

CMPT 983

Grounded Natural Language Understanding

April 12, 2021

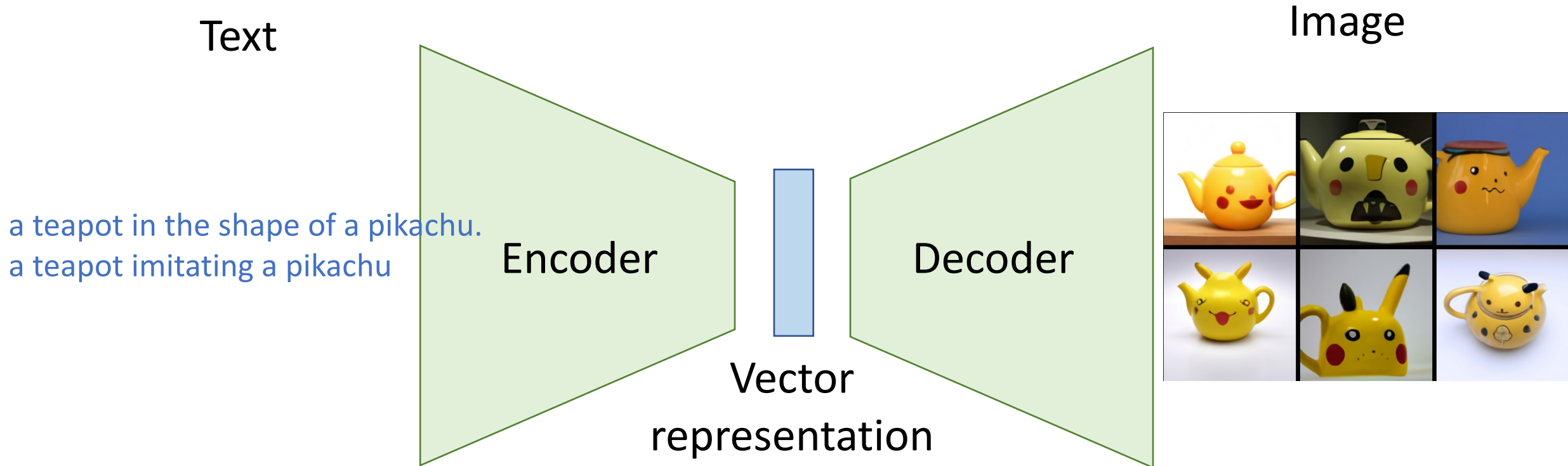
Content generation from language

Next time

- Thursday (4/15): Last day – project discussion and conclusion
 - Watch other group's project video before class
 - Project video due by 11:59pm 4/14
 - Project report due by 11:59pm 4/15

Content generation from
language

Translating across modalities



“Dall-e”

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

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:

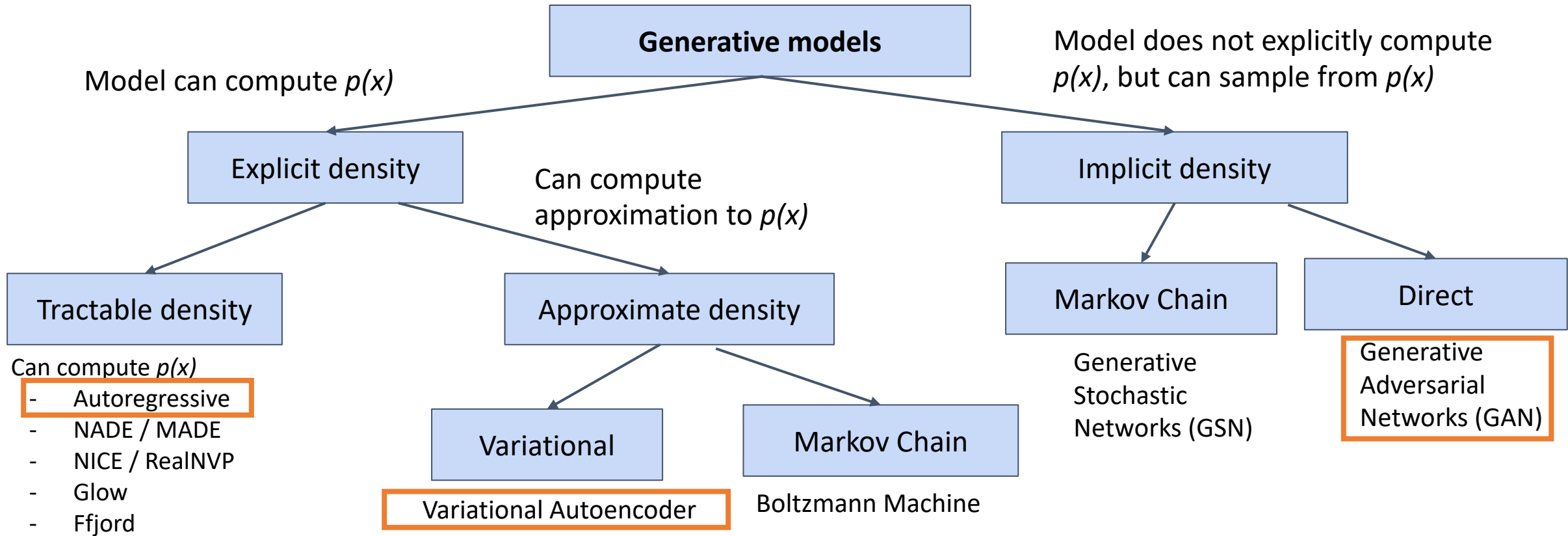
Learn $p(x|y)$



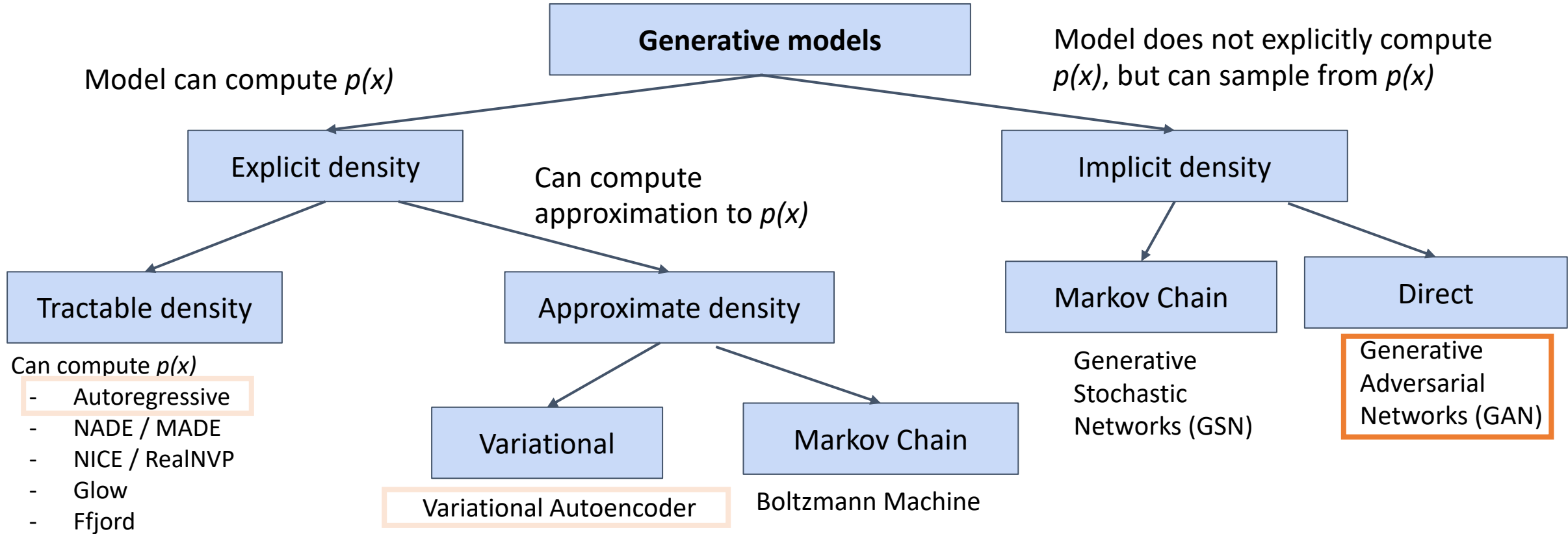
Assign labels, while rejecting outliers!

Generate new data conditioned on input labels

Taxonomy of generative models



Taxonomy of generative models



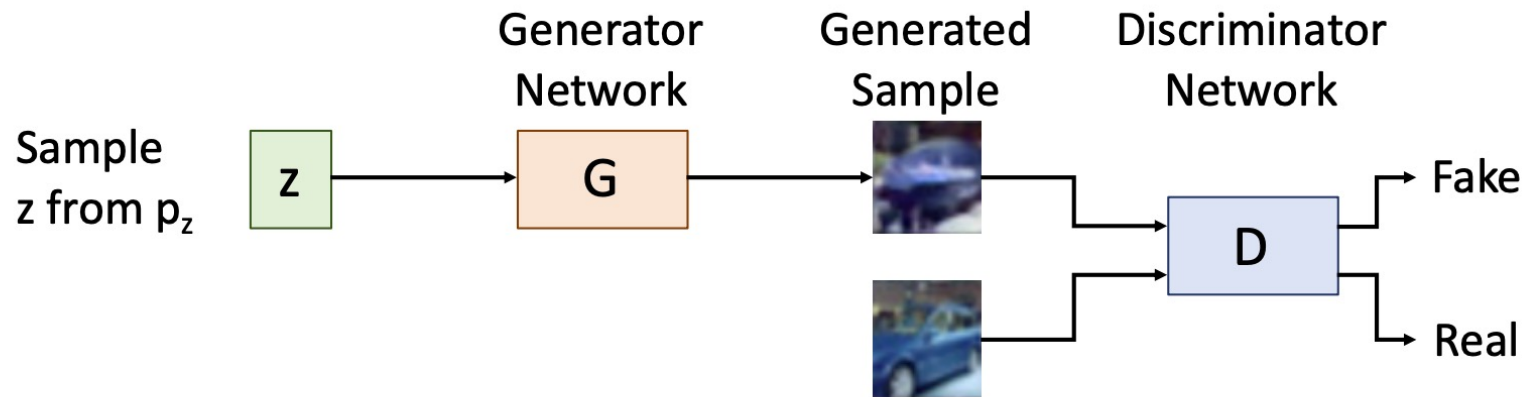
Generative Adversarial Networks (GAN)

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

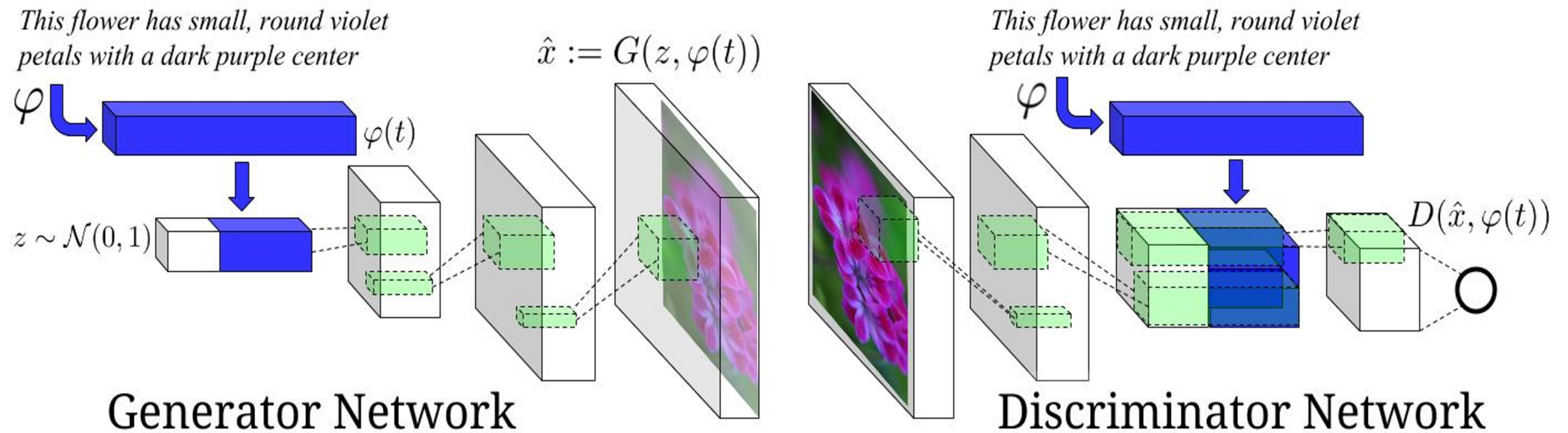
$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$



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

Text to image with GANs

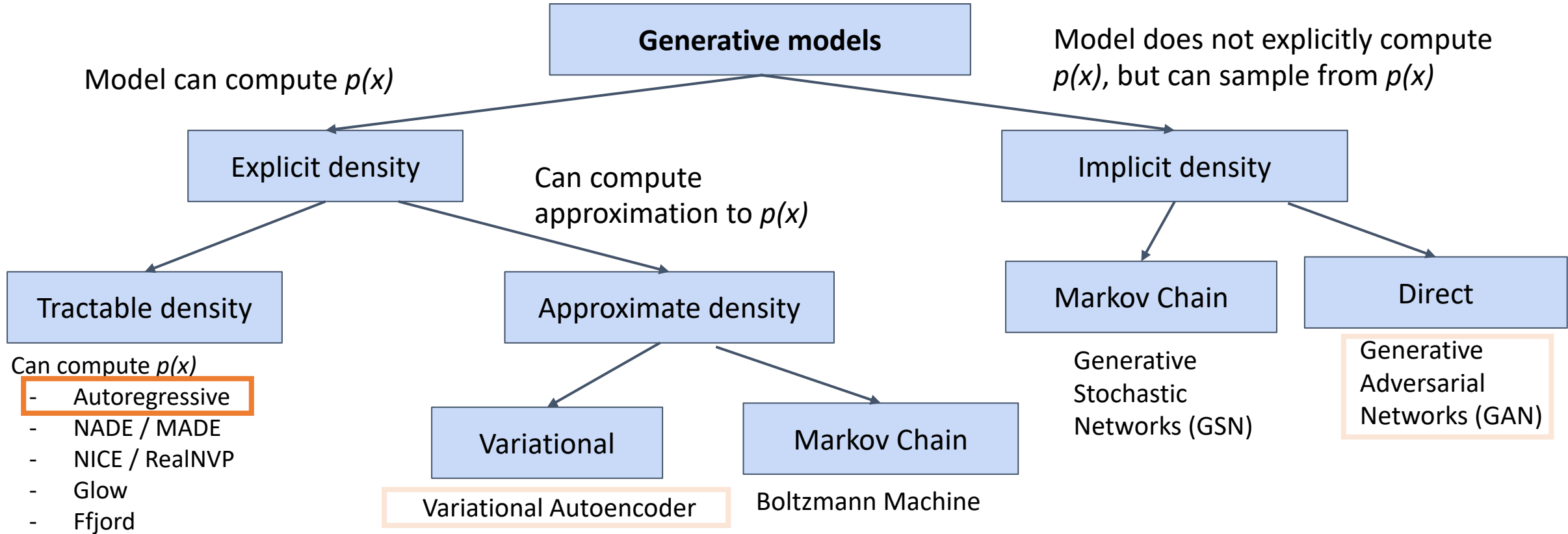
- Generator and Discriminators are trained alternately



GANs for text to image generation

- GAN+CLS+INT (Reed et al, ICML 2016)
 - Pre-train text (char-CNN-RNN) and image encoder (CNN) for joint-embedding
 - CLS: additional discriminator loss for if image/text match
 - INT: interpolated text for additional training data
- StackGAN (Zhang et al, ICCV 2017)
 - 2 level GANs stacked together for higher resolution
- StackGAN++ (Zhang et al, TPAMI 2018)
 - Generalized StackGAN (multiscale), trained end-to-end
 - Unconditional + conditional loss (similar to GAN+CLS)
- AttnGAN (Xu et al, CVPR 2018)
 - Series of GANs (like StackGAN++)
 - Attention based similarity (DAMSM loss) to encourage representations that align regions of images to words in the text
- Many more GAN papers: MirrorGAN, ControlGAN, DMGAN, DTGAN...

Taxonomy of generative models



Autoregressive models

- Explicit function for modeling $p(x) = f(x, W)$
- Assume x can be broken down into subparts and apply chain rule

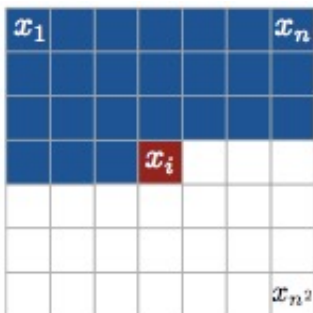
$$x = (x_1, x_2, \dots, x_T)$$

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

- Predict each part one after the other (autoregressive) using RNNs or Transformers

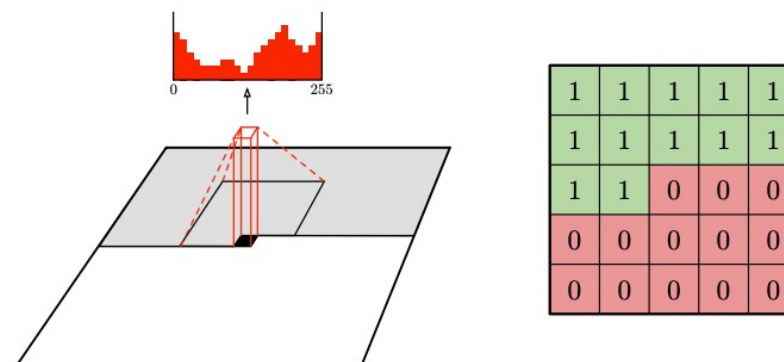
Autoregressive models

PixelRNN



PixelRNN, van der Oord et al, 2016

PixelCNN



PixelCNN, van der Oord et al, 2016

Image Transformer

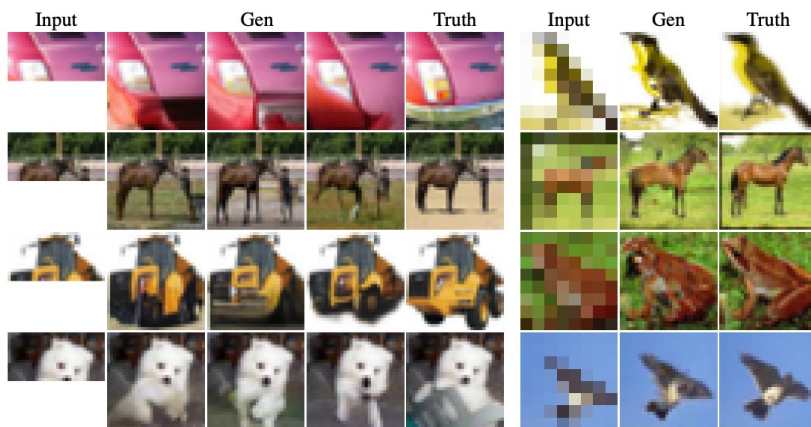
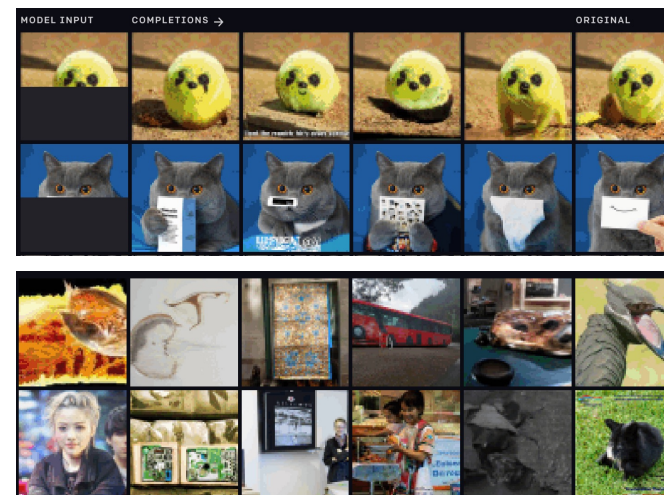


Image Transformer, Parmar et al, ICML 2018

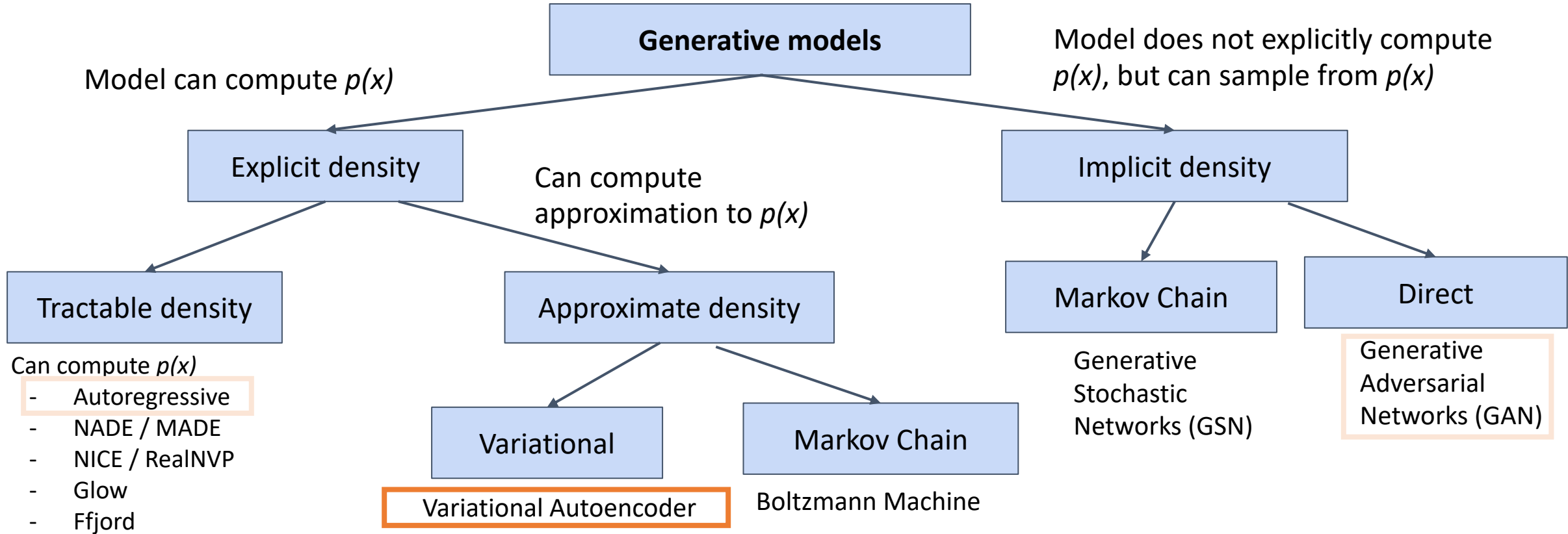
Image GPT



Generative Pretraining from Pixels, Chen et al, ICML 2020

<https://openai.com/blog/image-gpt/>

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?

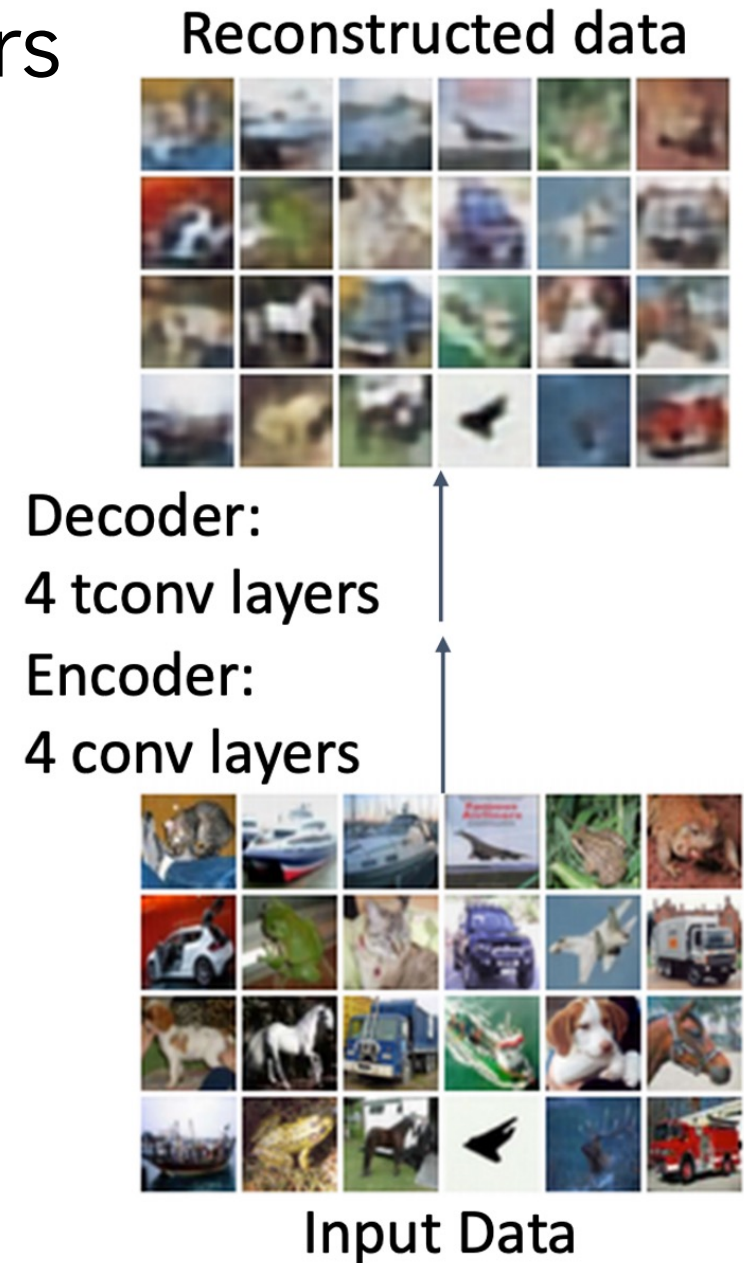
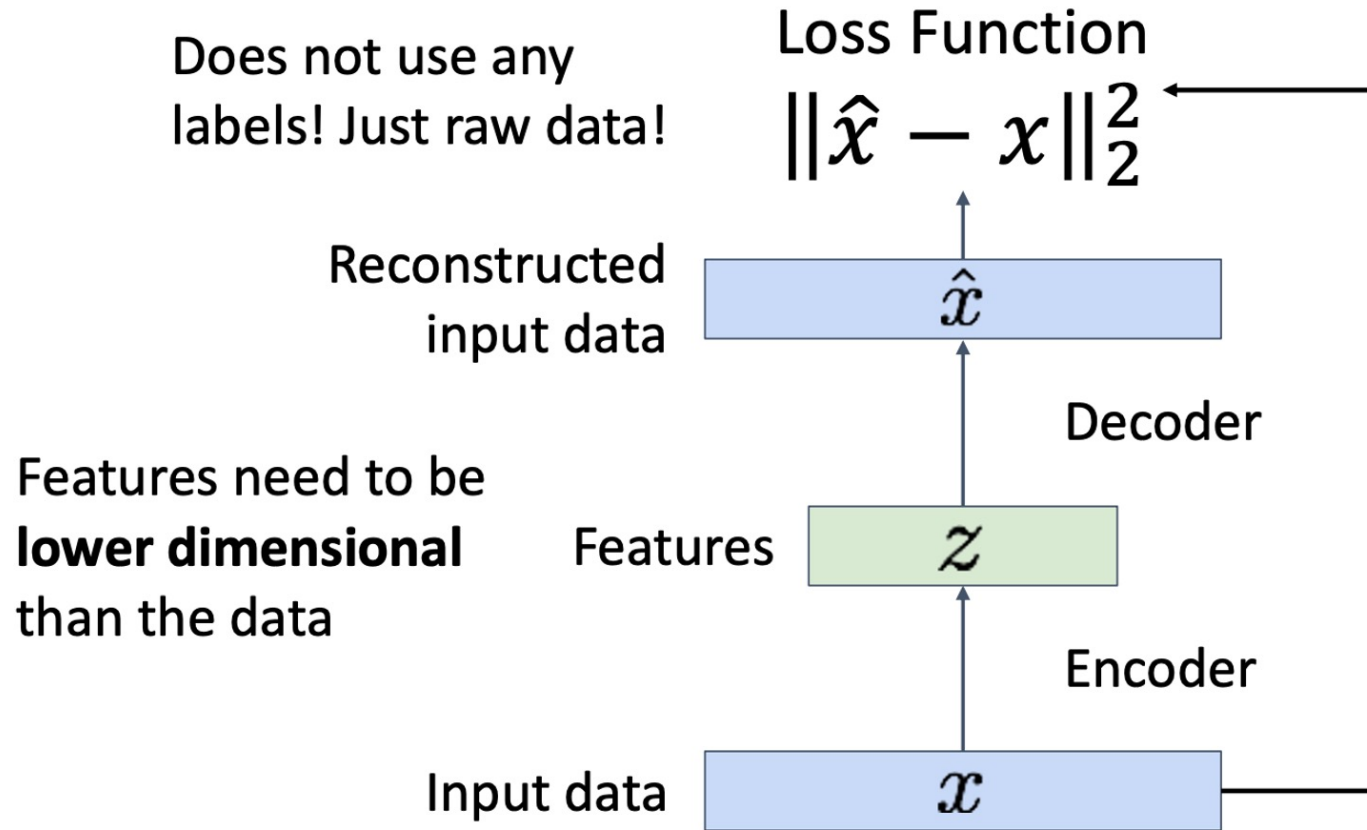
- Variational Autoencoders (VAE) use an **intractable density** that we cannot explicitly compute or optimize



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

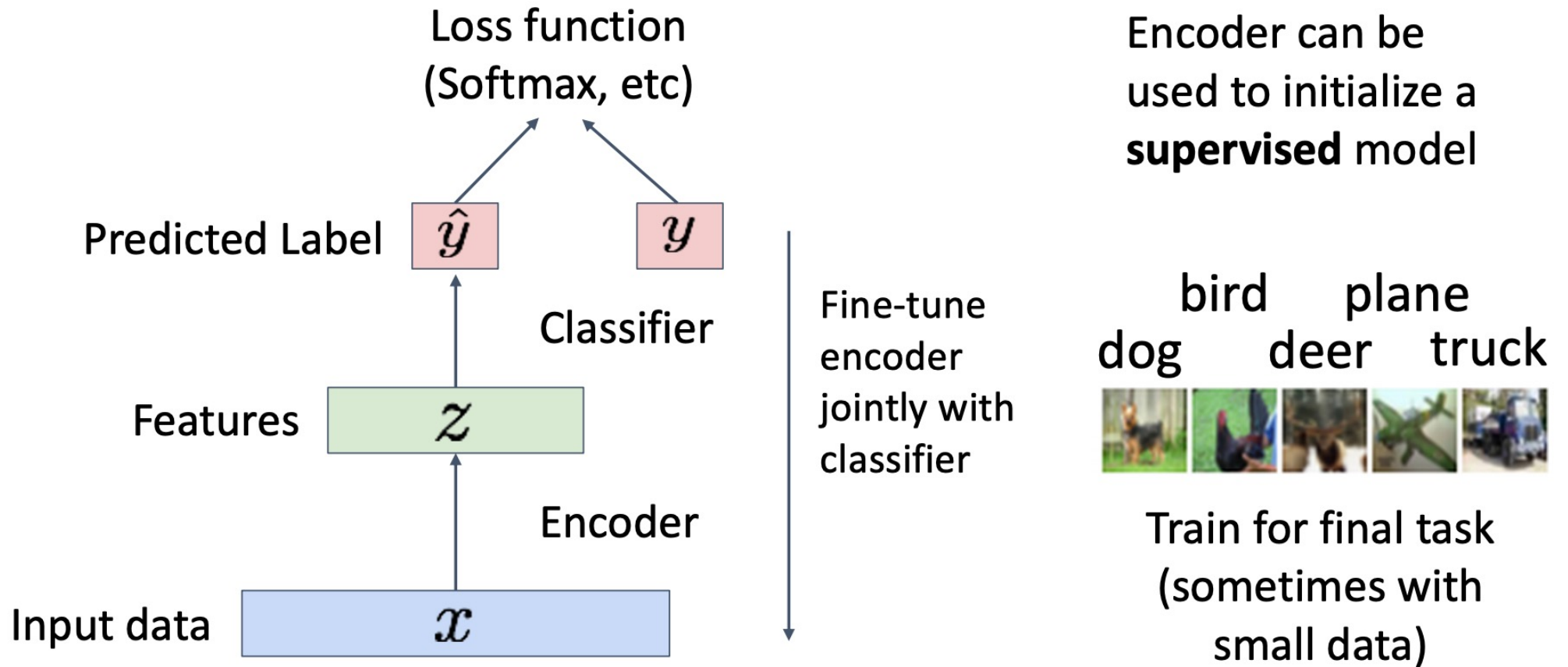
(Regular, non-variational) Autoencoders

Loss: L2 distance between input and reconstructed data.



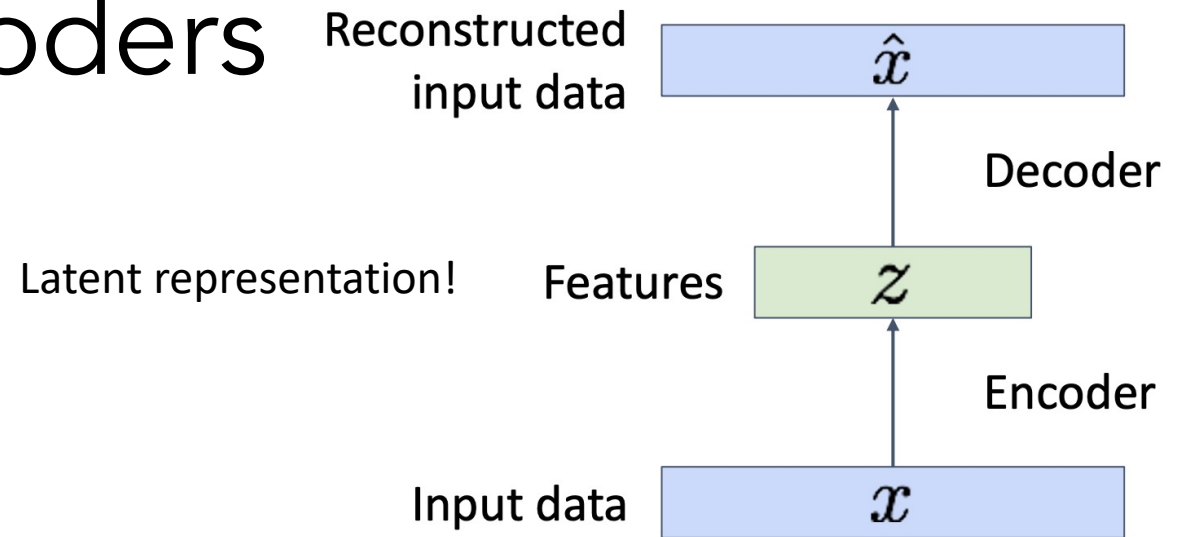
(Regular, non-variational) Autoencoders

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

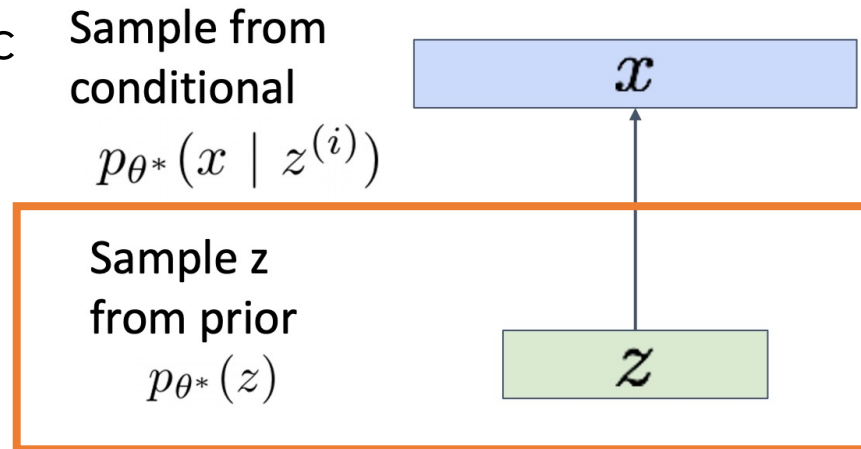


Variational Autoencoders

- Autoencoders
 - Not probabilistic
 - No sampling



- Variational
 - Probabilistic



How to sample?

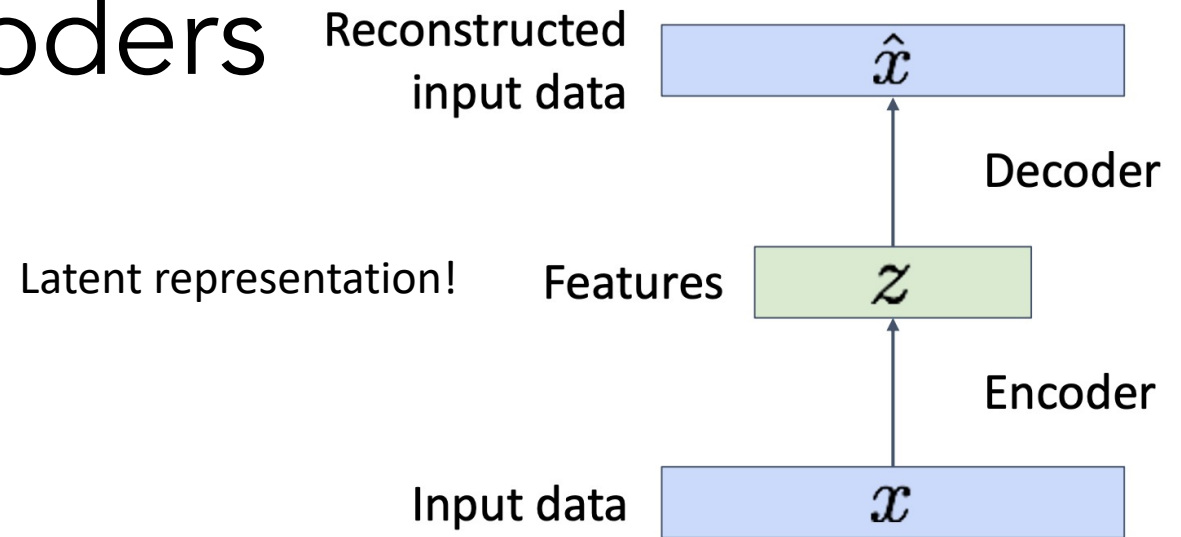
Assume simple prior $p(z)$, e.g. Gaussian with mean

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

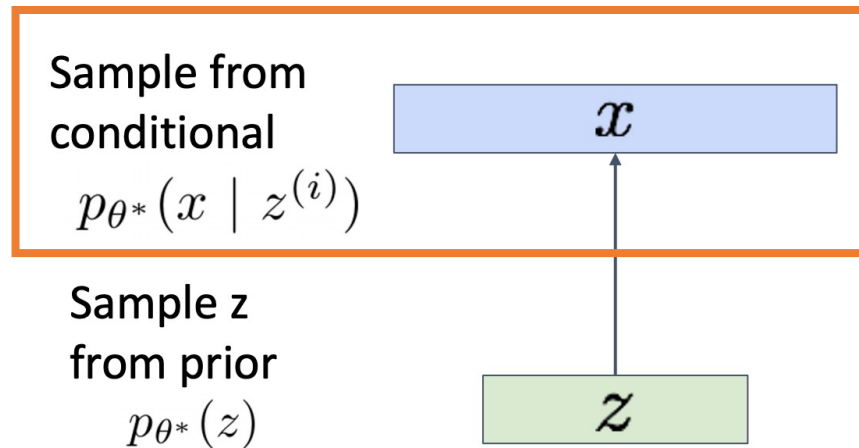
1. Learn latent representation
2. Sample to generate images

Variational Autoencoders

- Autoencoders
 - Not probabilistic
 - No sampling



- Variational
 - Probabilistic



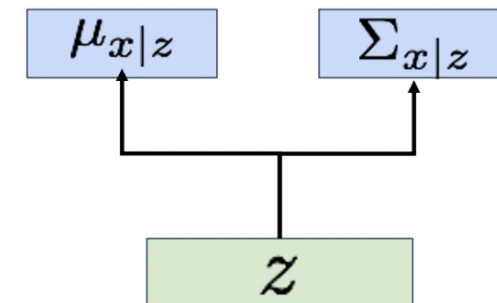
How to sample?

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?

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

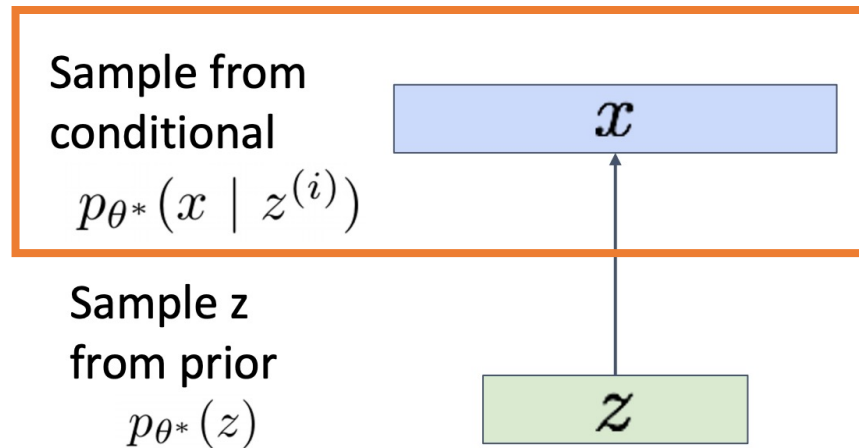
Problem: Impossible to integrate over all z!

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



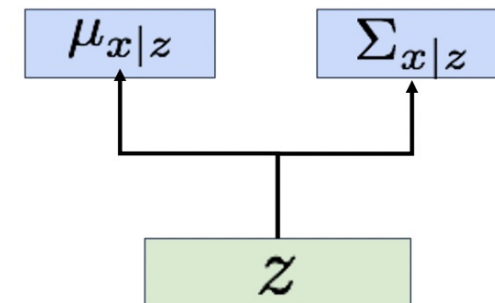
How to sample?

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!

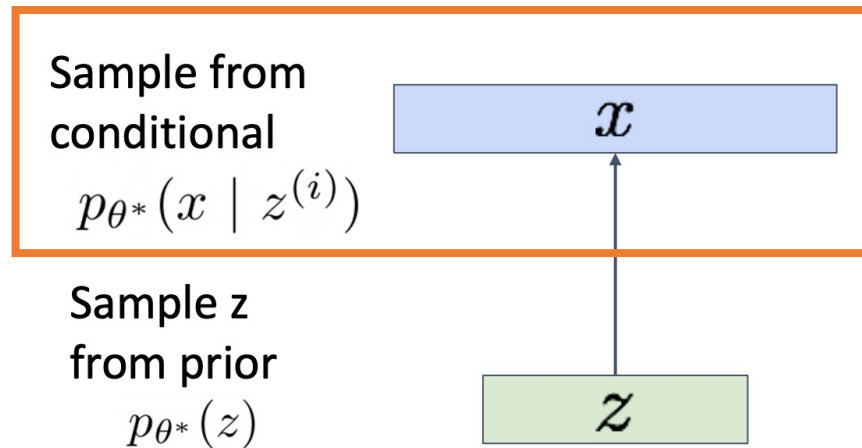
**Solution: Train
another network
(encoder) that learns**
 $q_{\phi}(z | x) \approx p_{\theta}(z | x)$

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



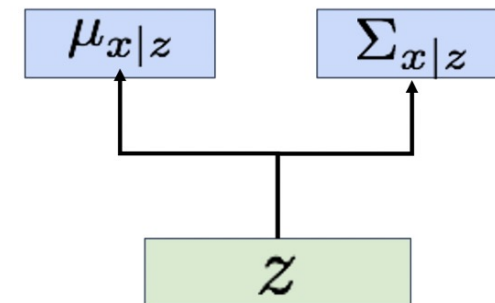
How to sample?

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

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



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

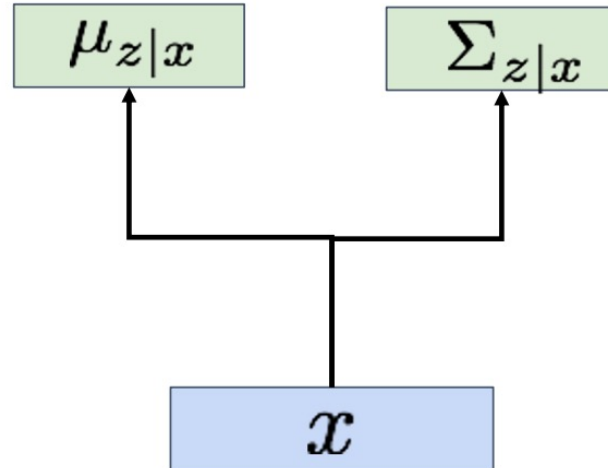
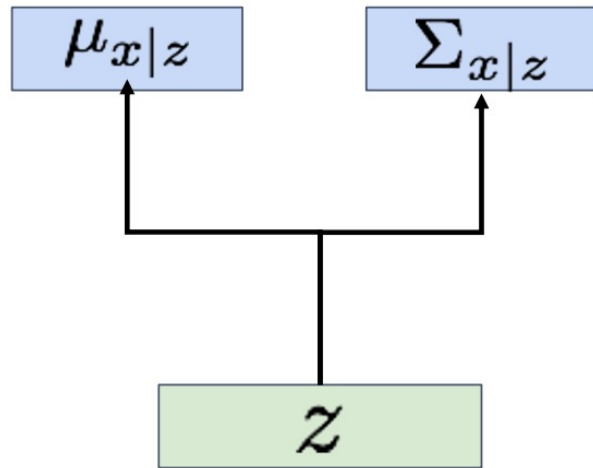
If we can ensure that $q_\phi(z | x) \approx p_\theta(z | x)$,

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

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

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:

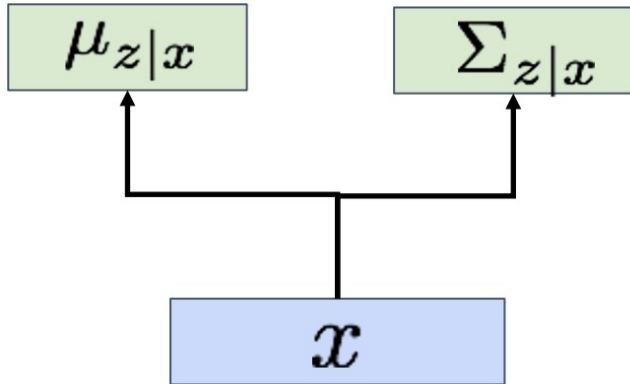
Variational Autoencoders (VAE)

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

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

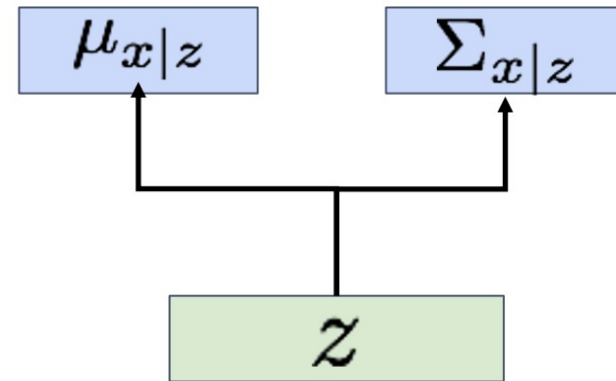
Encoder Network

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

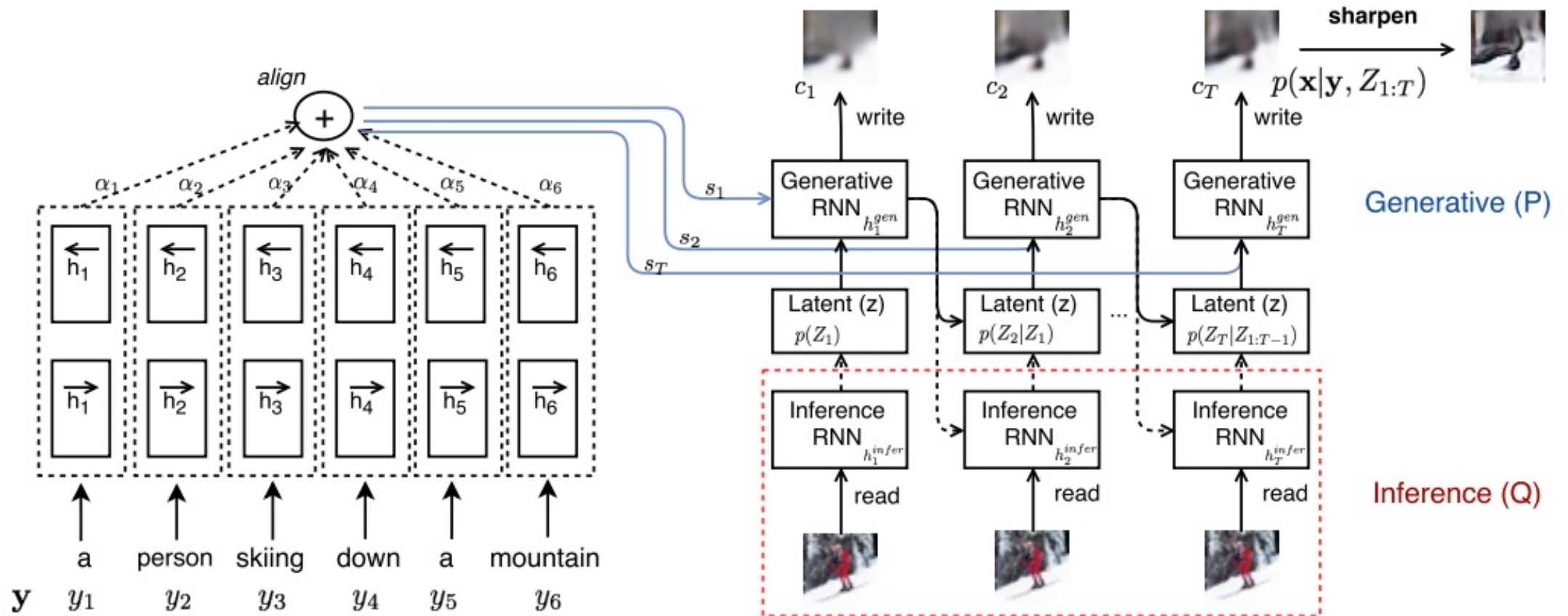


Decoder Network

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



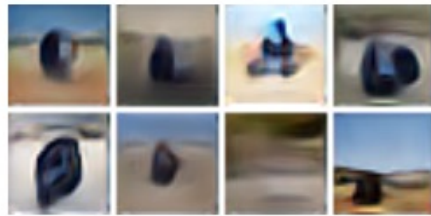
Text-based image generation with VAE



Generating Images from Captions with Attention

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

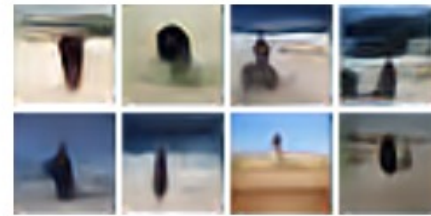
Text-based image generation with VAE



A rider on a blue motorcycle in the desert.



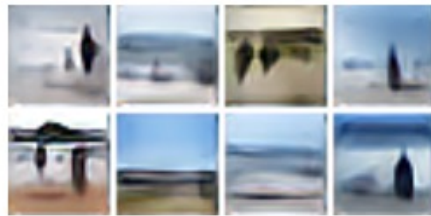
A rider on a blue motorcycle in the forest.



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

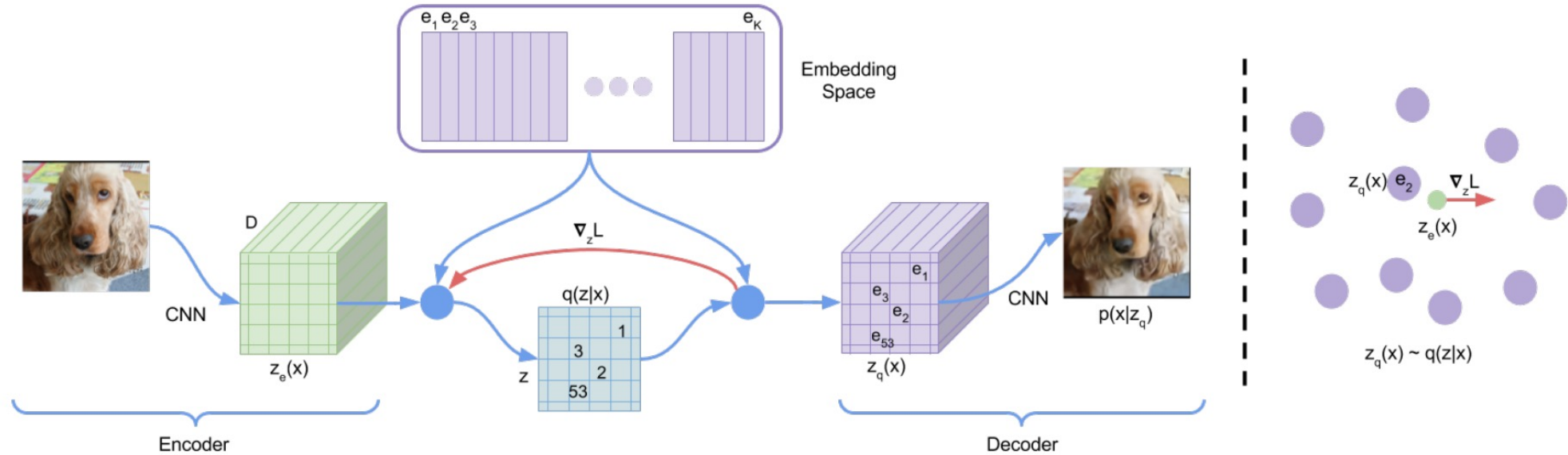
- Maximize lower-bound on $p(\text{data})$
- Generated images often blurry
- Very fast to generate images
- Learn rich latent codes

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

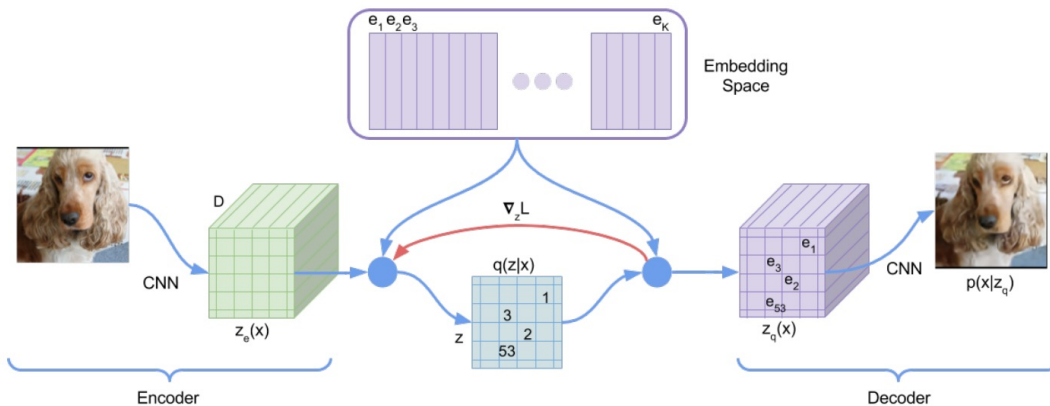


- Two-stage training process

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 $p(z)$

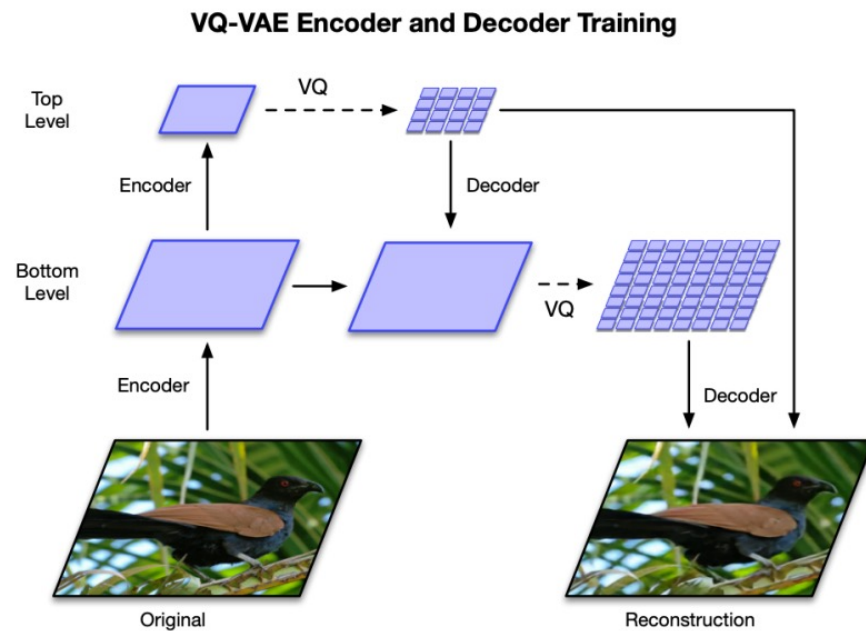
Neural Discrete Representation Learning

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

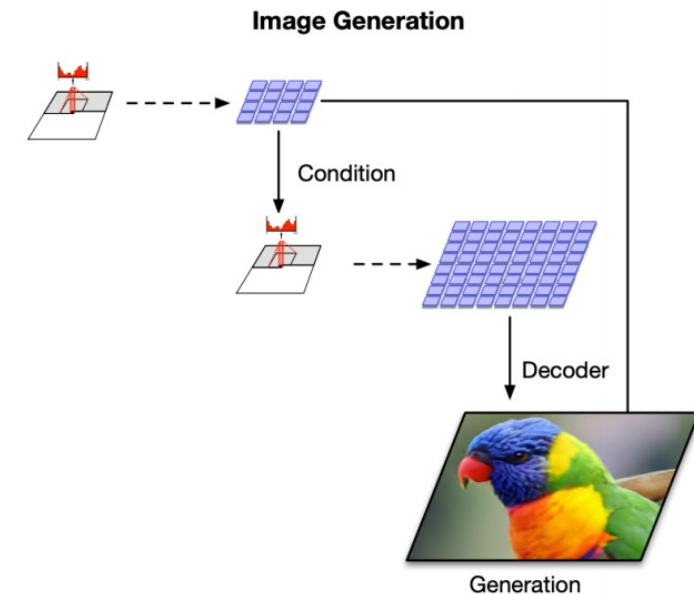
Combine VAE + Autoregressive Vector-Quantized Variational Autoencoder (VQ-VAE2)

- 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

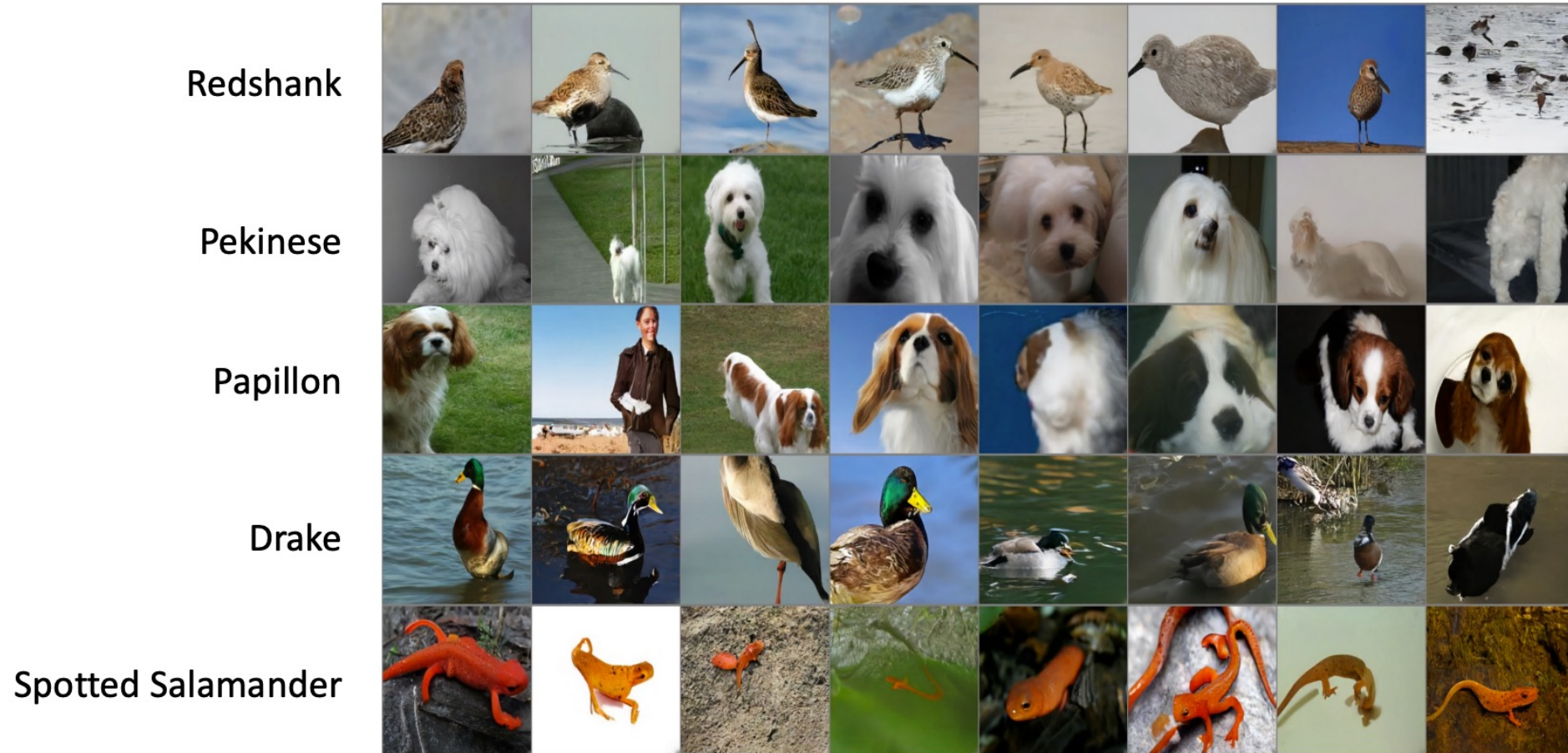
256 x 256 class-conditional samples, trained on ImageNet



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



<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



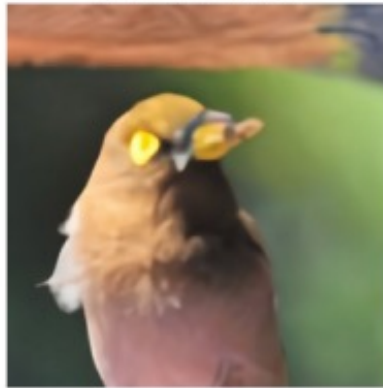
this is a small white bird with a yellow crown and a black eye ring and cheek patch and throat



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



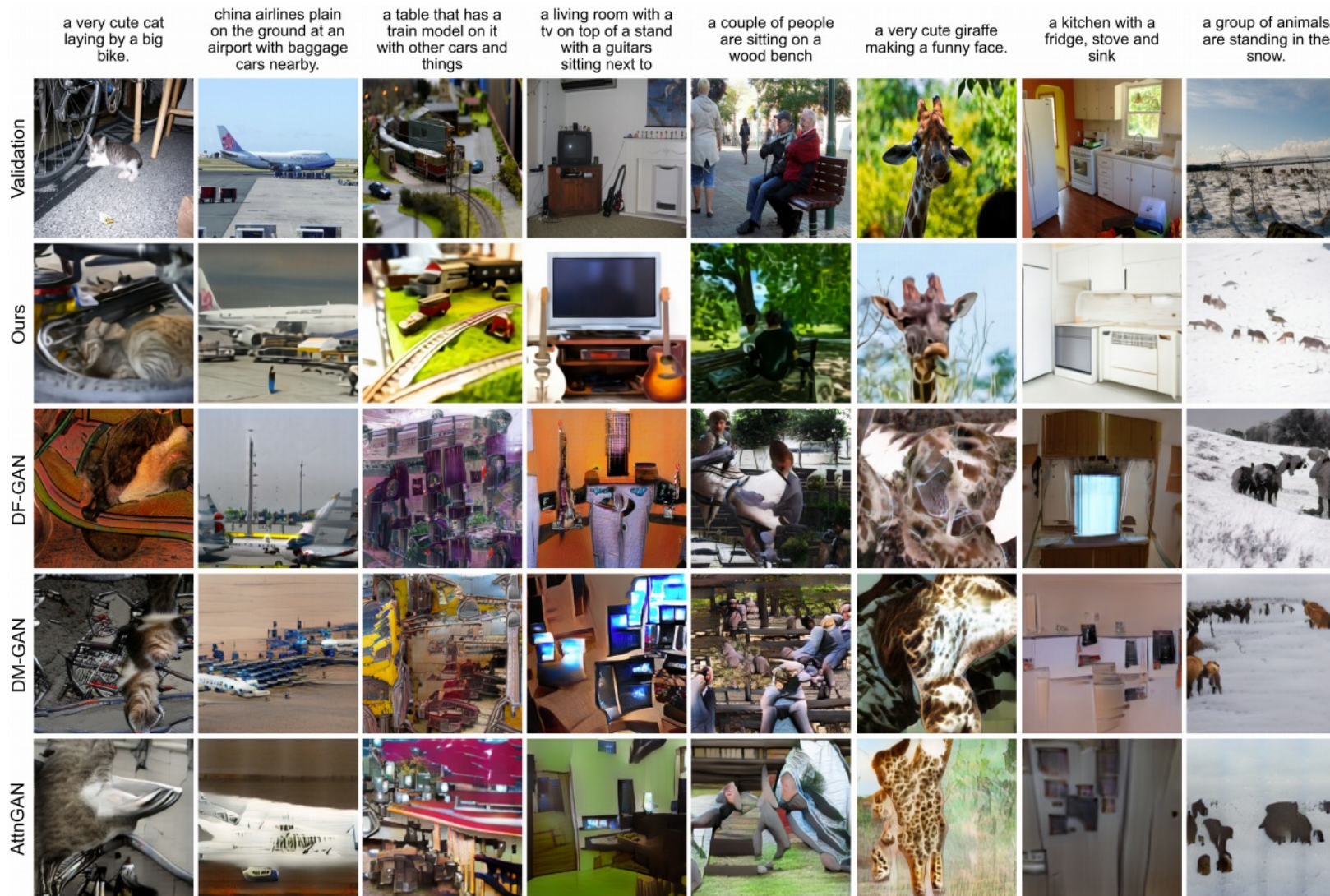
small bird with a pale yellow underside light brown crown and back gray tail and wing tips tip of tail feather bright yellow black eyes and black stripe over eyes



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



DALL-E: Results



Summary of Generative models

- Autoregressive models

$$p_{\theta}(x) = \prod_{i=1}^N p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

- PixelRNNs/CNNs, Image Transformers, ...
- Directly maximize likelihood of training data

- VAEs (Variational autoencoders)

- Probabilistic version of autoencoder to allow for sampling
- Introduces a latent z (assumed to be Gaussian) and maximizes a lower bound

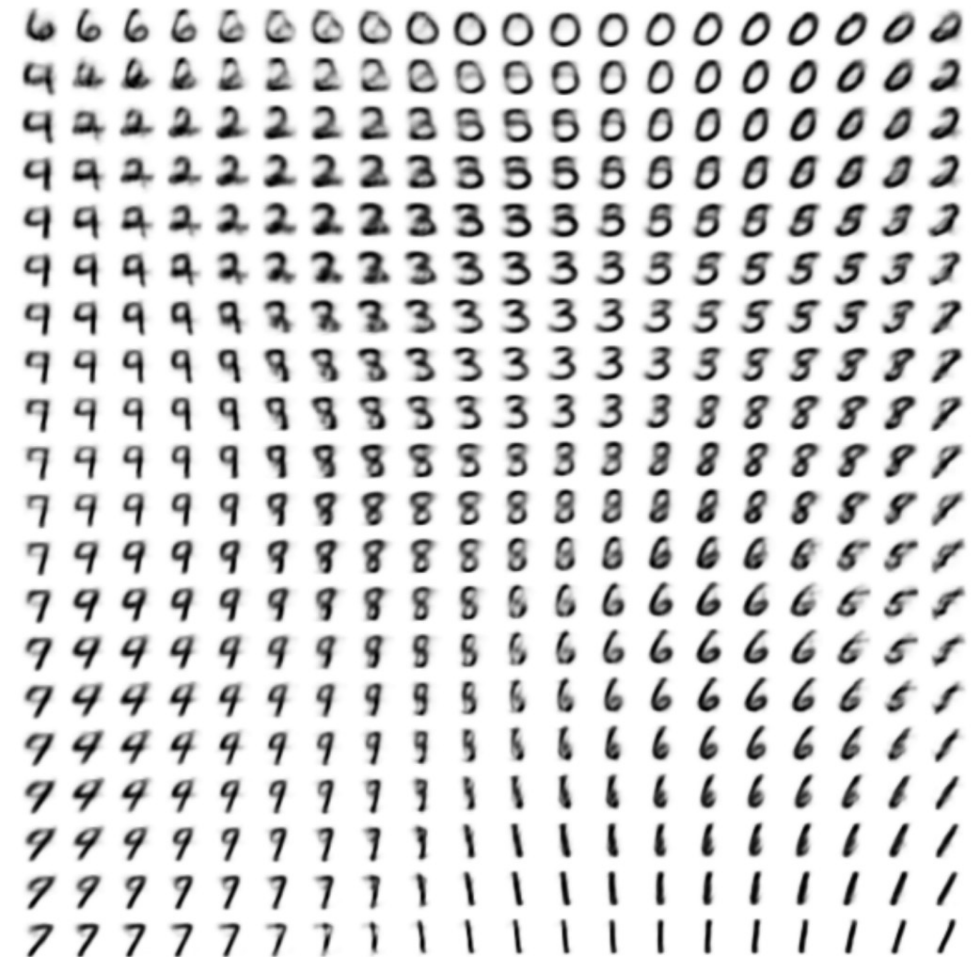
$$p_{\theta}(x) = \int_{\mathbf{z}} p_{\theta}(x|\mathbf{z})p(\mathbf{z})d\mathbf{z} \geq E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(x|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}), p(\mathbf{z}))$$

- GANs (Generative adversarial networks)

- Don't bother modeling $p(x)$, just try to sample from $p(x)$
- Generator + Discriminator (is it real or generated)

How good are these models?

- Is this model generating new and novel images? Or is it just retrieving from training data?
- Comparison of kNN retrieval vs images generated by the model.
- Showing of disentangled latent space, interpolation, and operations on the space.



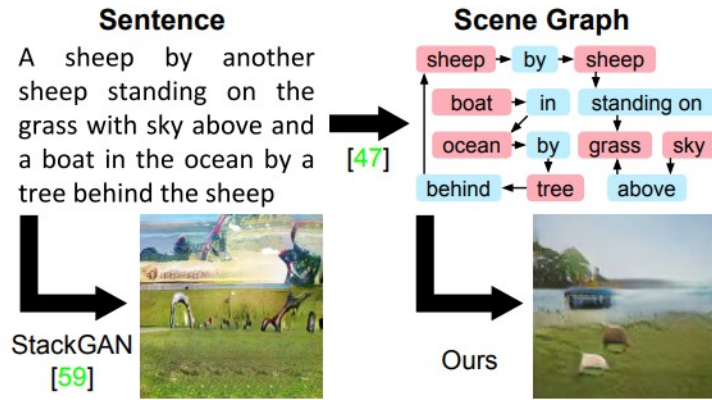
Auto-Encoding Variational Bayes

<https://arxiv.org/pdf/1312.6114.pdf>,

Kingma and Welling, ICLR 2014

Structured content
generation using language

Image generation from scene graphs



- Convert text to scene graph
- Layout objects to preserve relationships

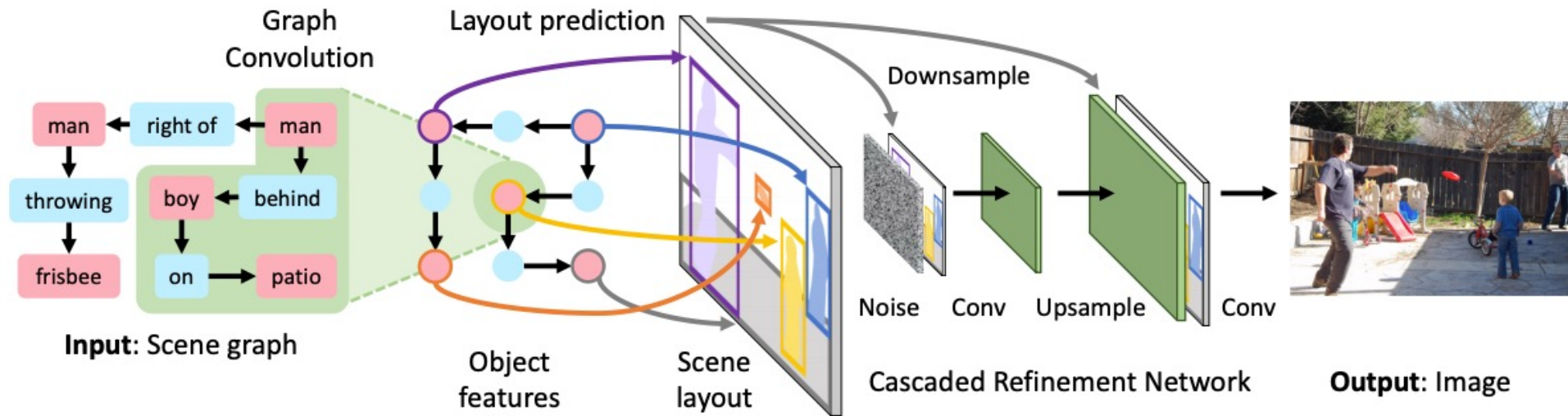


Image Generation from Scene Graphs

<https://arxiv.org/pdf/1804.01622.pdf>, Johnson et al, CVPR 2018

Image generation from scene graphs

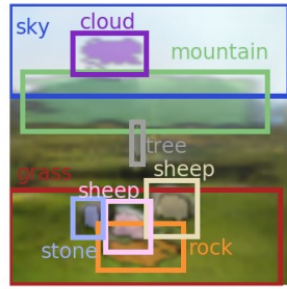
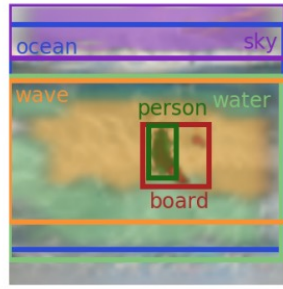

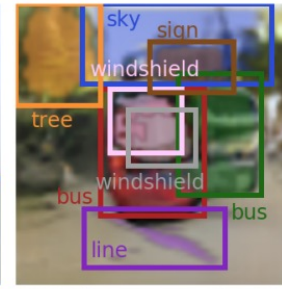
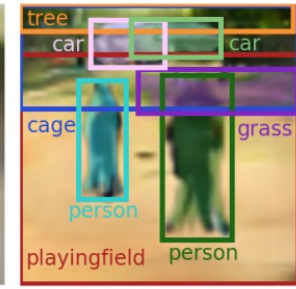
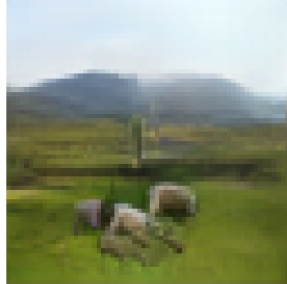
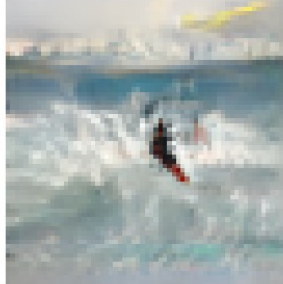
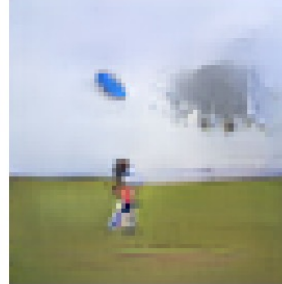


Graph	<pre> graph TD sky[sky] -- has --> cloud[cloud] cloud --> mountain[mountain] sheep1[sheep] -- eating --> grass1[grass] sheep2[sheep] -- eating --> grass2[grass] stone[stone] -- in front of --> tree[tree] rock[rock] -- behind --> tree </pre>	<pre> graph TD cloud[cloud] --> sky[sky] sky -- above --> water[water] person[person] -- by --> water riding1[riding] --> wave[wave] riding2[riding] --> board[board] background[background] --> edge[edge] </pre>	<pre> graph TD boy[boy] -- on top of --> grass[grass] looking[looking at] --> kite[kite] standing[standing on] --> brick[brick] sky[sky] --> field[field] field -- under --> mountain[mountain] </pre>	<pre> graph TD building[building] -- next to --> bus1[bus] bus1 -- behind --> tree[tree] bus2[bus] -- behind --> tree sign[sign] -- has --> windshield1[windshield] windshield1 --> windshield2[windshield] </pre>	<pre> graph TD car[car] -- left of --> car2[car] car2 -- above --> grass[grass] cage[cage] -- above --> person1[person] person1 -- left of --> grass person2[person] -- left of --> tree[tree] playingfield[playingfield] -- below --> person1 playingfield -- below --> person2 </pre>
Text	<p>Two sheep, one eating grass with a tree in front of a mountain; the sky has a cloud.</p>	<p>A person riding a wave and a board by the water with sky above.</p>	<p>A boy standing on grass looking at a kite and the sky with the field under a mountain</p>	<p>Two busses, one behind the other and a tree behind the second; both busses have windshields.</p>	<p>A person above a playingfield and left of another person left of grass, with a car left of a car above the grass.</p>
Layout					
Image					

Image Generation from Scene Graphs

<https://arxiv.org/pdf/1804.01622.pdf>, Johnson et al, CVPR 2018

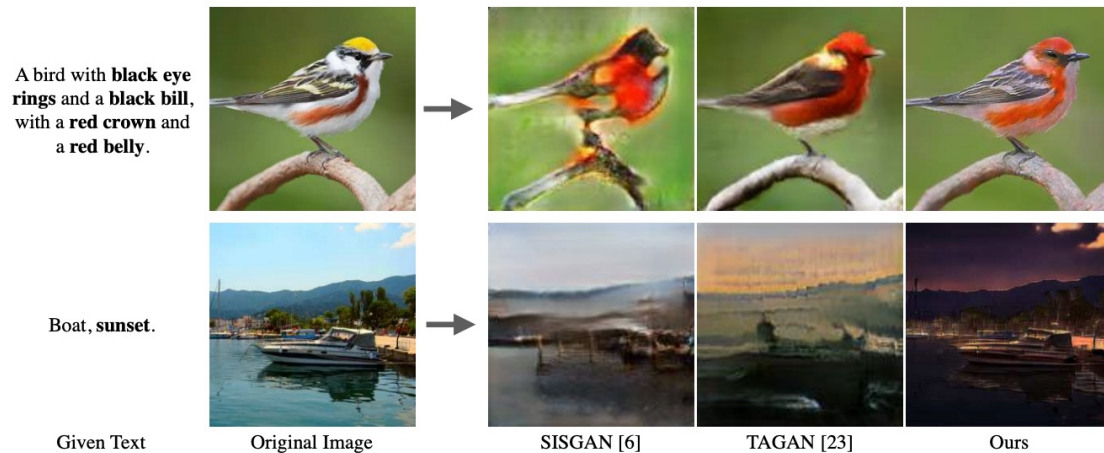
Content manipulation using language

Content manipulation

- Similarities to instruction following
- Structured representation + well defined operations
 - Semantic parsing!
- Unstructured representation + learned operations
 - Vector representation
 - Operation = some transformation on the encoded representation
 - Decode into image or shape

Content manipulation

ManiGAN



ManiGAN: Text-Guided Image Manipulation
<https://arxiv.org/pdf/1912.06203.pdf>
Li et al, 2020

RefineGAN

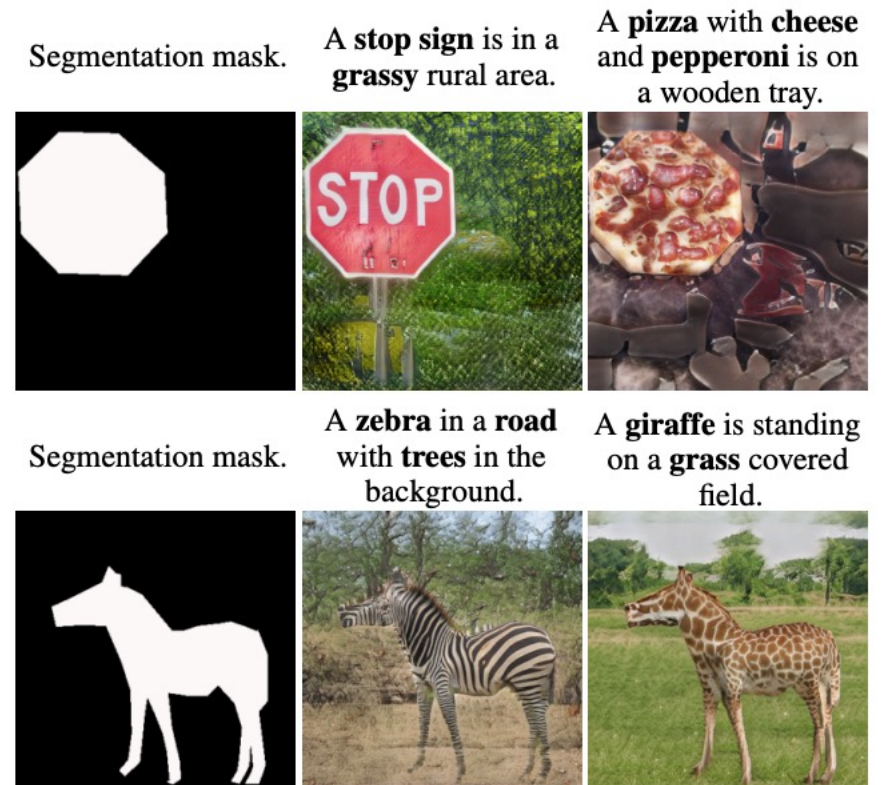


Image-to-Image Translation with Text Guidance
<https://arxiv.org/pdf/2002.05235.pdf>
Li et al, 2020

Text-to-image generation
with additional input

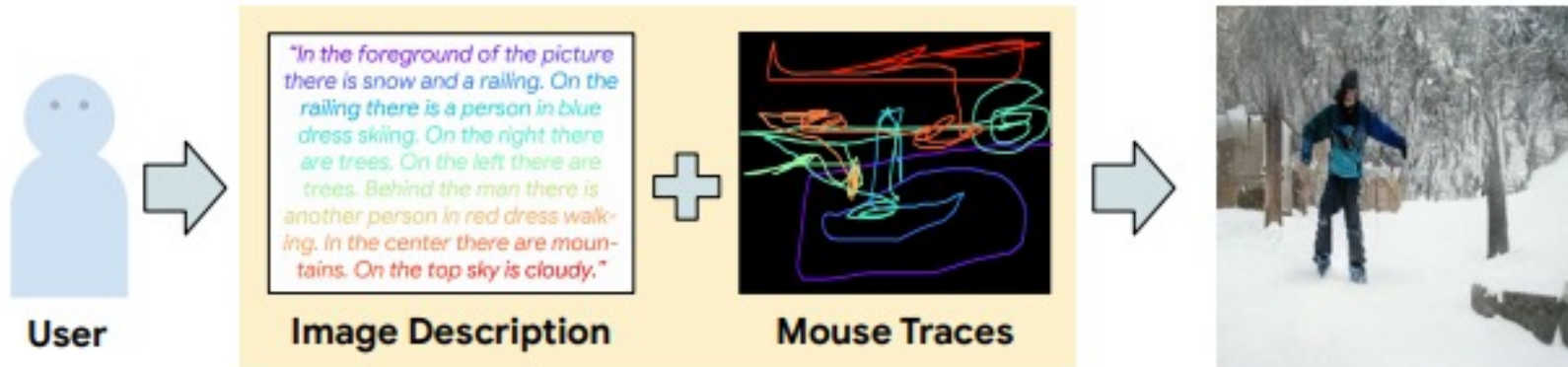
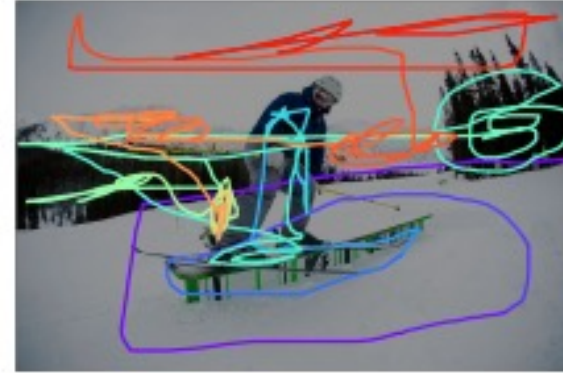
Text-to-Image generation with mouse traces

"A man skiing along a rail on the snowy hill."

Standard Caption

"In the foreground of the picture there is snow and a railing. On the railing there is a person in blue dress skiing. On the right there are trees. On the left there are trees. Behind the man there is another person in red dress walking. In the center there are mountains. On the top sky is cloudy."

Localized Narrative



Text-to-Image Generation Grounded by Fine-Grained User Attention
<https://arxiv.org/pdf/2011.03775.pdf>, Koh et al, 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	✓	✓	✓	✓	✓	✓	✓	✓	✓	

Text to 3D

State of text to 3D

- Much less work in text to 3D
- 3D generation less explored than 2D
 - Less 3D data → with more 3D data, 3D deep learning + 3D generation is more popular
 - Less developed methods
 - Choices of representation to use
- 2D or 3D, basic families of methods are the same
 - Details are tricky
- 3D generation broken down into two levels:
 - Scenes vs Shapes
- See survey paper for more on 3D generation

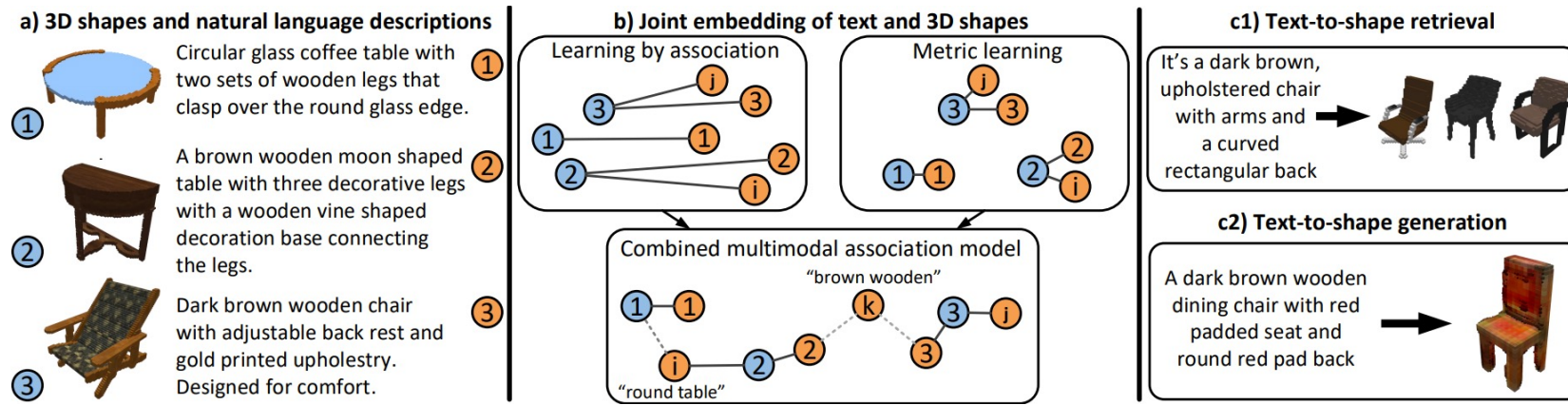
Shape vs Scene Generation

- Shape Generation typically treated as generating voxels/points/triangles directly
- Scene Generation typically treated as collection of objects (shapes): select/retrieve objects from DB and arrange them (position, rotation, scale)
- Blurry line
 - Can treat shape generation as generating + assembling parts
 - Grammar/program
 - Can treat scene generation as generating one big mesh (3D reconstruction)
- Not as emphasized in 2D (since output is just an image), but the two options also exist
 - Put together objects into an 2D scene or just generate pixels

Choice of output for shapes

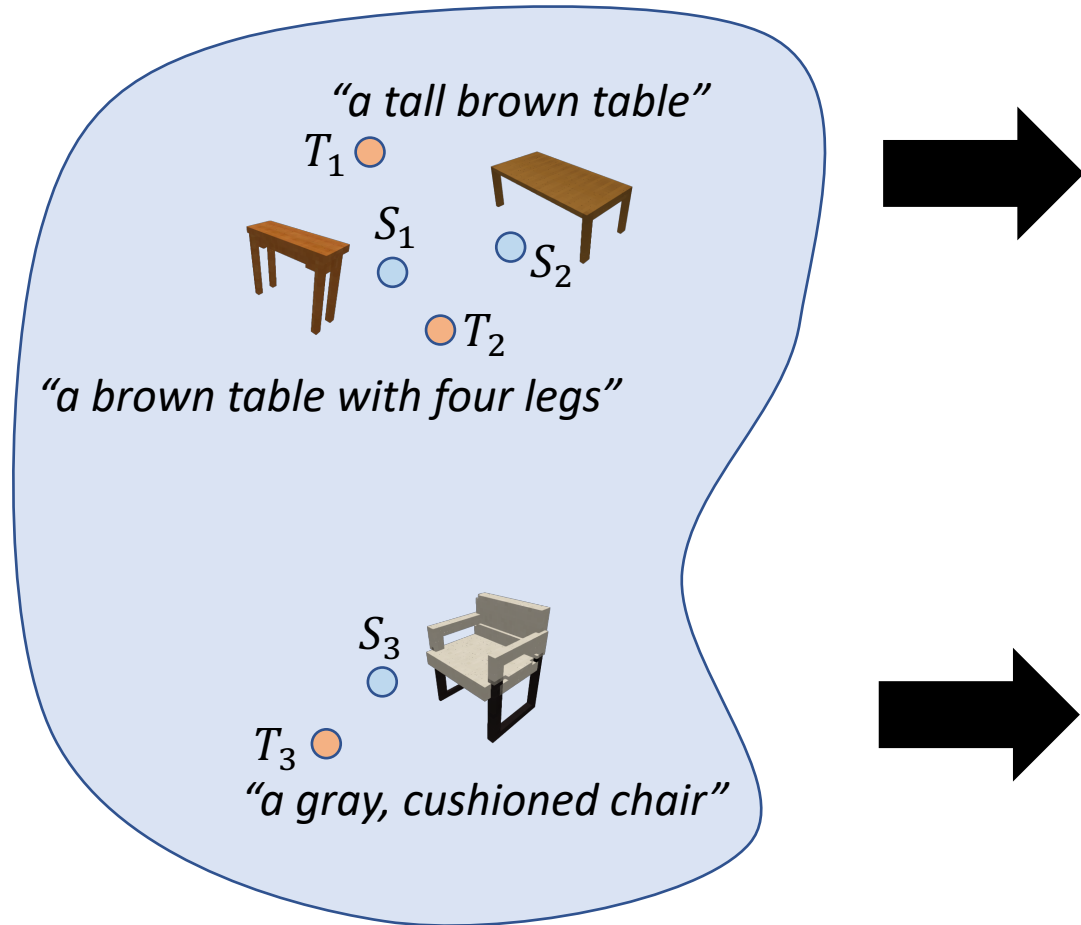
- Voxels: Direct analogue of 2D pixels, use convolutions
 - High resolution challenging
 - Dense voxels are expensive (N^3 vs N^2)
 - Work on sparse voxel representations (lots of empty space!)
- Point clouds
 - Does not capture topology
- Mesh: Traditionally what is used in graphics
 - Trickier to work work:
 - Use graph representation to predict triangle vertices and edges
 - Can go from Point clouds/Voxels to Mesh (use traditional methods)
 - Can also go from implicit surfaces to Mesh
 - For each point: predict if inside or outside of Mesh

Text2Shape: Generating Shapes from Natural Language by Learning Joint Embeddings

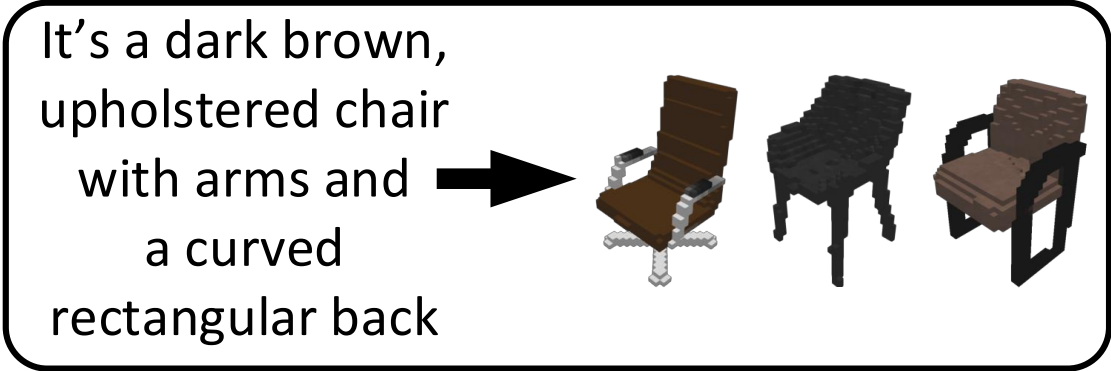


Kevin Chen, Christopher B. Choy, Manolis Savva, Angel Chang,
 Thomas Funkhouser, Silvio Savarese
 ACCV, 2018

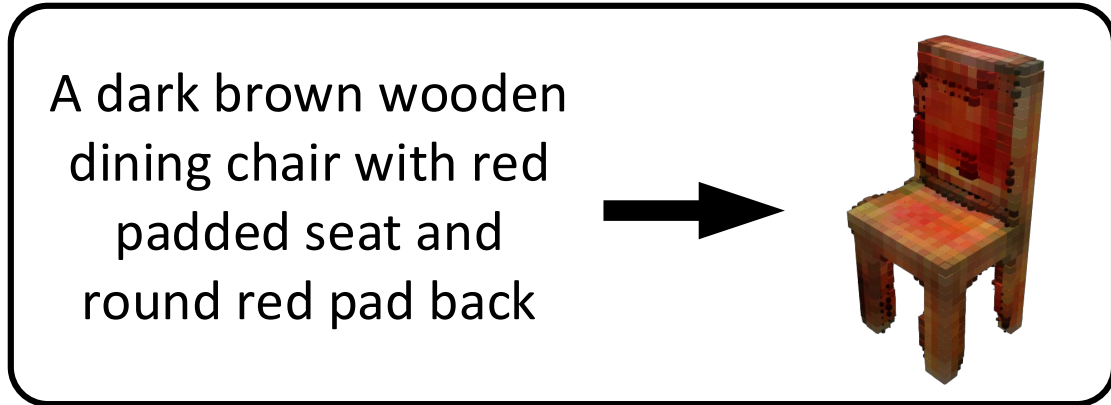
Text + Shape Joint Embedding



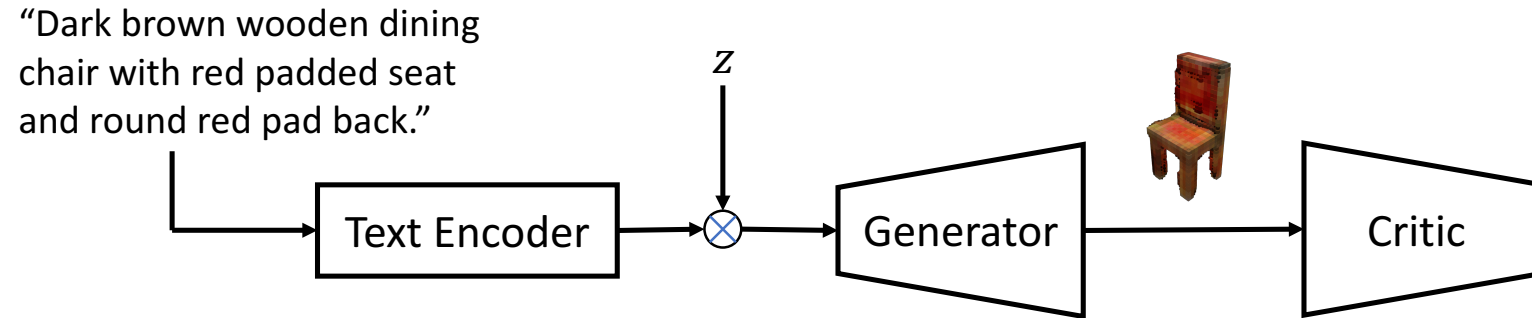
Text-to-shape retrieval



Text-to-shape generation



Text-to-shape generation



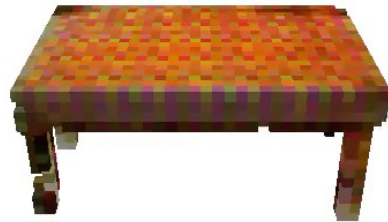
1. Encoder maps text description to latent space
2. Description embedding is concatenated with noise vector
3. Generator generates a plausible colored shape
4. Critic evaluates quality of generation

Text-to-shape generation

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



GAN-INT-CLS[1]



Ours CGAN



Ours CWGAN



Ground Truth

Text-to-shape generation

Input: Waiting room chair leather



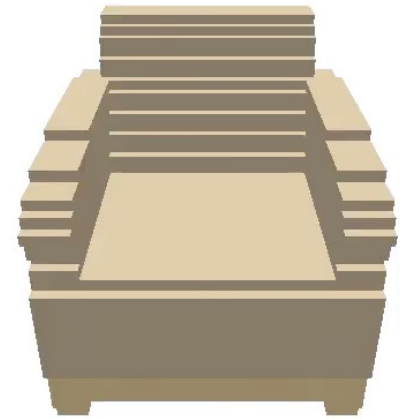
GAN-INT-CLS[1]



Ours CGAN



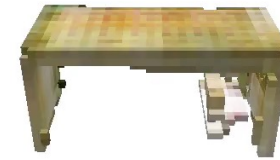
Ours CWGAN



Ground Truth

Shape manipulation

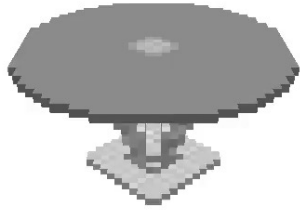
white table - white + brown =



gray table - gray + glass =



Shape manipulation



– round + rectangular =



– chair + table =



Text2Shape status

- Work from 2016-2017, published in 2018
- GAN based
- Voxels based
- Lots of improvements in generative models and shape generation since then!

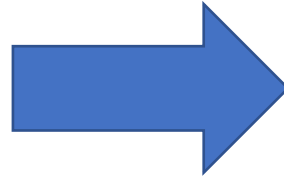
Text2Scene

There is a desk and a notepad on the desk.
A pen is next to the notepad.



How do we handle natural, underspecified language?

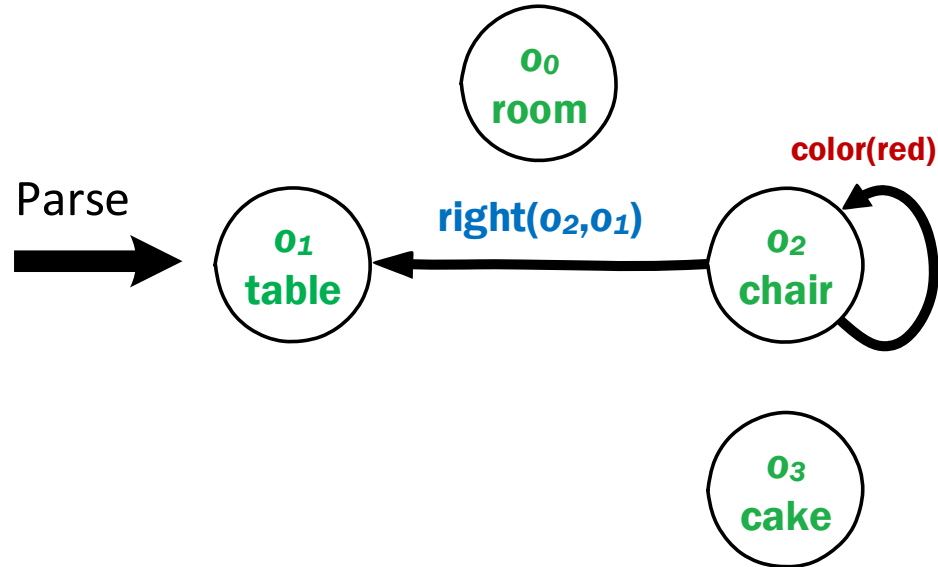
“Living room
with a red
couch”



- learn **common sense priors** on how objects are arranged in the real world
- view scene description as **constraints** on the scene

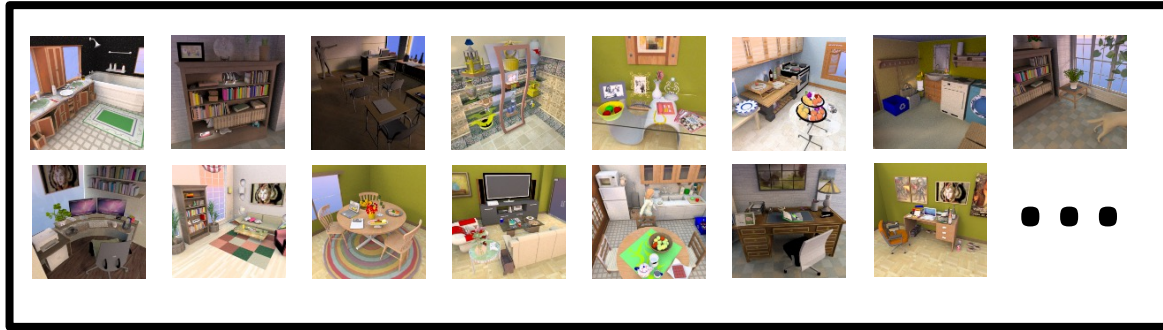
Language as constraint for 3D scene graphs

“There is a *room* with
a *table* and a *cake*.
There is a *red chair* to
the *right* of the *table*.”

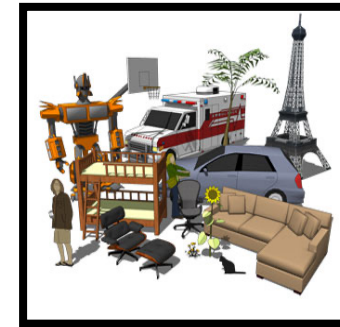


objects, attributes and relations

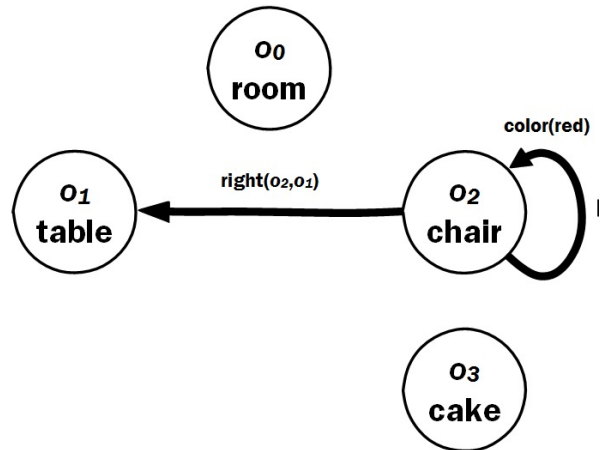
Scene Database



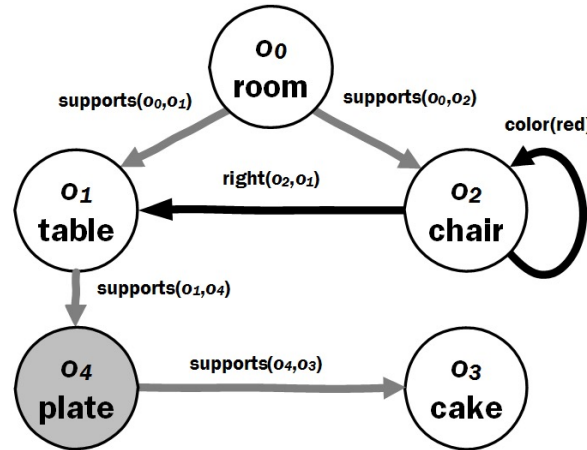
3D Models



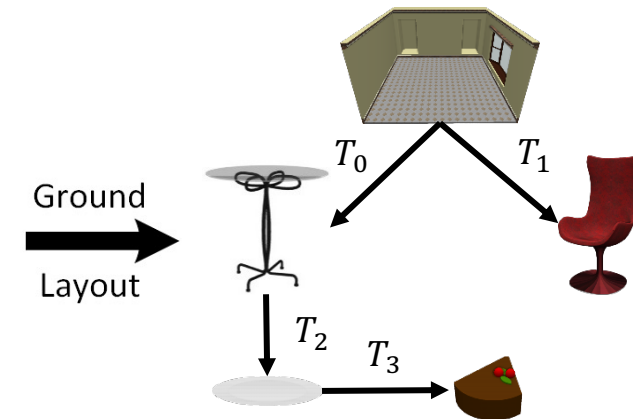
a) Explicit Constraints



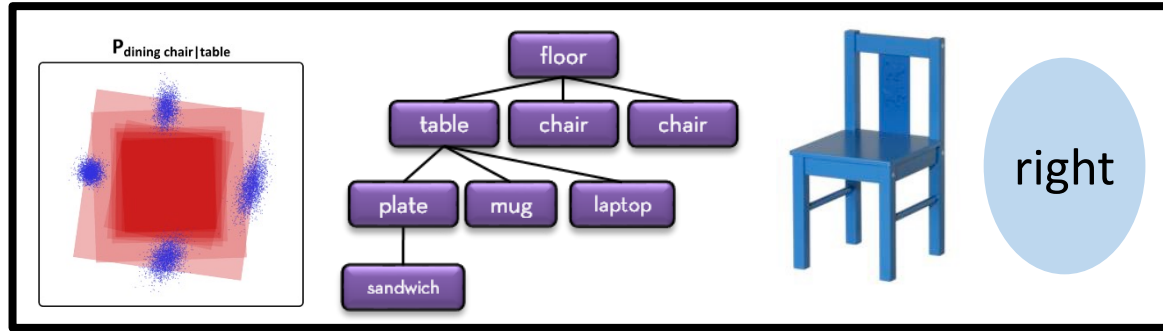
b) Inferred Scene Template



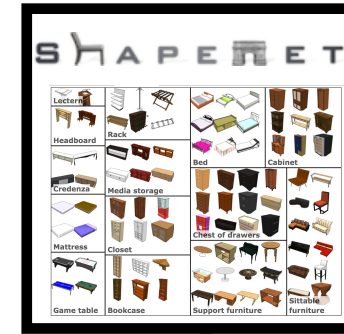
c) Geometric Scene



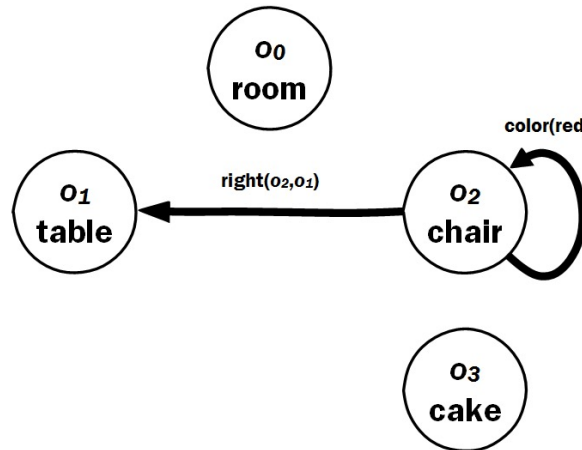
Spatial Knowledge Base



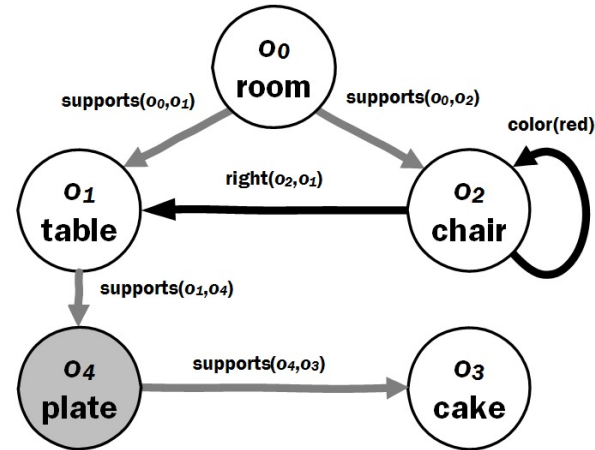
3D Models



a) Explicit Constraints



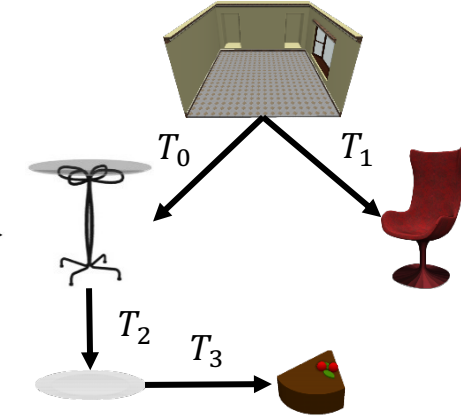
b) Inferred Scene Template



Infer

Ground
Layout

c) Geometric Scene

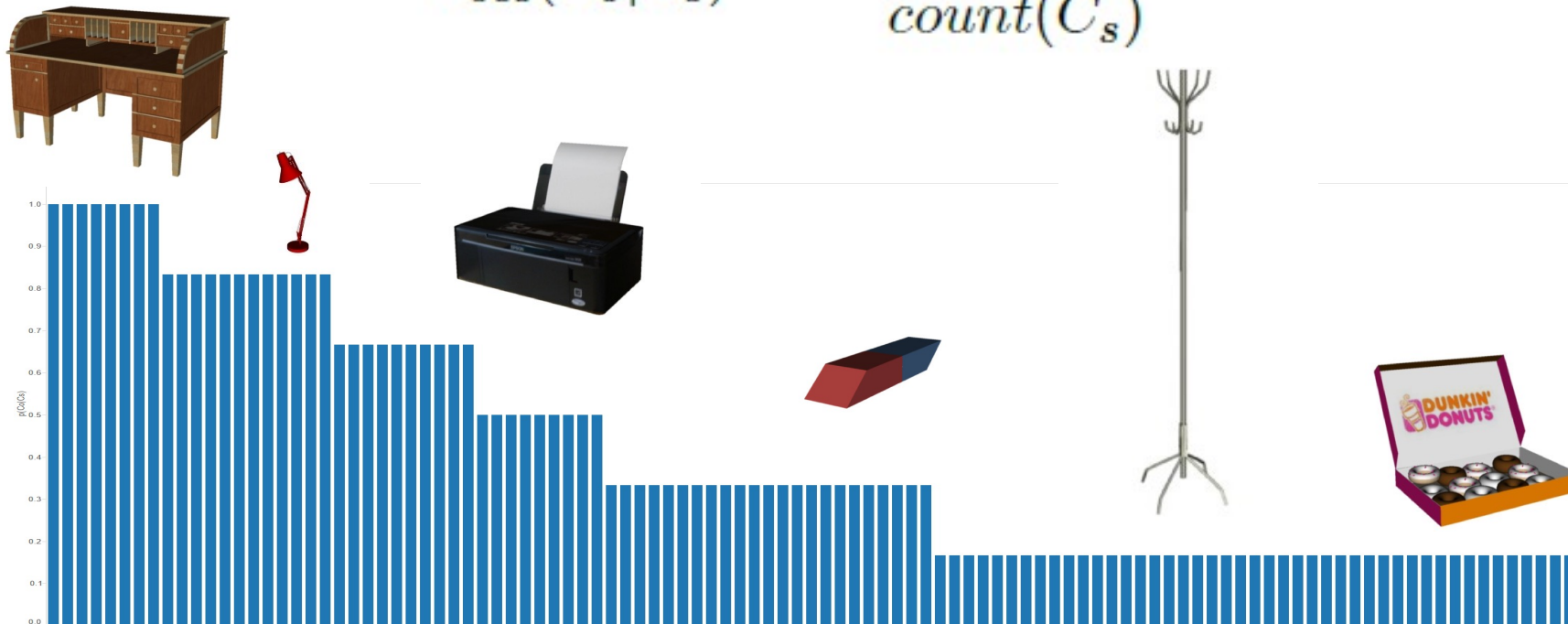


Object occurrences

What goes in an office?

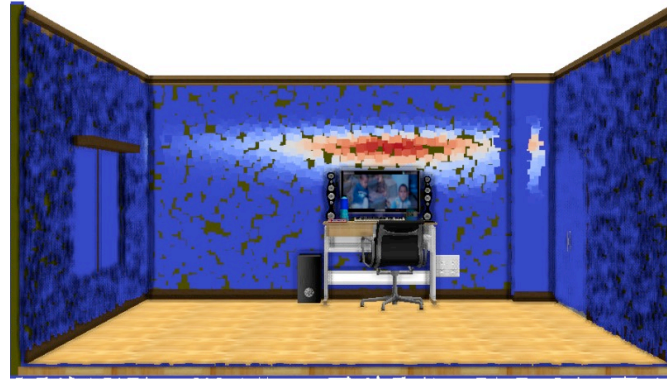
Probability that object of category C_o is found in scene type C_s

$$P_{occ}(C_o|C_s) = \frac{\text{count}(C_o \text{ in } C_s)}{\text{count}(C_s)}$$



Semantic queries – Where can X go?

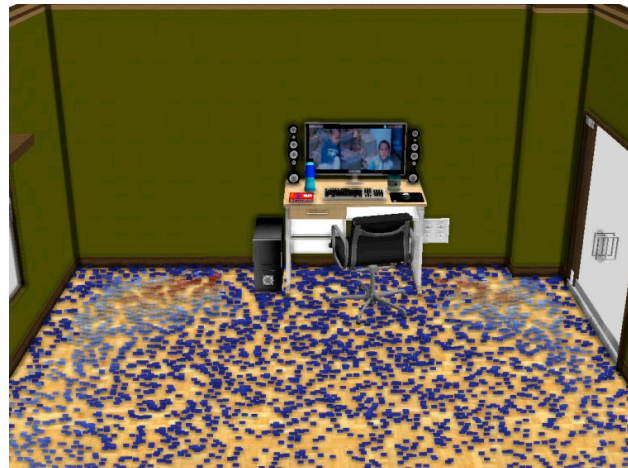
poster



rug



floor lamp



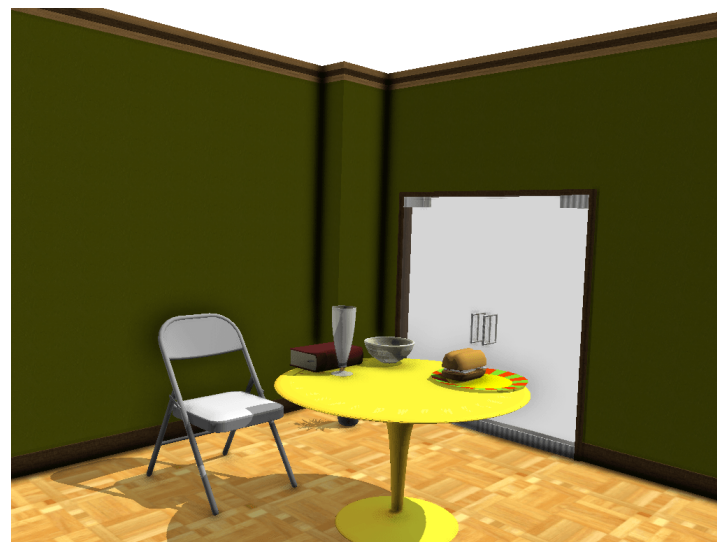
hat



There is a sandwich on a plate



There is a sandwich on a plate



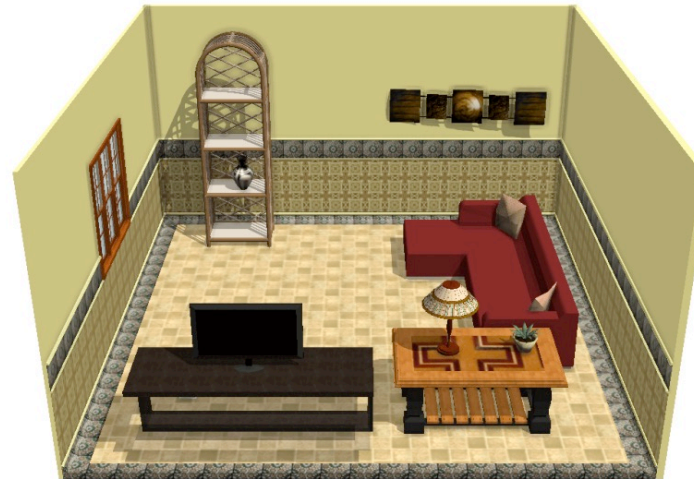
There is a desk and a chair



There is a computer desk with a red chair



“There is a living room with a red couch and a TV.”



Input Description

“There is a living room with a red couch and a TV”



Candidate Scenes



System Learning

...

More manipulation ...

Interactive Scene Manipulation



Manipulation

“Put a clock on the wall”

Scene refinement



“Put a cup on the bookcase.”

Scene refinement



“Put a clock on the wall.”

Scene refinement



“Put a painting on the wall.”

Scene refinement



“Look at the painting.”

Followup work

- Retrieve and edit approach

There is a desk with two monitors.



The desk is messy.



Replace the desk and monitors.

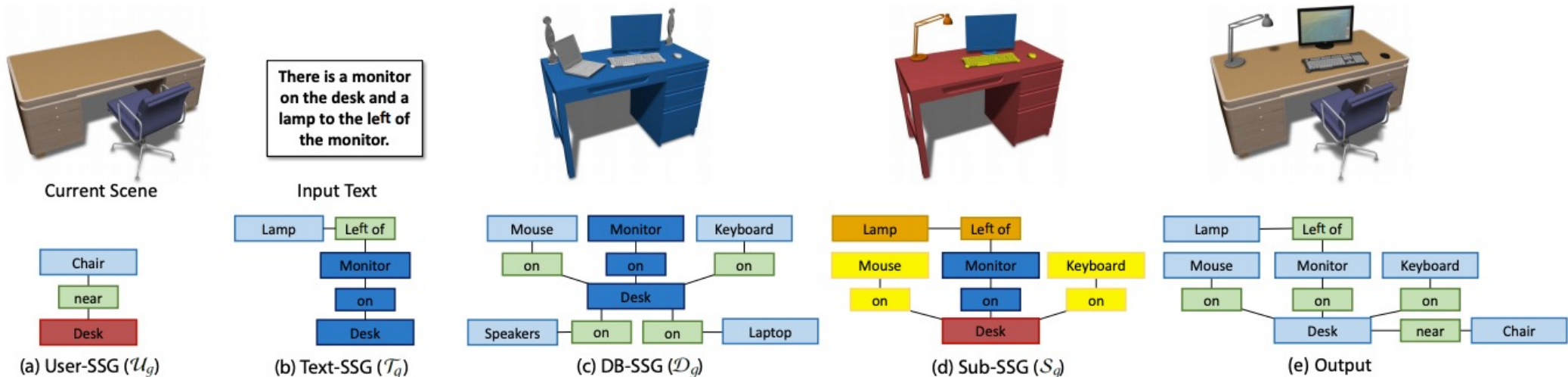


Language-Driven Synthesis of 3D Scenes from Scene Databases
https://manyili12345.github.io/Publication/2018/T2S/t2s_final.pdf

Ma et al, Siggraph Asia 2018

Followup work

- Retrieve and edit approach
- Also uses semantic scene graphs (extracted from larger datasets)
- Parse text, retrieve matching sub-scenes, edit to match desired semantic scene graph



Language-Driven Synthesis of 3D Scenes from Scene Databases
https://manyili12345.github.io/Publication/2018/T2S/t2s_final.pdf

Ma et al, Siggraph Asia 2018

Text to 3D scene status

- Work from 2013-2015, some work in 2018.
- Pre-deep learning, used probabilistic graphical models
- Lots of improvements to object retrieval, scene generation, text interpretation since then!
- Larger, more realistic 3D datasets
- Still, limited by the availability of text + 3D data

Next time

- Thursday (4/15): Last day – project discussion and conclusion
 - Watch other group's project video before class
 - Project video due by 11:59pm 4/14
 - Project report due by 11:59pm 4/15

