises \$101M

Kyle Wiggers @kyle_I_wiggers / 1:01 PM EDT • October 17, 2022

Comment



Stakeholder: Introduction

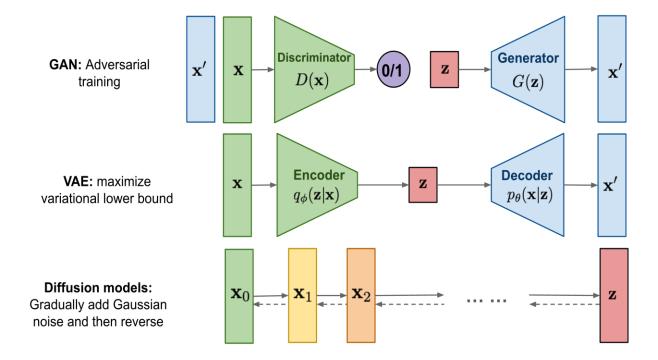
unCLIP: Two-stage model that generates image from text caption

- CLIP: Model for image representation
- Diffusion model: Model for image generation
- Combined: text-conditional generative model
 - Prior model: Given text caption, generates CLIP image embeddings (diffusion > autoregressive)
 - Decoder (invert CLIP encoder): Given image embedding, produces image



a propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese

Stakeholder: Related Work



Stakeholder: Method

Training data: (x,y) Embeddings: z_i, z_t

<u>**Prior:**</u> $P(z_i | y)$

Autoregressive: *z_i* converted into sequence of discrete codes and predicted autoregressively

- reduce dimensionality with PCA (transformer with causal attention mask)
- append to sequence: $y_i z_t, z_i \cdot z_t$

Diffusion: *z*_i modeled using Gaussian distribution model

- decoder-only Transformer with causal attention mask
- append to sequence: y, z_t, diffusion timestep embedding, noised z_i

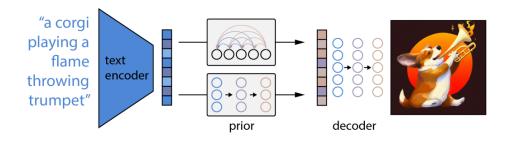
<u>**Decoder:**</u> $P(x \mid z_i, y)$

Generation (diffusion model): Modify GLIDE encoder to project CLIP embeddings 1) to existing timestep embedding and 2) into four tokens of context

- classifier-free guidance

High resolution (upsampler models): 64x64 > 256x256 > 1024x1024

- Gaussian blur | BSR degradation



Stakeholder: Image Manipulation

Image x encoded as (z_i, x_T)

Variation

Same content but vary in: shape and orientation

DDIM with η (=0 is deterministic)



Stakeholder: Image Manipulation

Interpolation: rotate embeddings with spherical interpolation



Stakeholder: Image Manipulation

Text Diff: Norm vector for text diff between two text captions

- rotate between embedding and text diff with spherical interpolation

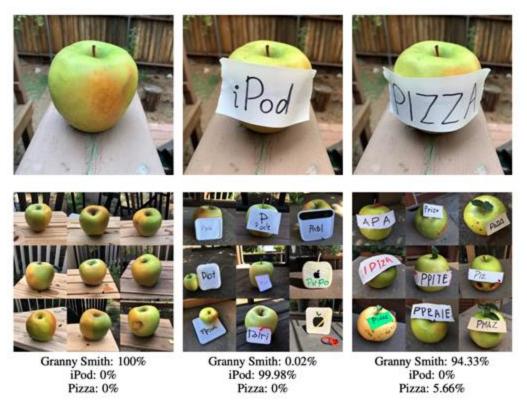


a photo of a victorian house \rightarrow a photo of a modern house



a photo of an adult lion \rightarrow a photo of lion cub

Probing the CLIP Latent space



₽: Isabel, Vicky

Decoder model allows us to visualize with the CLIP image encoder is seeing

Probing the CLIP Latent space



P: Isabel, Vicky
Use PCA on CLIP embeddings and progressively increase the dimension

Text to Image Generation - Human Evaluation

- Compare to GLIDE (SOTA in T2I Gen)
- Ask humans to give preference on the following:
 - Photorealism
 - Caption Similarity
 - Diversity

unCLIP Prior	Photorealism	Caption Similarity	Diversity
AR Diffusion	$\begin{array}{c} 47.1\% \pm 3.1\% \\ 48.9\% \pm 3.1\% \end{array}$	$\begin{array}{c} 41.1\% \pm 3.0\% \\ 45.3\% \pm 3.0\% \end{array}$	$\begin{array}{c} 62.6\% \pm 3.0\% \\ 70.5\% \pm 2.8\% \end{array}$

Probability humans preferred unCLIP over GLIDE

Improved Diversity-Fidelity Trade-off with Guidance



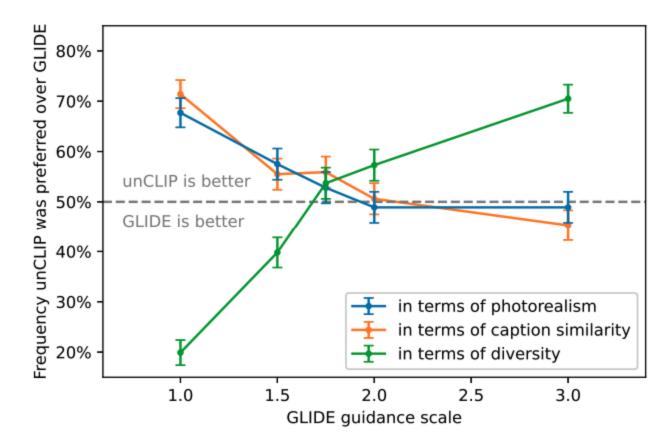
Guidance

₽: Isabel, Vicky

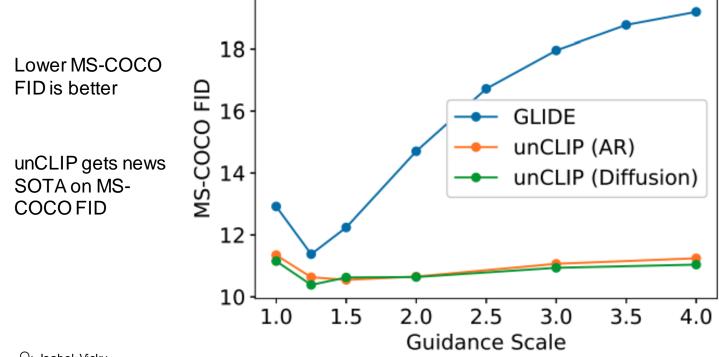
unCLIP

GLIDE

Improved Diversity-Fidelity Trade-off with Guidance

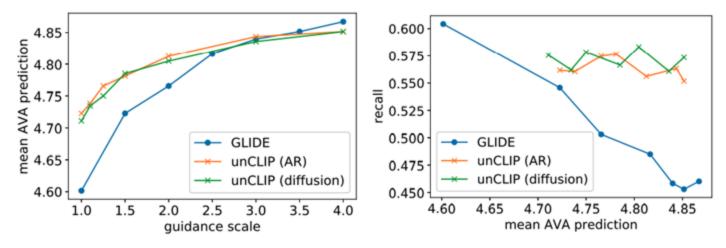


Improved Diversity-Fidelity Trade-off with Guidance



Aesthetic Quality Comparison

- 1. Use GPT-3 to sample artistic captions
- 2. Train a model on AVA dataset to predict aesthetic judgements
- 3. Use model to evaluate aesthetic quality



Limitations

unCLIP is worse than GLIDE at binding attributes to objects

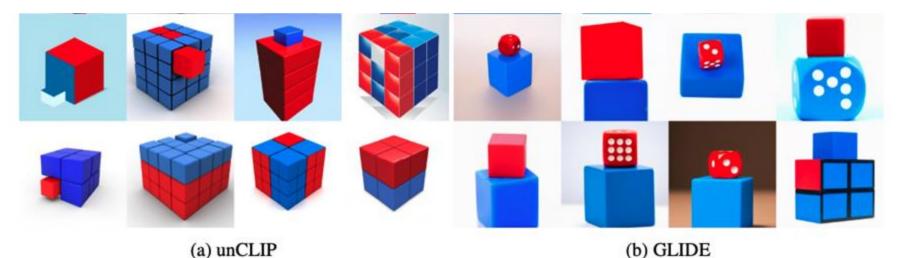


Figure 14: Samples from unCLIP and GLIDE for the prompt "a red cube on top of a blue cube".

Limitations

unCLIP struggles to produce coherent text



Figure 16: Samples from unCLIP for the prompt, "A sign that says deep learning."

Limitations

unCLIP produce low details for complex scenes



(b) A high quality photo of Times Square.

Review: Strengths

- 1. Importance of the Prior
 - 1. Modelling CLIP image embedding is integral to good performance. Alternative considered and tested
- 2. unCLIP latent space has nice properties like interpolation and controllable variation
- 3. Visualization of the latent space of CLIP a great tool for understanding the inner

processes of image encoder.



Granny Smith: 0.02% iPod: 99.98% Pizza: 0% Granny Smith: 94.33% iPod: 0% Pizza: 5.66%

Reviewer: Weaknesses

- 1. Too much unexplained concepts, hard to understand.
- 2. The model is weak with complex situations. e.g., producing coherent text and details in complex scenes.
- 3. Computation of the prior (AR & Diffusion) is poorly motivated:
 - 1. Autoregressive prediction of dimensionality-reduced CLIP image embeddings
 - 1. Assumptions of PCA met?
 - 2. Why autoregressive?
 - 2. Diffusion prior on sequence of [encoded text, CLIP text emb, emb for diffusion timestep, noised CLIP image emb, final emb to predict image embedding]

Reviewer: Weaknesses

- The label of each image is no longer a noun, but a sentence, so images that were forcibly divided into the same kind in the past have "infinitely fine-grained" labels.
 - Example: ImageNet labels images as "dogs," and with this paired example, one can learn the nuances of "dogs" in different environments and doing different things.
- The pairing of text and images is not strong enough. This is why the authors repeatedly stress collecting huge datasets, because they have to suppress the noise by means of big data.
- OpenAI has not released the data set, we only know that there are 500,000 queries, but we do not know what these queries are, such as simple descriptions or complex descriptions. If it is just a simple description, such as "black cat", "lecture room", etc.

- DALL-E 2 is the only open-source implementation of unclip (https://labs.openai.com)
- Can you guess which one is the real puppy?















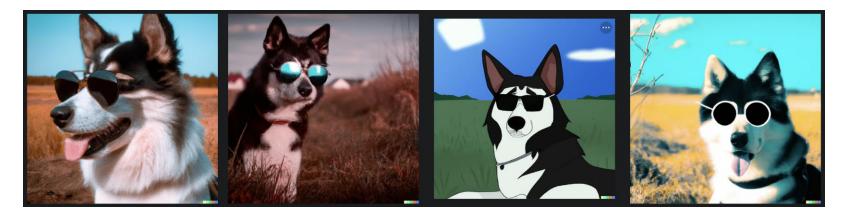
• Variations generated by the decoder diffusion model



• A cubist painting of a Yakutian Laika dog in a field with sunglasses



• A Yakutian Laika dog in a field with sunglasses in anime style



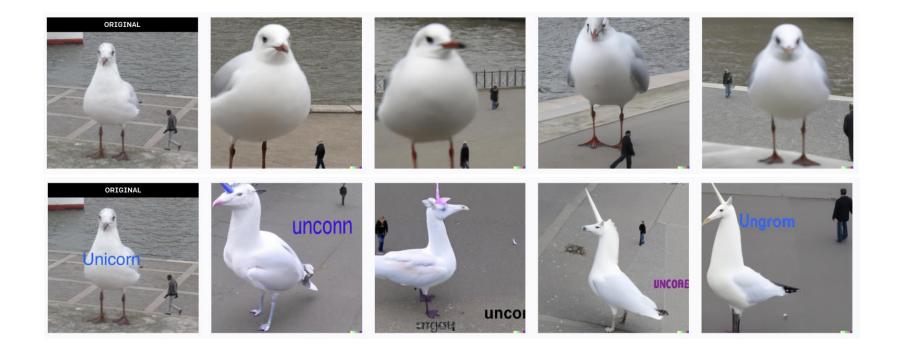
• A Yakutian Laika dog in a space suit sipping a piña colada in the desert











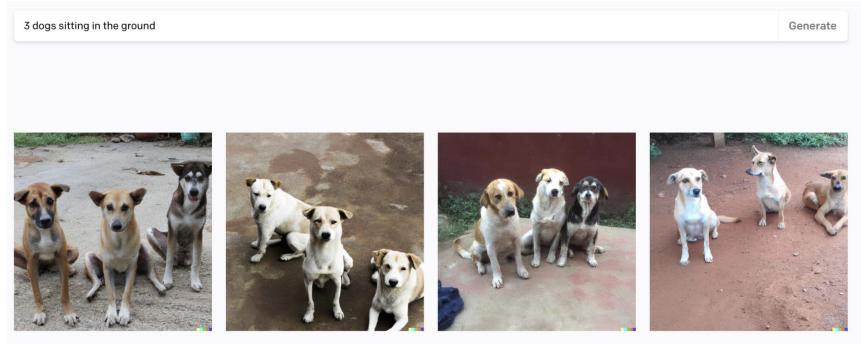
Limitation1: Number Understanding

Edit the detailed description	Surprise me	Upload	→I
A dog with the sign of 1024		Genera	ate

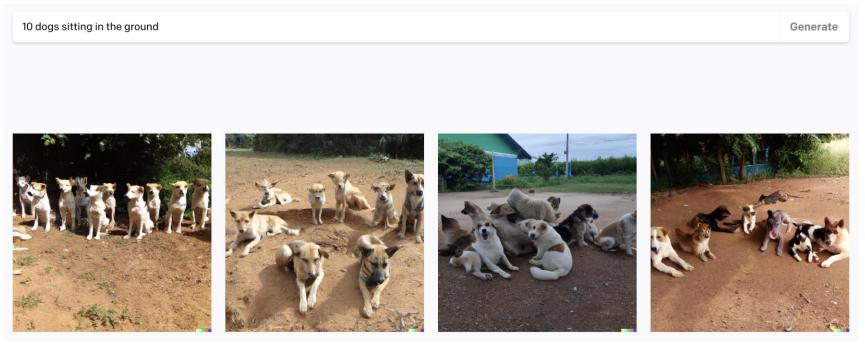




Limitation2: Counting



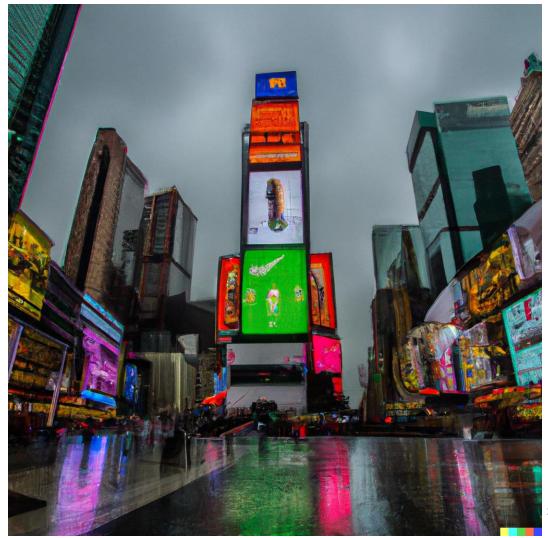
Limitation2: Counting





Limitation 2: Bad Details

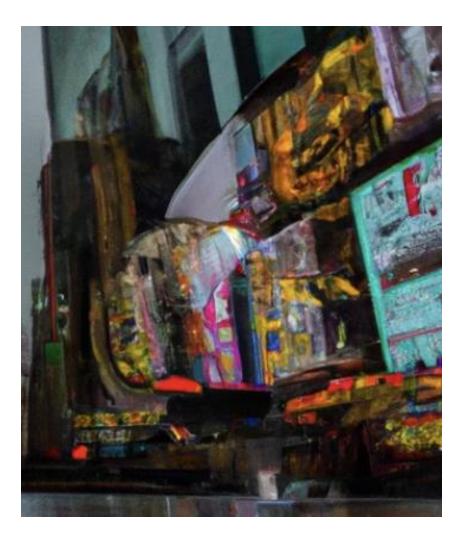
Example: A photo of the time square





Limitation 2: Bad Details

Example: A photo of the time square



Limitation 2: Bad Details

Example: A photo of the time square





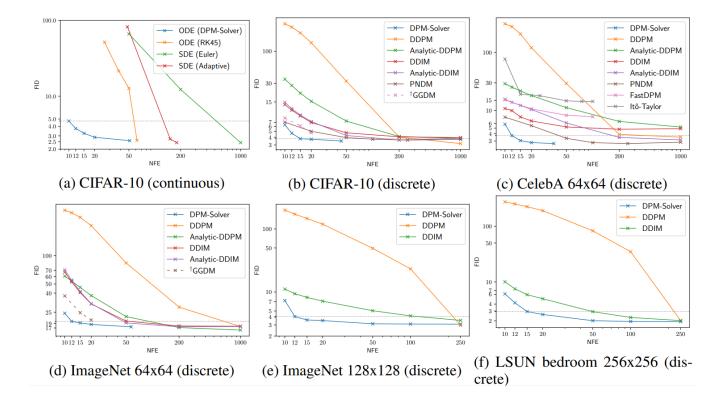
Visionary

• Practicality

• Theory

• Application

Practicality 1.Speed-up in sampling



(2)Speed-up in training?

A remaining open problem...

Relax Constraints?

- Replace Langevin dynamics with?
 - Replace gaussian noise with more complex perturbation?
 - Hierarchical VAE is more powerful -> but less stable and harder to train?

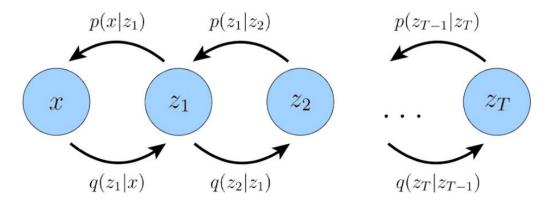
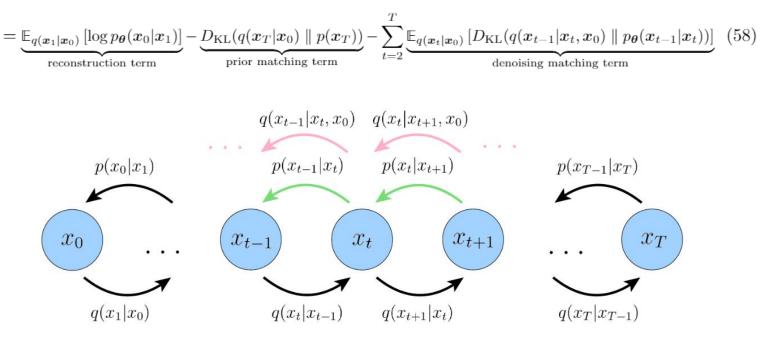


Figure 2: A Markovian Hierarchical Variational Autoencoder with T hierarchical latents. The generative process is modeled as a Markov chain, where each latent z_t is generated only from the previous latent z_{t+1} .

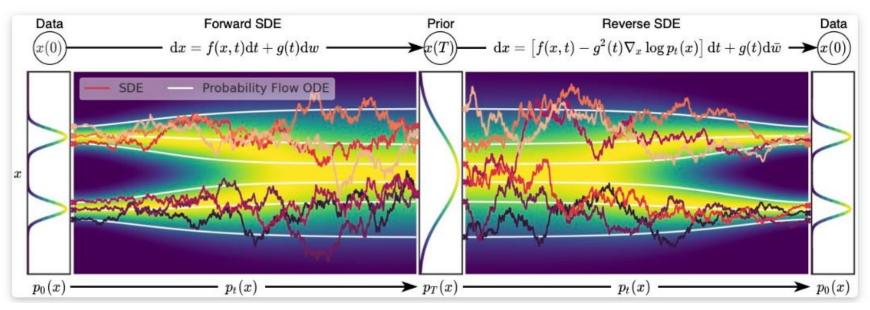
Relax Constraints?

- In variational diffusion model, xo's information is also used
- Can/should we relax the markov property further?

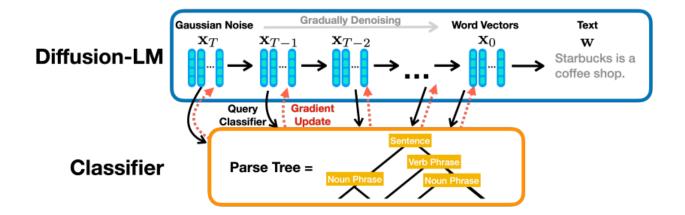


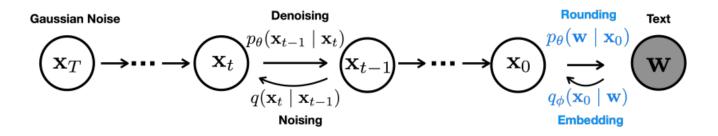
Other potential work

- Better sampling strategies
- Discrete -> Continuous process
 - o Stochastic differential equations (SDE)
 - o Probability flow ODE



Application-Diffusion for Language?





Application – Diffusion for Video?

- Video prediction
- Text-conditioned video generation



sawing a steel pipe, Technician concept

Aerial of horses on a pasture

Dramatic ocean sunset

Forest in Autumn



Berlin - Brandenburg Gate at

night



Traffic jam on 23 de Maio avenue, both directions, south of Sao Paulo,

Busy freeway at night



