

CMPT 983

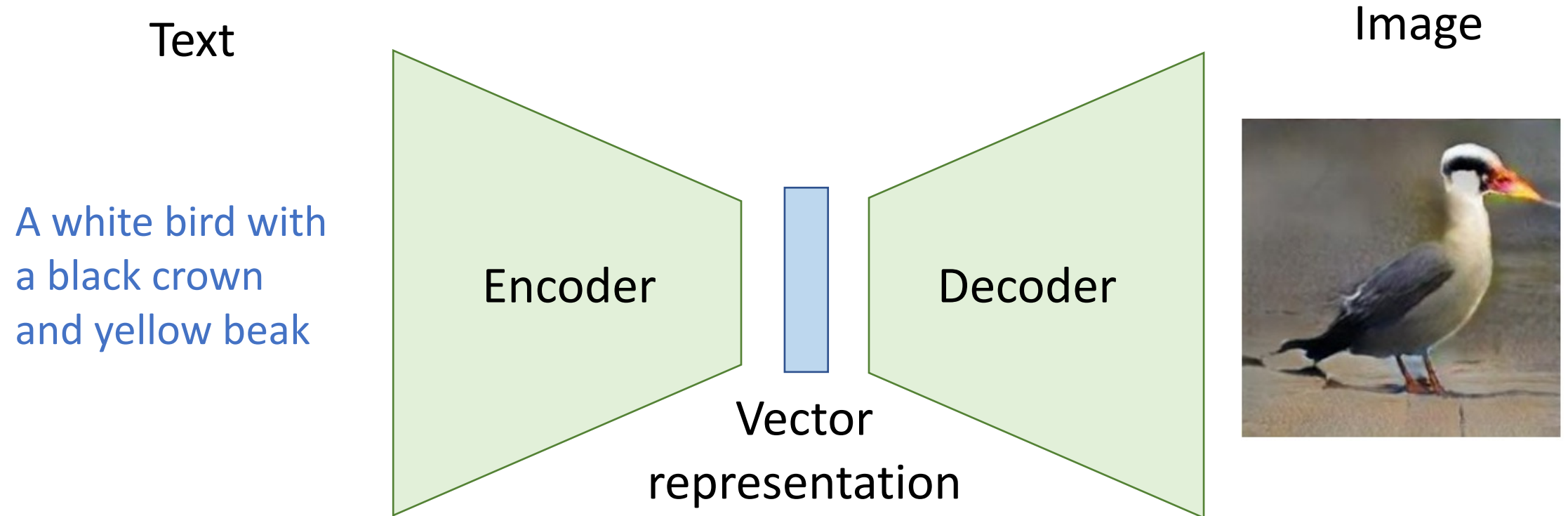
Grounded Natural Language Understanding

February 09, 2022

Content generation from language

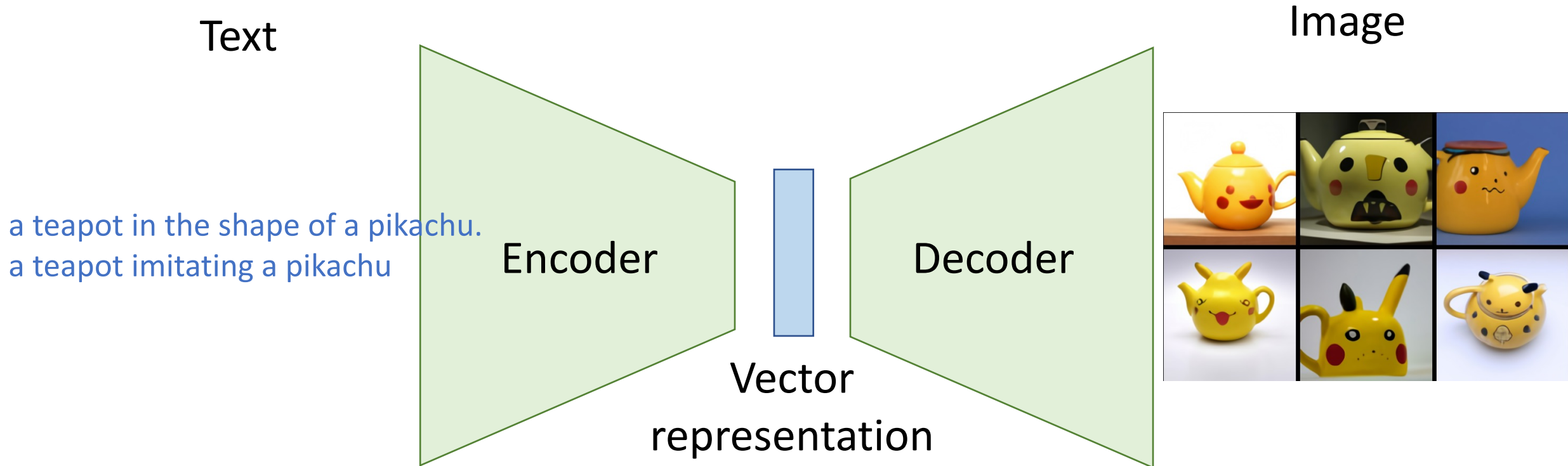
Content generation from
language

Translating across modalities



“StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks”
[Zhang et al, ICCV 2017]

Translating across modalities



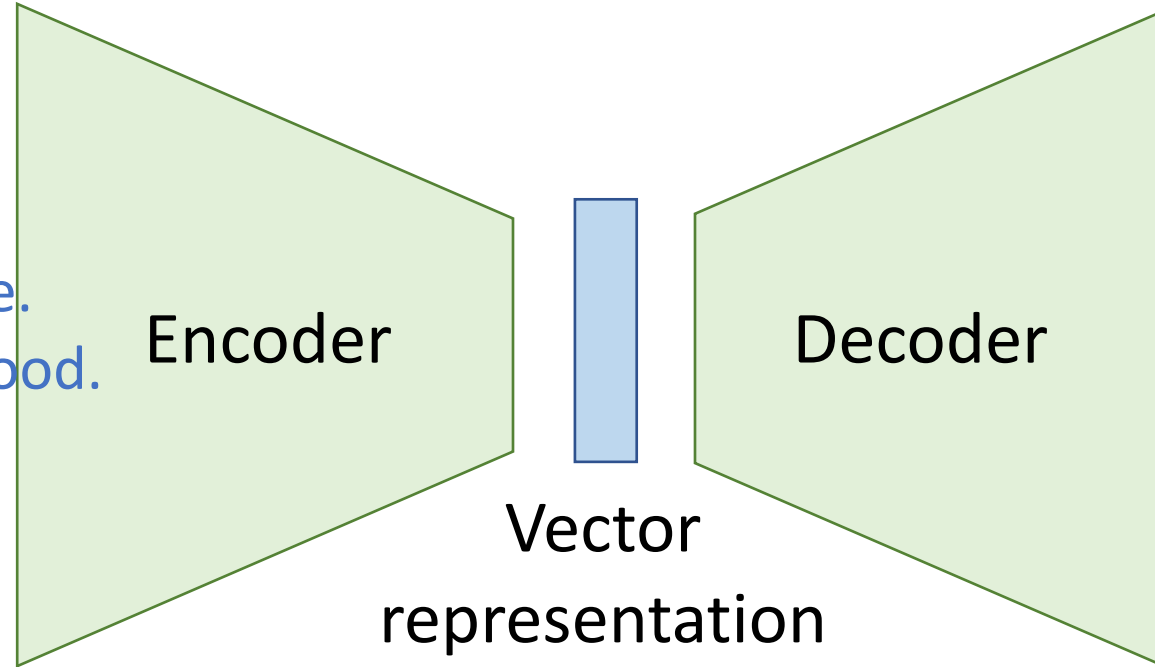
“Dall-e”

[Ramesh et al, <https://openai.com/blog/dall-e/>]

Translating across modalities

Text

Brown colored dining table.
It has four legs made of wood.



3D Shape



“Text2Shape: Generating Shapes from Natural Language by Learning Joint Embeddings”
[Chen et al, ACCV 2018]

How is generating *images*
and *shapes* different from
generating *text*?

Translating across modalities

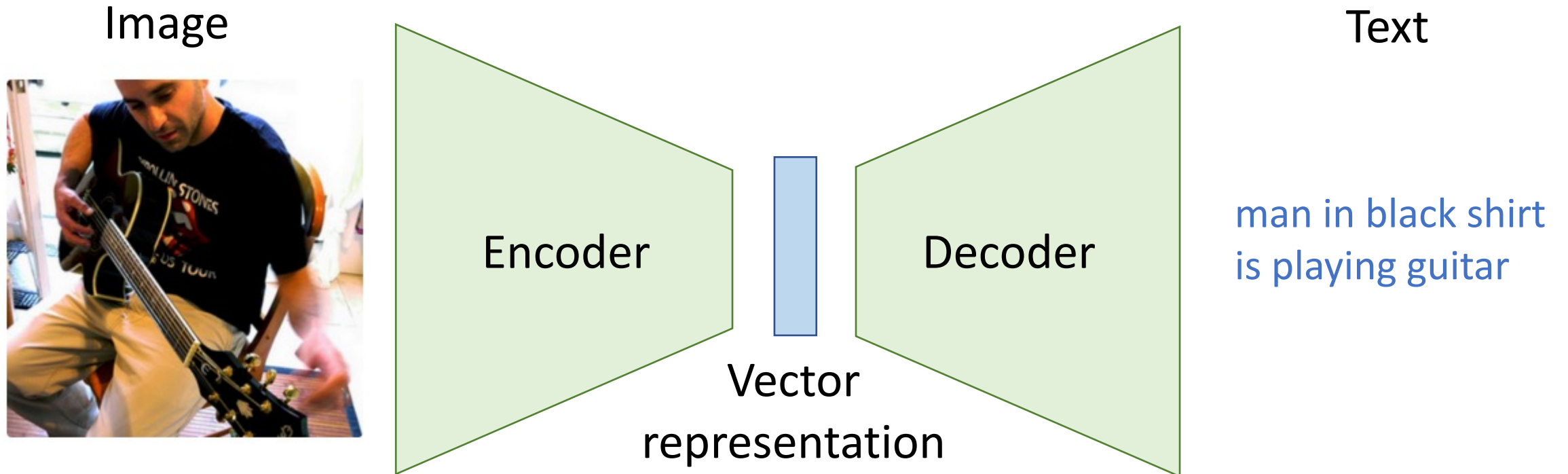


Image captioning

“Deep Visual-Semantic Alignments for Generating Image Descriptions”

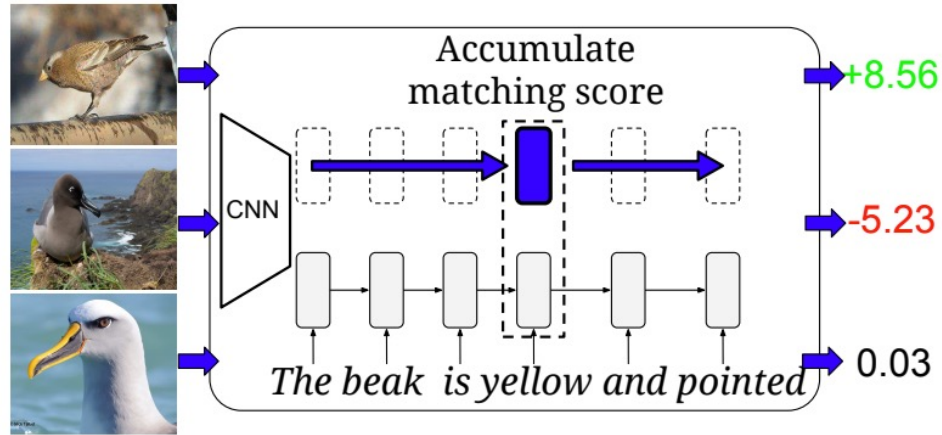
[Karpathy and Fei-Fei CVPR 2015]

Generating Content

- Note: [retrieval](#) as most basic form of generation

Generation as retrieval

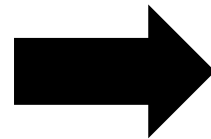
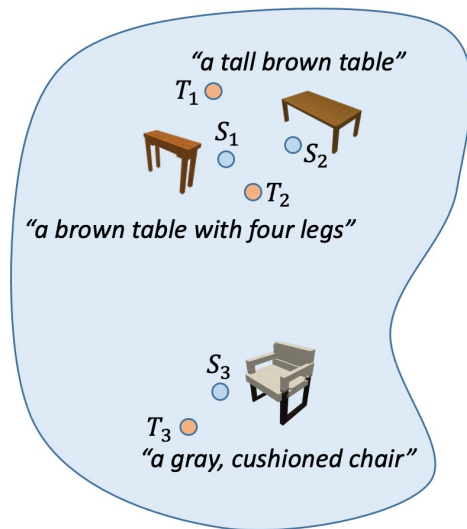
Learn joint embedding → Embed and retrieve



“This is a large black bird with a pointy black beak.”



“Learning Deep Representations of Fine-Grained Visual Descriptions” (Reed et al, CVPR 2016)



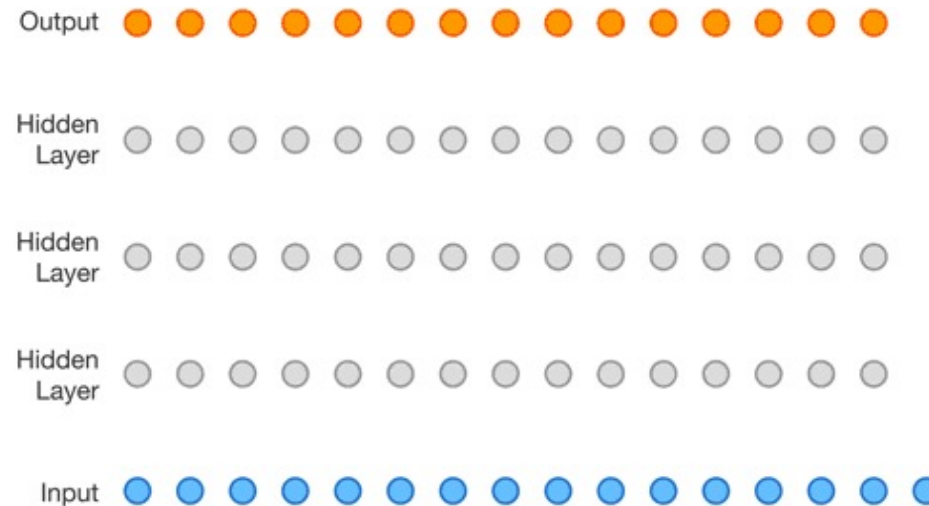
It’s a dark brown,
upholstered chair
with arms and
a curved
rectangular back



“Text2Shape: Generating Shapes from Natural Language by Learning Joint Embeddings” (Chen et al, ACCV 2018)

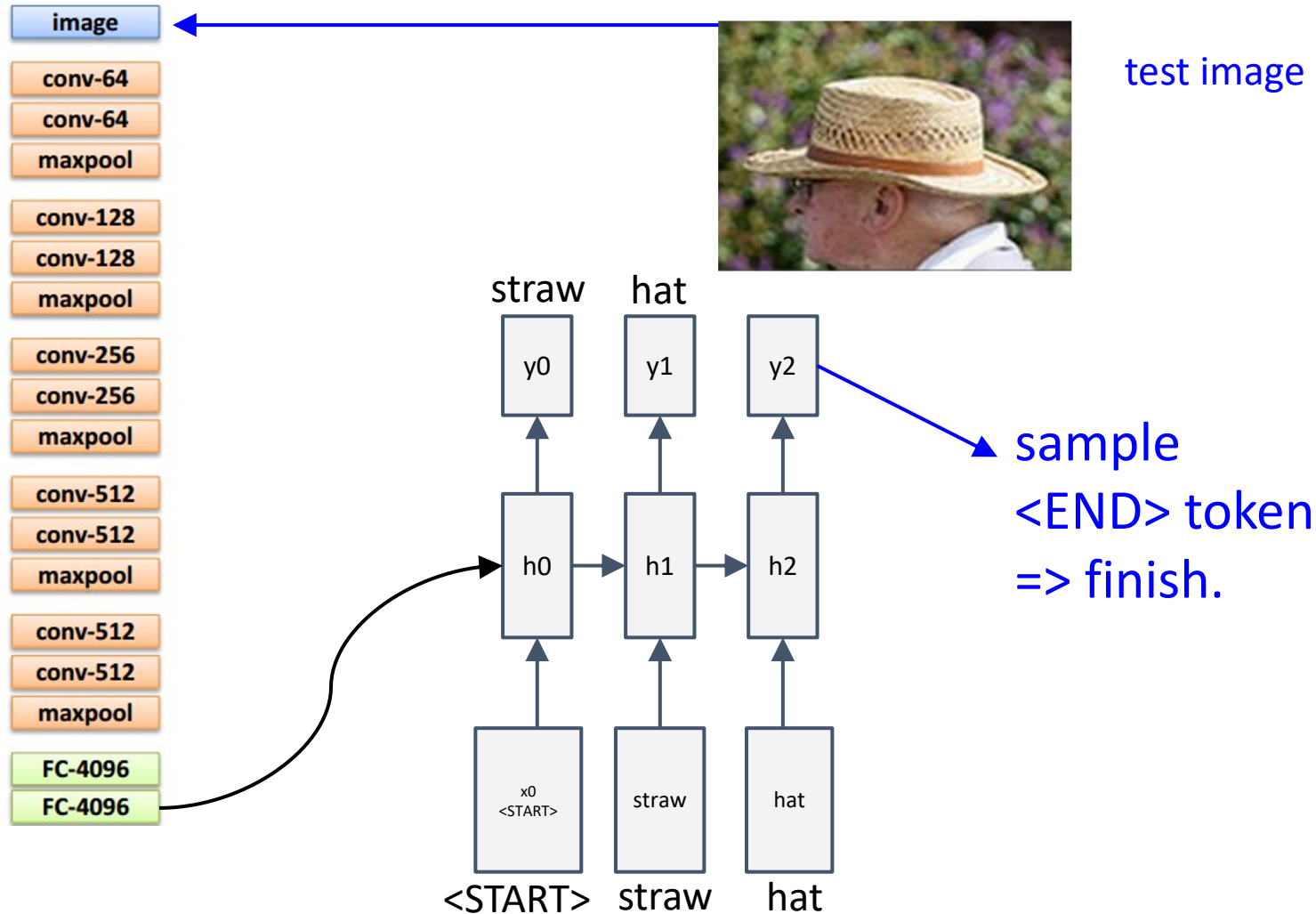
Generating Content

- Note: **retrieval** as most basic form of generation
 - Can also retrieve + edit
- Note : can model as output as a sequence and generate **autoregressively**



https://ml.berkeley.edu/blog/posts/AR_intro/

Autoregressive captioning



Output from previous step is fed as input into next

How to get different outputs?

Decoding strategies:

- Greedy decoding
 - Take $\operatorname{argmax} P_t(w)$
- Beam search
- Sampling
 - Basic sampling: sample from $P_t(w)$
 - Top- n sampling: restrict to top n words
 - Top- p sampling: restrict to top p proportion of words
- Temperature scaling (make distribution less spiky)

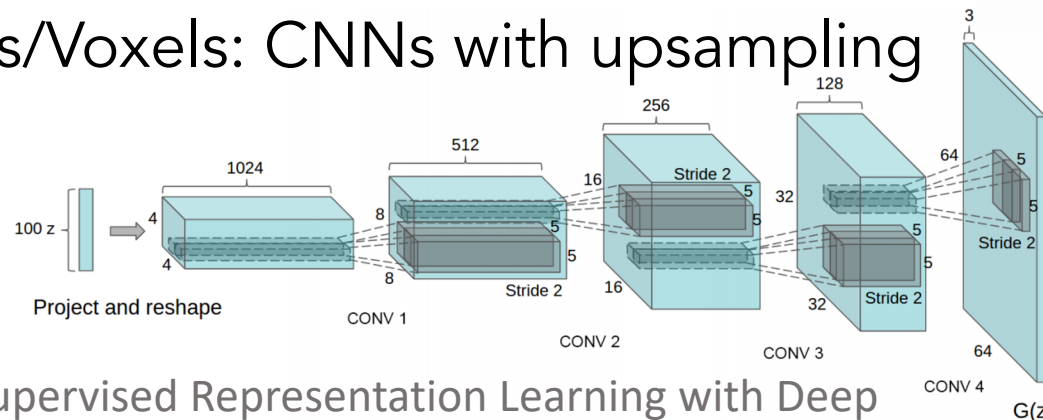
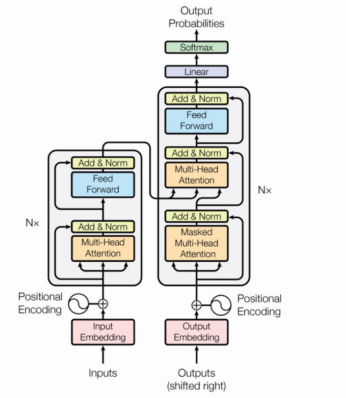
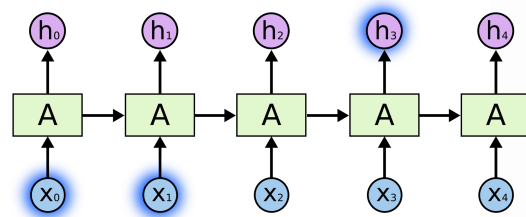
$$P_t(w) = \frac{\exp(s_w/\tau)}{\sum_{w' \in V} \exp(s_{w'}/\tau)}$$

Generating Content

- Note: **retrieval** as most basic form of generation
- Note : can model as output as a sequence and generate **autoregressively**

- Decoders:

- Language: RNNs/Transformers
- Images/Voxels: CNNs with upsampling



“wrongly called deconvolutions”

Taxonomy of machine learning models

Models different probability distributions

Discriminative models:

Learn $p(y|x)$



Assign labels to data

Feature learning (with labels)

Generative Model:

Learn $p(x)$



Detect outliers

Feature learning (without labels)

Sample to generate new data

Conditional Generative Model:

Learn $p(x|y)$



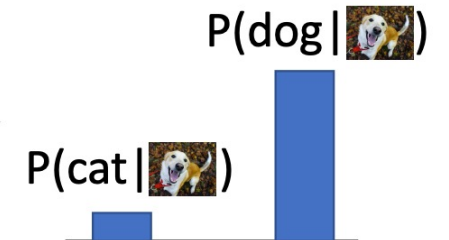
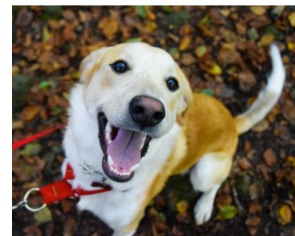
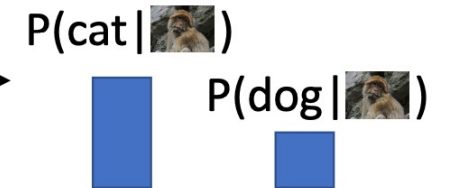
Assign labels, while rejecting outliers!

Generate new data conditioned on input labels

Taxonomy of machine learning models

Models different probability distributions

Discriminative models:
Learn $p(y|x)$

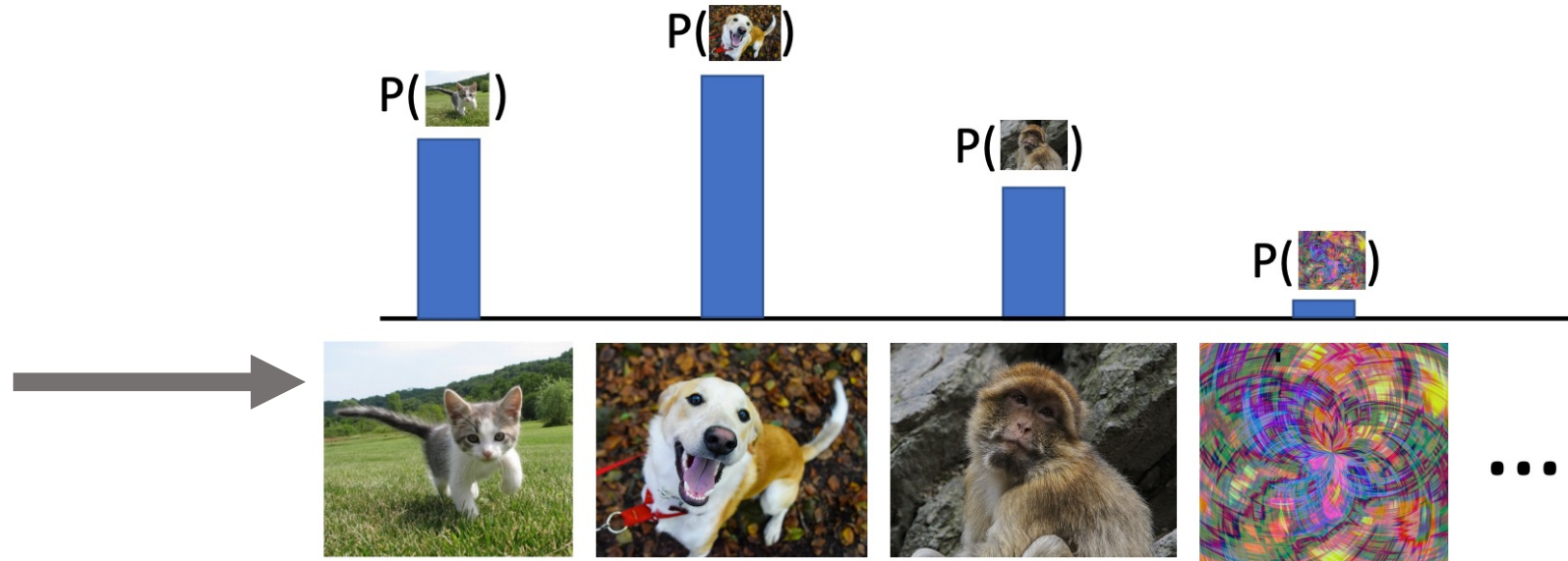


Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

Taxonomy of machine learning models

Models different probability distributions

Generative Model:
Learn $p(x)$



Generative model: All possible images compete with each other for probability mass

Model can “reject” unreasonable inputs by assigning them small values

Taxonomy of machine learning models

Models different probability distributions

Recall **Bayes' Rule**:

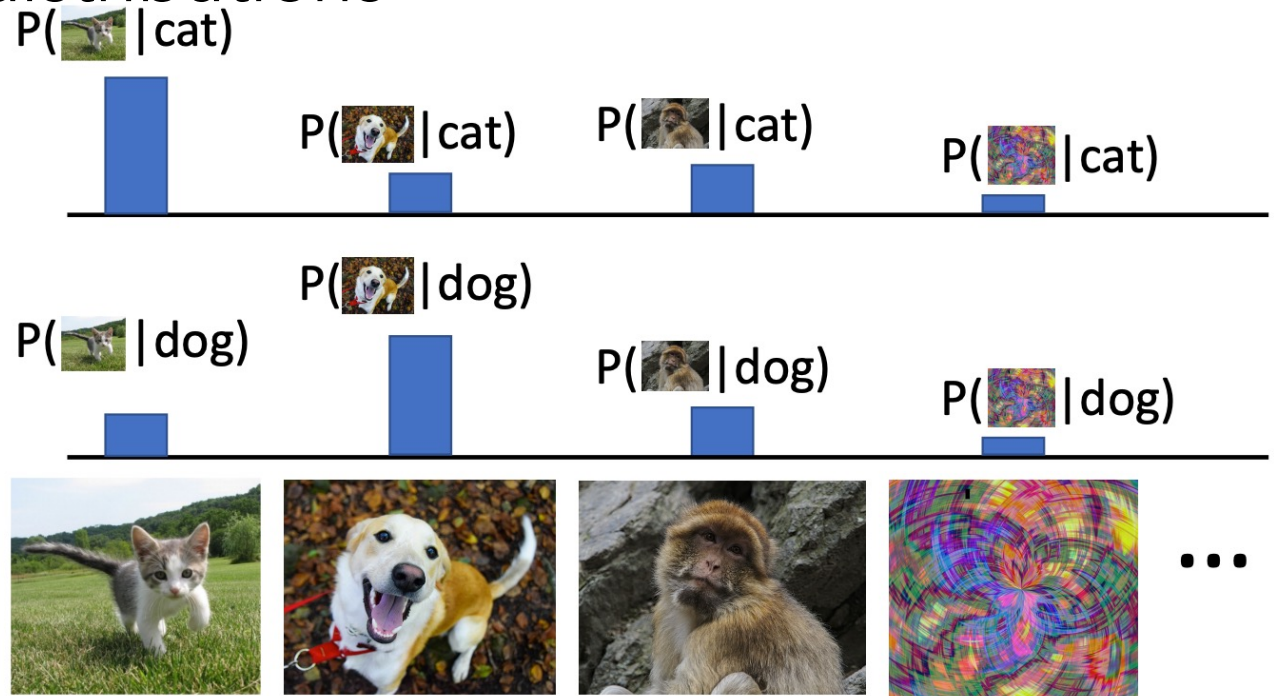
$$P(x | y) = \frac{P(y | x) P(x)}{P(y)}$$

Labels in the equation: $P(x | y)$ is a **Conditional Generative Model**; $P(y | x)$ is a **Discriminative Model**; $P(x)$ is an **(Unconditional) Generative Model**; $P(y)$ is a **Prior over labels**.

We can build a conditional generative model from other components!



Conditional Generative Model:
Learn $p(x|y)$



Conditional Generative Model: Each possible label induces a competition among all images

Taxonomy of machine learning models

Models different probability distributions

Discriminative models:

Learn $p(y|x)$



Assign labels to data

Feature learning (with labels)

Generative Model:

Learn $p(x)$



Detect outliers

Feature learning (without labels)

Sample to generate new data

Conditional Generative Model:

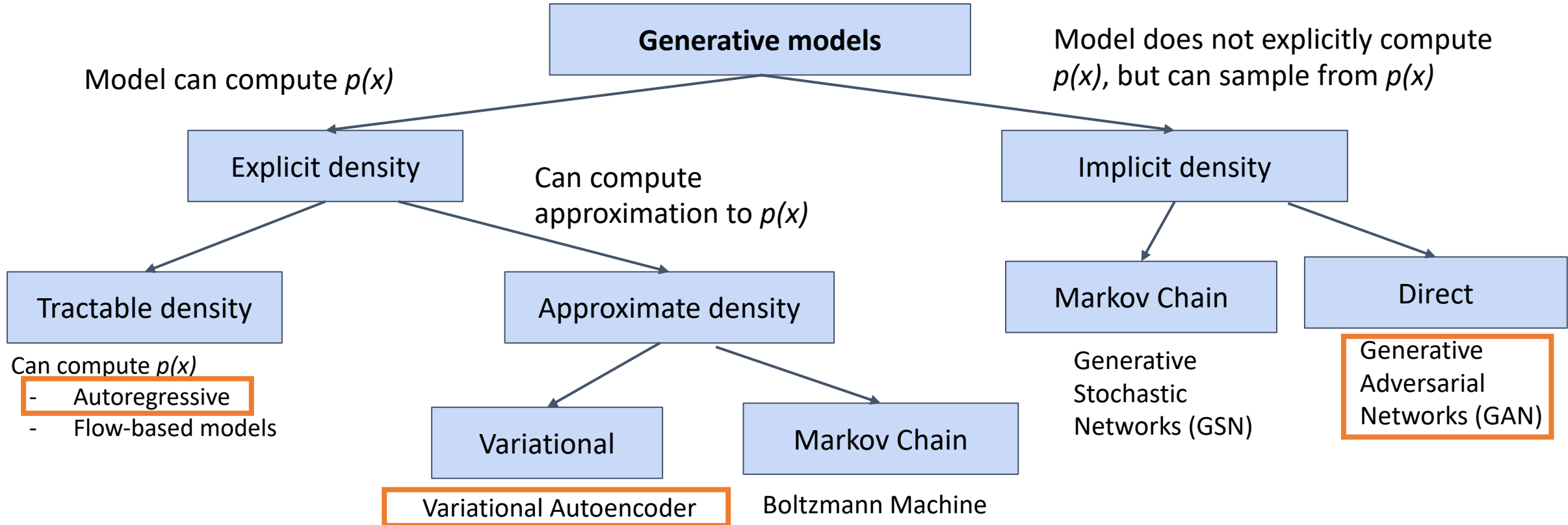
Learn $p(x|y)$



Assign labels, while rejecting outliers!

Generate new data conditioned on input labels

Taxonomy of generative models



Different types of generative models

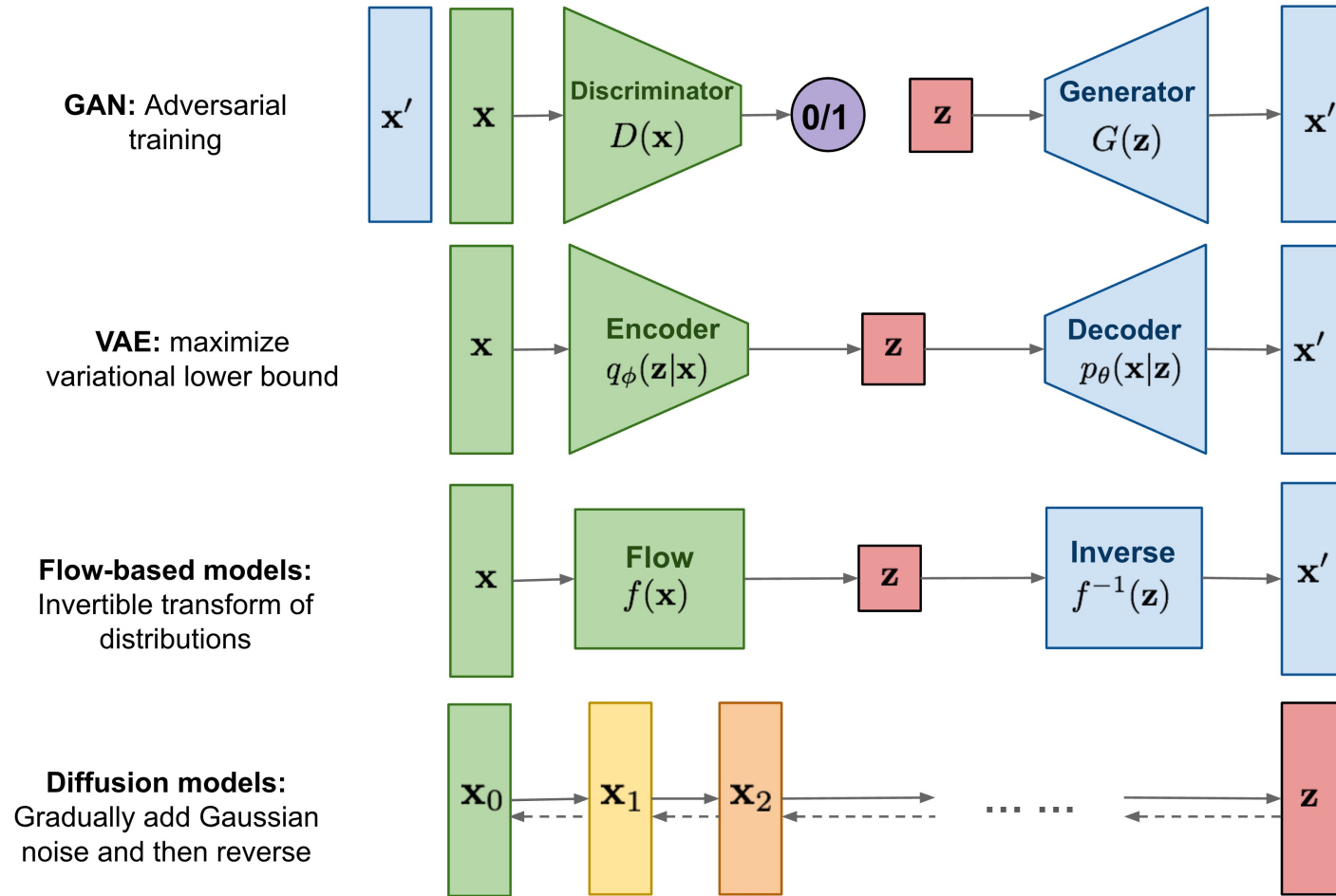
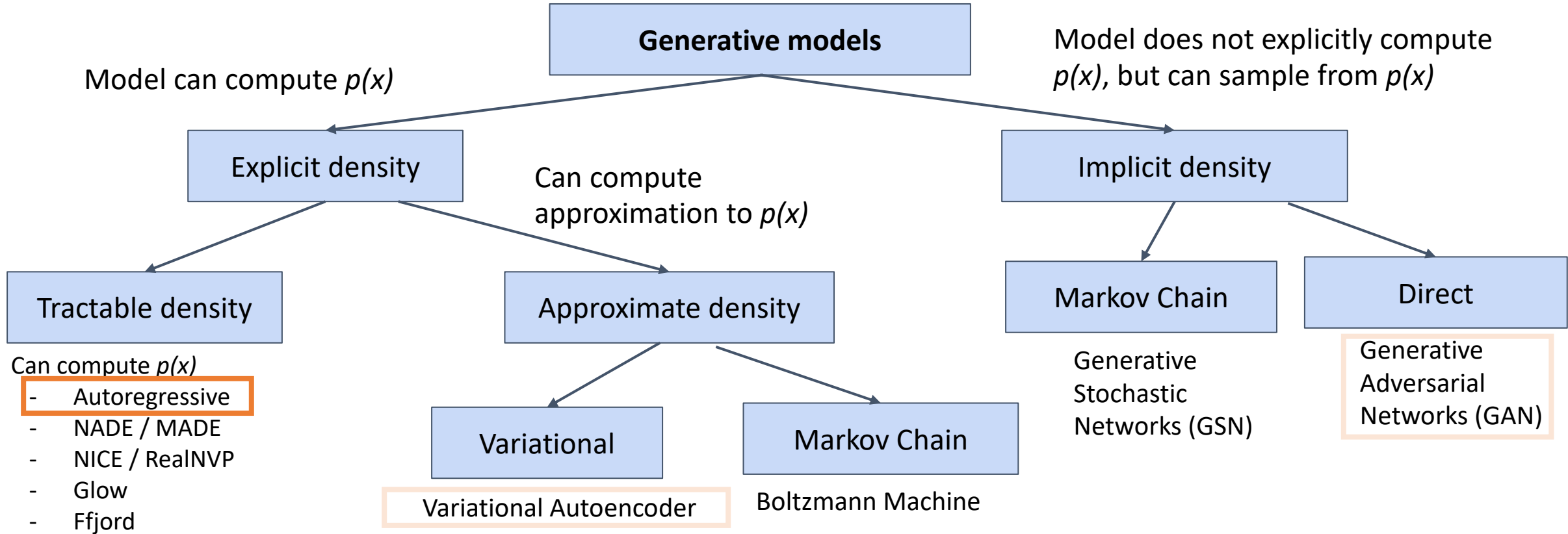


Figure credit: <https://lilianweng.github.io/lil-log/2021/07/11/diffusion-models.html>

Taxonomy of generative models



Explicit Density Estimation

Goal: Write down an explicit function for $p(x) = f(x, W)$

Given dataset $x^{(1)}, x^{(2)}, \dots, x^{(N)}$, train the model by solving:

$$\begin{aligned} W^* &= \arg \max_W \prod_i p(x^{(i)}) && \text{Maximize probability of training data} \\ &&& \text{(Maximum likelihood estimation)} \\ &= \arg \max_W \sum_i \log p(x^{(i)}) && \text{Log trick to exchange product for sum} \\ &= \arg \max_W \sum_i \log f(x^{(i)}, W) && \text{This will be our loss function!} \\ &&& \text{Train with gradient descent} \end{aligned}$$

Explicit Density: Autoregressive models

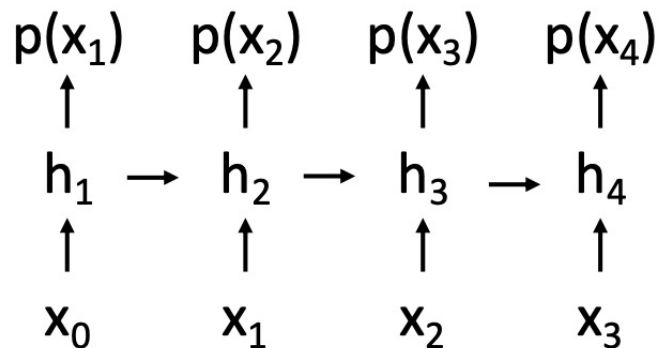
Goal: Write down an explicit function for $p(x) = f(x, W)$

Assume x consists of multiple subparts:

$$x = (x_1, x_2, x_3, \dots, x_T)$$

Break down probability using the chain rule:

$$\begin{aligned} p(x) &= p(x_1, x_2, x_3, \dots, x_T) \\ &= p(x_1)p(x_2 | x_1)p(x_3 | x_1, x_2) \dots \\ &= \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}) \end{aligned}$$



Probability of the next subpart given all the previous subparts

This is exactly what we had with the language modeling with RNNs and Transformers for captioning

PixelRNN

Generate image pixels one at a time, starting at the upper left corner

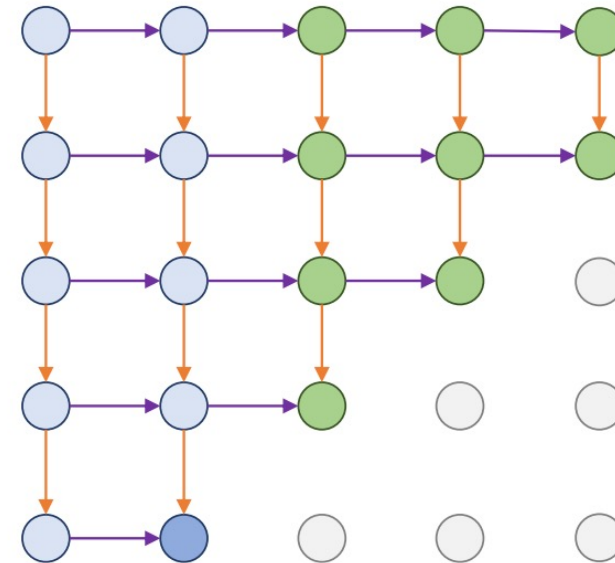
Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)

$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green: softmax over $[0, 1, \dots, 255]$

Each pixel depends **implicitly** on all pixels above and to the left:

Problem: Very slow during both training and testing; $N \times N$ image requires $2N-1$ sequential steps



Van den Oord et al, "Pixel Recurrent Neural Networks", ICML 2016

PixelCNN

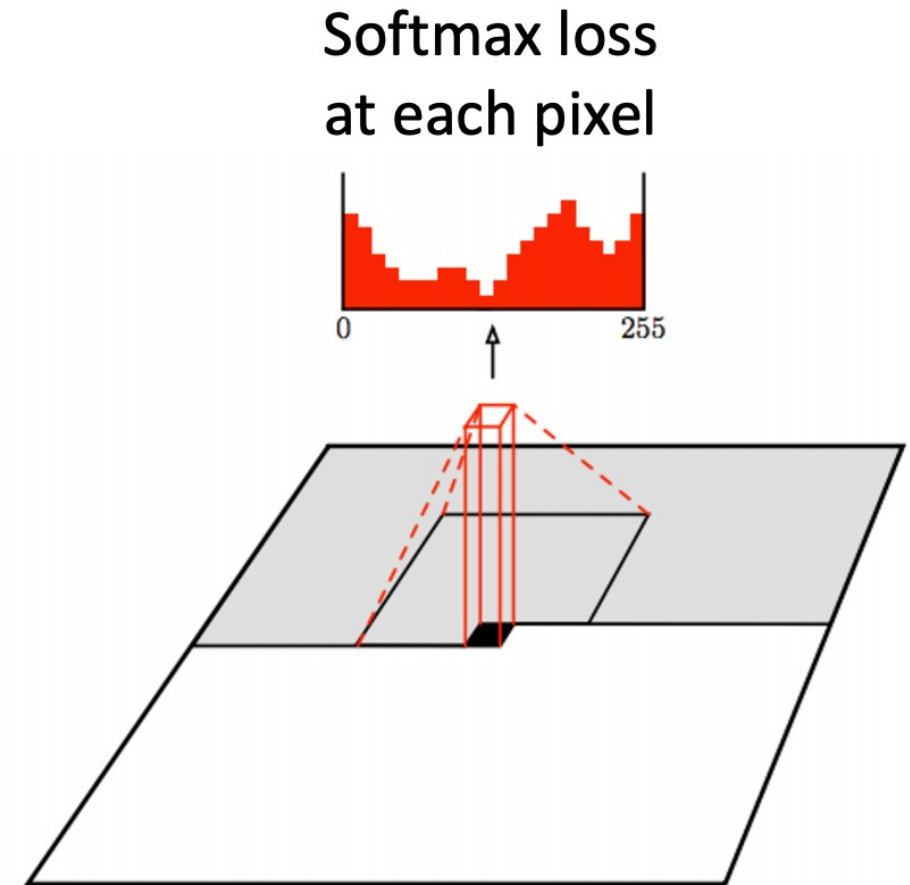
Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

Training is faster than PixelRNN
(can parallelize convolutions since context region values known from training images)

Generation must still proceed sequentially
=> still slow



Van den Oord et al, "Conditional Image Generation with PixelCNN Decoders", NeurIPS 2016

Slide credit: Justin Johnson (<https://web.eecs.umich.edu/~justincj/teaching/eecs498/FA2020/schedule.html>, L19,20)

Autoregressive models: PixelRNN and PixelCNN

Pros:

- Can explicitly compute likelihood $p(x)$
- Explicit likelihood of training data gives good evaluation metric
- Good samples

Con:

- Sequential generation => slow

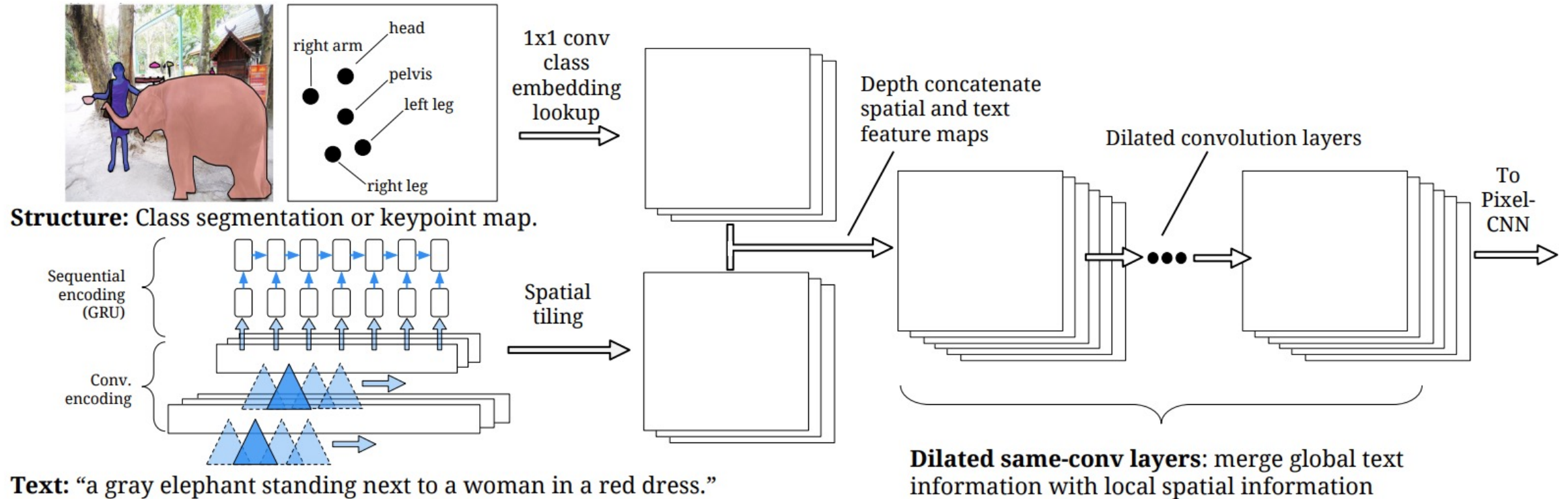
Improving PixelCNN performance

- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

See

- Van der Oord et al. NIPS 2016
- Salimans et al. 2017 (PixelCNN++)

Text-based image generation with PixelCNN



Text + segmentations

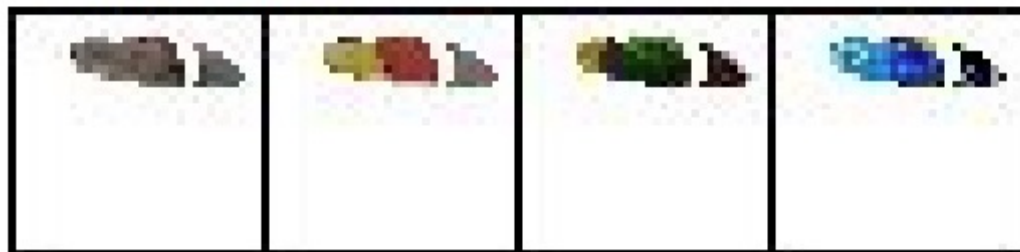
A person carrying their surfboard while walking along a beach.



Person



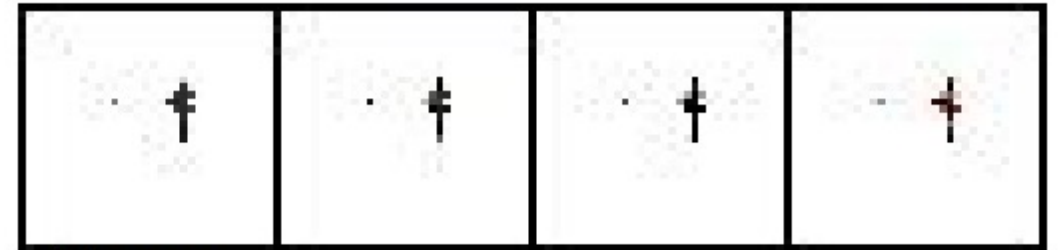
Surfboard



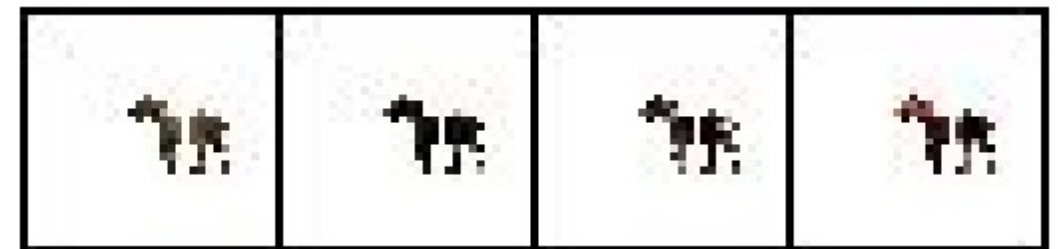
The woman is riding her horse on the beach by the water.



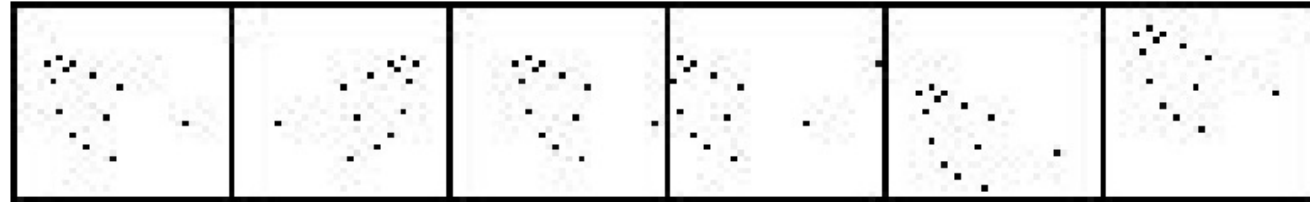
Person



Horse



Text + keypoints



This bird is bright yellow.



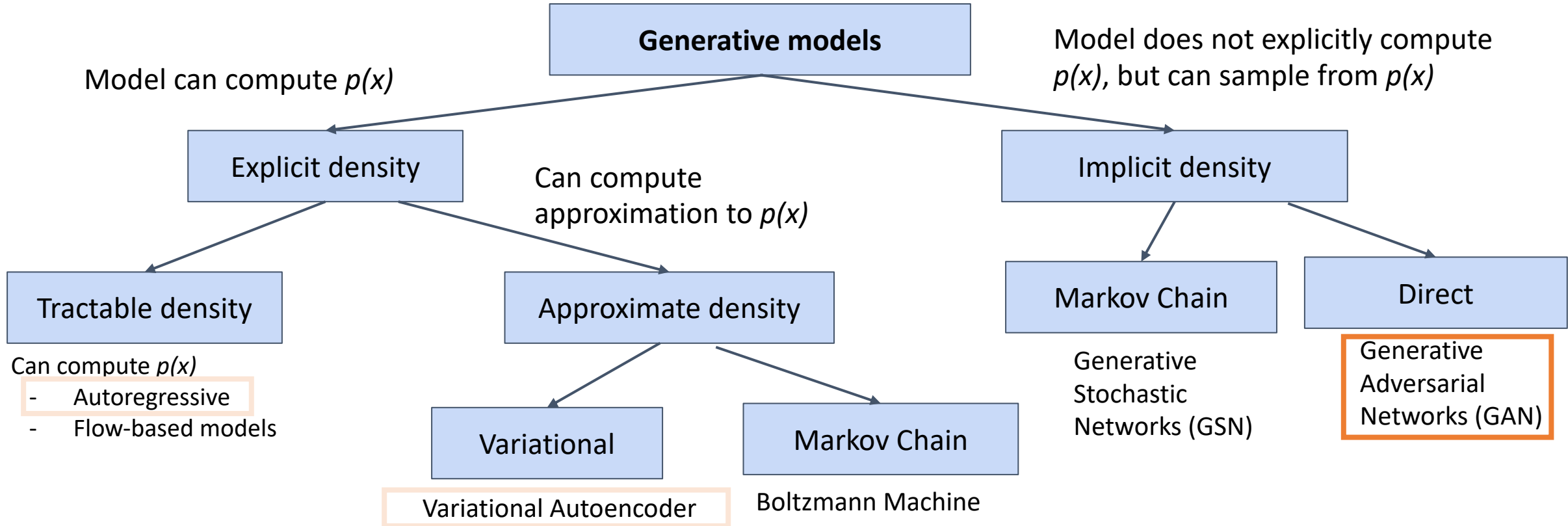
This bird is bright red.



This bird is bright blue.



Taxonomy of generative models



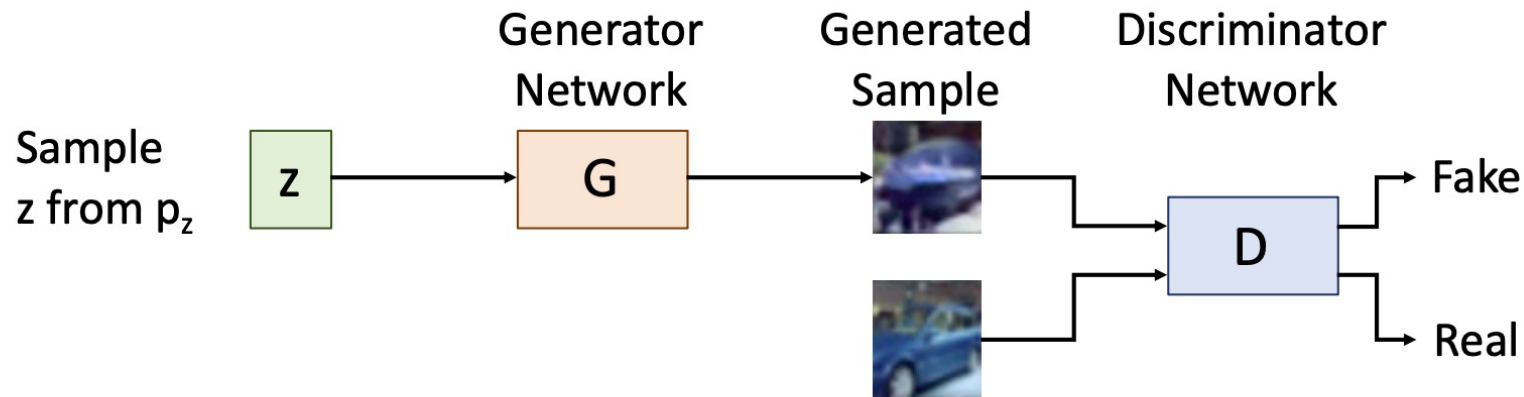
Generative Adversarial Networks (GAN)

Jointly train generator G and discriminator D with a **minimax game**

Discriminator wants
 $D(x) = 1$ for real data

Discriminator wants
 $D(x) = 0$ for fake data

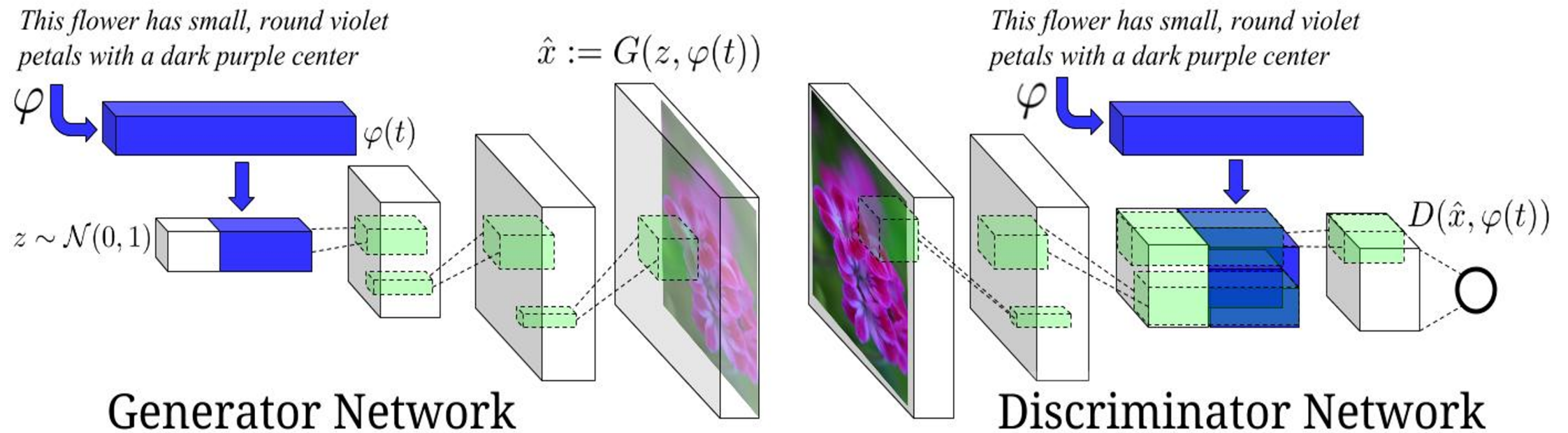
$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$



Generator wants
 $D(x) = 1$ for fake data

Text to image with GANs

- Generator and Discriminator are alternately trained



Text to image with GANs

- Image encoder (CNN ϕ) and text encoder (char-CNN-RNN φ) are pre-trained to produce a joint embedding where the embedded representations can be used to predict the class label of the image

$$\frac{1}{N} \sum_{n=1}^N \Delta(y_n, f_v(v_n)) + \Delta(y_n, f_t(t_n))$$

$$f_v(v) = \arg \max_{y \in \mathcal{Y}} \mathbb{E}_{t \sim \mathcal{T}(y)} [\phi(v)^T \varphi(t)]$$

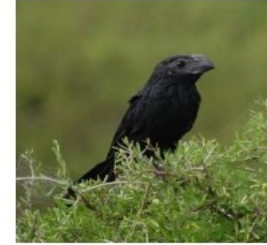
$$f_t(t) = \arg \max_{y \in \mathcal{Y}} \mathbb{E}_{v \sim \mathcal{V}(y)} [\phi(v)^T \varphi(t)]$$

Datasets

- CUB-200 (Birds)
 - 11,788 images of birds from 200 categories
- Oxford-102 (Flowers)
 - 8,189 images of flowers from 102 categories
- MSCOCO
 - 330K images
- 5 captions per image

Caltech-UCSD-Birds (CUB) 200

an all black bird with a distinct thick, rounded bill.



this small bird has a yellow breast, brown crown, and black superciliary



this flower is white and pink in color, with petals that have veins.



bright droopy yellow petals with burgundy streaks, and a yellow stigma.



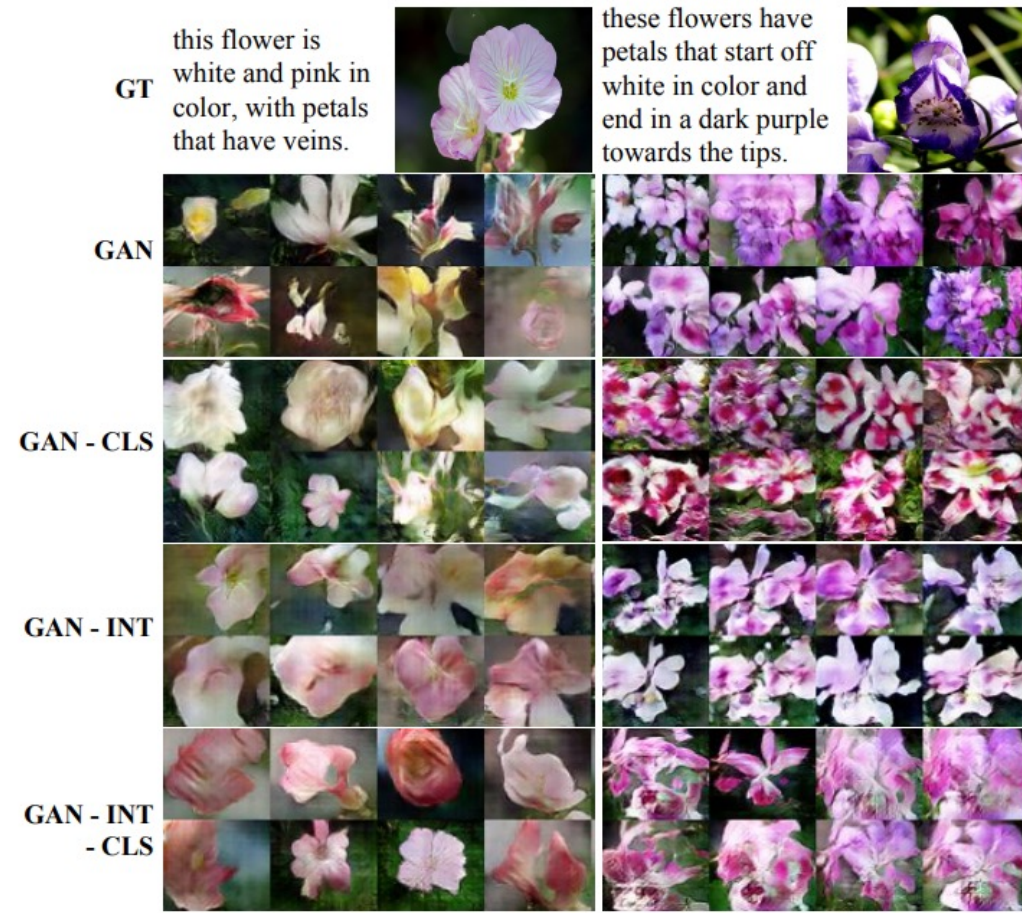
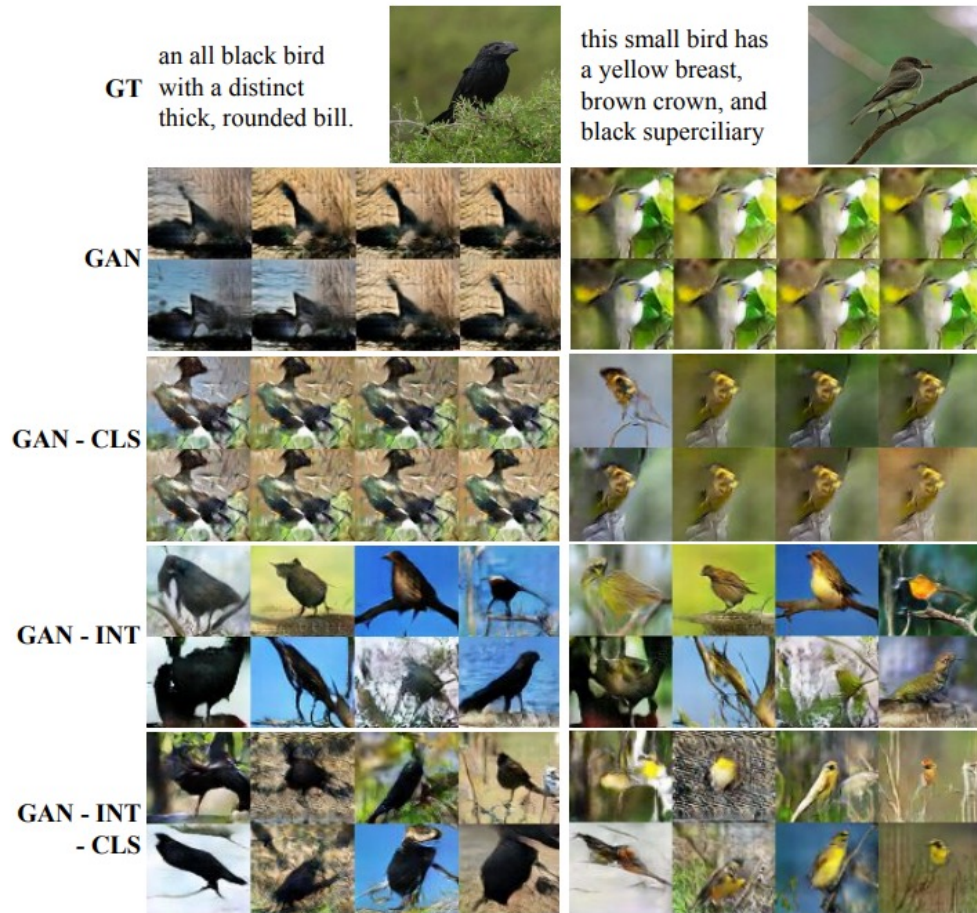
The man at bat readies to swing at the pitch while the umpire looks on.



















Bunk bed with a narrow shelf sitting underneath it.

Text to image with GANs: Results

- CLS: Add discriminator to distinguish if (image, text) match or not
(real image, right text), (real image, wrong text), (fake image, right text)
- INT: Add interpolated text embeddings (fake additional text embeddings)

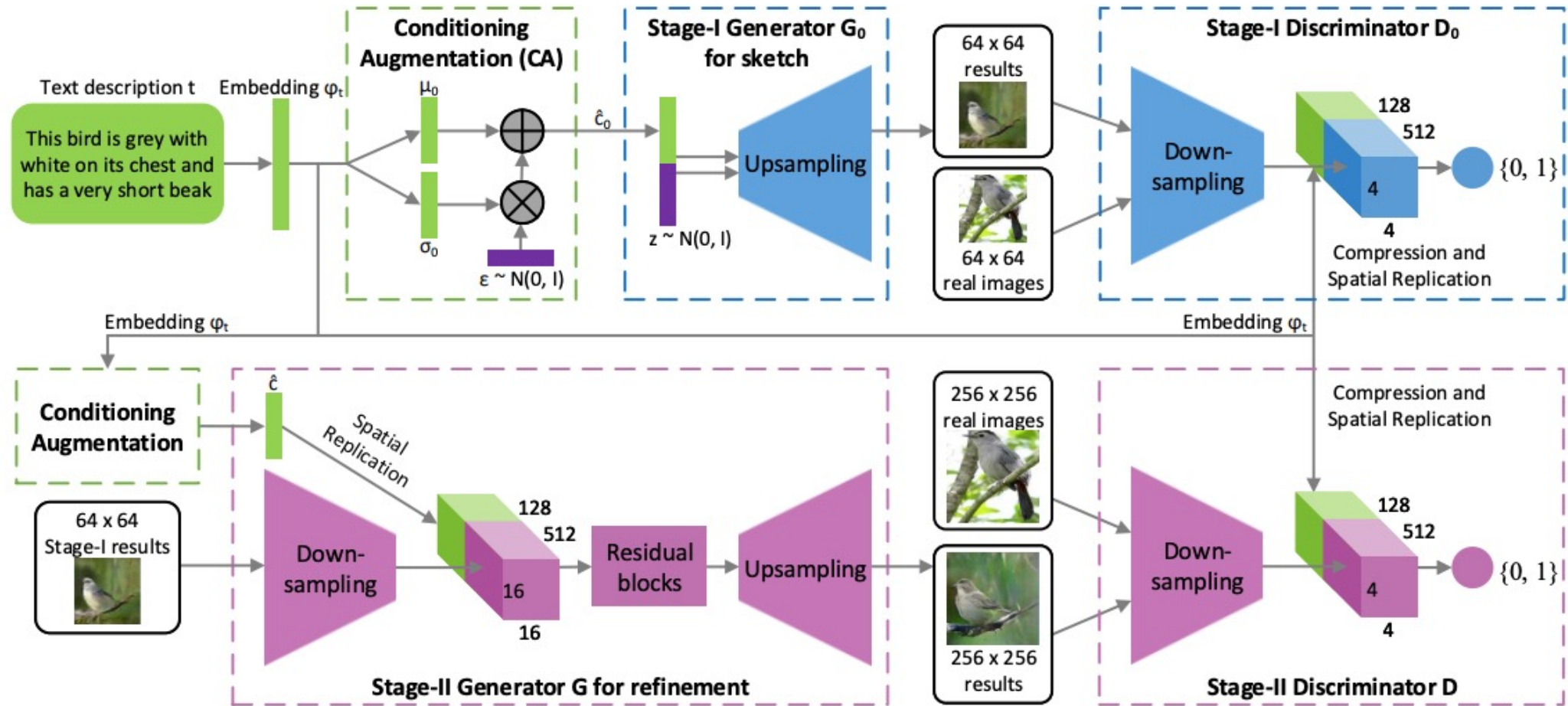


Text to image with GANS: Results

	GT	Ours		GT	Ours	
a group of people on skis stand on the snow.			a man in a wet suit riding a surfboard on a wave.			Very low res! 64 x 64 Follow up work: 128 x 128 Still low res!
a table with many plates of food and drinks			two plates of food that include beans, guacamole and rice.			
two giraffe standing next to each other in a forest.			a green plant that is growing out of the ground.			
a large blue octopus kite flies above the people having fun at the beach.			there is only one horse in the grassy field.			



StackGAN

Generate low resolution, and then pass through another GAN for improved resolution



















StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks
<https://arxiv.org/pdf/1612.03242.pdf>, Zhang et al, ICCV 2017

StackGAN: Results

Text description	This bird is red and brown in color, with a stubby beak	The bird is short and stubby with yellow on its body	A bird with a medium orange bill white body gray wings and webbed feet	This small black bird has a short, slightly curved bill and long legs	A small bird with varying shades of brown with white under the eyes	A small yellow bird with a black crown and a short black pointed beak	This small bird has a white breast, light grey head, and black wings and tail
64x64 GAN-INT-CLS							
128x128 GAWWN							
256x256 StackGAN							

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks
<https://arxiv.org/pdf/1612.03242.pdf>, Zhang et al, ICCV 2017

StackGAN: Results

Text description	This flower has a lot of small purple petals in a dome-like configuration	This flower is pink, white, and yellow in color, and has petals that are striped	This flower has petals that are dark pink with white edges and pink stamen	This flower is white and yellow in color, with petals that are wavy and smooth	A picture of a very clean living room	A group of people on skis stand in the snow	Eggs fruit candy nuts and meat served on white dish	A street sign on a stoplight pole in the middle of a day
64x64 GAN-INT-CLS								
256x256 StackGAN								

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks
<https://arxiv.org/pdf/1612.03242.pdf>, Zhang et al, ICCV 2017

StackGAN: Evaluation

- Inception Score: $I = \exp(\mathbb{E}_{\mathbf{x}} D_{KL}(p(y|\mathbf{x}) || p(y)))$
 - Use inception model to predict class y
 - Want good models to generate diverse but meaningful images
 - Large distance between marginal prior (of labels) and conditional prior
- Human rank images generated by models

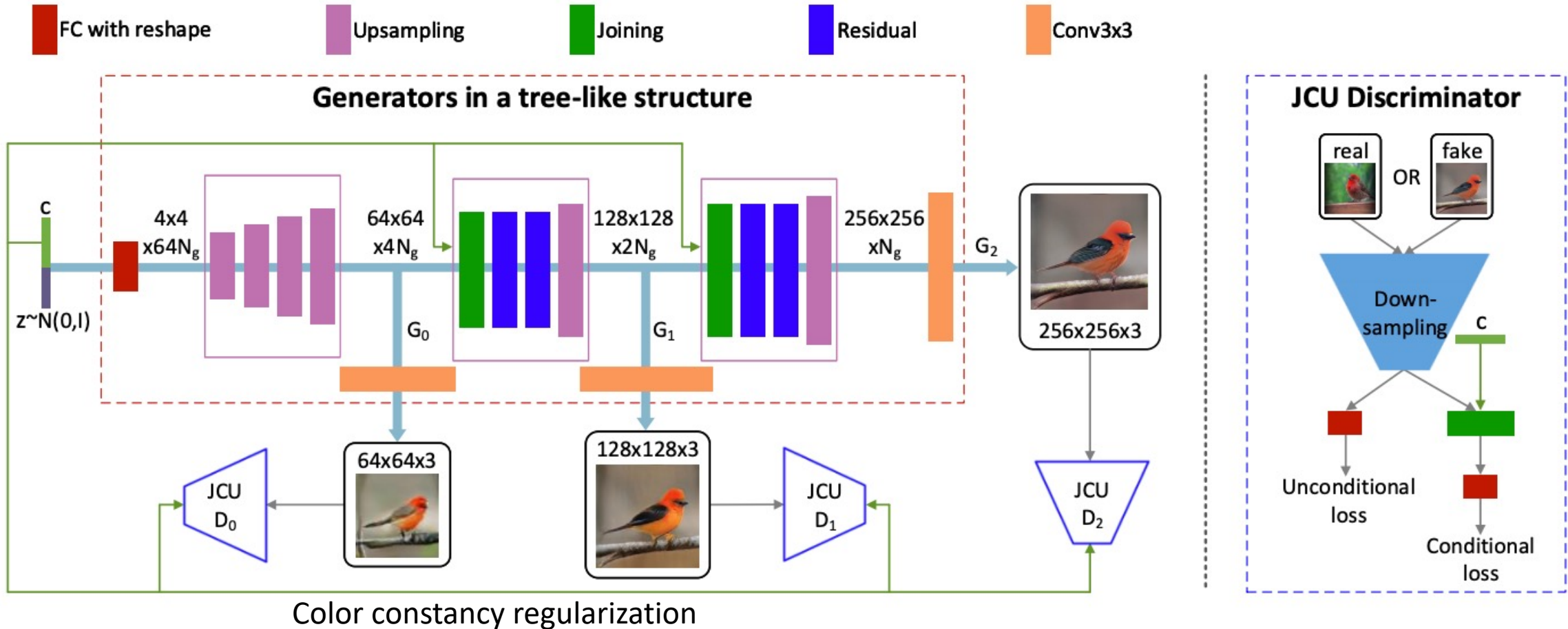
Metric	Dataset	GAN-INT-CLS	GAWWN	Our StackGAN
Inception score	CUB	2.88 ± .04	3.62 ± .07	3.70 ± .04
	Oxford	2.66 ± .03	/	3.20 ± .01
	COCO	7.88 ± .07	/	8.45 ± .03
Human rank	CUB	2.81 ± .03	1.99 ± .04	1.37 ± .02
	Oxford	1.87 ± .03	/	1.13 ± .03
	COCO	1.89 ± .04	/	1.11 ± .03

StackGAN++

Generalization of StackGAN (Multiscale)

Joint Discriminator

- if image is real/fake (unconditional loss)
- if text+image match (conditional loss)



StackGAN++: Realistic Image Synthesis with Stacked Generative Adversarial Networks

<https://arxiv.org/pdf/1710.10916.pdf>, Zhang et al, TPAMI 2018

StackGAN++

- Generalization of StackGAN (arbitrary number of Generators/Discriminators)
- Color constancy regularization
- Joint Discriminator (similar to +CLS from Reed et al)
 - if image is real/fake (unconditional loss)
 - if text+image match (conditional loss)
- Alternately train generator and discriminator

$$\mathcal{L}_{G_i} = \underbrace{-\frac{1}{2}\mathbb{E}_{\hat{x}_i \sim p_{G_i}} [\log(D_i(\hat{x}_i))]}_{\text{unconditional loss}} - \underbrace{\frac{1}{2}\mathbb{E}_{\hat{x}_i \sim p_{G_i}} [\log(D_i(\hat{x}_i, \bar{e}))]}_{\text{conditional loss}},$$

$$\mathcal{L}_{D_i} = \underbrace{-\frac{1}{2}\mathbb{E}_{x_i \sim p_{data_i}} [\log D_i(x_i)] - \frac{1}{2}\mathbb{E}_{\hat{x}_i \sim p_{G_i}} [\log(1 - D_i(\hat{x}_i))] +}_{\text{unconditional loss}} \underbrace{-\frac{1}{2}\mathbb{E}_{x_i \sim p_{data_i}} [\log D_i(x_i, \bar{e})] - \frac{1}{2}\mathbb{E}_{\hat{x}_i \sim p_{G_i}} [\log(1 - D_i(\hat{x}_i, \bar{e}))]}_{\text{conditional loss}},$$

StackGAN++: Realistic Image Synthesis with Stacked Generative Adversarial Networks

<https://arxiv.org/pdf/1710.10916.pdf>, Zhang et al, TPAMI 2018

StackGAN++: Results

- FID (Frechet Inception distance): measures distance between generated and real distribution

Metric	CUB			Oxford		COCO	
	GAN-INT-CLS	GAWWN	Our StackGAN-v1	GAN-INT-CLS	Our StackGAN-v1	GAN-INT-CLS	Our StackGAN-v1
FID ↓	68.79	67.22	51.89	79.55	55.28	60.62	74.05
FID* ↓	68.79	53.51	35.11	79.55	43.02	60.62	33.88
IS ↑	2.88 ± .04	3.62 ± .07	3.70 ± .04	2.66 ± .03	3.20 ± .01	7.88 ± .07	8.45 ± .03
IS* ↑	2.88 ± .04	3.10 ± .03	3.02 ± .03	2.66 ± .03	2.73 ± .03	7.88 ± .07	8.35 ± .11
HR ↓	2.76 ± .01	1.95 ± .02	1.29 ± .02	1.84 ± .02	1.16 ± .02	1.82 ± .03	1.18 ± .03

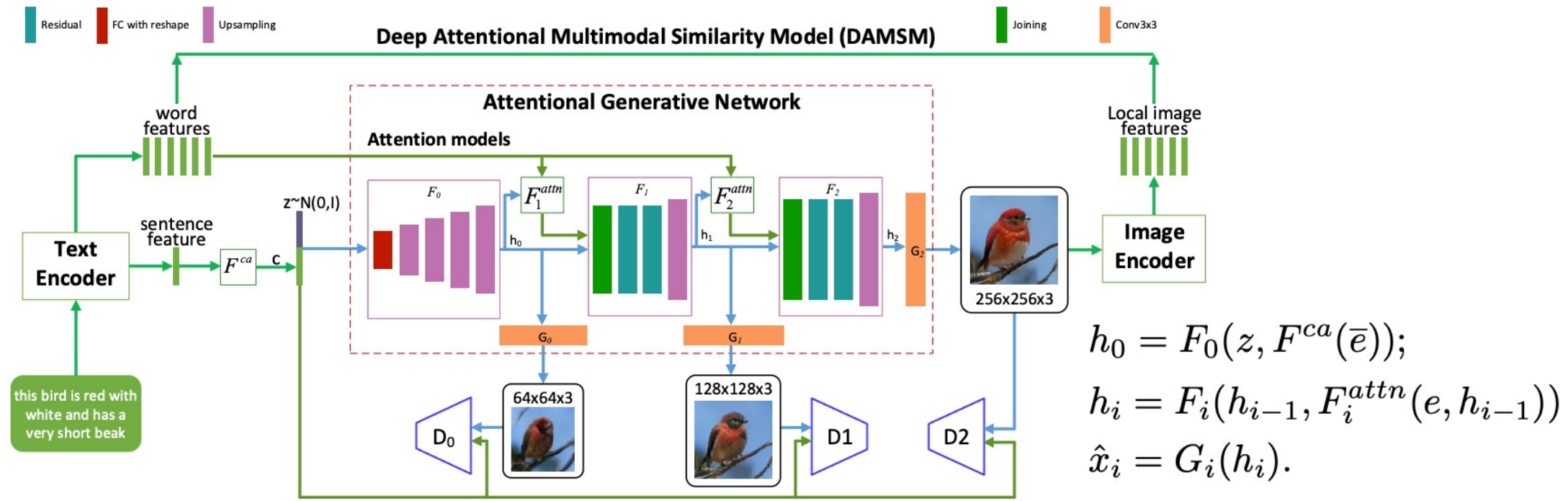
Dataset		CUB	Oxford-102	COCO	LSUN-bedroom	LSUN-church	ImageNet-dog	ImageNet-cat
FID ↓	StackGAN-v1	51.89	55.28	74.05	91.94	57.20	89.21	58.73
	StackGAN-v2	15.30	48.68	81.59	35.61	25.36	44.54	28.59
IS ↑	StackGAN-v1	3.70 ± .04	3.20 ± .01	8.45 ± .03	3.59 ± .05	2.87 ± .05	8.84 ± .08	4.77 ± .06
	StackGAN-v2	4.04 ± .05	3.26 ± .01	8.30 ± .10	3.02 ± .04	2.38 ± .03	9.55 ± .11	4.23 ± .05
HR ↓	StackGAN-v1	1.81 ± .02	1.70 ± .03	1.45 ± .04	1.95 ± .01	1.86 ± .02	1.90 ± .01	1.88 ± .02
	StackGAN-v2	1.19 ± .02	1.30 ± .03	1.55 ± .05	1.05 ± .01	1.14 ± .02	1.10 ± .01	1.12 ± .02

StackGAN++: Realistic Image Synthesis with Stacked Generative Adversarial Networks

<https://arxiv.org/pdf/1710.10916.pdf>, Zhang et al, TPAMI 2018

AttnGAN

- Attention based similarity matching of image and text that tries to align regions of the image to words in the text
- m generators (G_i), each taking hidden state h_i to produce image \hat{x}_i



AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks
<https://arxiv.org/pdf/1711.10485.pdf>, Xu et al, CVPR 2018

AttnGAN:

- Attention based similarity matching of image and text that tries to align regions of the image to words in the text

- m generators (G_i), each taking hidden state h_i to produce image \hat{x}_i

- Total Loss:
$$\mathcal{L} = \mathcal{L}_G + \lambda \mathcal{L}_{DAMSM}, \text{ where } \mathcal{L}_G = \sum_{i=0}^{m-1} \mathcal{L}_{G_i}.$$

- Main contribution:

- Semi-supervised training to match image regions to text
- Attention-based match score $R(Q, D)$ of image (Q) to text (D) based on attention-based match of words to regions in the image

- Train to optimize match based on words (w) and sentences (s)

- Estimate probability of text given image and vice versa
$$\mathcal{L}_{DAMSM} = \mathcal{L}_1^w + \mathcal{L}_2^w + \mathcal{L}_1^s + \mathcal{L}_2^s.$$

$$P(D_i|Q_i) = \frac{\exp(\gamma_3 R(Q_i, D_i))}{\sum_{j=1}^M \exp(\gamma_3 R(Q_i, D_j))} \quad \mathcal{L}_1^w = - \sum_{i=1}^M \log P(D_i|Q_i), \quad \mathcal{L}_2^w = - \sum_{i=1}^M \log P(Q_i|D_i),$$

AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks

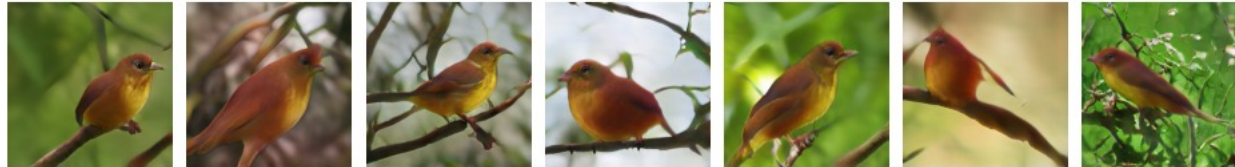
<https://arxiv.org/pdf/1711.10485.pdf>, Xu et al, CVPR 2018

AttnGAN: Results

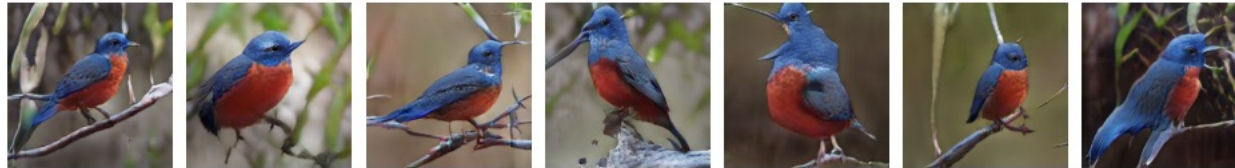
this bird has wings that are **black** and has a **white** belly



this bird has wings that are **red** and has a **yellow** belly



this bird has wings that are **blue** and has a **red** belly



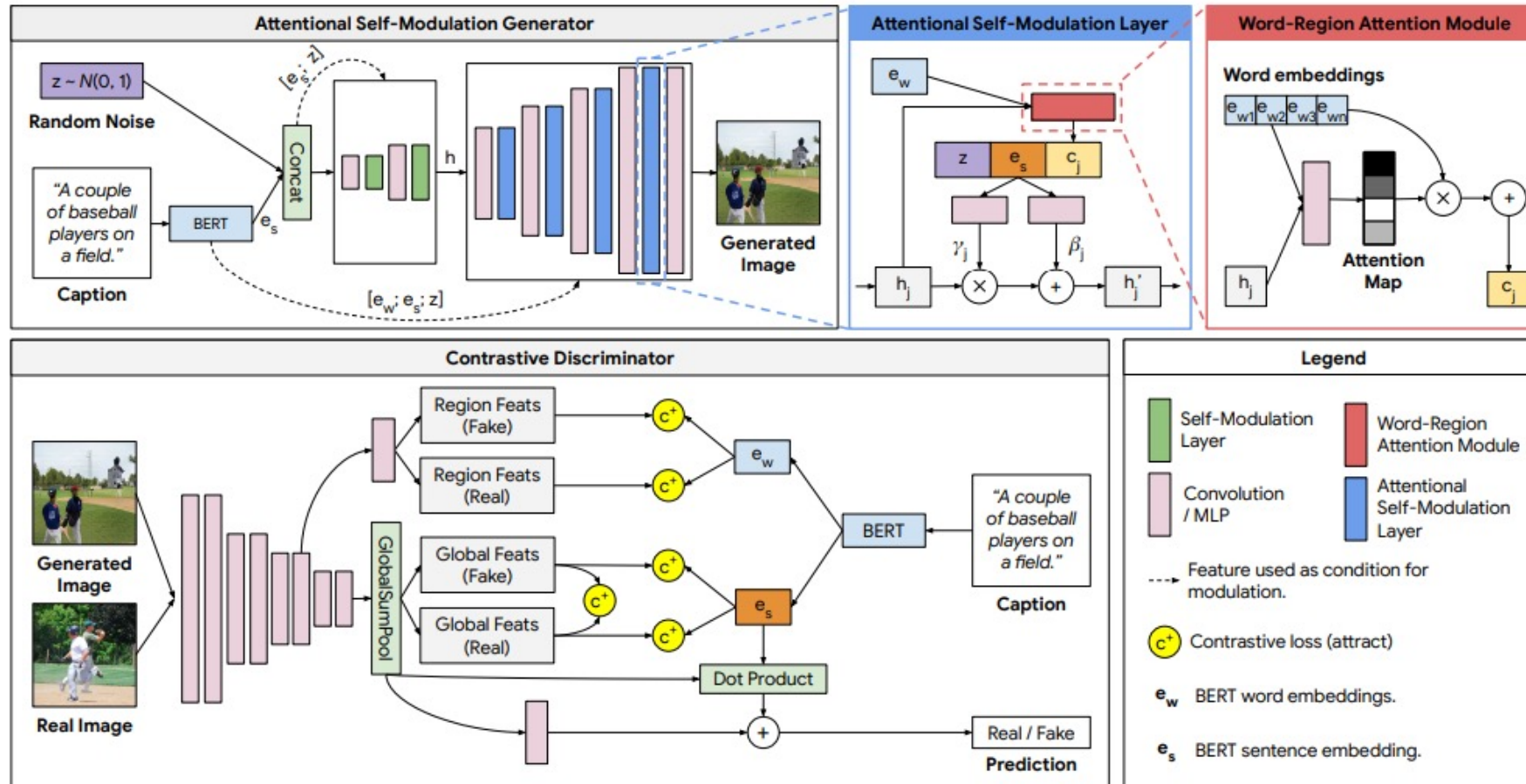
Inception Scores

Dataset	GAN-INT-CLS [20]	GAWWN [18]	StackGAN [31]	StackGAN-v2 [32]	PPGN [16]	Our AttnGAN
CUB	$2.88 \pm .04$	$3.62 \pm .07$	$3.70 \pm .04$	$3.82 \pm .06$	/	$4.36 \pm .03$
COCO	$7.88 \pm .07$	/	$8.45 \pm .03$	/	$9.58 \pm .21$	$25.89 \pm .47$

AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks
<https://arxiv.org/pdf/1711.10485.pdf>, Xu et al, CVPR 2018

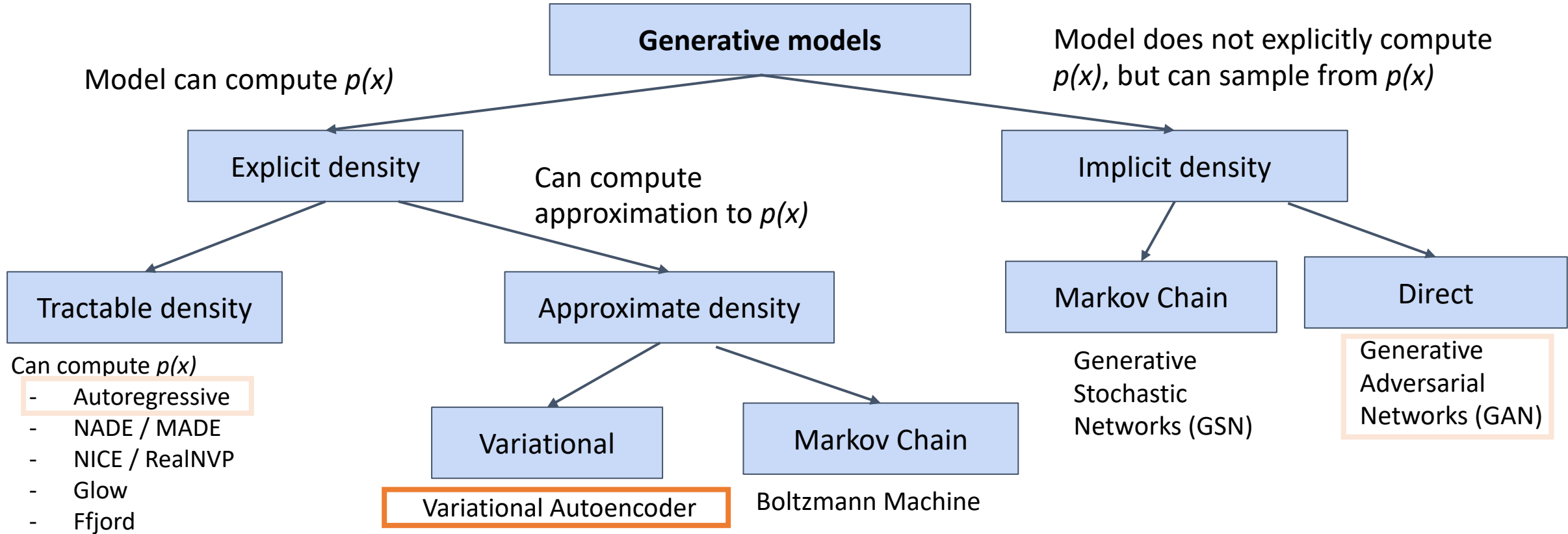
Cross-Modal Contrastive Learning

- GAN with contrastive losses



Cross-Modal Contrastive Learning for Text-to-Image Generation
<https://arxiv.org/pdf/2101.04702.pdf>, Zhang et al, CVPR 2021

Taxonomy of generative models



Variational Autoencoders

- PixelRNN/PixelCNN explicitly parameterizes density function with a neural network, so we can train to maximize likelihood of training data

$$p_{\theta}(x) = \prod_{t=1}^T p_{\theta}(x_t | x_1, \dots, x_{t-1})$$

Assume data can be broken into subparts!

What if we don't make this assumption?

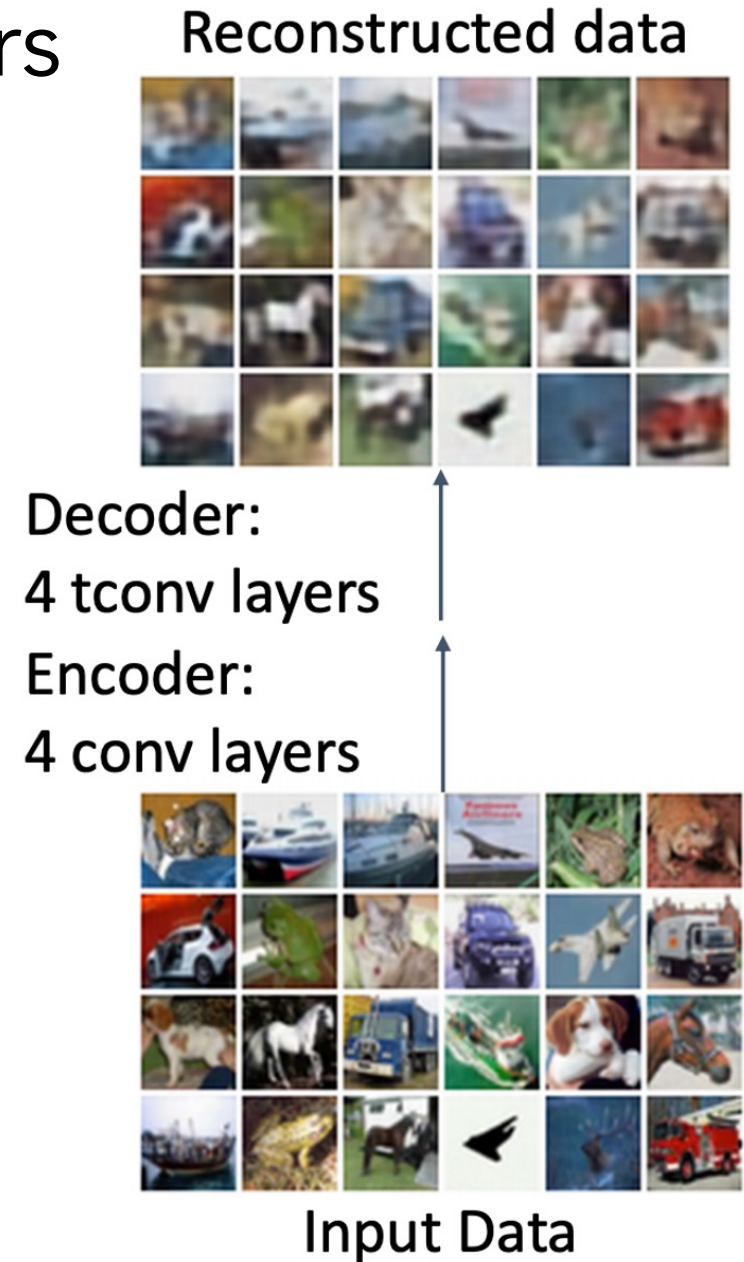
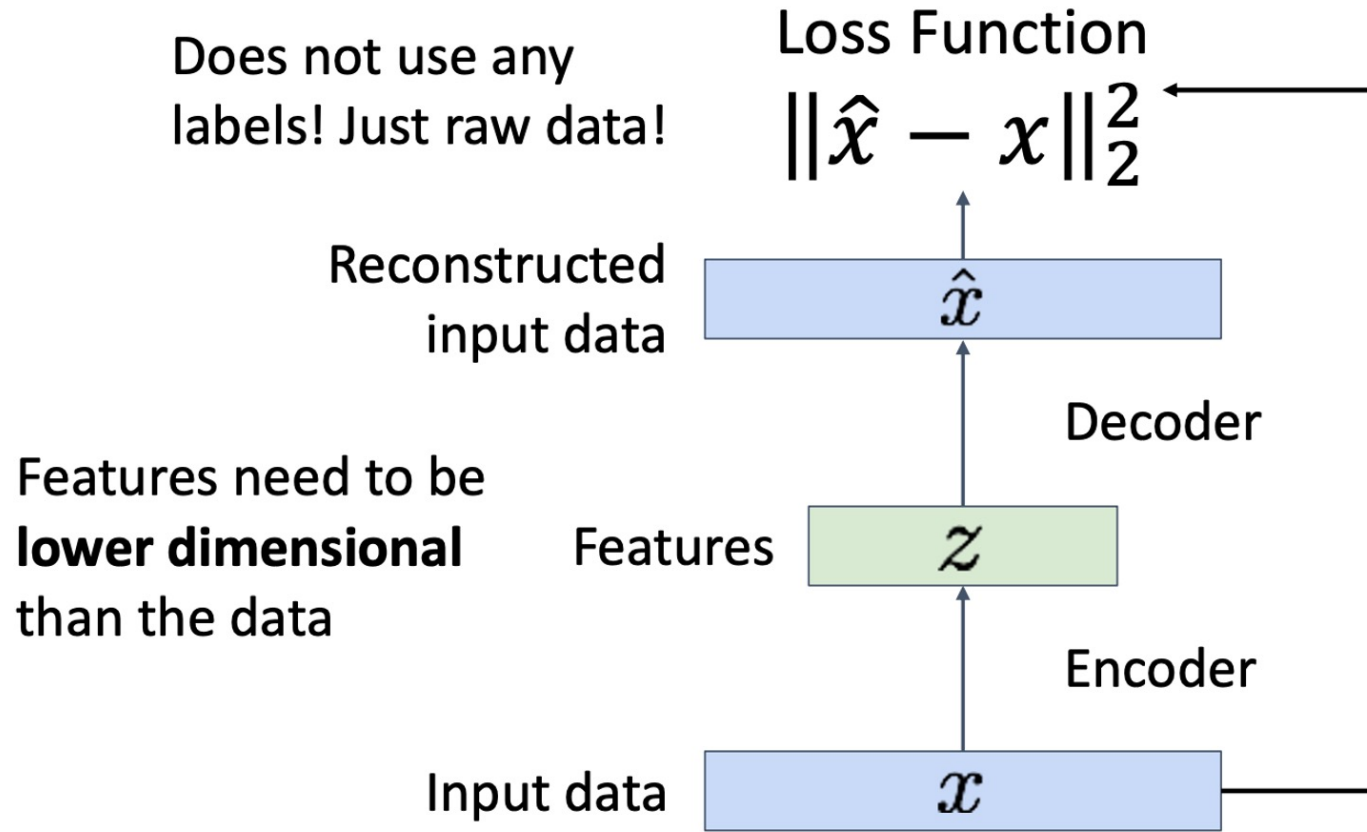
- Variational Autoencoders (VAE) use an **intractable density** that we cannot explicitly compute or optimize



- But we will be able to directly **optimize a lower bound** on the density

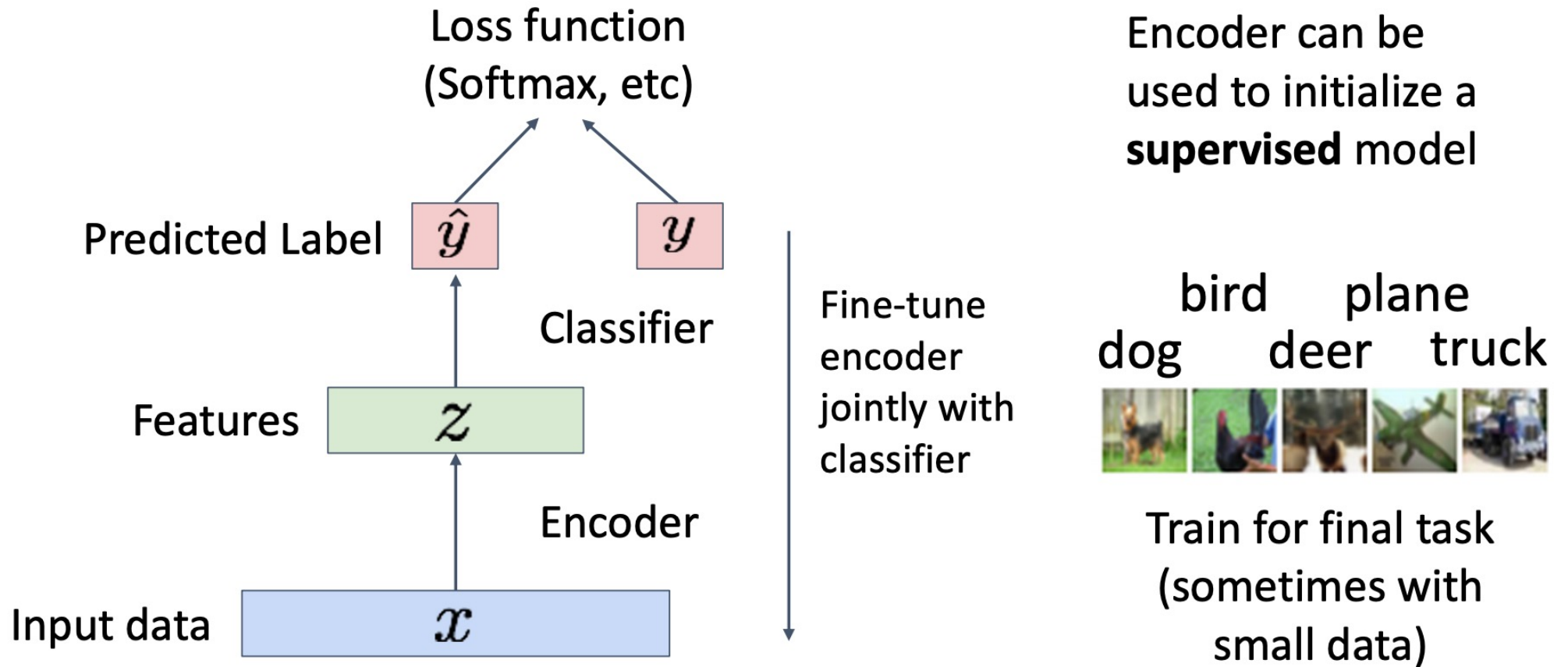
(Regular, non-variational) Autoencoders

Loss: L2 distance between input and reconstructed data.



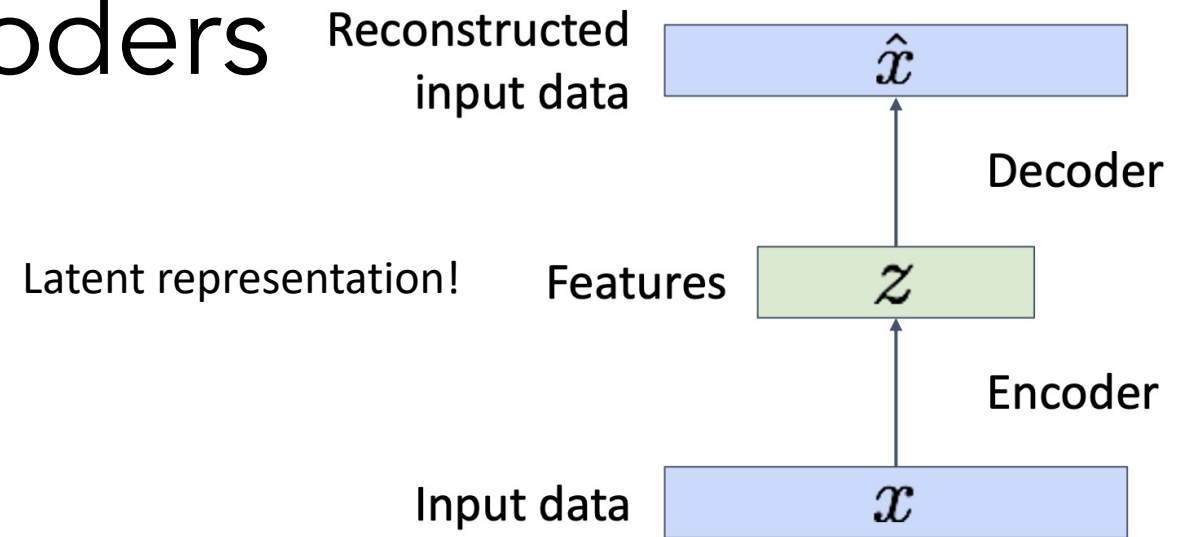
(Regular, non-variational) Autoencoders

After training, **throw away decoder** and use encoder for a downstream task

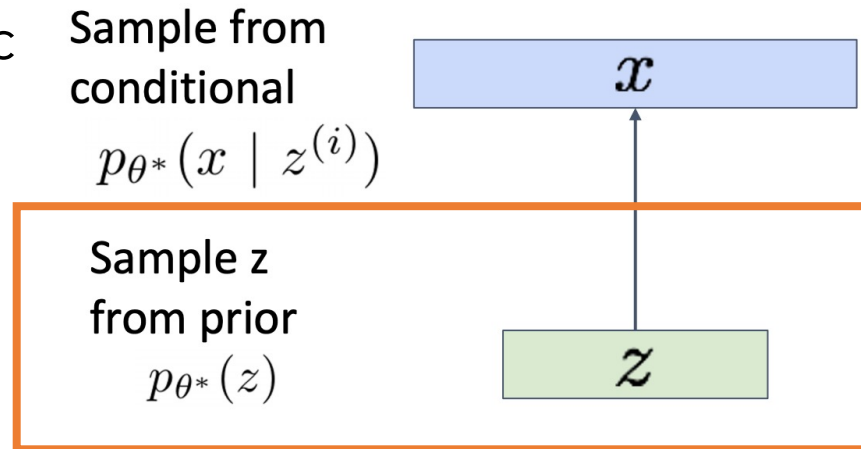


Variational Autoencoders

- Autoencoders
 - Not probabilistic
 - No sampling



- Variational
 - Probabilistic



How to sample?

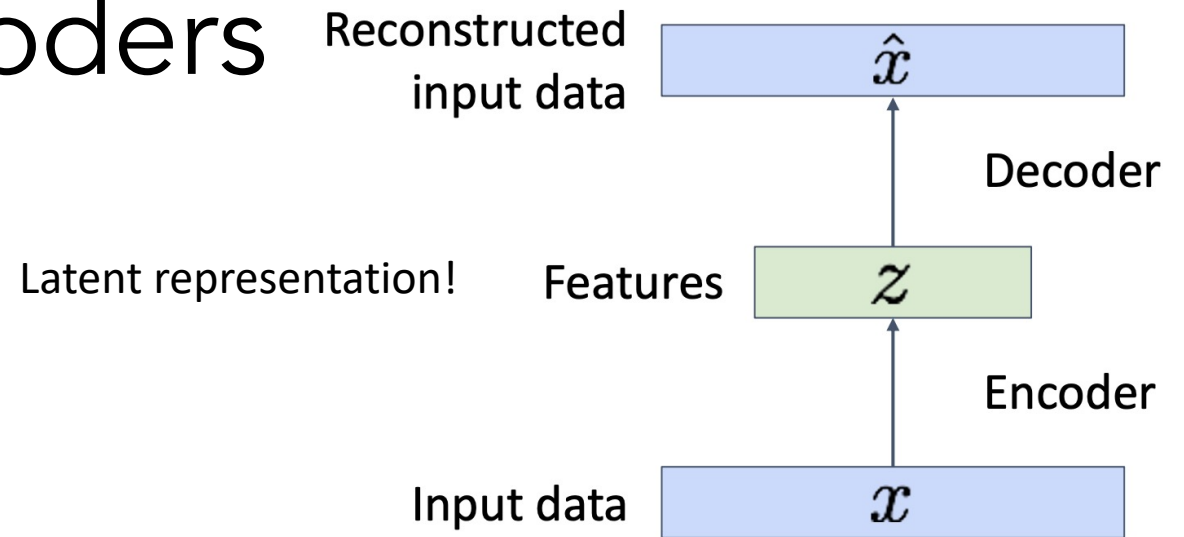
Assume simple prior $p(z)$, e.g. Gaussian with mean

Assume z is **latent representation** that we can **sample** from to generate image x .

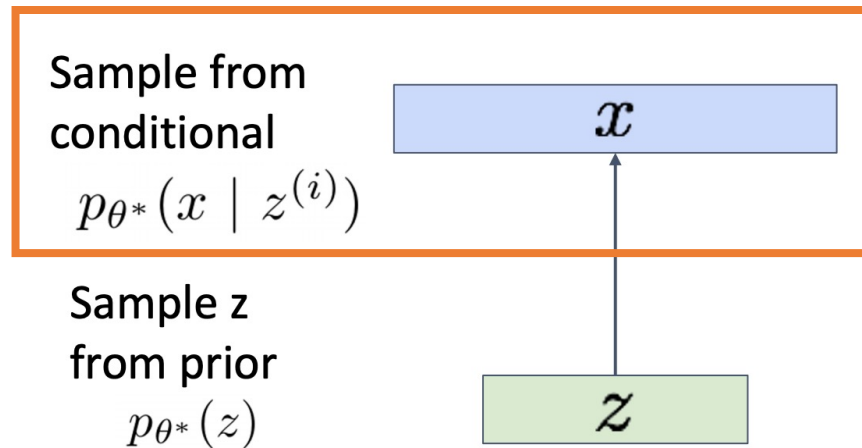
1. Learn latent representation
2. Sample to generate images

Variational Autoencoders

- Autoencoders
 - Not probabilistic
 - No sampling



- Variational
 - Probabilistic



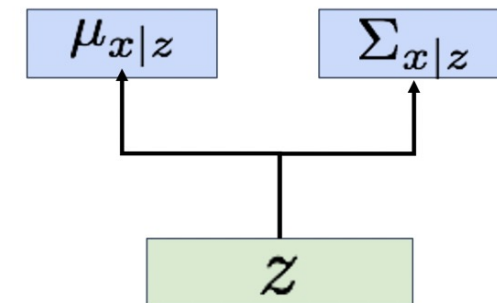
How to sample?

Assume simple prior $p(z)$, e.g. Gaussian with mean

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Decoder Network

$$p_{\theta}(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



Variational Autoencoders

- Let's maximize the likelihood of data! Need to compute $p_{\theta}(x)$

Marginalize?

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int p_{\theta}(x|z)p_{\theta}(z) dz$$

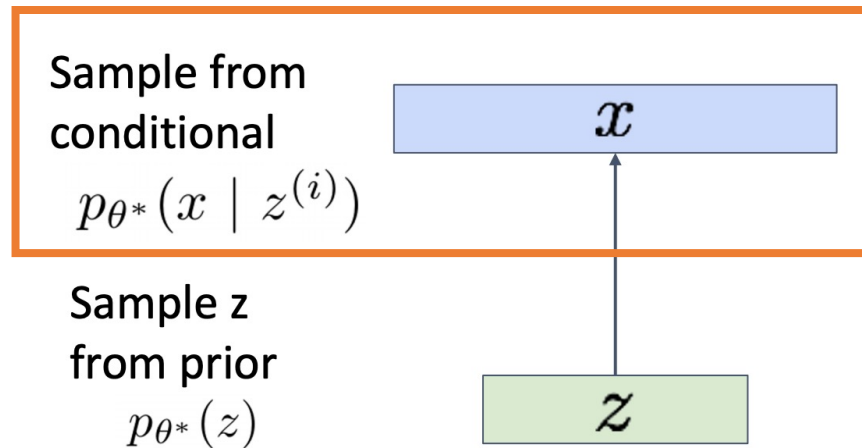
Problem: Impossible to integrate over all z!

Bayes Rule?

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Problem: No way to compute this!

- Variational
 - Probabilistic



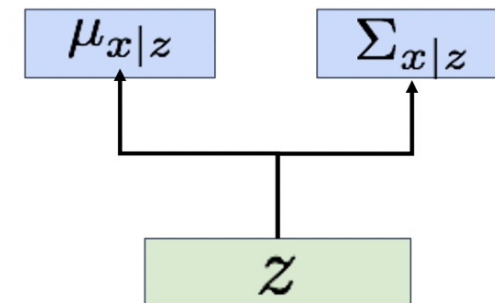
How to sample?

Assume simple prior $p(z)$, e.g. Gaussian with mean

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Decoder Network

$$p_{\theta}(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



Variational Autoencoders

- Let's maximize the likelihood of data! Need to compute $p_{\theta}(x)$

Let's train
encoder and
decoder jointly!

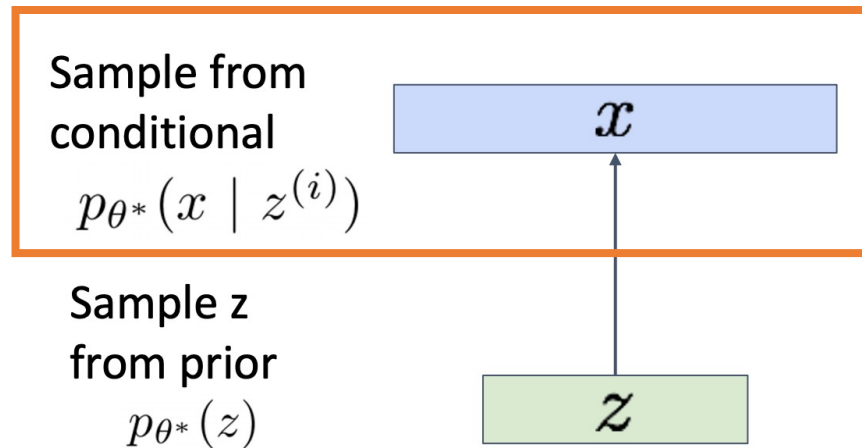
Solution: Train another network (encoder) that learns
 $q_{\phi}(z | x) \approx p_{\theta}(z | x)$

Bayes Rule?

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Problem: No way to compute this!

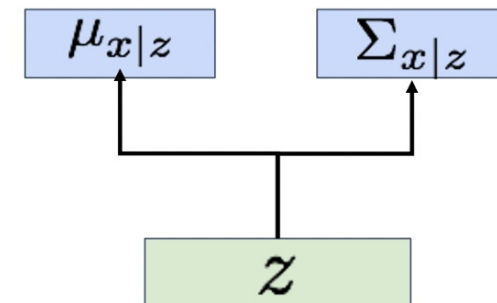
- Variational
 - Probabilistic



Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Decoder Network

$$p_{\theta}(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



How to sample?

Assume simple prior $p(z)$, e.g. Gaussian with mean

Adapted from slides by Justin Johnson

Variational Autoencoders (VAE)

Decoder network inputs latent code z , gives distribution over data x

Encoder network inputs data x , gives distribution over latent codes z

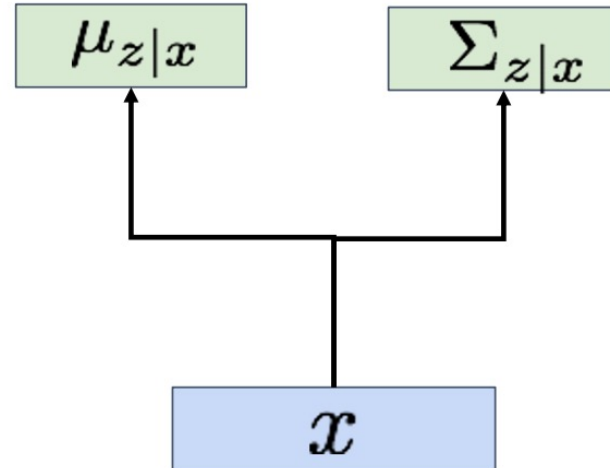
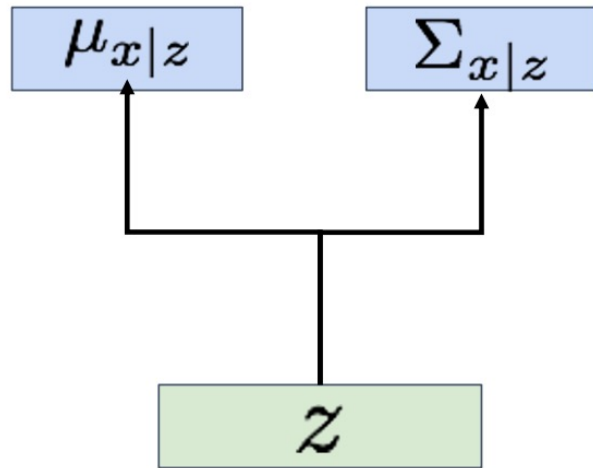
If we can ensure that $q_\phi(z | x) \approx p_\theta(z | x)$,

$$p_\theta(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$

$$q_\phi(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$$

then we can approximate

$$p_\theta(x) \approx \frac{p_\theta(x | z)p(z)}{q_\phi(z | x)}$$



Idea: Jointly train both encoder and decoder

Variational AutoEncoders (VAE)

Bunch of math to get a lower bound that we can optimize for!

$$\log p_{\theta}(x) = \log \frac{p_{\theta}(x|z)p(z)}{p_{\theta}(z|x)} = \log \frac{p_{\theta}(x|z)p(z)q_{\phi}(z|x)}{p_{\theta}(z|x)q_{\phi}(z|x)}$$

Variational AutoEncoders (VAE)

Bunch of math to get a lower bound that we can optimize for!

$$\begin{aligned}\log p_{\theta}(x) &= \log \frac{p_{\theta}(x|z)p(z)}{p_{\theta}(z|x)} = \log \frac{p_{\theta}(x|z)p(z)q_{\phi}(z|x)}{p_{\theta}(z|x)q_{\phi}(z|x)} \\ &= \log p_{\theta}(x|z) - \log \frac{q_{\phi}(z|x)}{p(z)} + \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)}\end{aligned}$$

Apply expectation (safely because x doesn't depend on z)

$$\log p_{\theta}(x) = E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x)]$$

$$= E_z [\log p_{\theta}(x|z)] - E_z \left[\log \frac{q_{\phi}(z|x)}{p(z)} \right] + E_z \left[\log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} \right]$$

Variational AutoEncoders (VAE)

Bunch of math to get a lower bound that we can optimize for!

$$\log p_{\theta}(x) = \log \frac{p_{\theta}(x|z)p(z)}{p_{\theta}(z|x)} = \log \frac{p_{\theta}(x|z)p(z)q_{\phi}(z|x)}{p_{\theta}(z|x)q_{\phi}(z|x)}$$

$$= E_z[\log p_{\theta}(x|z)] - E_z \left[\log \frac{q_{\phi}(z|x)}{p(z)} \right] + E_z \left[\log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} \right]$$

Data reconstruction

KL divergence between prior, and samples from the encoder network

KL divergence between encoder and posterior of decoder

$$= E_{z \sim q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - D_{KL}(q_{\phi}(z|x), p(z)) + D_{KL}(q_{\phi}(z|x), p_{\theta}(z|x))$$

KL is ≥ 0 , so dropping this term gives a **lower bound** on the data likelihood:

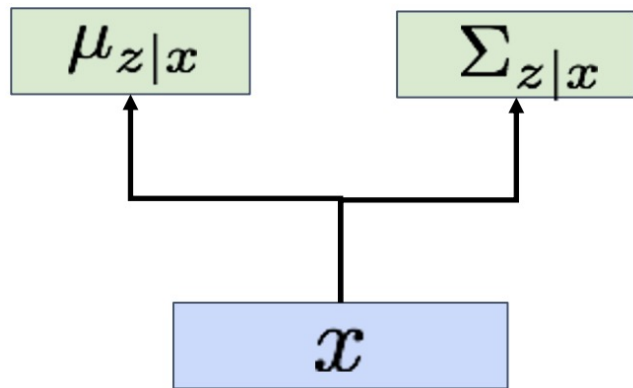
Variational Autoencoders (VAE)

Jointly train **encoder** q and **decoder** p to maximize the **variational lower bound** on the data likelihood

$$\log p_{\theta}(x) \geq E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x), p(z) \right)$$

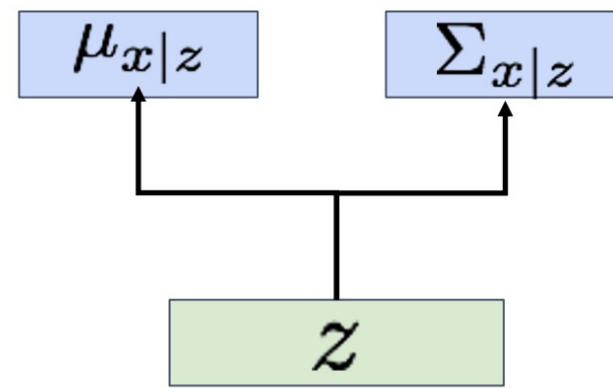
Encoder Network

$$q_{\phi}(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$$

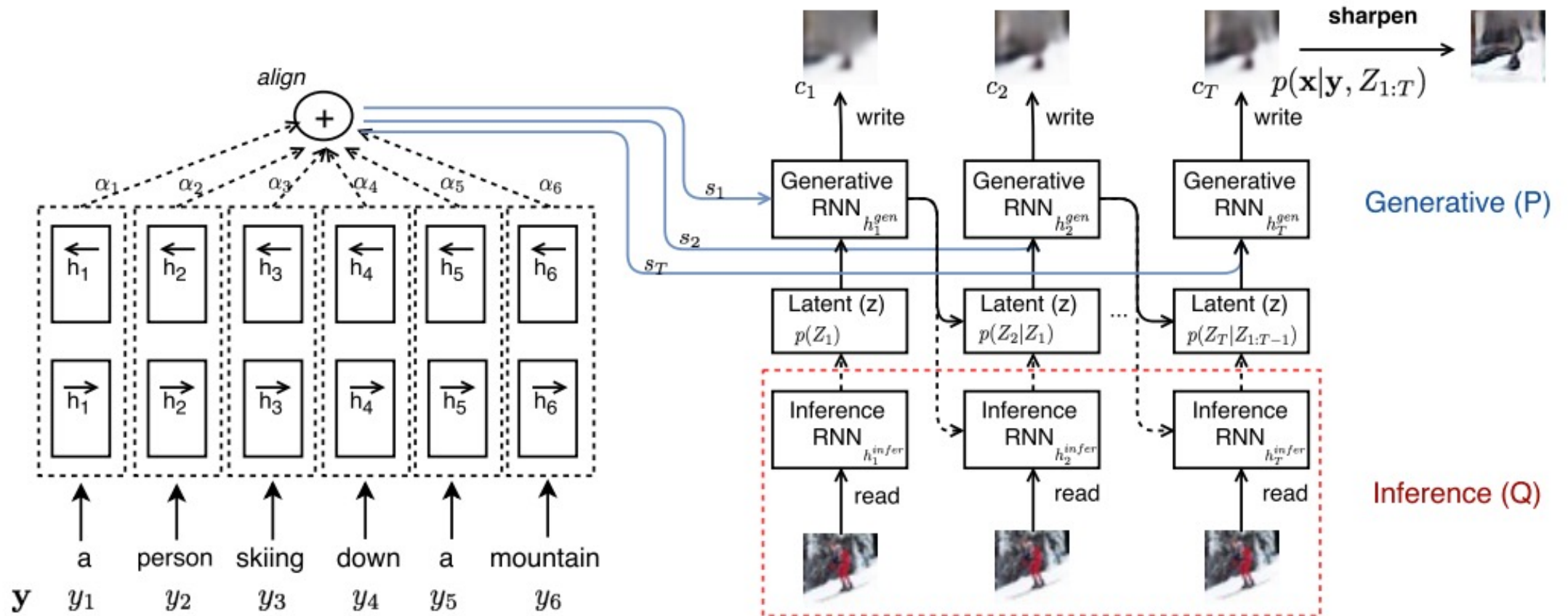


Decoder Network

$$p_{\theta}(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



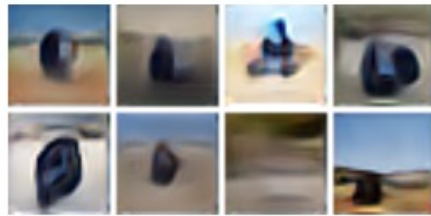
Text-based image generation with VAE



Generating Images from Captions with Attention

<https://arxiv.org/pdf/1511.02793.pdf>, Mansimov et al, ICLR 2016

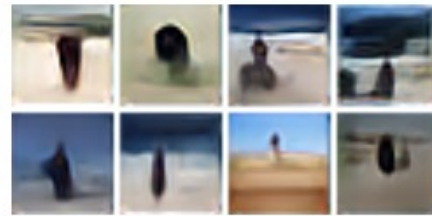
Text-based image generation with VAE



A rider on a blue motorcycle in the desert.



A rider on a blue motorcycle in the forest.



A surfer, a woman, and a child walk on the beach.



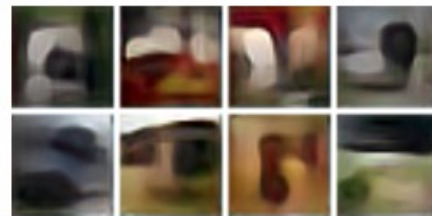
A surfer, a woman, and a child walk on the sun.



alignDRAW



LAPGAN



Conv-Deconv VAE



Fully-Conn VAE

Generating Images from Captions with Attention

<https://arxiv.org/pdf/1511.02793.pdf>, Mansimov et al, ICLR 2016

Compare AR and VAE models

Autoregressive models

- Directly maximize $p(\text{data})$
- High-quality generated images
- Slow to generate images
- No explicit latent codes

Variational models

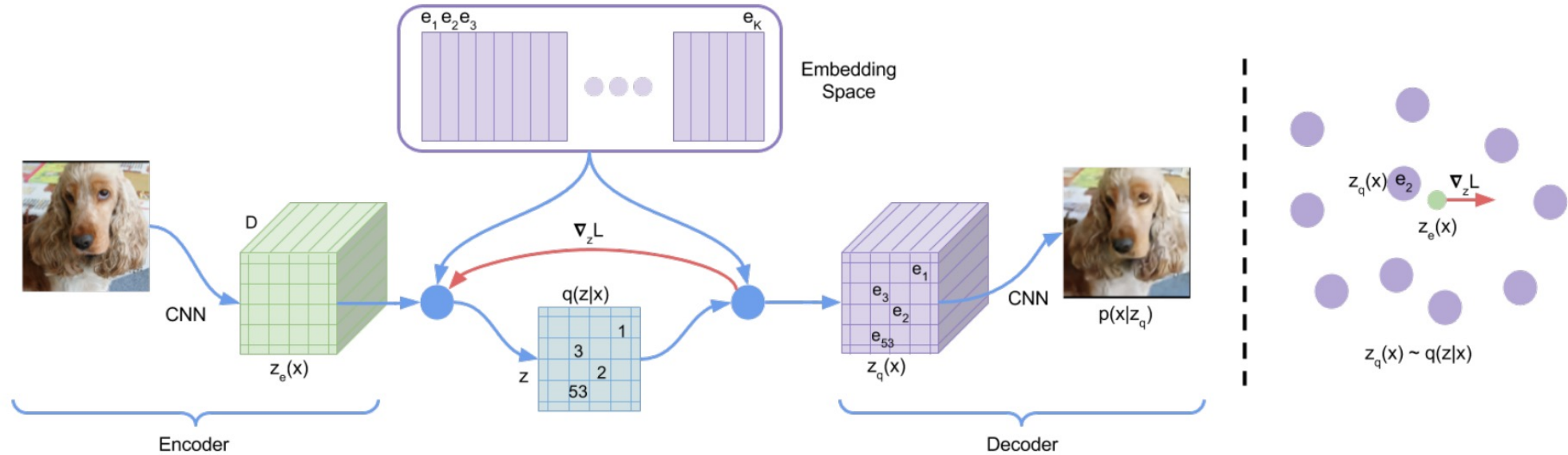
- Maximize lower-bound on $p(\text{data})$
- Generated images often blurry
- Very fast to generate images
- Learn rich latent codes

Can we combine them and get the best of both worlds?

Combine VAE + Autoregressive

Vector-Quantized Variational Autoencoder (VQ-VAE)

- Autoregressively model images
- But instead of directly on pixels, on image patches compressed into image "tokens" using VAE



- Two-stage training process

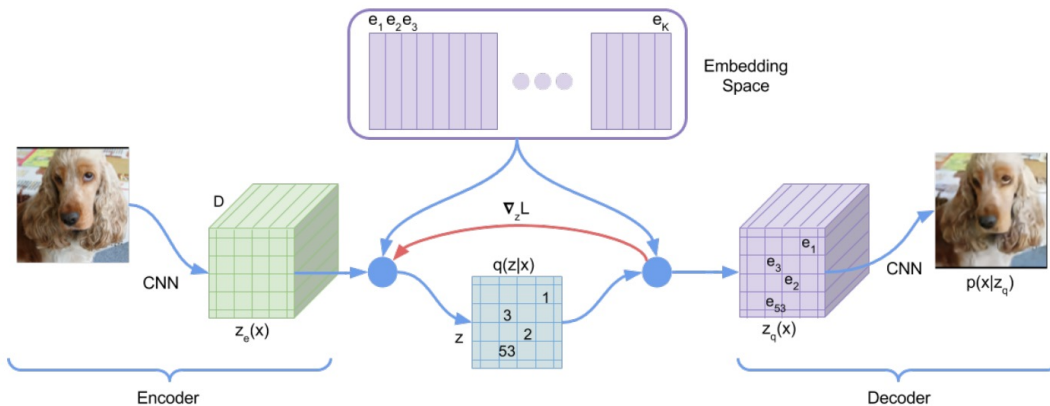
Neural Discrete Representation Learning

<https://arxiv.org/pdf/1711.00937.pdf>, Oord et al, NIPS 2017

Combine VAE + Autoregressive

Vector-Quantized Variational Autoencoder (VQ-VAE)

- Two-stage training process
- Use VAE to create a code book to encode image patch into latent quantized discrete vector



**128x128 class-conditional results
trained on ImageNet**



- Use autoregressive model (PixelCNN) to model latent prior $p(z)$

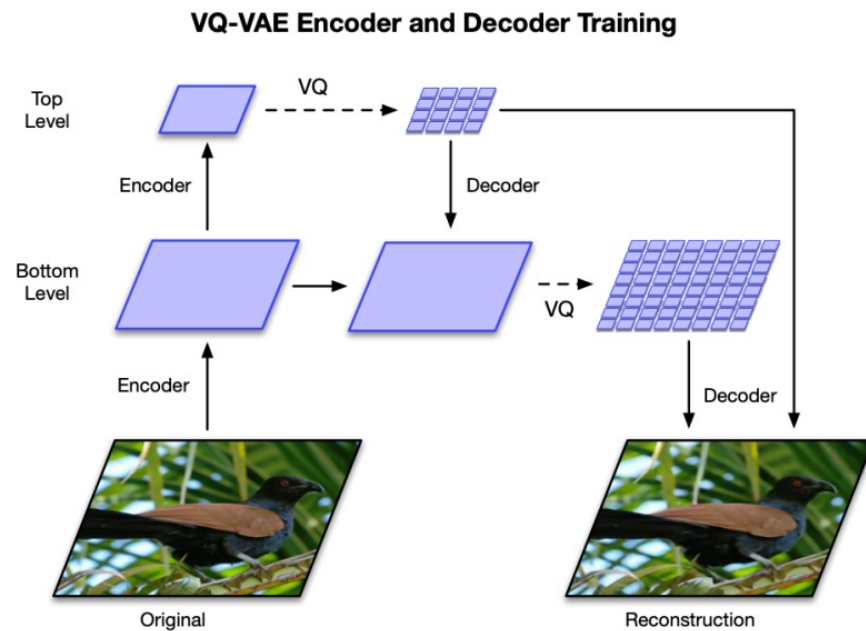
Neural Discrete Representation Learning

<https://arxiv.org/pdf/1711.00937.pdf>, Oord et al, NIPS 2017

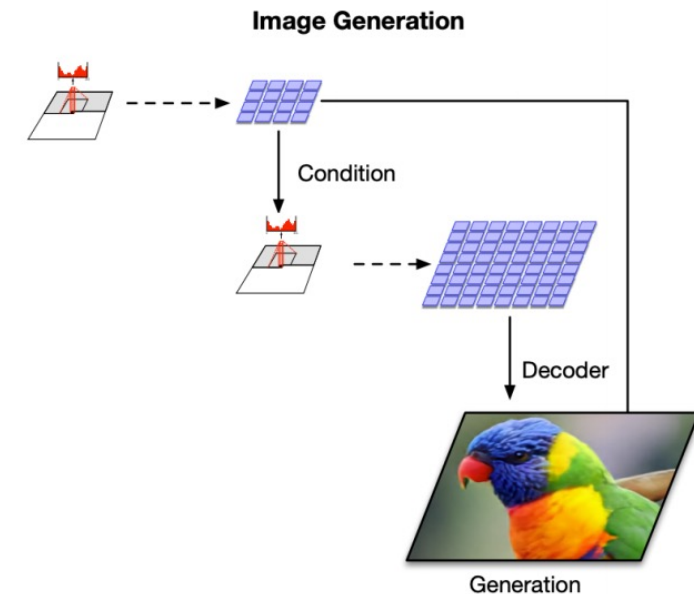
Combine VAE + Autoregressive Vector-Quantized Variational Autoencoder (VQ-VAE2)

- Hierarchical VQ-VAE

Train a VAE-like model to generate
multiscale grids of latent codes



Use a multiscale PixelCNN to
sample in latent code space



VQ-VAE2 Results

256 x 256 class-conditional samples, trained on ImageNet



Generating Diverse High-Fidelity Images with VQ-VAE-2
<https://arxiv.org/pdf/1906.00446.pdf>, Razavi et al, NeurIPS 2019

VQ-VAE2 Results

256 x 256 class-conditional samples, trained on ImageNet



Generating Diverse High-Fidelity Images with VQ-VAE-2

<https://arxiv.org/pdf/1906.00446.pdf>, Razavi et al, NeurIPS 2019

VQ-VAE2 Results

1024 x 1024 generated faces, trained on FFHQ



Generating Diverse High-Fidelity Images with VQ-VAE-2
<https://arxiv.org/pdf/1906.00446.pdf>, Razavi et al, NeurIPS 2019

DALL-E

- Like VQ-VAE2 but
 - Conditioned on text
 - Large network trained with tons of data
 - Used 3.3M text/image pairs (Conceptual Captions) for 1.2B parameter model
 - Used 120 text/image pairs (collected from Internet) for 12B parameter model
 - Uses autoregressive transformer vs PixelCNN
 - Uses CLIP to rerank generated images (vs classifier network trained on ImageNet)

“Dall-e”

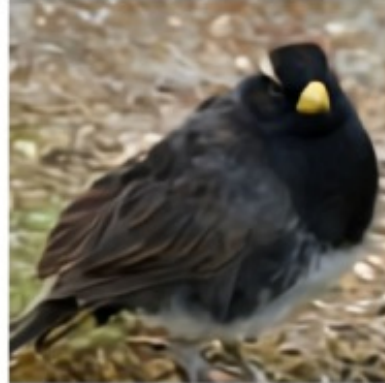
[Ramesh et al, <https://openai.com/blog/dall-e/>]

DALL-E: Results

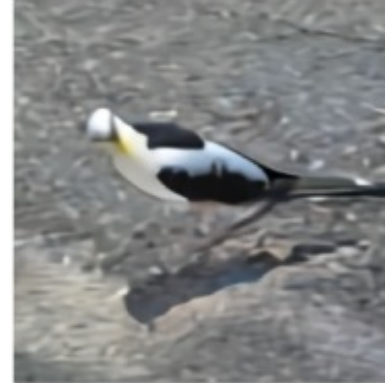
this gray bird has a pointed beak black wings with small white bars long thigh and tarsus and a long tail relative to its size



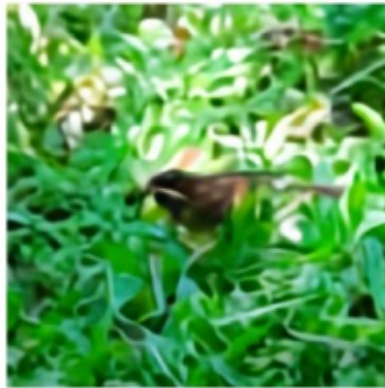
this rotund bird has a black tipped beak a black tail with a yellow tip and a black cheek patch



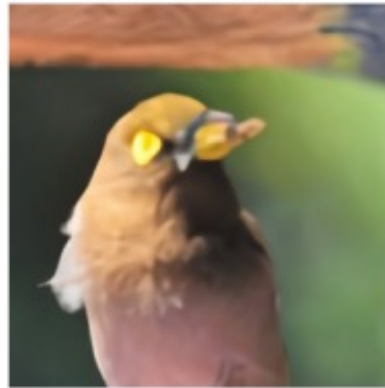
this is a small white bird with a yellow crown and a black eye ring and cheek patch and throat



the small bird has a dark brown head and light brown body



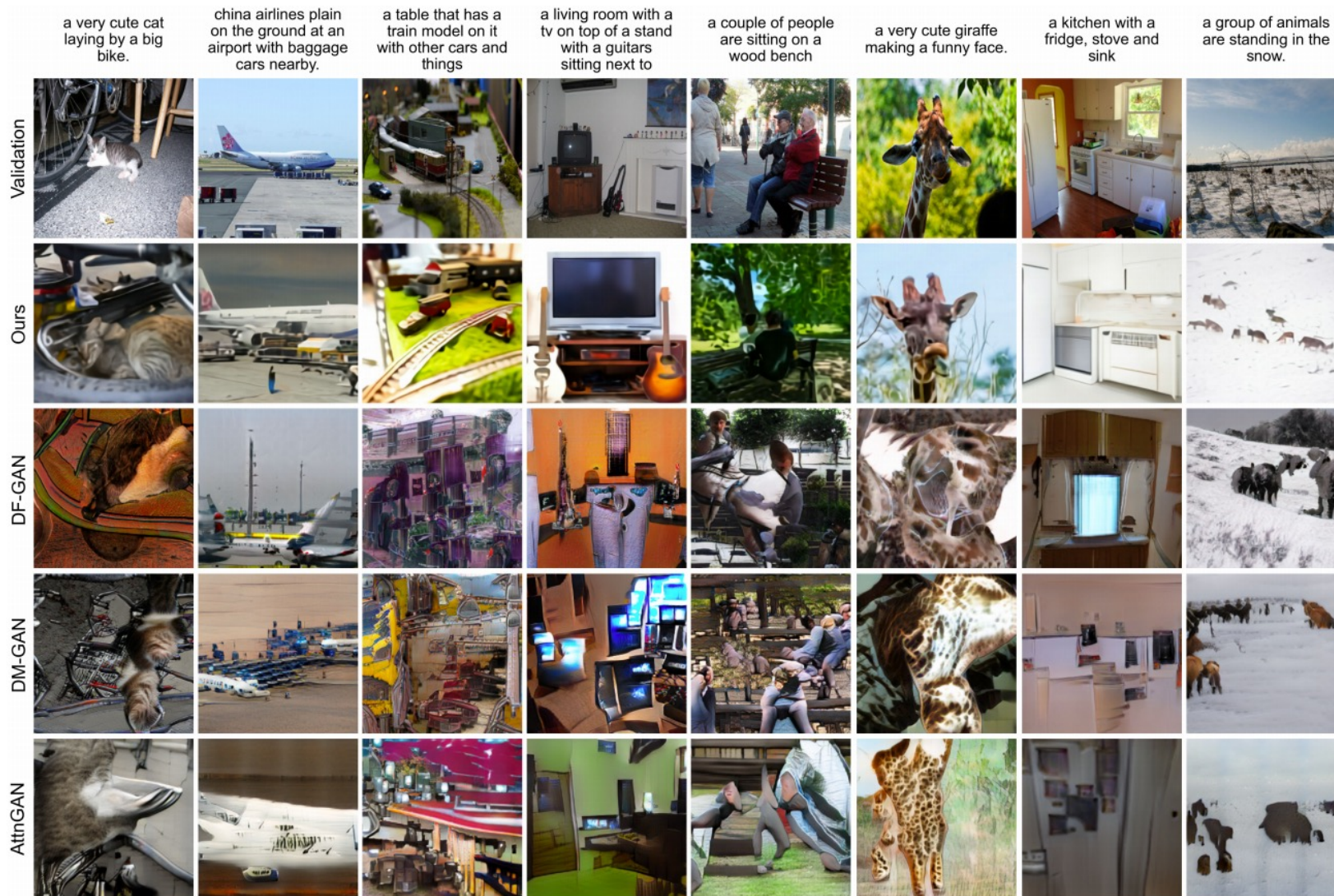
small bird with a pale yellow underside light brown crown and back gray tail and wing tips tip of tail feather bright yellow black eyes and black stripe over eyes



a small bird with a grey head and grey nape with grey black and white covering the rest of the body



DALL-E: Results



Diffusion models

- Define Markov chain of transitions from input to series of latent variables.

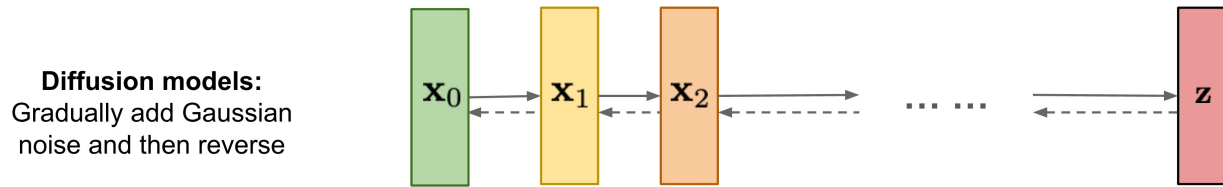
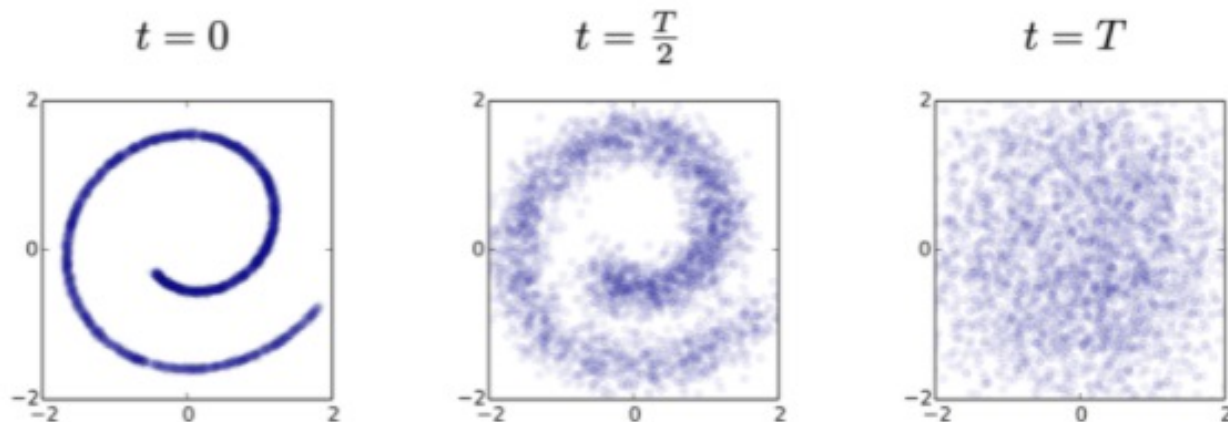


Figure credit: <https://lilianweng.github.io/lil-log/2021/07/11/diffusion-models.html>

- Forward process (diffusion process)

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

- Reverse diffusion: recreate sample (image) from latent (gaussian noise)



(Image source: [Sohl-Dickstein et al., 2015](#))

GLIDE: Diffusion Models

- Large diffusion model



“a hedgehog using a calculator”



“a corgi wearing a red bowtie and a purple party hat”



“robots meditating in a vipassana retreat”



“a fall landscape with a small cottage next to a lake”



“a surrealist dream-like oil painting by salvador dalí of a cat playing checkers”



“a professional photo of a sunset behind the grand canyon”



“a high-quality oil painting of a psychedelic hamster dragon”



“an illustration of albert einstein wearing a superhero costume”

GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models
<https://arxiv.org/pdf/2112.10741.pdf>, Nichol et al, arXiv 2021

Evaluating generated content

Evaluation

- Evaluation of these models are tricky!
- What makes for a good generation?
- General
 - Is the generated content high quality?
 - Does it match the distribution?
 - Is it diverse?
- For language conditioned generation:
 - Does the generated content match the language?
 - Are salient aspects of the language captured in the objects, appearance, and relationships?

GAN evaluation

- Inception Score: $I = \exp(\mathbb{E}_{\mathbf{x}} D_{KL}(p(y|\mathbf{x}) || p(y)))$
 - Use inception model to predict class y
 - Want good models to generate diverse but meaningful images
 - Large distance between marginal prior (of labels) and conditional prior
- FID (Frechet Inception distance): measures distance between generated and real distribution
- Human rank images generated by models

Metrics

- R-Precision (retrieval)
 - Randomly sample 99 other captions, where is the input caption ranked (using cosine similarity) compared to the rest (is it in the top r)?
- Visual similarity (VS)
 - how well does the encoded text and image match)
$$VS = \frac{f_t(t) \cdot f_x(x)}{\|f_t(t)\|_2 \cdot \|f_x(x)\|_2}$$
 - High variance, dependency on the specific encoders used
- Semantic Object Accuracy (SOA)
 - Use pretrained object detector to match words in text
- Captioning – generate caption and evaluate with original caption using standard captioning metrics

Metrics

Metric	Image Quality	Image Diversity	Object Fidelity	Text Relevance	Mentioned Objects	Numerical Alignment	Positional Alignment	Paraphrase Robustness	Explainable	Automatic
IS [130]	✓									✓
FID [131]	✓	✓								✓
SceneFID [103]			✓							✓
R-prec. [35]				✓						✓
VS [42]				✓						✓
SOA [108]				✓	✓					✓
Captioning				(✓)						✓
User Studies	✓	✓	✓	✓	✓	✓	✓	✓	✓	

Next time

- Monday: Paper presentations and discussions
 - (Tristan) Cross-Modal Contrastive Learning for Text-to-Image Generation
 - (Han-Hung) GLIDE: Toward Photorealist Image Generation
- Wednesday: Compositionality and structured representations