

Unsupervised Learning & Generative AI

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

1 Unsupervised Learning

- Text Models
- Image Models

2 ChatGPT

3 Autoencoders (AE)

- Manifold Learning*

4 Variational Autoencoders (VAE)

5 Flow-based Models

6 Diffusion Models

Outline

1 Unsupervised Learning

- Text Models
- Image Models

2 ChatGPT

3 Autoencoders (AE)

- Manifold Learning*

4 Variational Autoencoders (VAE)

5 Flow-based Models

6 Diffusion Models

Unsupervised Learning

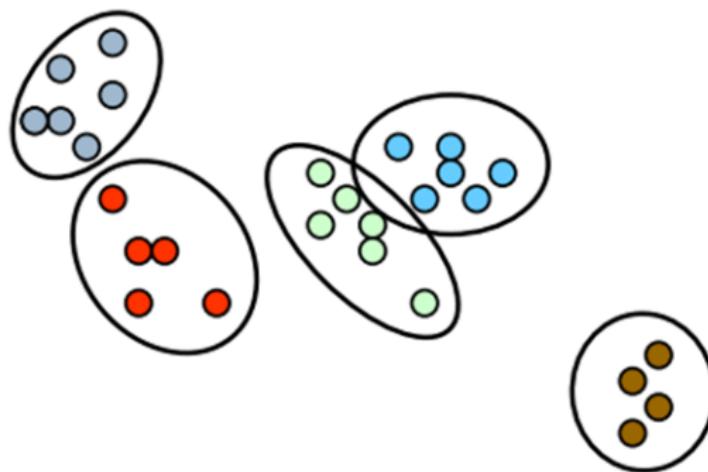
- Dataset: $\mathbb{X} = \{\mathbf{x}^{(i)}\}_i$, where $\mathbf{x}^{(i)}$'s are i.i.d. samples of \mathbf{x}
 - No supervision $\mathbf{y}^{(i)}$ (labels)
- What can we learn without labels?

Unsupervised Learning

- Dataset: $\mathbb{X} = \{\mathbf{x}^{(i)}\}_i$, where $\mathbf{x}^{(i)}$'s are i.i.d. samples of \mathbf{x}
 - No supervision $\mathbf{y}^{(i)}$ (labels)
- What can we learn without labels? The ***structures*** in \mathbb{X}
 - Inter-sample structures
 - Intra-sample structures

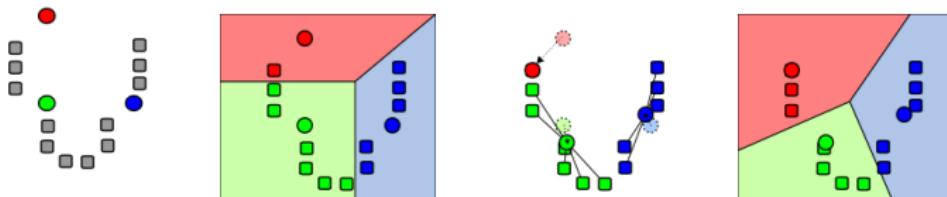
Clustering I

- Goal: to divide $x^{(i)}$'s into K groups/**clusters**
 - Based on some similarity/distance measure between $x^{(i)}$ and $x^{(j)}$



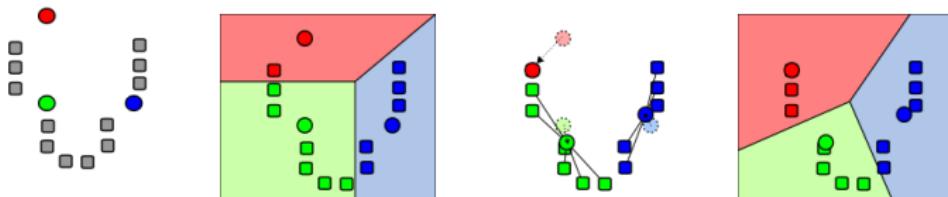
Clustering II

- K -means algorithm (K fixed): iteratively move clusterheads until convergence

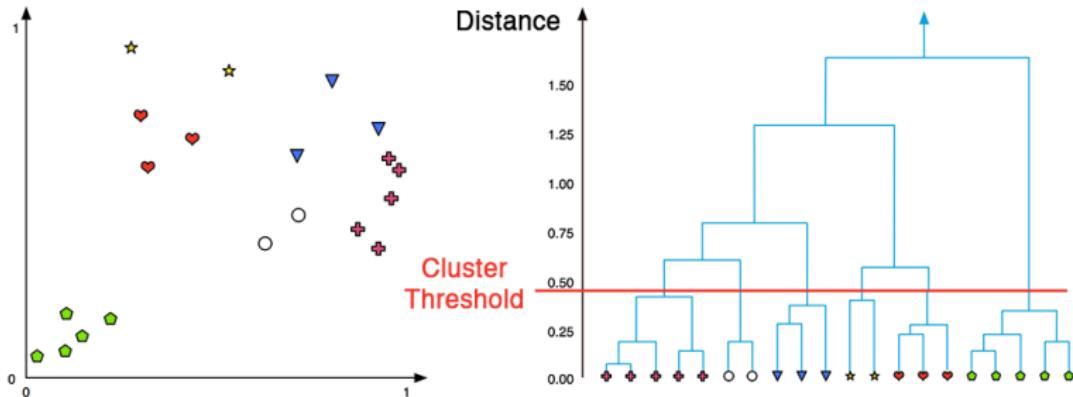


Clustering II

- K -means algorithm (K fixed): iteratively move clusterheads until convergence

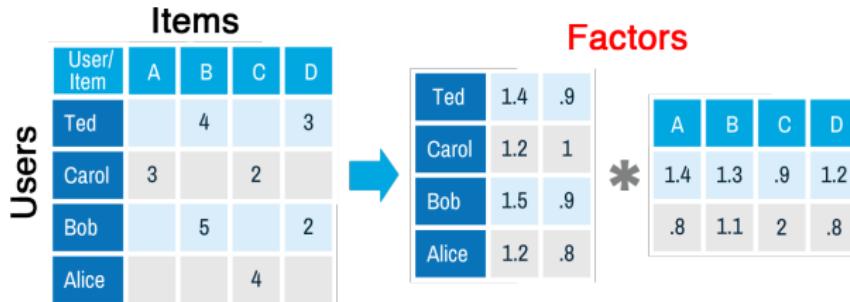


- Hierarchical clustering (variable K): iteratively merge two points/groups, then cut



Factorization for Tabular \mathbb{X}

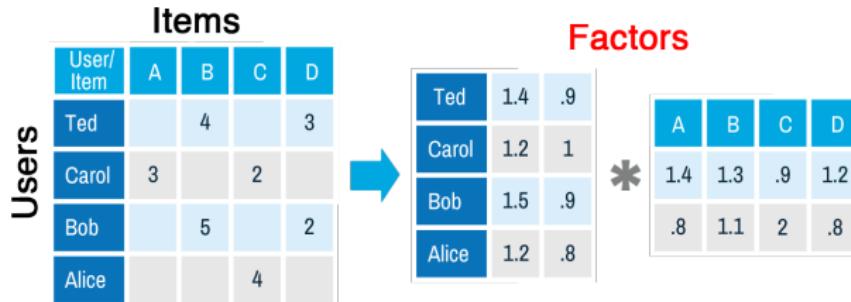
- Let X be a rating matrix where $X_{i,:} = \mathbf{x}^{(i)}$



- Goal: to approximate X with a dense matrix $\hat{X} = \mathbf{WH}$
 - To make recommendations based on $\hat{X}_{i,:}$ for user i , known as **collaborative filtering**
- How?

Factorization for Tabular \mathbb{X}

- Let X be a rating matrix where $X_{i,:} = \mathbf{x}^{(i)}$



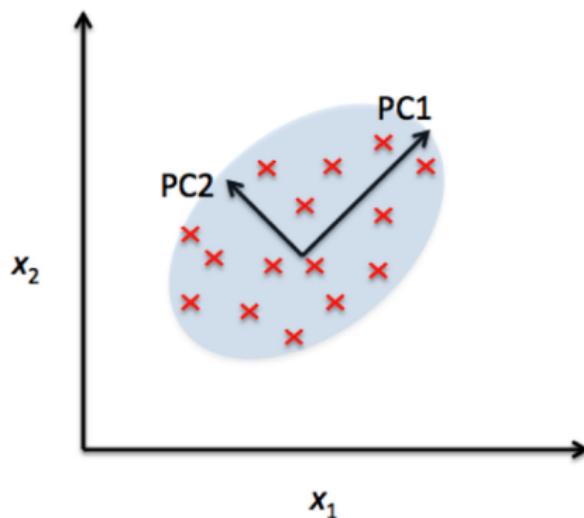
- Goal: to approximate X with a dense matrix $\hat{X} = \mathbf{WH}$
 - To make recommendations based on $\hat{X}_{i,:}$ for user i , known as **collaborative filtering**
- How? Non-negative matrix factorization (NMF) [11, 12]:

$$\arg \min_{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}} \|\mathbf{X} - \mathbf{WH}\|_F$$

- So, positive elements in \mathbf{W} and \mathbf{H} can be seen as **factor** degrees

Dimension Reduction

- Goal: to learn a low dimensional representation \mathbf{z} of \mathbf{x}
 - E.g., PCA



Self-Supervised Learning

- Goal: to learn a model that is able to “fill in the blanks”



Self-Supervised Learning

- Goal: to learn a model that is able to “fill in the blanks”
- Links unsupervised tasks with supervised models
 - Much more training data for models
 - Representation learning with deep networks



$x^{(i)}$

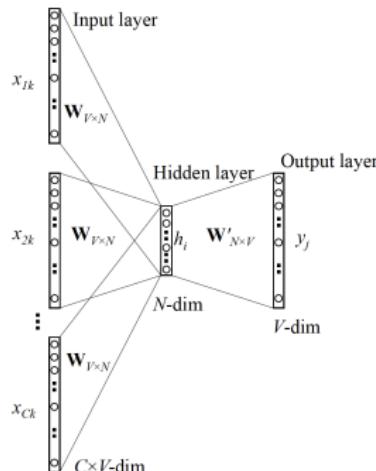
$y^{(i)}$

Example I: Word2Vec

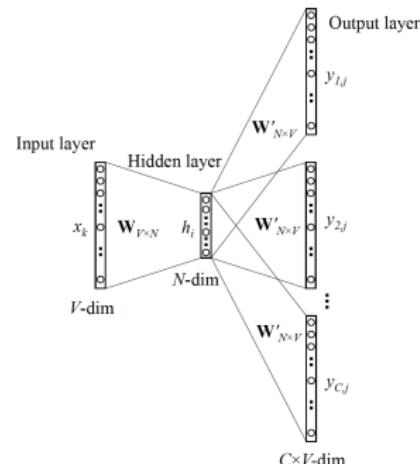
- Goal: to learn a model for blank filling

Example I: Word2Vec

- Goal: to learn a model for blank filling
 - E.g., word2vec [17, 16]: "... the cat sat on..."



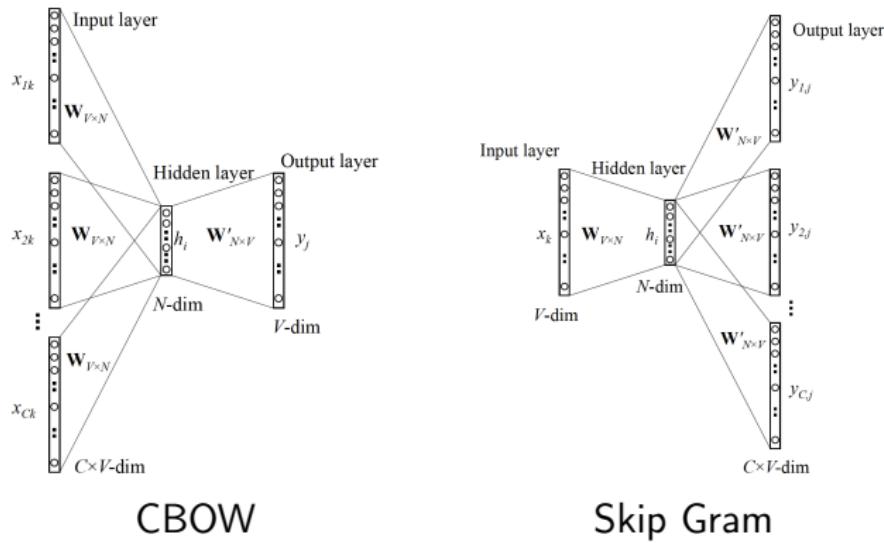
CBOW



Skip Gram

Example I: Word2Vec

- Goal: to learn a model for blank filling
 - E.g., word2vec [17, 16]: "... the cat sat on..."



- Latent representation \mathbf{h} encodes the ***semantics*** of a word
 - No need for synonym dictionary; big data tell that already

Example II: Doc2Vec

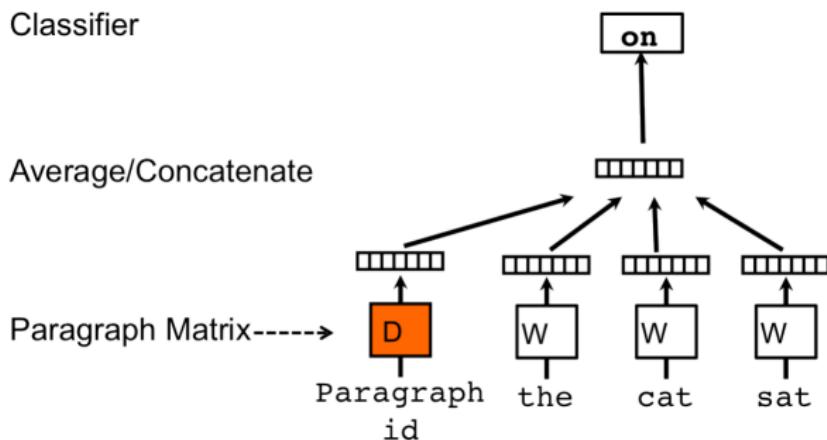
- How to encode a document?

Example II: Doc2Vec

- How to encode a document?
- Bag of words (TF-IDF), average word2vec, etc.
 - Do **not** capture the semantics due to sentence/paragraph/doc structure
 - "*John likes Mary*" \neq "*Mary likes John*"

Example II: Doc2Vec

- How to encode a document?
- Bag of words (TF-IDF), average word2vec, etc.
 - Do **not** capture the semantics due to sentence/paragraph/doc structure
 - “John likes Mary” ≠ “Mary likes John”
- Why not apply self-supervised learning to docs?
 - Doc2vec [10]: to capture the **context** not explained by words
 - **Transductive** rather than inductive; does not work with unseen docs



Generative Models

- Goal: to generate new samples of x
 - Can be conditioned on instructions (input)
- Largely based on self-supervised learning

The image is a collage of four screenshots of AI interfaces:

- DALL-E 2:** A dark-themed interface with the text "DALL-E 2" and "DALL-E is an AI system that can create realistic images and art from descriptions or natural language".
- ChatGPT:** A white-themed interface with the text "Welcome to ChatGPT" and "Log in with your OpenAI account to continue". It features a large "G" logo.
- Adobe Firefly:** A white-themed interface with the text "Text to image" and "Adobe Firefly". It shows examples of generated images like a car and a landscape.
- Bard:** A white-themed interface with the text "Bard can explain why large language models might make mistakes". It features a large "G" logo.

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

Generating Text

- With large (self-supervised) training data, *transformers* [33] have shown to perform better than RNNs and CNNs
 - Mainly due to the $O(1)$ point distance
- Two common text models based on transformer:
 - BERT [4]: non-autoregressive
 - GPT [21]: autoregressive

BERT [4]

- Massively trained *encoder* of the original transformer
 - Non-autoregressive
- Pre-training tasks:
 - Masked language model (“Cloze” task)
 - “A quick brown [MASK] jumps over the lazy dog” → “fox” 11%, “ant” 5%, ...
 - Next sentence prediction
 - “[MASK] go to store [SEP] to buy a [MASK] of milk” → True 93%, False 7%

One Pre-training, Multiple Fine-tuning Tasks

- Special input to identify downstream task: [CLS] token

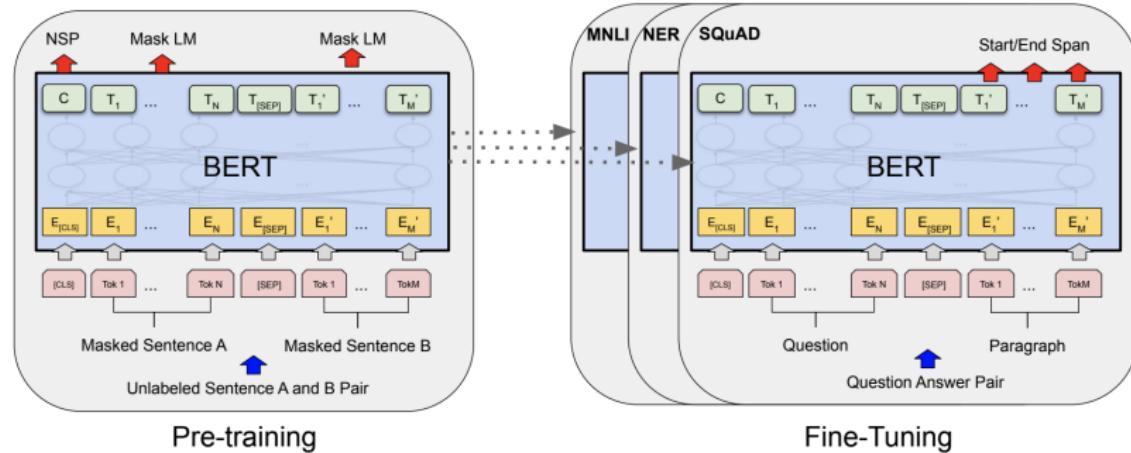
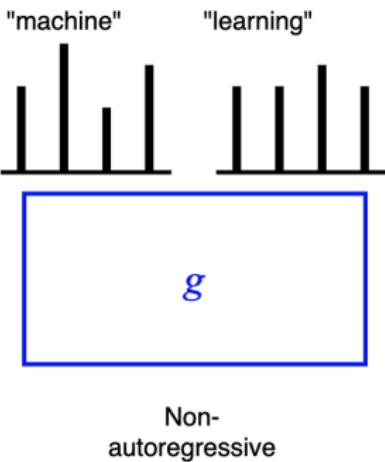
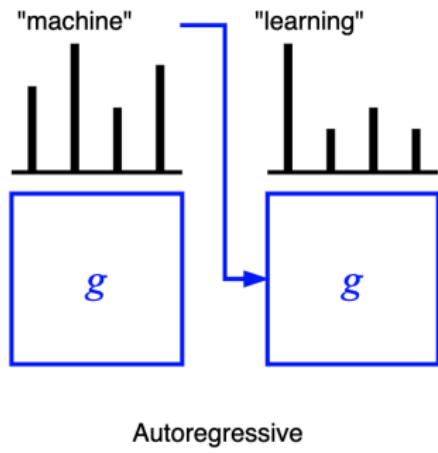


Figure 1: Overall pre-training and fine-tuning procedures for BERT. Apart from output layers, the same architectures are used in both pre-training and fine-tuning. The same pre-trained model parameters are used to initialize models for different down-stream tasks. During fine-tuning, all parameters are fine-tuned. [CLS] is a special symbol added in front of every input example, and [SEP] is a special separator token (e.g. separating questions/answers).

GPT [21]

- Massively trained **decoder** of the original transformer
 - Autoregressive
- Multitask pre-training by maximizing $\text{Pr}(\text{output}|\text{input}, \text{task})$:
 - Translation: $\text{Pr}(\text{french text}|\text{en text}, \text{translation})$
 - Question answering: $\text{Pr}(\text{answer}|\text{question}, \text{qa})$
 - Reading comprehension: $\text{Pr}(\text{answer}|\text{document}, \text{question}, \text{reading})$
- Usage for downstream task: fine-tuning or ***prompting***

Autoregressive or Not?



- It's easier for an autoregressive model to generate coherent text
- In this lecture, we focus on GPT and its variants

Outline

1 Unsupervised Learning

- Text Models
- Image Models

2 ChatGPT

3 Autoencoders (AE)

- Manifold Learning*

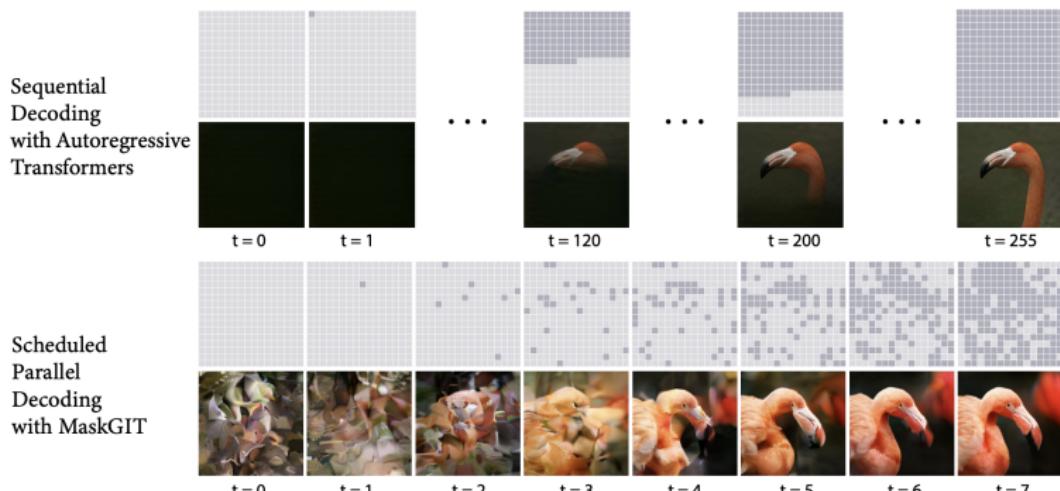
4 Variational Autoencoders (VAE)

5 Flow-based Models

6 Diffusion Models

Autoregressive or Not?

- Autoregressive models can still give good (if not better) performance
 - E.g., PixelRNN [32], PixelCNN [31], or MaskGIT [3]
- But unlike in text domain, generating images pixel-by-pixel is **very slow**
- Speed-up?
 - Scheduled **parallel-pixel** generation (e.g., MaskGIT)
 - **Stepwise** generation of **whole** images (e.g., flow-based and diffusion models)



Whole-Image Generation I

- Goal: given \mathbb{X} , to learn a *generator function* g such that $\hat{\mathbf{x}} = g(\mathbf{c}; \Theta_g)$ looks like a real image in \mathbb{X}
 - \mathbf{c} is a code or condition
 - Θ_g represents parameters of g
- Objective from Information Theory perspective:

$$\arg \min_g D(P_{\text{data}} \| P_g)$$

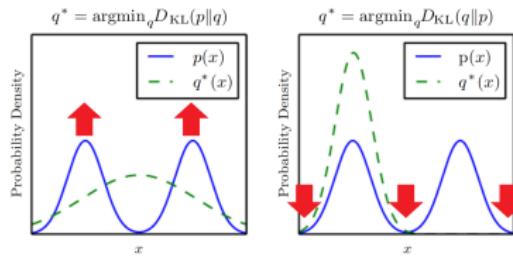
- P_{data} is the distribution of real data in the ground truth
- P_g is the distribution of generated data
- D is a divergence measure, e.g., D_{KL}

Whole-Image Generation I

- Goal: given \mathbb{X} , to learn a **generator function g** such that $\hat{\mathbf{x}} = g(\mathbf{c}; \Theta_g)$ looks like a real image in \mathbb{X}
 - \mathbf{c} is a code or condition
 - Θ_g represents parameters of g
- Objective from Information Theory perspective:

$$\arg \min_g D(P_{\text{data}} \| P_g)$$

- P_{data} is the distribution of real data in the ground truth
 - P_g is the distribution of generated data
 - D is a divergence measure, e.g., D_{KL}
- Why not $D_{\text{KL}}(P_g \| P_{\text{data}})$?



Whole-Image Generation II

- Minimizing $D_{KL}(P_{\text{data}} \| P_g)$ amounts to maximizing $P(\mathbb{X} | \Theta_g)$, the log likelihood of Θ_g :

$$\begin{aligned} g^* &= \arg \min_g D_{KL}(P_{\text{data}} \| P_g) \\ &= \arg \max_g E_{\mathbf{x} \sim P_{\text{data}}} [\log P_g(\mathbf{x})] + H(\mathbf{x} \sim P_{\text{data}}) \\ &= \arg \max_g E_{\mathbf{x} \sim P_{\text{data}}} [\log P_g(\mathbf{x})] \\ &\approx \arg \max_g \sum_{\mathbf{x}^{(i)} \in \mathbb{X}} \log P_g(\mathbf{x}^{(i)}) \\ &= \arg \max_g \log \prod_{\mathbf{x}^{(i)} \in \mathbb{X}} P_g(\mathbf{x}^{(i)}) \end{aligned}$$

$$\begin{aligned} \Theta_g^* &= \arg \max_{\Theta_g} \log \prod_{\mathbf{x}^{(i)} \in \mathbb{X}} P(\mathbf{x}^{(i)} | \Theta_g) \\ &= \arg \max_{\Theta_g} \log P(\mathbb{X} | \Theta_g) \end{aligned}$$

- Other divergence measures can lead to similar results

Common Whole-Image Generation Methods

- Autoencoder (single-step: $\hat{\mathbf{x}} = g(\mathbf{c})$)
 - Pros: easy
 - Cons: but no creativity, blurry images
- Variational Autoencoder (single-step: $\hat{\mathbf{x}} = g(\mathbf{c})$)
 - Pros: creative
 - Cons: only maximizes a lower bound of $P(\mathbb{X}|\Theta_g)$; blurry images
- Flow-based methods (multi-step: $\hat{\mathbf{x}} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$)
 - Pros: maximizes $P(\mathbb{X}|\Theta_g)$ directly
 - Cons: limited expressiveness of g for invertibility; slow training and inference
- GANs (single-step: $\hat{\mathbf{x}} = g(\mathbf{c})$)
 - Pros: good image quality (sharp and coherent)
 - Cons: difficult to train (convergence issue, mode collapse, vanishing gradients, etc.)
- Diffusion models (multi-step: $\hat{\mathbf{x}} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$)
 - Pros: good image quality; efficient & stable to train, can be made conditional easily
 - Cons: slow inference

Outline

1 Unsupervised Learning

- Text Models
- Image Models

2 ChatGPT

3 Autoencoders (AE)

- Manifold Learning*

4 Variational Autoencoders (VAE)

5 Flow-based Models

6 Diffusion Models

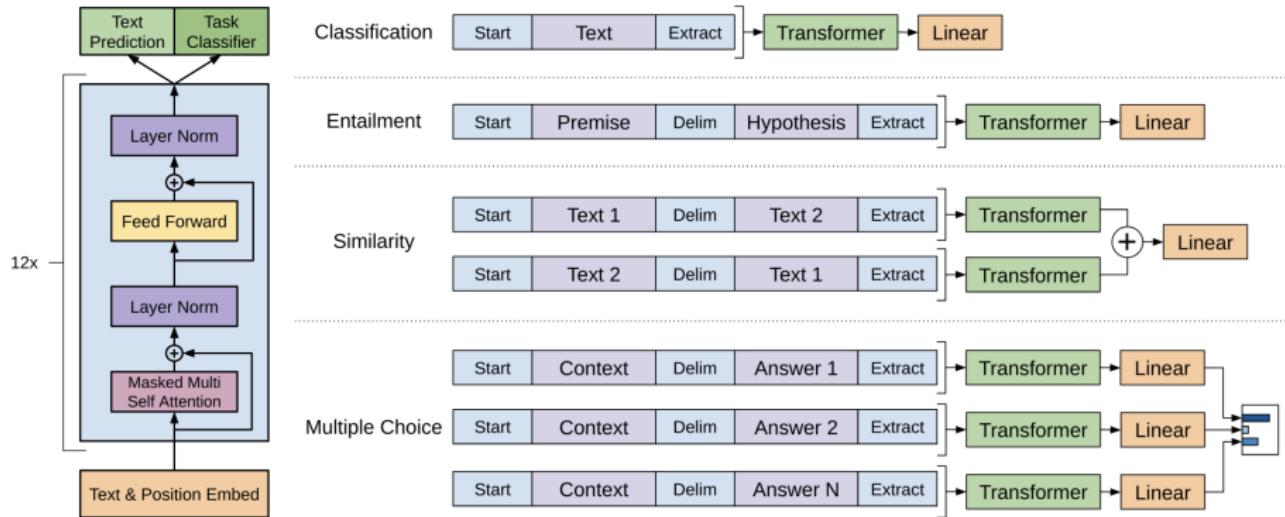
Evolutions

- 2018 GPTv1 [21]
 - Self-supervised pre-training
- 2019 GPTv2 [22]
 - **Multitask** pre-training
- 2020 GPTv3 [2]
 - **Few-shot & in-context** learning
- 2022 GPTv3.5 [20]
 - **Alignment** using (supervised) instruction tuning + reinforcement learning from human feedback (RLHF)
- 2023 GPT4: mixture of experts

GPTv1 [21]

- Aims at 2-step training process:
 - ① Self-supervised pre-training on unlabeled data
 - To predict next word in a sentence (language model)
 - ② Discriminative fine-tuning on labeled data in downstream tasks
 - Here, fine-tuning could also mean training a new model based on extracted features

Example Fine-tuning Tasks



Does Self-supervised Pre-training Help?

Table 5: Analysis of various model ablations on different tasks. Avg. score is a unweighted average of all the results. (*mc*= Mathews correlation, *acc*=Accuracy, *pc*=Pearson correlation)

Method	Avg. Score	CoLA (mc)	SST2 (acc)	MRPC (F1)	STS-B (pc)	QQP (F1)	MNLI (acc)	QNLI (acc)	RTE (acc)
Transformer w/ aux LM (full)	74.7	45.4	91.3	82.3	82.0	70.3	81.8	88.1	56.0
Transformer w/o pre-training	59.9	18.9	84.0	79.4	30.9	65.5	75.7	71.2	53.8
Transformer w/o aux LM	75.0	47.9	92.0	84.9	83.2	69.8	81.1	86.9	54.4
LSTM w/ aux LM	69.1	30.3	90.5	83.2	71.8	68.1	73.7	81.1	54.6

GPTv2 [22]

- Many NLP tasks can be formulated as the problem of maximizing $\Pr(\text{output}|\text{input}, \text{task})$
 - Translation: $\Pr(\text{french text}|\text{en text}, \text{translation})$
 - Question answering: $\Pr(\text{answer}|\text{question}, \text{qa})$
 - Reading comprehension: $\Pr(\text{answer}|\text{document}, \text{question}, \text{reading})$

GPTv3 [2]

- The model learns from few shots (examples), even *in context*
 - Enables **prompting** techniques for downstream tasks
 - E.g., few shots, chain of thought (CoT), or simply “Let’s work this out step-by-step to ensure the answer is correct”



Figure 1.1: Language model meta-learning. During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term “in-context learning” to describe the inner loop of this process, which occurs within the forward-pass upon each sequence. The sequences in this diagram are not intended to be representative of the data a model would see during pre-training, but are intended to show that there are sometimes repeated sub-tasks embedded within a single sequence.

Zero/Few Shot Prompting vs. Fine-tuning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



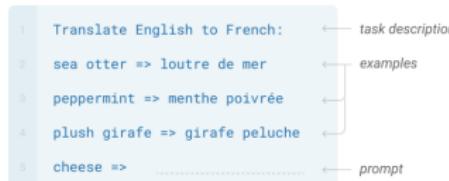
One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Performance vs. Model Size

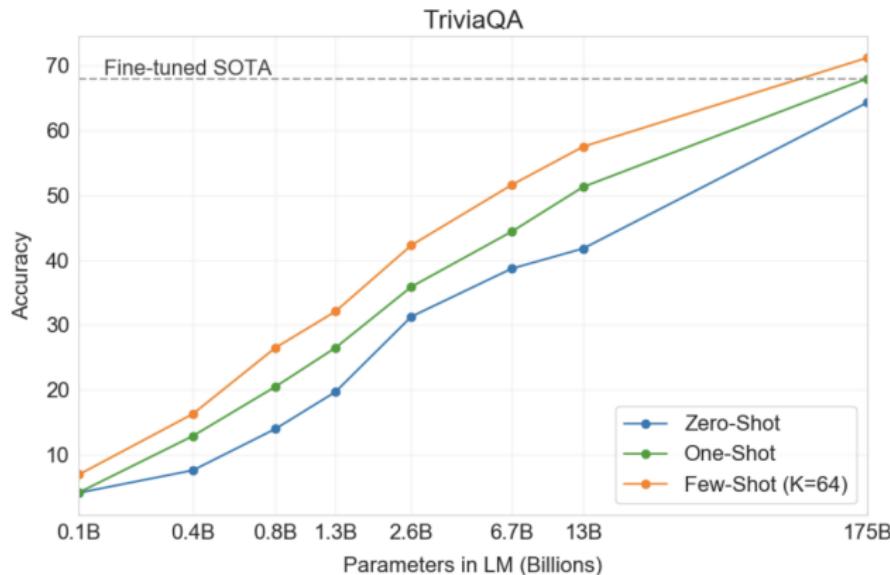


Figure 3.3: On TriviaQA GPT3's performance grows smoothly with model size, suggesting that language models continue to absorb knowledge as their capacity increases. One-shot and few-shot performance make significant gains over zero-shot behavior, matching and exceeding the performance of the SOTA fine-tuned open-domain model, RAG [LPP⁺20]

GPTv3.5 [20]

- From GPT to “ChatGPT” through *alignment*
 - Supervised (multi-task) instruction tuning
 - Reinforcement learning from human feedback (RLHF)

Step 1

Collect demonstration data, and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old



Some people went to the moon...

A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.

Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A Explain gravity... B Explain sat...

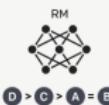
C Moon is natural satellite of... D People went to the moon...



D > C > A = B

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

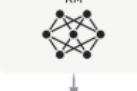
Write a story about frogs



The policy generates an output.



Once upon a time...

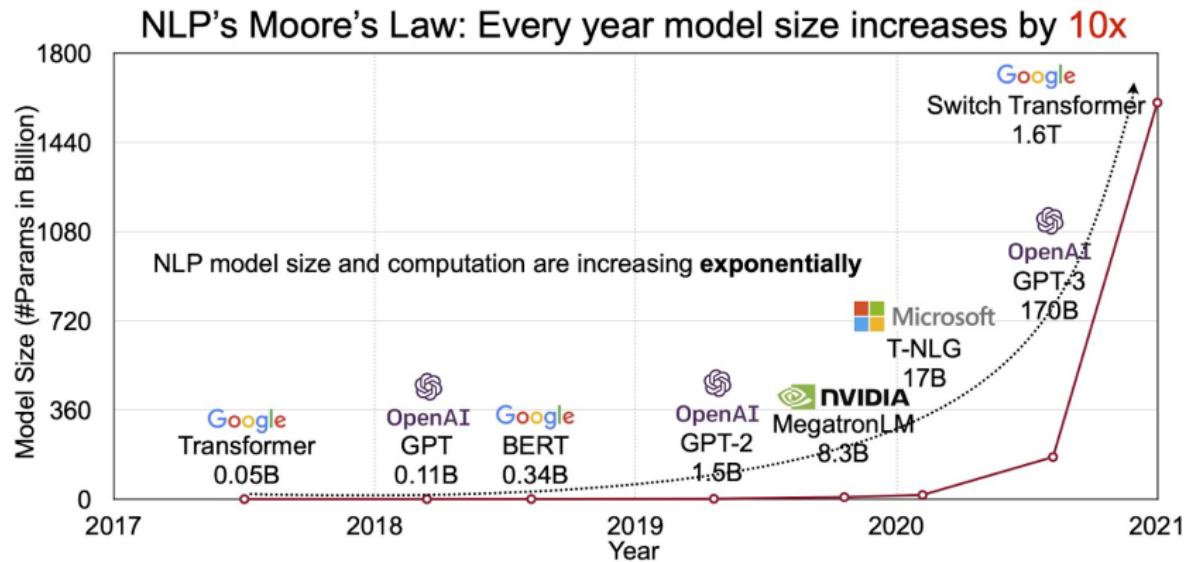


The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

r_k

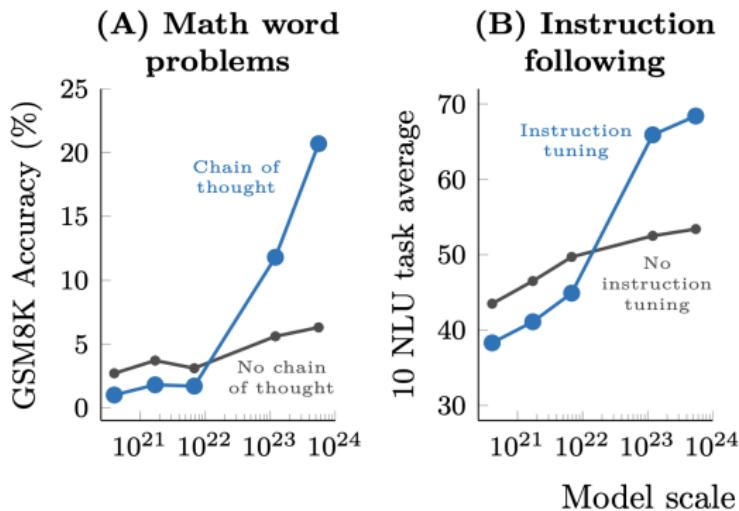
Sizes of Large Language Models (LLMs)



- Training costs [29]:
 - 110M params: \$2.5k–\$50k
 - 340M params: \$10k–\$200k
 - 1.5B param: \$80k–\$1.6m

Size Does Matter!

- Emerging abilities of LLMs [35]



- A balance: 70B parameters + 1.4T training tokens [6]

GPT Variants

- Low-rank adaptation (LoRA) [7]: to efficiently fine-tune LLMs
- WebGPT [19]: GPT that can search the web
- Retrieval-augmented generation (RAG) [13]: GPT to query external knowledge (vector) base
- GPTs that can reflect on their answers [1, 18]

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

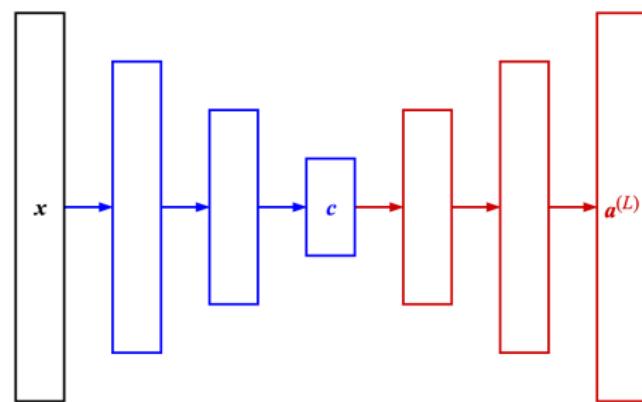
④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

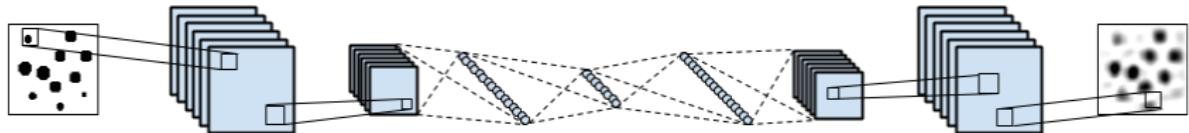
Autoencoders (AE)

- **Encoder**: to learn a low dimensional representation \mathbf{c} (called **code**) of input \mathbf{x}
- **Decoder**: to reconstruct \mathbf{x} from \mathbf{c}
- Objective: $\arg \max_{\Theta} \log P(\mathbb{X} | \Theta) = \arg \max_{\Theta} \sum_i \log P(\mathbf{x}^{(i)} | \Theta)$
- Assuming that $\mathbf{x} \sim \mathcal{N}(\mu, \cdot)$, we have linear output units
 $\mathbf{a}^{(L)} = \mathbf{z}^{(L)} = \hat{\mu}$
 - $\log P(\mathbf{x}^{(i)} | \Theta) \propto -\|\mathbf{x}^{(i)} - \mathbf{a}^{(i,L)}\|^2$
 - $\arg \max_{\Theta} \sum_i \log P(\mathbf{x}^{(i)} | \Theta) = \arg \min_{\Theta} \sum_i \|\mathbf{x}^{(i)} - \mathbf{a}^{(i,L)}\|^2$ (minimizing reconstruct error)



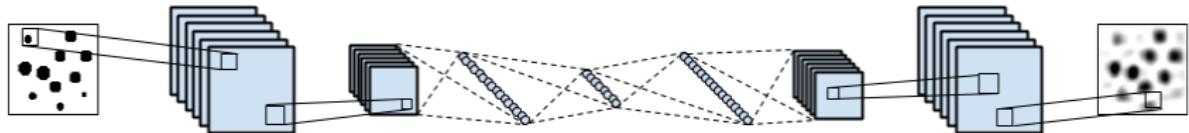
Convolutional Autoencoders

- Convolution + deconvolution layers:



Convolutional Autoencoders

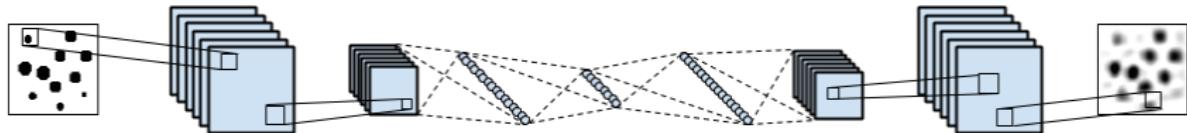
- Convolution + deconvolution layers:



- Decoder is a simplified DeconvNet [36] trained from scratch:

Convolutional Autoencoders

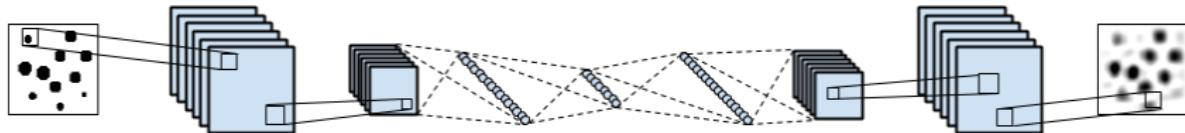
- Convolution + deconvolution layers:



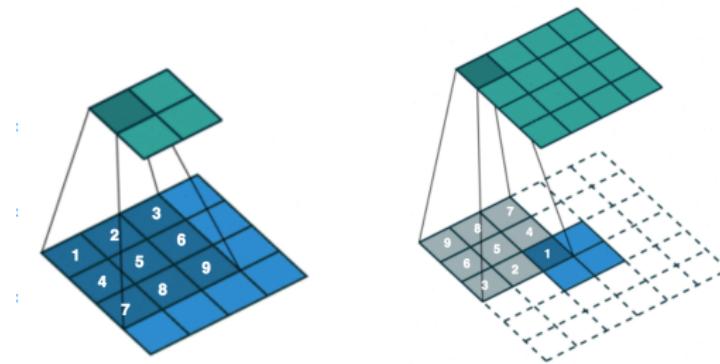
- Decoder is a simplified DeconvNet [36] trained from scratch:
 - Uppooling → upsampling (no need to remember max positions)

Convolutional Autoencoders

- Convolution + deconvolution layers:



- Decoder is a simplified DeconvNet [36] trained from scratch:
 - Uppooling → upsampling (no need to remember max positions)
 - Deconvolution → convolution



Codes & Reconstructed x

- A 32-bit code can roughly represents a 32×32 (1024 dimensional) MNIST image



Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

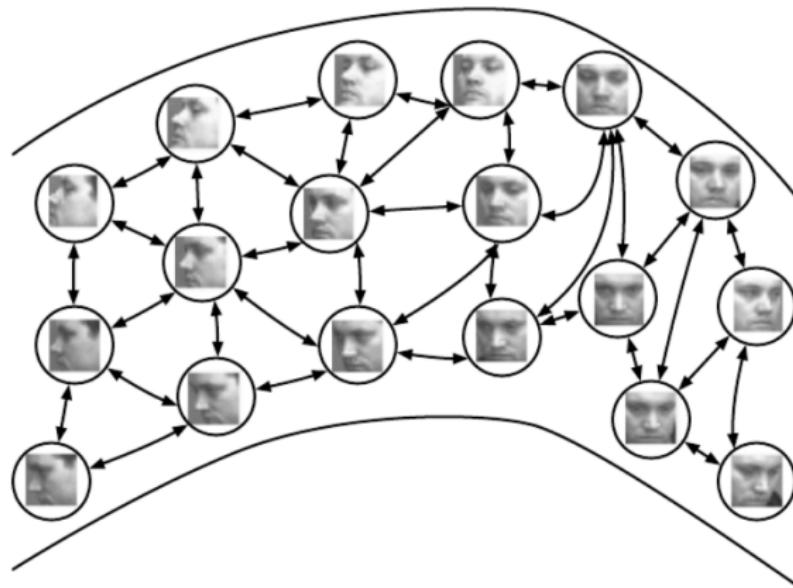
⑥ Diffusion Models

Manifolds I

- In many applications, data concentrate around one or more low-dimensional *manifolds*

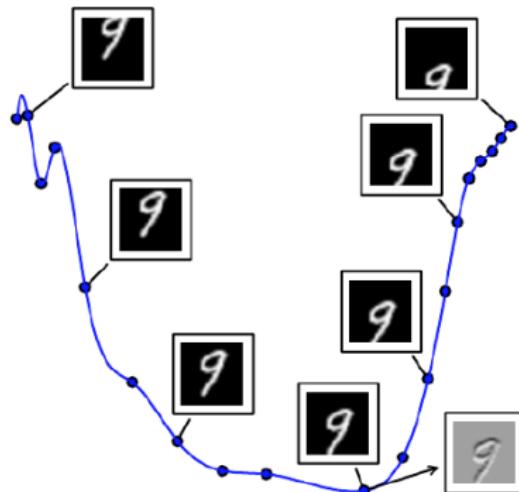
Manifolds I

- In many applications, data concentrate around one or more low-dimensional **manifolds**
- A manifold is a topological space that are **linear locally**



Manifolds II

- For each point x on a manifold, we have its **tangent space** spanned by **tangent vectors**
 - Local directions specify how one can change x infinitesimally while staying on the manifold



Learning Manifolds I

- How to make \mathbf{c} produced by autoencoders denote a *coordinate* of a dimensional manifold?

Learning Manifolds I

- How to make \mathbf{c} produced by autoencoders denote a *coordinate* of a dimensional manifold?
- Contractive autoencoder [24]: regularizes the code \mathbf{c} such that it is invariant to local changes of \mathbf{x} :

$$\Omega(\mathbf{c}) = \sum_i \left\| \frac{\partial \mathbf{c}^{(i)}}{\partial \mathbf{x}^{(i)}} \right\|_F^2$$

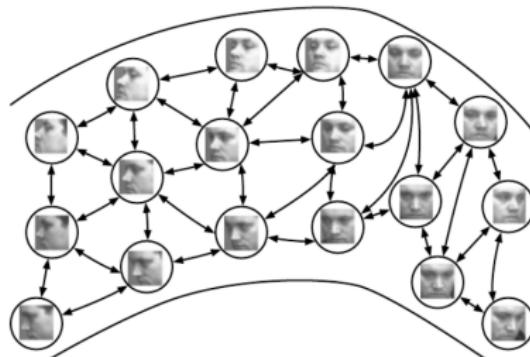
- $\partial \mathbf{c}^{(i)} / \partial \mathbf{x}^{(i)}$ is a Jacobian matrix

Learning Manifolds I

- How to make \mathbf{c} produced by autoencoders denote a *coordinate* of a dimensional manifold?
- Contractive autoencoder [24]: regularizes the code \mathbf{c} such that it is invariant to local changes of \mathbf{x} :

$$\Omega(\mathbf{c}) = \sum_i \left\| \frac{\partial \mathbf{c}^{(i)}}{\partial \mathbf{x}^{(i)}} \right\|_F^2$$

- $\partial \mathbf{c}^{(i)} / \partial \mathbf{x}^{(i)}$ is a Jacobian matrix
- Hence, \mathbf{c} represents only the variations needed to reconstruct \mathbf{x}
 - I.e., \mathbf{c} changes most along tangent vectors



Learning Manifolds II

- In practice, it is easier to train a denoising autoencoder [34]:



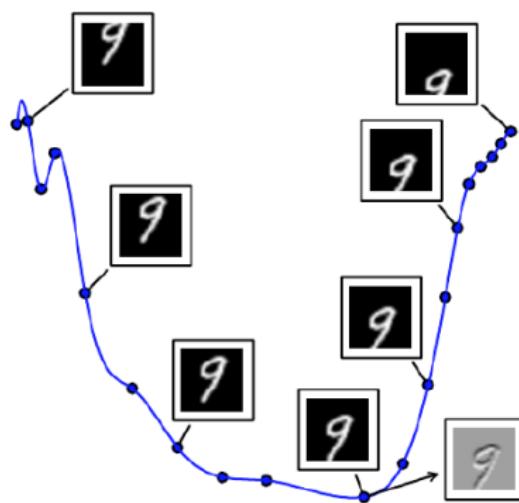
Learning Manifolds II

- In practice, it is easier to train a denoising autoencoder [34]:
 - Encoder: to encode \mathbf{x} **with** random noises
 - Decoder: to reconstruct \mathbf{x} **without** noises



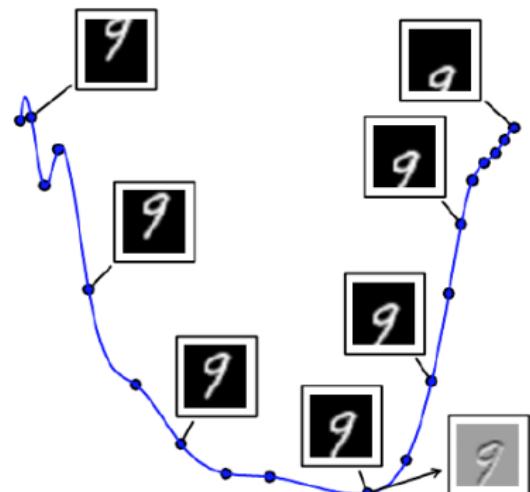
Getting Tangent Vectors I

- The code c represents a coordinate on a low dimensional manifold
 - E.g., the blue line
- How to get the tangent vectors of a given c ?



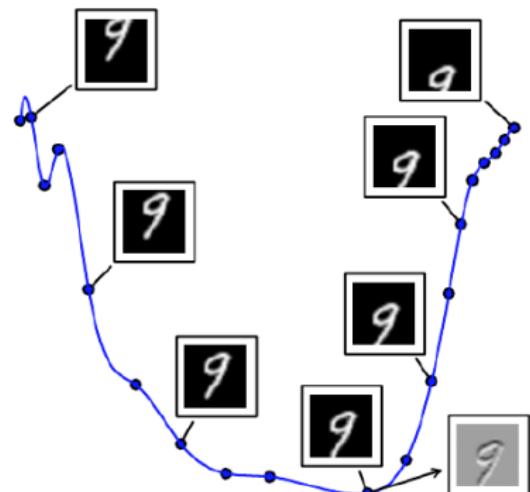
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors



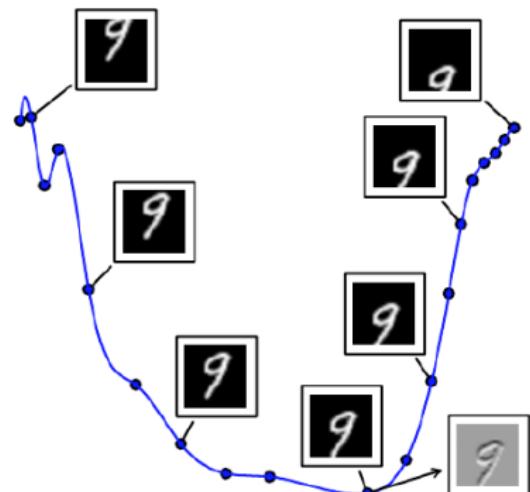
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors
- Given a point x , let c be the code of x and $J(x) = \frac{\partial c}{\partial x}$ be the Jacobian matrix of c at x



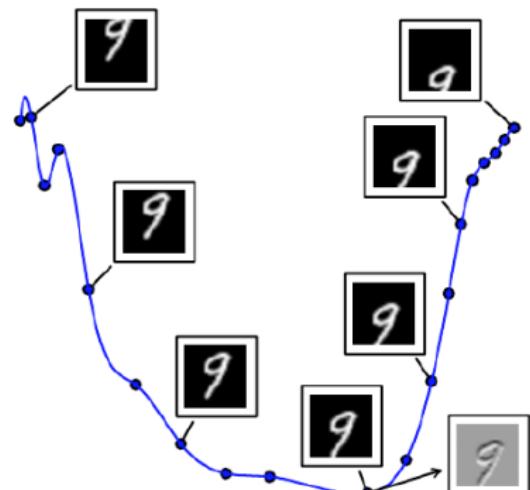
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors
- Given a point x , let c be the code of x and $J(x) = \frac{\partial c}{\partial x}$ be the Jacobian matrix of c at x
 - $J(x)$ summarizes how c changes in terms of x



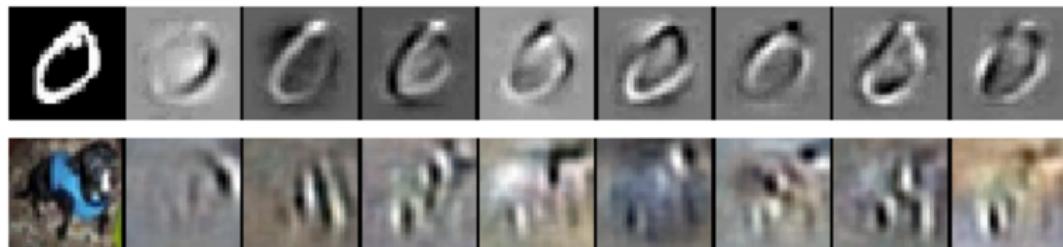
Getting Tangent Vectors II

- Recall: directions in the input space that *changes c most* should be tangent vectors
 - Given a point x , let c be the code of x and $J(x) = \frac{\partial c}{\partial x}$ be the Jacobian matrix of c at x
 - $J(x)$ summarizes how c changes in terms of x
- Decompose $J(x)$ using SVD such that $J(x) = UDV^\top$
 - Let tangent vectors be *rows of V corresponding to the largest singular values in D*



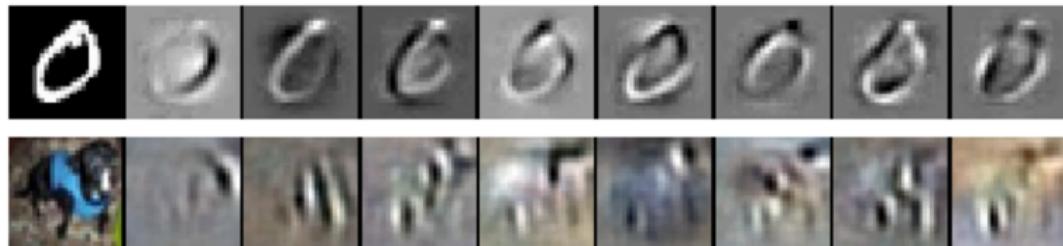
Getting Tangent Vectors III

- In practice, $J(x)$ usually has few large singular values
- Tangent vectors found by contractive/denoising autoencoders:



Getting Tangent Vectors III

- In practice, $J(\mathbf{x})$ usually has few large singular values
- Tangent vectors found by contractive/denoising autoencoders:



- Can be used by Tangent Prop [30]:
- Let $\{\mathbf{v}^{(i,j)}\}_j$ be tangent vectors of each example $\mathbf{x}^{(i)}$
- Trains an NN classifier f with cost penalty: $\Omega[f] = \sum_{i,j} \nabla_{\mathbf{x}} f(\mathbf{x}^{(i)})^\top \mathbf{v}^{(i,j)}$
 - Points in the same manifold share the same label

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

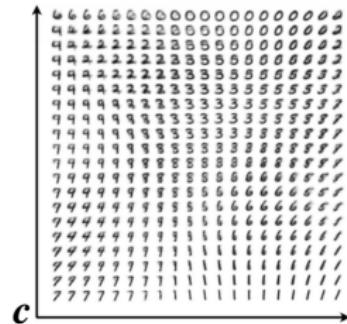
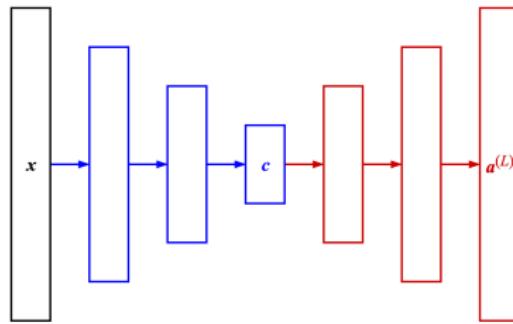
④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

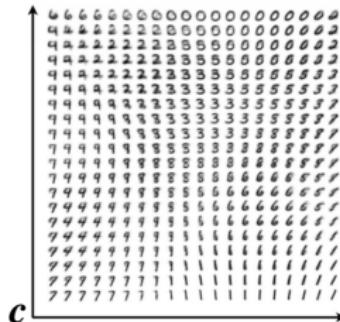
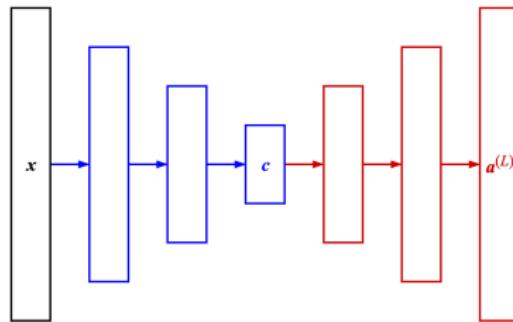
Problems of Autoencoder I

- Ideally, the decoder of an autoencoder can be used to generate images even with *synthetic codes*



Problems of Autoencoder I

- Ideally, the decoder of an autoencoder can be used to generate images even with *synthetic codes*



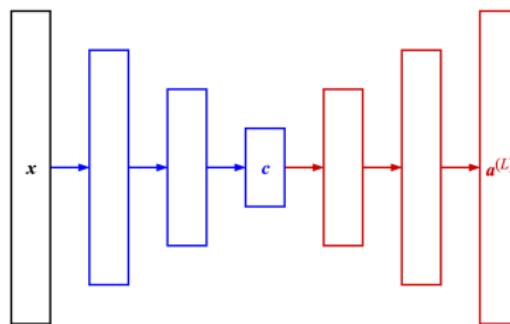
- In reality, the learnt c in code space has many “holes” that fail to map to images
 - More complex image patterns
 - Training data are never enough

Problems of Autoencoder II

- Blurry images: the objective

$$\arg \max_{\Theta} \sum_i \log P(\mathbf{x}^{(i)} | \Theta) = \arg \min_{\Theta} \sum_i \|\mathbf{x}^{(i)} - \mathbf{a}^{(i,L)}\|^2$$

does not penalize Gaussian pixel noises in $\mathbf{a}^{(i,L)}$



Variational Autoencoders (VAE) [8]

- Encoder $f(\cdot; \Theta_f)$: maps each sample of \mathbf{x} (i.e., $\mathbf{x}^{(i)} \in \mathbb{X}$) to an ***axis-aligned normal distribution*** $\mathcal{N}(\mu, \sigma)$
 - Each code dimension is independent with each other
 - $f(\mathbf{x}) = (\mu, \sigma)$
- Decoder $g(\cdot; \Theta_g)$: same as that of AE
- How to minimize the objective $\arg \max_{\Theta_f, \Theta_g} \log P(\mathbb{X} | \Theta_f, \Theta_g))$?

VAE Objective

- Objective:

$$\arg \max_{\Theta_f, \Theta_g} \log P(\mathbb{X} | \Theta_f, \Theta_g) = \arg \max_{\Theta_f, \Theta_g} \sum_i \log P(x^{(i)} | \Theta_f, \Theta_g)$$

- Considering $\log P(\mathbf{x})$ for any sample \mathbf{x} , we have

$$\begin{aligned}\log P(\mathbf{x}) &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log P(\mathbf{x}) d\mathbf{c} \quad // \text{Q can be any distribution} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{P(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} = \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \frac{Q(\mathbf{c}|\mathbf{x})}{P(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} + \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{Q(\mathbf{c}|\mathbf{x})}{P(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} + D_{\text{KL}}(Q(\mathbf{c}|\mathbf{x}) \| P(\mathbf{c}|\mathbf{x})) \\ &\geq \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{c}, \mathbf{x})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \quad // \text{lower bound} \\ &= \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}|\mathbf{c})P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c}\end{aligned}$$

- VAE lets $Q(\cdot|\mathbf{x}) = \mathcal{N}(f(\mathbf{x}; \Theta_f))$ and $P(\cdot|\mathbf{c}) = \mathcal{N}(g(\mathbf{c}; \Theta_g))$
 - So $Q(\mathbf{c}|\mathbf{x}) = Q(\mathbf{c}|\mathbf{x}, \Theta_f)$ and $P(\mathbf{x}|\mathbf{c}) = P(\mathbf{x}|\mathbf{c}, \Theta_g)$
- New objective: finds Θ_f and Θ_g that maximize the lower bound
 - Not necessarily maximize $\log P(\mathbf{x}|\Theta_f, \Theta_g)$

Maximizing Lower Bound

$$\begin{aligned} & \int_c Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}|\mathbf{c})P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} \\ &= \int_c Q(\mathbf{c}|\mathbf{x}, \Theta_f) \log \left(\frac{P(\mathbf{x}|\mathbf{c}, \Theta_g)P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x}, \Theta_f)} \right) d\mathbf{c} \\ &= \int_c Q(\mathbf{c}|\mathbf{x}, \Theta_f) \log \left(\frac{P(\mathbf{c})}{Q(\mathbf{c}|\mathbf{x}, \Theta_f)} \right) d\mathbf{c} + \int_c Q(\mathbf{c}|\mathbf{x}, \Theta_f) \log P(\mathbf{x}|\mathbf{c}, \Theta_g) d\mathbf{c} \\ &= -D_{KL}(Q(\mathbf{c}|\mathbf{x}, \Theta_f) \| P(\mathbf{c})) + E_{(\mathbf{c}|\mathbf{x}, \Theta_f) \sim Q} [\log P(\mathbf{x}|\mathbf{c}, \Theta_g)] \end{aligned}$$

- For $P(\mathbf{c}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$:
- To minimize the first term, the encoder $f(\mathbf{x}; \Theta_f) = (\mu, \sigma)$ has a loss term [8]

$$\exp(\sigma) - (1 + \sigma) + \|\mu\|^2$$

- To maximize the second term, the decoder
 - ① Samples \mathbf{c} from $\mathcal{N}(\mu, \sigma)$ to get $g(\mathbf{c}; \Theta_g)$, the mean of output \mathcal{N}
 - ② Minimizes a loss term $\|\mathbf{x} - \mathbf{a}^{(L)}\|^2$ as in AE

Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

Problems of VAE

- Only maximize a lower bound of the likelihood $P(\mathbb{X} | \Theta_f, \Theta_g)$
 - Θ_f and Θ_g are encoder and decoder weights, respectively
- Still blurry images
 - The decoder's loss is the same with that of AE

Flow-based Models

- Idea: let the decoder $g(\cdot; \Theta_g)$ be a **deterministic invertible** function
 - Given $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we have $\mathbf{x} = g(\mathbf{c}; \Theta_g)$ of complex distribution
 - Conversely, given \mathbf{x} , we have $\mathbf{c} = g^{-1}(\mathbf{x}; \Theta_g)$
- The likelihood can be maximize directly:

$$\begin{aligned} & \arg \max_g \log P_g(\mathbb{X}) \\ &= \arg \max_g \sum_i \log P_{\mathbf{x}}(g(\mathbf{c}^{(i)})) \\ &= \arg \max_g \sum_i \log [P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) | \det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)})) |] \\ &= \arg \max_g \sum_i \log P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) + \log |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))| \end{aligned}$$

- First term: finds $g^{-1}(\Theta_g)$ that maps all $\mathbf{x}^{(i)}$ to 0 (mean of \mathcal{N})
- Second term: prevents $g^{-1}(\Theta_g)$ from mapping all $\mathbf{x}^{(i)}$ to only 0
 - As $\log |\det(\mathbf{O})| = -\infty$

Flow-based Models

- Idea: let the decoder $g(\cdot; \Theta_g)$ be a **deterministic invertible** function
 - Given $\mathbf{c} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we have $\mathbf{x} = g(\mathbf{c}; \Theta_g)$ of complex distribution
 - Conversely, given \mathbf{x} , we have $\mathbf{c} = g^{-1}(\mathbf{x}; \Theta_g)$
- The likelihood can be maximize directly:

$$\begin{aligned} & \arg \max_g \log P_g(\mathbb{X}) \\ &= \arg \max_g \sum_i \log P_{\mathbf{x}}(g(\mathbf{c}^{(i)})) \\ &= \arg \max_g \sum_i \log [P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))|] \\ &= \arg \max_g \sum_i \log P_{\mathbf{c}}(g^{-1}(\mathbf{x}^{(i)})) + \log |\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)}))| \end{aligned}$$

- First term: finds $g^{-1}(\Theta_g)$ that maps all $\mathbf{x}^{(i)}$ to 0 (mean of \mathcal{N})
- Second term: prevents $g^{-1}(\Theta_g)$ from mapping all $\mathbf{x}^{(i)}$ to only 0
 - As $\log |\det(\mathbf{O})| = -\infty$
- But how to ensure the followings during training?
 - g is invertible
 - $\det(\mathbf{J}(g^{-1})(\cdot))$ can be easily computed

Ensuring Invertibility

- Glow [9]: make g^{-1} an 1×1 convolution layer

- $\mathbf{x}_{i,j,:} = \begin{bmatrix} x_{i,j,1} \\ x_{i,j,2} \\ x_{i,j,3} \end{bmatrix}, \quad g^{-1} = \mathbf{W}_{3 \times 3},$

$$g^{-1}(\mathbf{x}_{i,j,:}) = \mathbf{W}_{3 \times 3} \begin{bmatrix} x_{i,j,1} \\ x_{i,j,2} \\ x_{i,j,3} \end{bmatrix} = \begin{bmatrix} c_{i,j,1} \\ c_{i,j,2} \\ c_{i,j,3} \end{bmatrix}$$

- At training time, initialize $\mathbf{W}_{3 \times 3}$ as an invertible matrix
 - $g^{-1} = \mathbf{W}_{3 \times 3}$ is likely to be invertible after SGD updates
 - Determinant is easy to compute: $\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)})) = \det(\mathbf{W}_{3 \times 3})^{W \times H}$
- At inference time, use $g = \mathbf{W}_{3 \times 3}^{-1}$ to generate images

Ensuring Invertibility

- Glow [9]: make g^{-1} an 1×1 convolution layer

- $\mathbf{x}_{i,j,:} = \begin{bmatrix} x_{i,j,1} \\ x_{i,j,2} \\ x_{i,j,3} \end{bmatrix}, \quad g^{-1} = \mathbf{W}_{3 \times 3},$

$$g^{-1}(\mathbf{x}_{i,j,:}) = \mathbf{W}_{3 \times 3} \begin{bmatrix} x_{i,j,1} \\ x_{i,j,2} \\ x_{i,j,3} \end{bmatrix} = \begin{bmatrix} c_{i,j,1} \\ c_{i,j,2} \\ c_{i,j,3} \end{bmatrix}$$

- At training time, initialize $\mathbf{W}_{3 \times 3}$ as an invertible matrix
 - $g^{-1} = \mathbf{W}_{3 \times 3}$ is likely to be invertible after SGD updates
 - Determinant is easy to compute: $\det(\mathbf{J}(g^{-1})(\mathbf{x}^{(i)})) = \det(\mathbf{W}_{3 \times 3})^{W \times H}$
- At inference time, use $g = \mathbf{W}_{3 \times 3}^{-1}$ to generate images
- Problem: g has limited expressiveness

Step-wise Generation

- Cascade multiple invertible g to have a more complex one:

$$\mathbf{x} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$$

Step-wise Generation

- Cascade multiple invertible g to have a more complex one:

$$\mathbf{x} = g^{(T)}(\cdots g^{(2)}(g^{(1)}(\mathbf{c})))$$

- Result: sharp images



Outline

① Unsupervised Learning

- Text Models
- Image Models

② ChatGPT

③ Autoencoders (AE)

- Manifold Learning*

④ Variational Autoencoders (VAE)

⑤ Flow-based Models

⑥ Diffusion Models

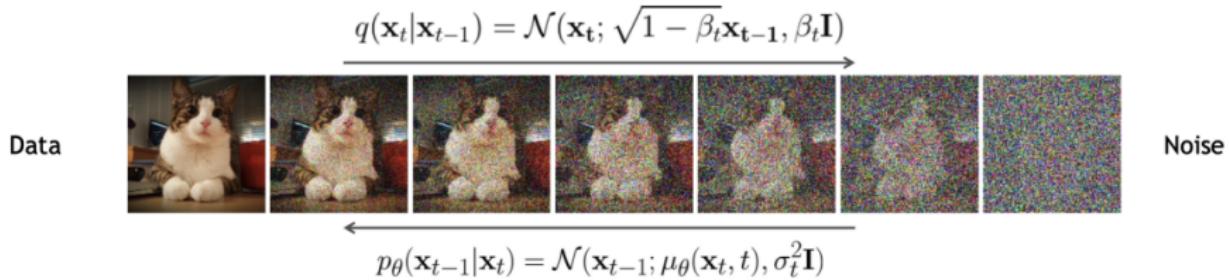
Problems of Flow-based Models

- Limited model expressiveness
- Slow training and inference
 - #steps can be large

Denoising Diffusion Probabilistic Models (DDPM)

[5]

- Borrow some good ideas from previous works
 - Probabilistic formulation of VAE that models encoder in the objective
 - Step-wise encoding/decoding in generative flows
- But, unlike flows, the encoding steps
 - Are **predefined**, no parameter to learn
 - Can be simplified to **one encoding step**



Encoding: multi-steps = single-step

- Predefined encoding functions $f^{(t)}(\cdot; \beta^{(t)})$, $t = 1, \dots, T$:
 - $\beta^{(t)}$, $\forall t$, are hyperparameters; no learning needed
- Let each sample $\mathbf{x} = \mathbf{x}^{(0)}$
- $\mathbf{x}^{(1)} = f^{(1)}(\mathbf{x}^{(0)}; \beta^{(1)}) = \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{\beta^{(1)}} \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - So, $(\mathbf{x}^{(1)} | \mathbf{x}^{(0)} = \mathbf{x}^{(0)}) \sim \mathcal{N}(\sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)}, \beta^{(1)} \mathbf{I})$

Encoding: multi-steps = single-step

- Predefined encoding functions $f^{(t)}(\cdot; \beta^{(t)})$, $t = 1, \dots, T$:
 - $\beta^{(t)}$, $\forall t$, are hyperparameters; no learning needed
- Let each sample $\mathbf{x} = \mathbf{x}^{(0)}$
- $\mathbf{x}^{(1)} = f^{(1)}(\mathbf{x}^{(0)}; \beta^{(1)}) = \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{\beta^{(1)}} \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - So, $(\mathbf{x}^{(1)} | \mathbf{x}^{(0)} = \mathbf{x}^{(0)}) \sim \mathcal{N}(\sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)}, \beta^{(1)} \mathbf{I})$
- $$\begin{aligned}\mathbf{x}^{(2)} &= f^{(2)}(\mathbf{x}^{(1)}; \beta^{(2)}) = \sqrt{1 - \beta^{(2)}} \mathbf{x}^{(1)} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \\ &= \sqrt{1 - \beta^{(2)}} \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - (1 - \beta^{(2)})(1 - \beta^{(1)})} \boldsymbol{\varepsilon} \\ &= \sqrt{\alpha^{(2)} \alpha^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - \alpha^{(2)} \alpha^{(1)}} \boldsymbol{\varepsilon}\end{aligned}$$
 - $\sqrt{1 - \beta^{(1)}} \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, ((1 - \beta^{(1)}) \beta^{(2)} + \beta^{(2)}) \mathbf{I})$
 - Let $\alpha^{(t)} = 1 - \beta^{(t)}$ (derived hyperparameter)
 - Only one $\boldsymbol{\varepsilon}$ is added

Encoding: multi-steps = single-step

- Predefined encoding functions $f^{(t)}(\cdot; \beta^{(t)})$, $t = 1, \dots, T$:
 - $\beta^{(t)}$, $\forall t$, are hyperparameters; no learning needed
- Let each sample $\mathbf{x} = \mathbf{x}^{(0)}$
- $\mathbf{x}^{(1)} = f^{(1)}(\mathbf{x}^{(0)}; \beta^{(1)}) = \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{\beta^{(1)}} \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - So, $(\mathbf{x}^{(1)} | \mathbf{x}^{(0)} = \mathbf{x}^{(0)}) \sim \mathcal{N}(\sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)}, \beta^{(1)} \mathbf{I})$
- $$\begin{aligned}\mathbf{x}^{(2)} &= f^{(2)}(\mathbf{x}^{(1)}; \beta^{(2)}) = \sqrt{1 - \beta^{(2)}} \mathbf{x}^{(1)} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \\ &= \sqrt{1 - \beta^{(2)}} \sqrt{1 - \beta^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - (1 - \beta^{(2)})(1 - \beta^{(1)})} \boldsymbol{\varepsilon} \\ &= \sqrt{\alpha^{(2)} \alpha^{(1)}} \mathbf{x}^{(0)} + \sqrt{1 - \alpha^{(2)} \alpha^{(1)}} \boldsymbol{\varepsilon}\end{aligned}$$
 - $\sqrt{1 - \beta^{(1)}} \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} + \sqrt{\beta^{(2)}} \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, ((1 - \beta^{(1)}) \beta^{(2)} + \beta^{(2)}) \mathbf{I})$
 - Let $\bar{\alpha}^{(t)} = 1 - \beta^{(t)}$ (derived hyperparameter)
 - Only one $\boldsymbol{\varepsilon}$ is added
- $\mathbf{x}^{(t)} = \sqrt{\bar{\alpha}^{(t)}} \mathbf{x}^{(0)} + \sqrt{1 - \bar{\alpha}^{(t)}} \boldsymbol{\varepsilon}$
 - Let $\bar{\alpha}^{(t)} = \alpha^{(t)} \alpha^{(t-1)} \dots \alpha^{(1)}$ (derived hyperparameter)

Objective

- As VAE, DDPM maximizes a lower bound of $P(\mathbb{X})$
- VAE: for each \mathbf{x} , where $P(\mathbf{x}) = \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log P(\mathbf{x}) d\mathbf{c}$, maximize:

$$\int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}, \mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} = E_{\mathbf{c}|\mathbf{x} \sim Q} \left[\log \left(\frac{P(\mathbf{x}, \mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) \right] \leq P(\mathbf{x})$$

Objective

- As VAE, DDPM maximizes a lower bound of $P(\mathbb{X})$
- VAE: for each \mathbf{x} , where $P(\mathbf{x}) = \int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log P(\mathbf{x}) d\mathbf{c}$, maximize:

$$\int_{\mathbf{c}} Q(\mathbf{c}|\mathbf{x}) \log \left(\frac{P(\mathbf{x}, \mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) d\mathbf{c} = E_{\mathbf{c}|\mathbf{x} \sim Q} \left[\log \left(\frac{P(\mathbf{x}, \mathbf{c})}{Q(\mathbf{c}|\mathbf{x})} \right) \right] \leq P(\mathbf{x})$$

- DDPM: for each \mathbf{x} where $P(\mathbf{x}^{(0)}) = \int_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}} Q(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)} | \mathbf{x}^{(0)}) \log P(\mathbf{x}^{(0)}) d\mathbf{x}^{(1)} \dots d\mathbf{x}^{(T)}$, maximize:

$$E_{(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)} | \mathbf{x}^{(0)}) \sim Q} \left[\log \left(\frac{P(\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)})}{Q(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)} | \mathbf{x}^{(0)})} \right) \right],$$

which can be simplified to [15]:

$$-D_{KL}(Q(\mathbf{x}^{(T)} | \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(T)})) + E_{\mathbf{x}^{(1)} | \mathbf{x}^{(0)} \sim Q} [\log P(\mathbf{x}^{(0)} | \mathbf{x}^{(1)})] + \\ - \sum_{t=2}^T E_{\mathbf{x}^{(t)} | \mathbf{x}^{(0)} \sim Q} [D_{KL}(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}))]$$

- First term controlled by encoding process (predefined)
- Second & third term controlled by denoising process (learnable)

Denoising I

- For simplicity, we focus on maximizing the third term:

$$-\sum_{t=2}^T \mathbb{E}_{\mathbf{x}^{(t)} | \mathbf{x}^{(0)} \sim Q} \left[D_{KL} \left(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right) \right]$$

- Goal: for each observed $\mathbf{x}^{(t)}$, minimize

$$D_{KL} \left(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right)$$

- Note that

$$\begin{aligned} Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) &= \frac{Q(\mathbf{x}^{(t-1)}, \mathbf{x}^{(t)}, \mathbf{x}^{(0)})}{Q(\mathbf{x}^{(t)}, \mathbf{x}^{(0)})} \\ &= \frac{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)}) Q(\mathbf{x}^{(0)})}{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)}) Q(\mathbf{x}^{(0)})} \\ &= \frac{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}) Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)})}{Q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)})} \end{aligned}$$

- Since $Q(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$ & $Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(0)})$ are Gaussian, we have [15]:

$$Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) = \mathcal{N} \left(\frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}}, \frac{1-\bar{\alpha}^{(t-1)}}{1-\bar{\alpha}^{(t)}}\beta^{(t)}\mathbf{I} \right)$$

Denoising II

- Goal: for each observed $\mathbf{x}^{(t)}$, minimize

$$D_{KL} \left(Q(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \mathbf{x}^{(0)}) \| P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) \right),$$

where $Q(\dots) = \mathcal{N} \left(\frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}}, \dots \right)$ is **fixed**

- DDPM finds Θ that move the mean of $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$ (also Gaussian) toward $Q(\dots)$'s mean:

$$\frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}}$$

Noise Predictor

- Θ moves $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$'s mean toward

$$\begin{aligned}& \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}} \\&= \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\frac{\mathbf{x}^{(t)} - \sqrt{1-\bar{\alpha}^{(t)}}\varepsilon}{\sqrt{\bar{\alpha}^{(t)}}}}{1-\bar{\alpha}^{(t)}} \\&= \frac{1}{\sqrt{\alpha^{(t)}}} \left(\mathbf{x}^{(t)} - \frac{1-\alpha^{(t)}}{\sqrt{1-\bar{\alpha}^{(t)}}} \varepsilon \right)\end{aligned}$$

- What's its corresponding network?

Noise Predictor

- Θ moves $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$'s mean toward

$$\begin{aligned}& \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}} \\&= \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\frac{\mathbf{x}^{(t)} - \sqrt{1-\bar{\alpha}^{(t)}}\boldsymbol{\varepsilon}}{\sqrt{\bar{\alpha}^{(t)}}}}{1-\bar{\alpha}^{(t)}} \\&= \frac{1}{\sqrt{\alpha^{(t)}}} \left(\mathbf{x}^{(t)} - \frac{1-\alpha^{(t)}}{\sqrt{1-\bar{\alpha}^{(t)}}} \boldsymbol{\varepsilon} \right)\end{aligned}$$

- What's its corresponding network?
- By definition: let Θ parametrize a network outputting $\mathbf{x}^{(t-1)}$ given $\mathbf{x}^{(t)}$
 - But the input $\mathbf{x}^{(t)}$ also resides in the output target

Noise Predictor

- Θ moves $P(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}, \Theta)$'s mean toward

$$\begin{aligned}& \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\mathbf{x}^{(0)}}{1-\bar{\alpha}^{(t)}} \\&= \frac{\sqrt{\alpha^{(t)}}(1-\alpha^{(t-1)})\mathbf{x}^{(t)} + \sqrt{\bar{\alpha}^{(t-1)}}\beta^{(t)}\frac{\mathbf{x}^{(t)} - \sqrt{1-\bar{\alpha}^{(t)}}\varepsilon}{\sqrt{\bar{\alpha}^{(t)}}}}{1-\bar{\alpha}^{(t)}} \\&= \frac{1}{\sqrt{\alpha^{(t)}}} \left(\mathbf{x}^{(t)} - \frac{1-\alpha^{(t)}}{\sqrt{1-\bar{\alpha}^{(t)}}} \varepsilon \right)\end{aligned}$$

- What's its corresponding network?
- By definition: let Θ parametrize a network outputting $\mathbf{x}^{(t-1)}$ given $\mathbf{x}^{(t)}$
 - But the input $\mathbf{x}^{(t)}$ also resides in the output target
- DDPM: let Θ parametrize a **noise predictor** outputting ε given $\mathbf{x}^{(t)}$
 - Objective: $\arg \min_{\Theta} \|\varepsilon - e(\mathbf{x}^{(t)}, t; \Theta)\|^2$
 - Of U-Net [26] architecture
 - Shared between all t

Training & Inference Algorithms

- One-step encoding during training time
- Multi-step inference (sampling), with each intermediate decoding step t , $t > 1$, comes with extra noise $\sigma^{(t)} \mathbf{z}$
 - Similar to output token sampling in GPT
 - Improves performance empirically

Algorithm 1 Training

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
     
$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$

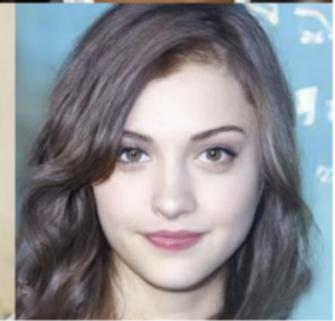
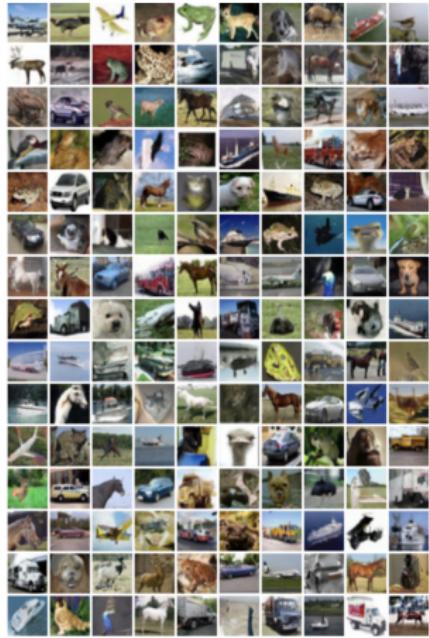
6: until converged
```

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

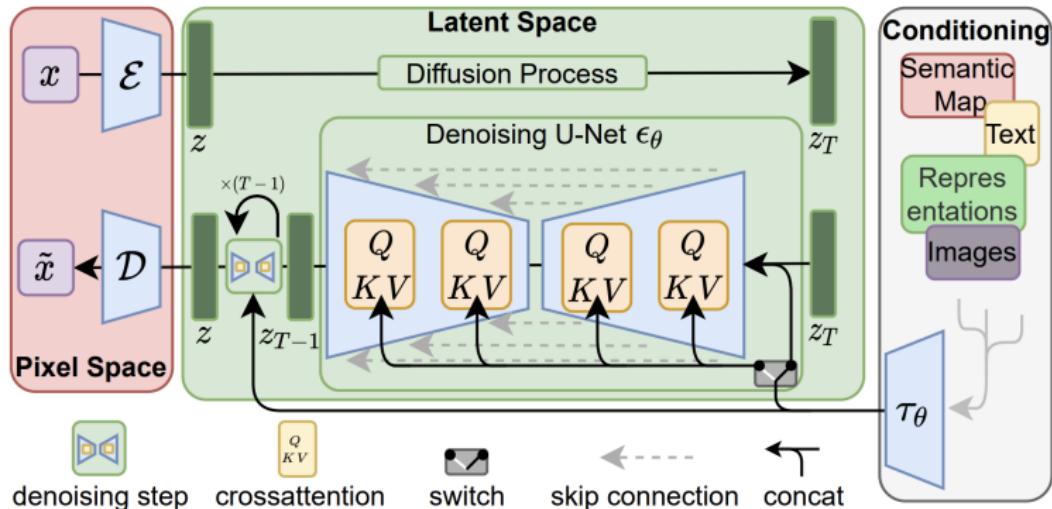
Results

- Sharp and coherent



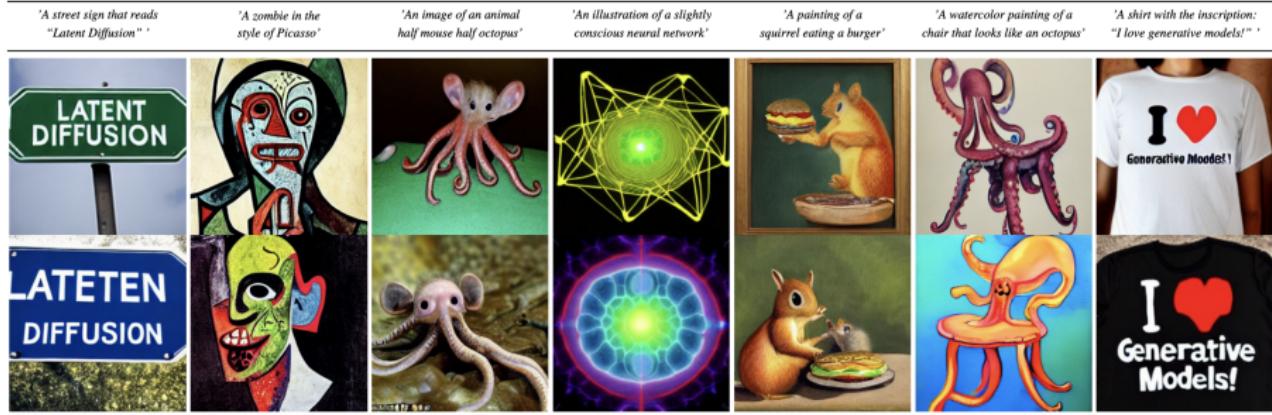
Conditioning & Scaling

- Stable Diffusion [25]: separately train 3 networks
 - Text/condition embedding: only needs text data
 - Diffusion: needs paired (text-image) data but works at **latent** (low dimensional) space
 - Image encoder/decoder: only needs image data
- Text embedding as extra input for denoising net
- Other models like DALL-E [23] also use similar strategies



Results

Text-to-Image Synthesis on LAION. 1.45B Model.



- Even the smallest text-image training set (LAION [28, 27]) has >400M samples!
- MS COCO [14] “only” has 328K images

Reference I

- [1] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al.
Constitutional ai: Harmlessness from ai feedback.
arXiv preprint arXiv:2212.08073, 2022.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al.
Language models are few-shot learners.
Advances in neural information processing systems, 33:1877–1901, 2020.

Reference II

- [3] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman.
Maskgit: Masked generative image transformer.
In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova.
Bert: Pre-training of deep bidirectional transformers for language understanding.
arXiv preprint arXiv:1810.04805, 2018.
- [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel.
Denoising diffusion probabilistic models.
Advances in neural information processing systems, 33:6840–6851, 2020.

Reference III

- [6] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
Training compute-optimal large language models.
arXiv preprint arXiv:2203.15556, 2022.
- [7] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
Lora: Low-rank adaptation of large language models.
arXiv preprint arXiv:2106.09685, 2021.
- [8] Diederik P Kingma and Max Welling.
Auto-encoding variational bayes.
arXiv preprint arXiv:1312.6114, 2013.
- [9] Durk P Kingma and Prafulla Dhariwal.
Glow: Generative flow with invertible 1x1 convolutions.
Advances in neural information processing systems, 31, 2018.

Reference IV

- [10] Quoc V Le and Tomas Mikolov.
Distributed representations of sentences and documents.
In *ICML*, volume 14, pages 1188–1196, 2014.
- [11] Daniel D Lee and H Sebastian Seung.
Learning the parts of objects by non-negative matrix factorization.
Nature, 401(6755):788–791, 1999.
- [12] Daniel D Lee and H Sebastian Seung.
Algorithms for non-negative matrix factorization.
In *Advances in neural information processing systems*, pages 556–562, 2001.

Reference V

- [13] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al.
Retrieval-augmented generation for knowledge-intensive nlp tasks.
Advances in Neural Information Processing Systems, 33:9459–9474, 2020.
- [14] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick.
Microsoft coco: Common objects in context.
In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- [15] Calvin Luo.
Understanding diffusion models: A unified perspective.
arXiv preprint arXiv:2208.11970, 2022.

Reference VI

- [16] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean.
Efficient estimation of word representations in vector space.
arXiv preprint arXiv:1301.3781, 2013.
- [17] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean.
Distributed representations of words and phrases and their compositionality.
In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [18] Varun Nair, Elliot Schumacher, Geoffrey Tso, and Anitha Kannan.
Dera: enhancing large language model completions with dialog-enabled resolving agents.
arXiv preprint arXiv:2303.17071, 2023.

Reference VII

- [19] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al.
Webgpt: Browser-assisted question-answering with human feedback.
arXiv preprint arXiv:2112.09332, 2021.
- [20] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al.
Training language models to follow instructions with human feedback.
Advances in Neural Information Processing Systems, 35:27730–27744, 2022.
- [21] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al.
Improving language understanding by generative pre-training.
OpenAI blog, 2018.

Reference VIII

- [22] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
Language models are unsupervised multitask learners.
OpenAI blog, 2019.
- [23] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen.
Hierarchical text-conditional image generation with clip latents.
arXiv preprint arXiv:2204.06125, 2022.
- [24] Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio.
Contractive auto-encoders: Explicit invariance during feature extraction.
In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 833–840, 2011.

Reference IX

- [25] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
High-resolution image synthesis with latent diffusion models.
In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox.
U-net: Convolutional networks for biomedical image segmentation.
In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.

Reference X

- [27] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al.
Laion-5b: An open large-scale dataset for training next generation image-text models.
Advances in Neural Information Processing Systems, 35:25278–25294, 2022.
- [28] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki.
Laion-400m: Open dataset of clip-filtered 400 million image-text pairs.
arXiv preprint arXiv:2111.02114, 2021.

Reference XI

- [29] Or Sharir, Barak Peleg, and Yoav Shoham.
The cost of training nlp models: A concise overview.
arXiv preprint arXiv:2004.08900, 2020.
- [30] Patrice Simard, Bernard Victorri, Yann LeCun, and John S Denker.
Tangent prop-a formalism for specifying selected invariances in an adaptive network.
In *NIPS*, volume 91, pages 895–903, 1991.
- [31] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al.
Conditional image generation with pixelcnn decoders.
In *Advances in Neural Information Processing Systems*, pages 4790–4798, 2016.

Reference XII

- [32] Aaron Van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu.
Pixel recurrent neural networks.
In *International Conference on Machine Learning*, pages 1747–1756, 2016.
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin.
Attention is all you need.
Advances in neural information processing systems, 30, 2017.
- [34] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol.
Extracting and composing robust features with denoising autoencoders.
In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.

Reference XIII

- [35] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al.
Emergent abilities of large language models.
arXiv preprint arXiv:2206.07682, 2022.
- [36] Matthew D Zeiler and Rob Fergus.
Visualizing and understanding convolutional networks.
In *European conference on computer vision*, pages 818–833. Springer, 2014.