

GANs

Applications

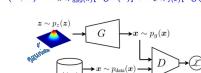
- Labels to street scenes
- How to make
- Labels to code
- How to crop
- How to move
- Labels to photo

Generative Adversarial Network GANs

- ↳ An unsupervised architecture
- ↳ Uses 2 neural networks, competing against each other in order to generate new synthetic instances of data
- ↳ They can learn to mimic the distribution of data
- ↳ Can be used for:
 - ↳ Image generation
 - ↳ Voice generation
 - ↳ Video generation
- ↳ Steps
 - ↳ The generator takes in random numbers
 - ↳ The discriminator is fed:
 - ↳ Real images
 - ↳ Fakes generated by the generator
 - ↳ A stream of images will be generated, some fake images and others previously representing \xrightarrow{G} a Predictor of differentials
 - ↳ To fool

$$\min_{\theta_G} \max_{\theta_D} V(D, G)$$

$$V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(D(x))] + \mathbb{E}_{x \sim p_G(x)} [\log(1 - D(G(x)))]$$



QUESTION

1. GANs are hard to train → Div. convergence hard to detect → how realistic do I want it to be?
2. Hard to meet convergence criteria
3. Vanilla GANs suffer from vanishing gradient problem
4. Overwhelming adversary
5. Can't count objects
6. Sensitive to choice of hyperparameters

as it makes learning O stand because O is hard to learn. Can't be done.

Mode Collapse Problem in GAN

- The mode collapse problem is a common problem in GAN models.
- The generator learns to map several different z values to the same generated data point x .
- Mode collapse usually happens in GAN when the distribution of training data, $p_{\text{data}}(x)$, has multiple modes.

SOLUTIONS

- Minibatch discrimination**
 - resolve the mode collapse problem
 - When mode collapses, all images created look similar.
 - Feed real images and generated images into the discriminator separately in different batches and compute the similarity of the image x with images in the same batch.
 - append the similarity $\omega(x)$ in one of the dense layers in the discriminator to classify whether this image is real or generated.
 - If the mode starts to collapse, the similarity of generated images increases.
 - The discriminator can use this score to detect generated images and penalize the generator if mode is collapsing.

Alternative Strategies to Prevent Mode Collapse

- Uncoupled GANs:** They allow the generator to update itself using a copy of the discriminator from several steps ahead, preventing the generator from overfitting to the current discriminator.
- Experience Replay:** This involves keeping a memory bank of past generated images and occasionally showing them to the discriminator, which helps maintain diversity in the generator's output.
- Modified Training Objectives:** Alternative loss functions, like Wasserstein loss or square losses, provide more stable gradients and can help prevent mode collapse.
- Regularization:** Adding noise to inputs or labels, or using dropout in the discriminator, can prevent it from making overly confident decisions, which in turn pressures the generator to be more diverse.
- Architecture Tweaks:** Adjusting the network architecture, such as adding more layers or changing activation functions, can help the generator explore a wider range of outputs.
- Feature Matching:** The generator is trained to match the statistical features of the real data in some intermediate layer of the discriminator, promoting diversity in the generated data.
- Penalizing the Discriminator:** Introducing penalties for the discriminator when it gets too confident can prevent it from overpowering the generator, leading to a more varied generation.
- Two-Time-Scale Update Rule (TTUR):** Using different learning rates for the generator and the discriminator can help balance their training and prevent the generator from collapsing to two modes.

Application: Image to Image Translation by GAN

- Interactive GAN (IGAN)
- PatchGAN
- CycleGAN (Cycle-Consistent Generative Adversarial Networks)
- Simulated Gan (SimGAN)
- DeepFaceDrawing



Transforming zebra to horse and vice versa



Generating a facial image from a facial sketch.

Training Generative Adversarial Networks

$$\max_{\theta_G} V(D, G) = \max_{\theta_D} \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x))] + p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$x = G(z) \Rightarrow z = G^{-1}(x) \Rightarrow d\pi(z) = d\pi(G^{-1}(x)) dx$$

$$\Rightarrow p_G(x) = p_G(G^{-1}(x)) \cdot |\det(G^{-1}(x))|$$

$$\int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(G(x)))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x)) dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(x))]$$

Conditional GAN

- Assume that the dataset with which GAN is trained has c number of classes.
- The original GAN generates points from any class, and we do not have control to generate a point from a specific class.
- Conditional GAN, also called the conditional adversarial network, gives the user the opportunity to choose the class of generation of points.

For the dataset $\{(x_i \in \mathbb{R}^{d_x})_{i=1}^n, \{y_i \in \mathbb{R}^{c_y}\}_{i=1}^n\}$, let the one-hot encoded class labels be $\{y_i \in \mathbb{R}^{c_y}\}_{i=1}^n$. In conditional GAN, we can use the following loss function instead:

$$\min_{\theta_G} \max_{\theta_D} V_c(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(D(x|y))] + \mathbb{E}_{x \sim p_G(x)} [\log(1 - D(x|y))]$$

$$\min_{\theta_G} \max_{\theta_D} V_c(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(D(x|y))] + \mathbb{E}_{x \sim p_G(x)} [\log(1 - D(G(x|y)))]$$

the discriminator and generator are both conditioned on the labels

Inference from Conditional GAN

- In the inference phase, the user chooses the desired class label and the generator generates a new point from that class.

Application: Text to Image Generation

- Image is generated from some descriptive text.



StackGAN

text-to-image generation.

MixMatch, FineGAN

Designed to combine attributes (like background, texture, shape) from different images or labels.

Other Applications

- In some structures of GAN, learned latent space is meaningful and we can do vector arithmetic in the latent space.
- For example in Deep Convolutional GAN (DCGAN)
- DCGAN uses an all-convolutional network for both generator and discriminator.
- **Inpainting**
- GAN learns to inpaint the lost part based on the available pixels in the image.
- **Medical application**
- Generating histopathology images which can give insight into cancer diagnosis from pathology whole slide images
- **NLP**
- **Speech processing**
- **Network embedding**
- **Logic**
- **Sketch retrieval**

Understanding the objective function

$$\max_{\theta_G} V(D, G) = \max_{\theta_D} \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x))] + p_G(x) [\log(1 - D(x))] dx$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) [\log(D(x))] dx + \int_{\mathcal{X}} p_G(x) [\log(1 - D(x))] dx$$

$$\frac{\partial}{\partial D(x)} (\rho_{\text{data}}(x) \log(D(x)) + p_G(x) \log(1 - D(x))) = 0$$

$$\Rightarrow \frac{\rho_{\text{data}}(x) - p_G(x)}{1 - D(x)} = 0$$

$$\Rightarrow D(x) = \frac{\rho_{\text{data}}(x)}{\rho_{\text{data}}(x) + p_G(x)}$$

$$= \int_{\mathcal{X}} p_{\text{data}}(x) \log(\frac{\rho_{\text{data}}(x)}{\rho_{\text{data}}(x) + p_G(x)}) dx + p_G(x) \log(\frac{p_G(x)}{\rho_{\text{data}}(x) + p_G(x)}) dx - \log(4)$$

$$C(G) = \int_{\mathcal{X}} p_{\text{data}}(x) \log(\frac{\rho_{\text{data}}(x)}{\rho_{\text{data}}(x) + p_G(x)}) dx + \int_{\mathcal{X}} p_G(x) \log(\frac{p_G(x)}{\rho_{\text{data}}(x) + p_G(x)}) dx - \log(4) \geq 0$$

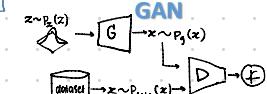
$$\min_{\theta_G} C(G) = 0 + 0 - \log(4) = -\log(4)$$

$$KL(p_{\text{data}}(x) || \frac{\rho_{\text{data}}(x) + p_G(x)}{2}) = 0$$

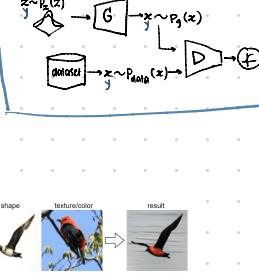
$$\text{when } \rho_{\text{data}}(x) = \frac{\rho_{\text{data}}(x) + p_G(x)}{2} \Rightarrow \rho_{\text{data}}(x) = p_G(x)$$

Implementation of Conditional GAN

- Concatenate the one-hot encoded label y to the point x for the input to the discriminator.
- Concatenate the one-hot encoded label y to the noise z for the input to the generator.
- For these, the input layers of discriminator and generator are enlarged to accept the concatenated inputs



Conditional GAN



Application: Mixing Image Characteristics

- FineGAN
- An unsupervised GAN model which disentangles the features of the generated image to background, shape, and color/texture.
- FineGAN generates an image hierarchically. It starts with generating the background.

MixMatch

- built upon FineGAN
- It gives the user the opportunity to choose the background, shape, and colour/texture from three pictures and it generates an image with the chosen characteristics.

generating an image by borrowing its characteristics from three images using MixMatch

generating an image by borrowing its characteristics from different domains using improved MixMatch.

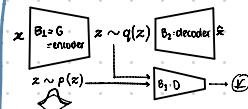


Vector arithmetic in the latent space on images by DCGAN.

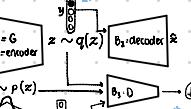
Adversarial Autoencoder (AAE)

- In contrast to variational autoencoder which uses KL divergence and evidence lower bound, AAE uses adversarial learning for imposing a specific distribution on the latent variable in its coding layer.

Unsupervised AAE

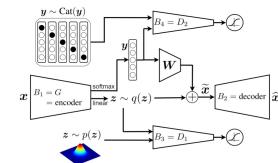


Supervised AAE



Dimensionality Reduction with AAE

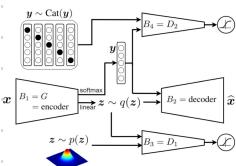
The low-dimensional representation of.



Semi-Supervised AAE

- Consider a partially labelled dataset.
- The labelled part of data has n number of classes.
- AAE can be used for semi-supervised learning with the partially labelled dataset.
- We can use the same structure for Semi-Supervised.
- But rather than the clusters, we assume we have n number of classes.
- We have a partially labelled part of the dataset.

Semi-Supervised AAE

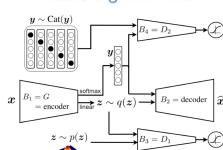


- If the point x has a label use it. Labels are one-hot vectors.
- If the point x does not have any label, randomly sample a label $y \in \mathbb{R}^n$ from a categorical distribution, i.e., $y \sim \text{Cat}(y)$
- This categorical distribution gives a one-hot encoded vector where the prior probability of each class is estimated by the proportion of the class's population to the total number of labelled points.

Clustering with AAE

- Assume we have c number of clusters.
- All points are unlabeled.
- The cluster indices are sampled randomly by the categorical distribution.
- The cluster labels and the latent code are both trained.

Clustering with AAE



Comparisons

GAN:

Think of a forger (G) and a detective (D). G tries to fool D with fake images.

VAE:

Learns a smooth, interpretable "latent space" (compressed version of data). Generates new samples by decoding latent vectors.

Loss is based on how well it reconstructs and how close latent vectors are to a normal distribution. Tends to produce blurrier images but better control.

AAE:

Combines VAE's structure (encoder-decoder) with GAN's loss in the latent space. Instead of using KL divergence (like VAE), uses a GAN to make sure latent codes look like the prior (e.g., normal distribution).

Better at clustering and generating structured latent representations.

Conditional GAN:

Like GAN but with labels added as input.

For example, "generate a 7" instead of any digit.

Helps solve the lack of control in vanilla GAN.

Model	Pros	Cons
GAN	Realistic generation, no labels needed	Training instability, mode collapse, no inference
AAE	Generative + encoding, flexible latent priors	Less sharp images, adversarial tuning needed
Supervised AAE	Combines generation & classification	Needs labeled data, may overfit
Unsupervised AAE	No labels, useful for clustering	No semantic control, weaker generation quality
Conditional GAN	Controlled generation by labels	Needs labels, can overfit, GAN issues remain
Semi-Supervised AAE	Uses limited labels effectively	Complex to train, balancing losses is tricky
Clustering AAE	Finds groups without labels, useful latent representations	May need to predefine cluster count, cluster quality depends on prior

When to Use

1. GANs (Generative Adversarial Networks)

Use when: You want to generate realistic data (e.g., images, speech) from noise.

Input: Random noise z

Output: Synthetic sample

Labels needed? No

Use cases: Image generation, super-resolution, style transfer, image inpainting

2. AAE (Adversarial Autoencoder)

Use when: You want to impose a prior distribution on latent space and reconstruct input data.

Combines: Autoencoder + GAN

Input: Data sample x

Output: Reconstructed x + latent vector z

Labels needed? No

Use cases: Dimensionality reduction, clustering, anomaly detection

3. Supervised AAE

Use when: You have labeled data and want the latent space to reflect class structure.

Labels needed? Yes

Learn: Both reconstruction and class-discriminative latent codes

Use cases: Semi-supervised learning, structured latent spaces

4. Unsupervised AAE

Use when: You want to learn structure in data without labels

Labels needed? No

Learn: Prior matching of latent space + reconstruction

Use cases: Clustering, representation learning, generative modeling

5. Conditional GAN (cGAN)

Use when: You want to generate data conditioned on labels or attributes

Input: Noise z + condition y (e.g., class label)

Output: Sample matching the label

Labels needed? Yes (for training)

Use cases: Text-to-image, image-to-image translation, label-controlled generation

6. Semi-Supervised AAE

Use when: You have both labeled and unlabeled data

Use: A small amount of labeled data to guide training; unlabeled data to improve generalization

Labels needed? Partially

Use cases: When labeled data is scarce but unlabeled data is abundant (e.g., medical imaging, fraud detection)

7. Clustering with AAE

Use when: You want to automatically cluster data in latent space

Technique: Sample latent codes from a mixture of distributions (e.g., Gaussian + Categorical)

Use cases: Discovering groups in data, unsupervised classification, data exploration

B) Assume that we want to classify images in 3 categories of Apple, Orange and Banana for which the prior probabilities are $P(A) = 0.3$, $P(O) = 0.4$, $P(B) = 0.2$. We see a new image x from which we extract a feature vector x . The likelihoods are as:

$$P(x | A) = 0.24$$

$$P(x | O) = 0.15$$

$$P(x | B) = 0.38$$

What should image x be classified as? An apple, an orange, or a banana? Why?

Write down the complete derivations. (7 points)

$$P(A | x) = p(x | A) p(A) / p(x) = 0.24 * 0.3 / p(x) = 0.72 / p(x)$$

$$P(O | x) = p(x | O) p(O) / p(x) = 0.15 * 0.4 / p(x) = 0.75 / p(x)$$

$$P(B | x) = p(x | B) p(B) / p(x) = 0.38 * 0.2 / p(x) = 0.76 / p(x)$$

$$\Rightarrow P(B | x) > P(O | x) > P(A | x)$$

\Rightarrow It has to be classified as **Banana**

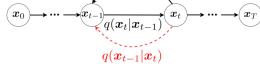
7. Diffusion Models

Diffusion Models

Diffusion models are generative models that:
Gradually add noise to an image (forward process).
Then learn to reverse the process and remove the noise to generate a clean image (reverse process).
• Diffusion models use a Markov chain to do it

Denoising Diffusion Probabilistic Models (DDPMs)

So, Training = learn this
 $p_{\theta}(x_{T-1}|x_T)$
reverse process.



Joint Distribution over Latent States in DDPMs

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

on total initial state
steps from time 1 to T

$$q(x_t|x_{t-1}) \propto \mathcal{N}(x_t; \mu_t = \theta_t(x_{t-1}), \sigma_t^2 = \beta_t)$$

Forward Pass
Noise

$$p_{\theta}(x_t|x_0) \sim \mathcal{N}(x_t; \mu_t = \theta_t(x_0), \sigma_t^2 = \beta_t)$$

Backward Pass
Noise

$$x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_T$$

Data
Noise
 $x_0 \leftarrow x_1 \leftarrow \dots \leftarrow x_T$

*Identity matrix is used to simplify calculation
 $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ - covariance

How Works DDPMs

1. Forward Diffusion Process

Start with a real image x_0 .
At each timestep t , add a small amount of Gaussian noise:

$$\epsilon \sim \mathcal{N}(0, I)$$

where $\epsilon \sim \mathcal{N}(0, I)$ and β_t is a small noise variance.

After many steps (say 1000), the image becomes nearly random noise.

2. Reverse Process (Generation)

Train a neural network (usually a U-Net) to predict the noise ϵ added at each step, given the noisy image x_t and the timestep t .
The model learns to denoise the image one step at a time, using the learned noise estimate: $x_{t-1} = x_t - \theta_t(\epsilon_t)$.

*Starting from pure noise x_T , apply the denoising steps iteratively to get a new, realistic sample x_0 .

1. What is the key difference between a standard autoencoder and a VAE?

A standard autoencoder learns a deterministic mapping (input \rightarrow compressed code \rightarrow output).

A VAE learns a probabilistic mapping: it maps input to a distribution (usually Gaussian), samples from it, and then decodes.

2. Why do we need the KL divergence term in the VAE loss function?

It ensures that the latent distribution stays close to a standard normal distribution ($\mathcal{N}(0, 1)$). Without this, sampling from the latent space would generate unrealistic or untrained outputs.

3. What is the parameterization trick, and why is it used in VAEs?

Trick: Instead of sampling $\epsilon \sim \mathcal{N}(0, I)$ we write,

$$z = \mu + \epsilon \sim \mathcal{N}(\mu, \sigma^2)$$

It's used to make sampling differentiable, so we can backpropagate through it.

4. In a VAE, what are μ and σ used for?

They define the mean and standard deviation of the latent distribution for a given input.
They're used to sample a latent vector z that gets decoded.

5. What does it mean to "sample" from the latent space in a VAE?

It means generating a value z from the learned Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$, rather than using a fixed code.

6. Can a VAE generate new data without any input image? Explain how.

Yes! You can randomly sample $z \sim \mathcal{N}(0, 1)$ and feed it into the decoder to generate a new image.

7. Why do diffusion models add noise to data during training?

To simulate a process where clean data slowly becomes noise.

This helps the model learn how to denoise during generation.

8. What is the difference between the forward process and reverse process in diffusion models?

Forward: Gradually add Gaussian noise: $x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_T$

Reverse: Learn to remove noise: $x_T \rightarrow x_{T-1} \rightarrow \dots \rightarrow x_0$

9. In DDPMs, what does the model learn during training?

The model learns to predict the noise added at each step, or equivalently, how to denoise the image.

10. How are DDPMs different from VAEs in terms of latent representation?

VAEs use one latent variable z .

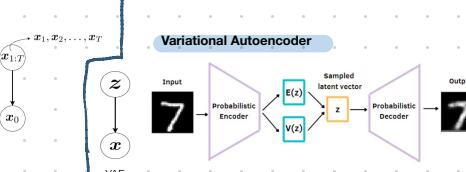
DDPMs use a sequence of noisy variables x_1, x_2, \dots, x_T .

DDPMs' latent space is temporal, evolving step by step.

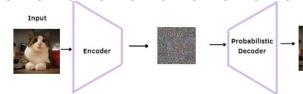
11. Why do Diffusion Models produce better images than VAEs?

VAEs optimize a tradeoff between reconstruction and regularization \rightarrow can produce blurry results.

Diffusion models focus on denoising at each step, allowing sharper and more high-fidelity images, at the cost of slower sampling.



VAEs (Variational Autoencoders) are a type of generative model in machine learning used to learn the underlying structure of data and generate new data that is similar to the original. They are especially popular in unsupervised learning and are used in applications like image generation, anomaly detection, and representation learning.



how it works VAE

Input to Encoder: An image (or other data) is fed into the encoder.

Latent Distribution: Instead of producing a single vector, the encoder outputs two vectors:

Mean vector (μ)

Standard deviation vector (σ)

These define a normal distribution: $z \sim \mathcal{N}(\mu, \sigma^2)$.

Sampling (Reparameterization Trick): A latent vector z is sampled from this distribution. To make it differentiable (for backpropagation), a trick is used:

$$z = \mu + \sigma \epsilon$$

4. Decoder Reconstruction Output: The decoder takes z and tries to reconstruct the original input.

Loss Functions

Reconstruction Loss: Measures how close the output is to the original input (e.g., using MSE or binary cross-entropy).

-KL Divergence: Encourages the latent distribution to be close to a standard normal distribution ($\mathcal{N}(0, 1)$), helping regularize the latent space.

Total Loss = Reconstruction Loss + KL Divergence

VAE SAMPLING

$$\theta: \mu = 0, \sigma = 0.5, \epsilon = 0.3$$

$$\begin{aligned} z &= \mu + \sigma \epsilon \\ &= 0 + 0.5(0.3) \\ &= 0.15 \end{aligned}$$

Diffusion Model vs VAEs

In DPM

- Observed variable x
- Hidden Variable z
- Latent variable z , direct mapping
- Multiple noisy steps — a full sequence, not just one latent variable.

Objective Functions

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

$$p_{\theta}(x_0) = \int p_{\theta}(x_0, x_1, \dots, x_T) dx_1 \dots dx_T$$

Variational Lower Bounds

$$p_{\theta}(x) \geq \mathbb{E}_{q(z|x)} [\log p_{\theta}(x|z)] - D_{KL}(q(z|x) || p_{\theta}(z))$$

YES

Yes, but optimized differently

$$p_{\theta}(x_0) \geq \mathbb{E}_{q(x_1:T|x_0)} [\log p_{\theta}(x_0|x_1:T)] - D_{KL}(q(x_1:T|x_0) || p_{\theta}(x_1:T))$$

LOSS

KL divergence + reconstruction

Target

$$\text{maximize ELBO}$$

ELBO

ELBO ensures that latent variable z follows a desired distribution (like Gaussian).

Encoder

$$\text{Maps } n \rightarrow 2$$

Decoder

$$\text{Maps } 2 \rightarrow n$$

Process Type

$$1 \text{ step}$$

Similar principle, but over many time steps.
Loss is decomposed for each timestep and simplified

Forward Pass
Adds noise: $x_0 \rightarrow x_1$

Backward Pass

Learns denoising: $x_T \rightarrow x_0$

Multistep (Markov Chain)

Diversion w/o Gaussians

$$\theta: \mu = 1, \sigma = 1$$

$$KL(N(\mu, \sigma) || N(0, 1)) = -\frac{1}{2} (1 + \log(\sigma^2) - \mu^2 - 1)$$

$$\therefore \sigma^2 = 1$$

$$\therefore \log(\sigma^2) = 0$$

$$\therefore KL = \frac{1}{2} (1 + 0 - 1 - 1)$$

$$\therefore 0.5$$

Diffusion Forward Process calculation

$$\theta: x_0 = 0, x_{T-1} = 1, \epsilon = 0.5$$

$$x_T = \sqrt{\epsilon} x_{T-1} + \sqrt{1-\epsilon} (\epsilon)$$

$$\therefore \sqrt{0.5} (0.5) + \sqrt{1-0.5} (0.5)$$

$$\therefore 0.911$$

Diffusion Backward Process calculation

$$\theta: x_T = 2, \epsilon = 0.9$$

$$\therefore 0.7$$

1. ELBO (Evidence Lower Bound)

ELBO is the function VAEs maximize during training.
It is a lower bound on the log-likelihood of the data ($\log p(x)$).
VAEs can't directly maximize $\log p(x)$ (too complex), so they optimize ELBO instead.

Variational Lower Bound

$$\begin{aligned} & \mathbb{E}_{q(x_t | x_0)} [\log p_{\theta}(x_t | x_0)] - D_{KL}(q(x_t | x_0) || p_{\theta}(x_t | x_0)) \\ &= \mathbb{E}_{q(x_t | x_0)} [\log p_{\theta}(x_0 | x_t)] - \int q(x_t | x_0) \log \frac{p_{\theta}(x_t | x_0)}{p_{\theta}(x_0 | x_t)} dx_0 \\ &= \mathbb{E}_{q(x_t | x_0)} [\log p_{\theta}(x_0 | x_t)] - \mathbb{E}_{q(x_t | x_0)} \left[\log \frac{p_{\theta}(x_t | x_0)}{p_{\theta}(x_0 | x_t)} \right] \\ &= \mathbb{E}_{q(x_t | x_0)} \left[\log p_{\theta}(x_0 | x_t) - \log \frac{p_{\theta}(x_t | x_0)}{p_{\theta}(x_0 | x_t)} \right] \\ &= \mathbb{E}_{q(x_t | x_0)} \left[\log p_{\theta}(x_0 | x_t) + \log \frac{p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] \\ &= \mathbb{E}_{q(x_t | x_0)} \left[\log \frac{p_{\theta}(x_0 | x_t) p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] = \mathbb{E}_{q(x_t | x_0)} \left[\log \frac{p_{\theta}(x_t | x_0)}{q(x_t | x_0)} \right] \\ &= \mathbb{E}_q \left[\log \frac{p_{\theta}(x_0 | x_t)}{q(x_t | x_0)} \right] \\ &= \mathbb{E}_q \left[\log \frac{p(x_T)}{\prod_{t=1}^T p(x_t | x_{t-1})} \right] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q [\log p(x_t | x_{t-1})] - \mathbb{E}_q [\log q(x_t | x_{t-1})] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q [\log p_{\theta}(x_t | x_{t-1})] - \mathbb{E}_q \left[\log \prod_{t=1}^T q(x_t | x_{t-1}) \right] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q [\log p_{\theta}(x_t | x_{t-1})] - \sum_{t=1}^T \mathbb{E}_q [\log q(x_t | x_{t-1})] \\ &= \mathbb{E}_q [\log p(x_T)] + \sum_{t=1}^T \mathbb{E}_q \left[\log \frac{p_{\theta}(x_t | x_{t-1})}{q(x_t | x_{t-1})} \right] \end{aligned}$$

2. Decomposing the ELBO

ELBO has two main parts:

ELBO = Reconstruction Term - KL Divergence

Reconstruction Term: Measures how well the decoder reconstructs input x from latent code z . It encourages accurate generation.

KL Divergence: Regularizes the latent distribution to be close to a standard normal ($N(0, 1)$), preventing overfitting.

Decomposing the ELBO

$\mathbb{E}_q[\log p(x_T)]$: Validates the model's calibration to the true noise at the last diffusion step, ensuring the reverse process starts accurately.

$\sum_{t=1}^T \mathbb{E}_q \left[\log \frac{p_{\theta}(x_t | x_{t-1})}{q(x_t | x_{t-1})} \right]$: Acts as a regularizer by comparing the model's backward predictions to the known forward diffusion, guiding the model to correct any discrepancies.

$$\mathcal{L}(x_0) = \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p(x_T) + \sum_{t=1}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right]$$

By the Markov property:

$$q(x_t | x_{t-1}) = q(x_t | x_{t-1}, x_0),$$

and by Bayes' rule:

$$q(x_t | x_{t-1}, x_0) = \frac{q(x_t | x_1, x_0) q(x_1 | x_0)}{q(x_{t-1} | x_0)}.$$

Plugging this equation into the ELBO, we get

$$\begin{aligned} \mathcal{L}(x_0) &= \mathbb{E}_{q(x_{1:T}|x_0)} \left[\log p(x_T) + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \right] \\ &+ \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t)}{q(x_t | x_0)} = \log p_{\theta}(x_0 | x_1) \end{aligned}$$

* = $-\log q(x_T | x_0) + \log q(x_1 | x_0)$

Hence the negative ELBO (variational upper bound) becomes

$$\begin{aligned} & -\mathbb{E}_q \left[\log \frac{p(x_T)}{q(x_T | x_0)} + \sum_{t=2}^T \log \frac{p_{\theta}(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} + \log p_{\theta}(x_0 | x_1) \right] \\ &= \underbrace{D_{KL}(q(x_T | x_0) || p(x_T))}_{L_T(\epsilon_T)} \\ &+ \sum_{t=2}^T \underbrace{\mathbb{E}_{q(x_{1:T}|x_0)} D_{KL}(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t))}_{L_{t-1}(\epsilon_{t-1})} \\ &- \underbrace{\mathbb{E}_{q(x_{1:T}|x_0)} \log p_{\theta}(x_0 | x_1)}_{L_0(\epsilon_0)} \end{aligned}$$

3. Reverse Process (in Diffusion Models)

Starts with pure noise and gradually denoises it using a trained model.

At each step, the model predicts the noise and removes it to recover the original data.

It's the generation phase of diffusion models.

Revers Process

We need to learn a neural network to approximate the conditioned probability distributions in the reverse diffusion process.

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \Sigma_{\theta}(\mathbf{x}_t, t)).$$

$$q(\mathbf{x}_t | \mathbf{x}_0)$$

Recall that the forward process is defined by the equation:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \sqrt{1-\alpha_t} \mathbf{x}_{t-1}, \beta_t I)$$

We can sample \mathbf{x}_t at any time step t in a closed form.

Let $\alpha_t = 1 - \beta_t$ and $\tilde{\alpha}_t = \prod_{i=1}^t \alpha_i$:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \sqrt{\tilde{\alpha}_t} \mathbf{x}_0, (1 - \tilde{\alpha}_t)I)$$

4. Mean (μ) and Variance (σ^2) in VAEs

The encoder outputs two vectors: mean (μ) and log-variance ($\log \sigma^2$). They define a Gaussian distribution in latent space:

$\mathbf{z} \sim N(\mu, \sigma^2)$

Sampling from this gives diverse and continuous latent representations.

Mean of $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$

The reverse conditional probability is tractable when conditioned on \mathbf{x}_0 :

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \hat{\mu}(\mathbf{x}_t, \mathbf{x}_0), \hat{\Sigma}(t))$$

Using Bayes' rule, we have:

$$\begin{aligned} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) &= q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \\ &\propto \exp \left(-\frac{1}{2} \left(\frac{(\mathbf{x}_t - \sqrt{\tilde{\alpha}_{t-1}} \mathbf{x}_{t-1})^2}{\beta_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{1-\alpha_{t-1}} \mathbf{x}_0)^2}{1-\alpha_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\tilde{\alpha}_{t-1}} \mathbf{x}_0)^2}{1-\tilde{\alpha}_{t-1}} \right) \right) \\ &= \exp \left(-\frac{1}{2} \left(\frac{\alpha_t}{\beta_t} \mathbf{x}_t^2 - 2\sqrt{\tilde{\alpha}_{t-1}} \mathbf{x}_{t-1} \mathbf{x}_t + \frac{\alpha_{t-1}^2}{1-\alpha_{t-1}} \mathbf{x}_{t-1}^2 - 2\sqrt{1-\alpha_{t-1}} \mathbf{x}_{t-1} \mathbf{x}_0 + \alpha_{t-1} \mathbf{x}_0^2 }{1-\tilde{\alpha}_{t-1}} \right) \right) \\ &= \exp \left(-\frac{1}{2} \left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\alpha_{t-1}} \right) \mathbf{x}_t^2 - 2\frac{\sqrt{\tilde{\alpha}_{t-1}}}{\sqrt{1-\alpha_{t-1}}} \mathbf{x}_{t-1} \mathbf{x}_t + \left(\frac{2\sqrt{\tilde{\alpha}_{t-1}}}{\sqrt{1-\alpha_{t-1}}} \mathbf{x}_{t-1} + \alpha_{t-1} \mathbf{x}_0 \right) \mathbf{x}_t + C(\mathbf{x}_t, \mathbf{x}_0) \right) \right) \end{aligned}$$

where $C(\mathbf{x}_t, \mathbf{x}_0)$ is some function not involving \mathbf{x}_{t-1} .

We can express the mean and variance in the following manner: (recall that $\alpha_t = 1 - \beta_t$ and $\tilde{\alpha}_t = \prod_{i=1}^t \alpha_i$):

$$\begin{aligned} \hat{\beta}_t &= 1 / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\tilde{\alpha}_{t-1}} \right) = 1 / \left(\frac{\alpha_t - \tilde{\alpha}_t + \beta_t}{\beta_t (1 - \tilde{\alpha}_{t-1})} \right) = \frac{1 - \tilde{\alpha}_{t-1}}{1 - \alpha_t} \cdot \beta_t \\ \hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) &= \left(\frac{\sqrt{\tilde{\alpha}_t}}{\beta_t} \mathbf{x}_t + \frac{\sqrt{1-\alpha_{t-1}}}{1-\alpha_{t-1}} \mathbf{x}_0 \right) / \left(\frac{\alpha_t}{\beta_t} + \frac{1}{1-\alpha_{t-1}} \right) \\ &= \left(\frac{\sqrt{\tilde{\alpha}_t}}{\beta_t} \mathbf{x}_t + \frac{\sqrt{1-\alpha_{t-1}}}{1-\alpha_{t-1}} \mathbf{x}_0 \right) \frac{1}{1-\alpha_t} \cdot \beta_t \\ &= \frac{\sqrt{\tilde{\alpha}_t} (1 - \alpha_{t-1})}{1 - \alpha_t} \mathbf{x}_t + \frac{\sqrt{1-\alpha_{t-1}}}{1 - \alpha_t} \mathbf{x}_0 \end{aligned}$$

Recall $\mathbf{x}_0 = \frac{1}{\sqrt{\tilde{\alpha}_t}} (\mathbf{x}_t - \sqrt{1-\tilde{\alpha}_t} \mathbf{x}_0)$

$$\begin{aligned} \hat{\mu}_t &= \frac{\sqrt{\tilde{\alpha}_t} (1 - \alpha_{t-1})}{1 - \alpha_t} \mathbf{x}_t + \frac{\sqrt{1-\alpha_{t-1}}}{1 - \alpha_t} \frac{1}{\sqrt{\tilde{\alpha}_t}} (\mathbf{x}_t - \sqrt{1-\tilde{\alpha}_t} \mathbf{x}_0) \\ &= \frac{1}{\sqrt{\tilde{\alpha}_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \mathbf{x}_0 \right) \end{aligned}$$

5. KL Divergence

Measures how much one distribution (e.g., $q(\mathbf{x}|\mathbf{x}_0)$) differs from another (e.g., $N(0, 1)$).

In VAEs, it keeps the latent space well-behaved and prevents overfitting.

Lower KL = closer to normal distribution.

KL for two Gaussian

$$D_{KL}(q(x_{t-1} | x_t, x_0) || p_{\theta}(x_{t-1} | x_t))$$

$$q(x) = \mathcal{N}(x; \mu_q, \Sigma_q) \quad \text{and} \quad p(x) = \mathcal{N}(x; \mu_p, \Sigma_p),$$

$$D_{KL}(q || p) = \frac{1}{2} \left(\text{tr}(\Sigma_p^{-1} \Sigma_q) + (\mu_p - \mu_q)^T \Sigma_p^{-1} (\mu_p - \mu_q) - k + \ln \left(\frac{\det \Sigma_p}{\det \Sigma_q} \right) \right)$$

for q

$$\hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\tilde{\alpha}_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \tilde{\alpha}_t}} \mathbf{x}_0 \right)$$

for p

$$\mu_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right)$$

6. VAE Loss Function

Total loss = Reconstruction Loss + KL Divergence

Reconstruction Loss ensures the output matches the input.

KL Divergence regularizes the latent space.

Loss function

$$\begin{aligned} L_t &= E_{\mathbf{x}_0, \epsilon} \left[\frac{1}{2 \| \hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) \|^2} \| \hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_{\theta}(\mathbf{x}_t, t) \|^2 \right] \\ &= E_{\mathbf{x}_0, \epsilon} \left[\frac{1}{2 \| \hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) \|^2} \left\| \frac{1}{\sqrt{\tilde{\alpha}_t}} (\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} \epsilon) - \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t)) \right\|^2 \right] \\ &= E_{\mathbf{x}_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \tilde{\alpha}_t) \| \hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) \|^2} \| \epsilon_t - \epsilon_{\theta}(\mathbf{x}_t, t) \|^2 \right] \\ &= E_{\mathbf{x}_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \tilde{\alpha}_t) \| \hat{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) \|^2} \| \epsilon_t - \epsilon_{\theta}(\sqrt{\tilde{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \tilde{\alpha}_t} \epsilon_{\theta}(\mathbf{x}_t, t)) \|^2 \right] \end{aligned}$$

• Ho et al. (2020) discovered empirically that the diffusion model produces higher-quality images when using a simplified objective, omitting the weighting term

$$\begin{aligned} L_{\text{simple}} &= E_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} [\| \epsilon_t - \epsilon_{\theta}(\mathbf{x}_t, t) \|^2] \\ &= E_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} [\| \epsilon_t - \epsilon_{\theta}(\sqrt{\tilde{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \tilde{\alpha}_t} \epsilon_{\theta}(\mathbf{x}_t, t)) \|^2] \end{aligned}$$

Transformer

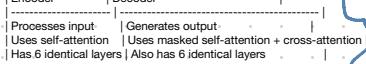
What is a Transformer?

A Transformer is a deep learning architecture

Unlike RNNs, it:

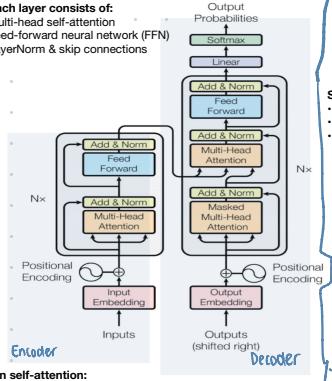
Processes all words in parallel

Uses self-attention to understand relationships between words is the foundation for models like BERT, GPT, T5, etc.



Each layer consists of:

- Multi-head self-attention
- Feed-forward neural network (FFN)
- LayerNorm & skip connections



In self-attention:

Every word compares itself to every other word.
It answers: "How important is word X to word Y?"

Final Output Layer

Linear Projection:

• Primary Role: Adjusting dimensionality.

• The linear layer serves to change the dimensionality of the feedforward network's output to match the size of the vocabulary.

• This ensures that the output has a dimension corresponding to every word in the dictionary.

Softmax Activation:

• This function transforms the linear layer's output into probabilities.

• It's applied to predict the next word.

1. Why do Transformers need Positional Encoding?

Transformers don't use RNNs or CNNs, so they have no built-in sense of order.

Positional Encoding adds information about the position of each word/token in the sequence to the input embeddings so the model can still understand sequence structure.

2. What is the difference between Self-Attention and Cross-Attention?

Self-Attention

Input = Output = Same seq,

| Input × Output sequences
Used in Encoder/Decoder | Used only in Decoder

Each token attends to all others in the same sequence | Each token in the decoder attends to the encoder output |

Cross-Attention

Input = Output = Same seq,

| Input × Output sequences
Used only in Decoder

Each token attends to all others in the same sequence | Each token in the decoder attends to the encoder output |

II. Self Supervised Learning → BERT, GPT

Bi-directional Encoder Representations from Transformers (BERT)

- Pre-training of Deep Bi-directional Transformers for Language Understanding
- BERT is built using transformer encoder blocks.

BERT: Bi-directional language model

- Masks words in the input and asks the model to predict the missing word.
- Additional task: Given two sentences (A and B), is B likely to be the sentence that follows A, or not?

- BERT is designed to pretrain bidirectional representations from unlabeled text.
- Jointly conditioning on both left and right context.
- The pre-trained BERT model can be finetuned with just one
- It creates state-of-the-art models for a wide range of tasks, such as question answering and language inference.
- additional output layer.

[CLS] Token in BERT

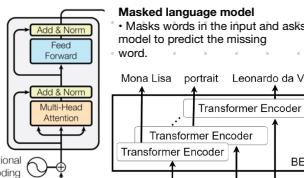
- The [CLS] token is prepended to the input text and travels through the Transformer layers alongside other tokens.
- All tokens, including [CLS], gather contextual information from the entire sequence due to the self-attention mechanism.
- For sentence-level tasks, the final hidden state of the [CLS] token is used as the sentence representation.
- During fine-tuning on a specific task, the model learns to imbue the [CLS] token with a meaningful representation of the entire sentence, optimized for that task.
- Example Usage: In classification tasks, the [CLS] token representation is fed into a classifier to determine the sentence's class.

1. Multilingual BERT:

- Trained on 104 different languages, capable of "zero-shot" adapting to a new language domain.

2. Domain Specific BERT Variants:

- BiBERT: Retrained on a biomedical corpus.
- SciBERT: Trained on over one million published articles.
- BERTweet: A Roberta model trained on 850 million tweets.
- FinBERT: Adapted to the financial domain.



BERT is basically a trained Transformer Encoder stack

1. Encoder Layers: 6
2. FFNN Hidden Layer Units: 512 3. Attention Heads: 8
3. BERT Base
1. Encoder Layers: 12
2. FFNN Hidden Layer Units: 768 3. Attention Heads: 12
4. Total Parameters: 110 million
3. BERT Large
1. Encoder Layers: 24
2. FFNN Hidden Layer Units: 1024 3. Attention Heads: 16
4. Total Parameters: 340 million

RoBERTa:

- Optimizes BERT's training process by using more data, larger batch sizes, and longer training times, resulting in improved performance on NLP tasks.

TinyBERT:

- A smaller and faster version of BERT designed for resource-constrained environments, retaining competitive performance with significantly fewer parameters.

Pretraining Objectives

- Supervised training conducted on downstream tasks provided by GLUE and SuperGLUE benchmarks, reformulated into text-to-text tasks.
- Self-supervised training employs corrupted tokens: 15% of tokens are randomly removed and replaced with sentinel tokens.

Generative Pre-trained Transformer (GPT)

- Improving Language Understanding by Generative Pre-Training
- The GPT is built using transformer decoder blocks.

Predict the next word, given all of the previous words

Architecture:

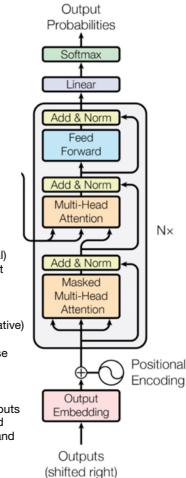
- Stack of transformer decoder blocks
- No encoder, hence no cross-attention module.
- Components: Positional encoding, masked multihead self-attention, and feedforward network.
- Directionality:
- Only considers previous (left) words in attention, not bidirectional like BERT.
- Utilizes masked multihead self-attention for this purpose.

GPT 3

- Released: 2018
- Parameters: 117 Million
- Layers: 12
- Training Data: Books1 Corpus
- Focus: Unsupervised pre-training, Transformer architecture, large-scale language modeling

GPT 4

- Released: 2019
- Parameters: ~1.5 Billion
- Layers: 48
- Training Data: 40GB (English)
- Focus: Transformer architecture, self-attention mechanism
- Release: Not Yet
- Parameters: ~100 Trillion (speculative)
- Layers: Unknown
- Training Data: Larger, more diverse (speculative)
- Focus: GPT-4 is known to be a multimodal model, capable of processing both text and image inputs to generate text outputs. Advanced few-shot learning, improved NLU and NLG, reasoning and inference



T5 - Text-to-Text Transfer Transformer.

- Introduced as a unified framework for addressing NLP tasks by converting them into a text-to-text format.
- Key Focus:** Transfer learning - Pre-training on data-rich tasks followed by fine-tuning on downstream tasks.

- Operates on an encoder-decoder model.
- Pretrained on a multi-task mixture of both unsupervised and supervised tasks, each converted into a text-to-text format.
- Utilizes relative scalar embeddings; encoder input padding can be performed on both left and right.

Task Adaptation

- T5 adapts to various tasks by prepending task-specific prefixes to the input, e.g., for translation: "translate English to German: ...", for summarization: "summarize: ...".
- Achieves very good results on many benchmarks including summarization, question answering, and text classification.

The Challenge:

GPT just predicts what comes next – but it doesn't inherently follow instructions.

The Fix: RLHF (Reinforcement Learning with Human Feedback)

3-Step Process:

Supervised Fine-Tuning (SFT): Train with human-curated answers.

Reward Model (RM): Learn what answers humans prefer

Reinforcement Learning (RLHF): Optimize GPT to produce better aligned outputs

Core Challenge with GPT Models:

- GPT models predict the next token based on historical data, lacking an innate ability to follow instructions.
- GPT (Generative Pre-trained Transformer) models, at their core, predict the next word or token in a sequence based on the probabilities derived from pre-training on extensive text corpora.
- They don't inherently "understand" instructions or follow commands but generate what's statistically likely to come next, given their training.

The Goal of Alignment:

- The aim is to align GPT's responses with specific user instructions and ethical standards, beyond just generating probable text.
- The primary objective is to bridge the gap between these statistical predictions and meaningful adherence to instructions provided by users.
- This involves ensuring that the AI's responses are not just contextually appropriate or conversationally relevant but also aligned with the specific intentions, ethical expectations, and task-oriented goals of the user.

Key Areas of Focus in AI Alignment

- Learning from Human Feedback:
 - Tailoring AI through human interaction.
 - Training to Follow Instructions.
 - Teaching AI specific task adherence.
 - Evaluating AI Harms:
 - Identifying risks in AI outputs.
 - Modifying Behavior to Mitigate Harms:
 - Adjusting AI to prevent negative impacts.

1. Main difference between self-supervised and supervised learning?

Self-supervised learning uses unlabeled data and creates learning objectives from the data itself (e.g., predicting a missing word). Supervised learning requires labeled datasets with human-annotated answers.

2. Why does BERT use bidirectional attention but GPT does not?

BERT is trained to understand context from both sides of a word to capture deeper meaning (e.g., in masked word prediction). GPT is trained for generating text, so it must only look at previous tokens, not future ones – hence unidirectional.

3. How does Masked Language Modeling work in BERT?

Random words in a sentence are replaced with a [MASK] token.

The model learns to predict the masked word based on its left and right context.

4. What is the purpose of the [CLS] token in BERT?

The [CLS] token is prepended to every input. After going through the Transformer, it represents the whole sentence. Used for sentence-level tasks like classification.

5. Function of masked multi-head self-attention in GPT?

It ensures that each word only attends to previous words, not future ones – important for generating text sequentially.

6. Why can GPT generate text but BERT cannot?

GPT is trained to predict the next word, making it naturally suited for text generation. BERT is bidirectional and masked, so it lacks a natural left-to-right generation structure.

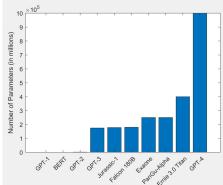
7. What does NSP (Next Sentence Prediction) teach BERT?

NSP trains BERT to understand sentence relationships – whether one sentence follows another in natural text.

8. 3 steps in GPT alignment with RLHF?

- Supervised Fine-Tuning (SFT)
- Training a Reward Model (RM)
- Reinforcement Learning from Human Feedback (RLHF)

Number of Parameters in Various Language Models



GPT 3 Applications

- Create fiction so close to human level
- Write Poem
- Chat similar to humans
- Write computer codes
- Design applications
- Summarize texts
- Answer questions
- Design application (Figma): Enter the description of an application. Then the Software can design an app for that description.
- Write code (): Enter the description of an application. Then the Software will generate react code for described app
- Ella a table: Make an empty table with some column names. Then the model fills the table by predicting what the entries of the table should be
- Semantic search: Go to a web page like Wikipedia and ask a question. Then the model can show you the paragraph where the answer of your question can be found there.

12. CLIP

1. Introduction

Motivated by NLP's success with task-agnostic web-scale pretraining (e.g., GPT-3).

Goal: Can we do the same for vision, by training on natural language supervision instead of human-labeled datasets?

CLIP achieves zero-shot performance (without training on the target task).

Caveat: Align images and texts in a joint embedding space using contrastive learning.

2. Approach

2.1 Natural Language Supervision

Uses image-text pairs from the internet.

Treats language as a rich source of supervision, more scalable than labels.

Leverages contrastive learning to associate correct (image, text) pairs and separate mismatched ones.

2.2 Dataset Construction

400 million image-text pairs scraped from the internet.

Different from prior datasets like ImageNet – unfiltered and noisy but massive in size.

No manual labeling.

2.3 Model Architecture

Dual-encoder architecture:

Image Encoder: ResNet or Vision Transformer (ViT)

Text Encoder: Transformer (e.g., 12-layer, 8-head, 63M params)

Each encoder maps its input to a 512-dimensional embedding.

Uses dot product similarity in the shared embedding space.

3. Experiments

3.1 Zero-Shot Transfer

Instead of evaluating representation quality (like in self-supervised learning), they measure task generalization.

Zero-shot: CLIP predicts by choosing the class text whose embedding is closest to the image.

3.2 Robustness

CLIP models are more robust to distribution shifts than supervised baselines. Shows better generalization to datasets outside training distribution.

Figure 2: Shows CLIP's 3x to 4x greater efficiency compared to Bag-of-Words and language-modeling baselines.

3.4 Results

Evaluated on 30+ datasets: OCR, action recognition, geo-localization, emotion recognition.

CLIP outperforms supervised models on many.

Example: Matches ResNet-50 on ImageNet without any ImageNet training.

Few-shot learning is possible via fitting linear classifiers on frozen CLIP features.

Figure 15: Shows few-shot performance drops compared to zero-shot due to lack of sample efficiency.

Table 2: Human vs CLIP accuracy on Oxford Pets. CLIP zero-shot beats human zero-shot.

3.5 Data Overlap Analysis

Checks if evaluation datasets leaked into the training set.

Built a duplicate detector using ResNet + custom augmentation.

3.6 Bias and Fairness (FairFace Analysis)

Evaluated CLIP on bias across race, gender, and age.

CLIP shows significant biases (especially in zero-shot mode).

Example: Underperforms on non-white FairFace categories compared to white.

Table 3 & 4: CLIP's accuracy on race/gender/age for white and non-white categories.

3.7 Broader Impacts

CLIP can classify arbitrary categories (e.g., shoplifters) which raises ethical concerns.

Promotes flexibility but also risk of misuse.

Large models like CLIP require societal accountability.

3.8 Related Work

Builds on contrastive learning, weakly supervised learning, image-text retrieval (e.g., ViTEx, ConVIRT, etc.)

Distinct from works that rely on complex attention-based architectures like ViLBERT or VisualBERT.

3.9 Conclusion

CLIP demonstrates the feasibility of language-supervised pretraining for vision.

Offers strong zero-shot generalization and task-agnostic transfer.

Promising direction, but has biases, poor few-shot performance, and data limitations.

4 PRACTICE QUESTIONS & ANSWERS

Q1: What is CLIP and what problem does it solve?

A: CLIP is a vision-language model that learns to associate images and text using contrastive learning on a web-scale dataset. It enables zero-shot learning for image classification tasks without needing task-specific training.

Q2: How does CLIP perform zero-shot classification?

A: By embedding a set of text prompts (like "A photo of a cat") and selecting the one whose embedding is closest (dot product similarity) to the image embedding.

Q3: What are the two types of encoders used in CLIP?

A: Image encoder (ResNet or ViT) and Text encoder (Transformer).

Q4: What loss function does CLIP use?

A: Symmetric InfoNCE loss (cross-entropy over pairwise cosine similarities).

5 Technical / Architecture Questions

Q5: What is the dimensionality of the embedding space?

A: 512-dimensional shared embedding space.

Q6: How is the temperature parameter τ used?

A: Scales the logits before softmax to control sharpness in the similarity distribution.

Q7: Why does few-shot performance drop compared to zero-shot in CLIP?

A: Because CLIP isn't optimized for few-shot settings, unlike humans who benefit significantly from even a single example.

Q8: What are the image encoder variants used in CLIP?

A: ResNet-50 family (e.g., RN50x4, RN50x16, RN50x64).

Vision Transformer (ViT-B/32, ViT-B/16, ViT-L/14).

ViTs are found to be 3x more compute efficient for CLIP's objectives

Q9: How is attention used in CLIP's ResNet variant?

A: Replaces global average pooling with a single-layer multi-head attention pooling. This improves performance by letting the network learn which regions to emphasize for the image representation.

Q10: Why is the temperature τ learned during training?

A: It controls the scale of logits before softmax in the contrastive loss. Learning it avoids manual tuning and helps balance similarity scores across batches.

6 Dataset and Training

Q11: How many image-text pairs are used in CLIP's training?

A: 400 million unfiltered internet (image, text) pairs.

Q12: What was the largest model trained?

A: ViT-L/14@336px and RN50x64, taking up to 592 GPUs and 18 days.

Q13: What is the loss function used in CLIP, and how does it work?

A: CLIP uses the InfoNCE loss, a symmetric contrastive loss. For a batch of image-text pairs:

It maximizes cosine similarity for matching (image, text) pairs. Minimizes it for all non-matching pairs in the batch (N-N negative pairs).

Implements it as a cross-entropy loss over similarity logits, scaled by a learnable temperature parameter τ .

Q14: What augmentations are applied to images during training?

A: Only a random square crop from resized images. The training is kept minimal and straightforward to avoid data-specific tuning.

Q15: How does CLIP differ from prior contrastive-image-text models like Visual N-Grams?

A: CLIP trains on a dataset 10x larger, uses 100x more compute, leverages transformers, and shows a 95% top-5 ImageNet accuracy, outperforming Visual N-Grams on multiple benchmarks.

7 BIAS, FAIRNESS & ETHICAL IMPLICATIONS

Q16: What benchmark was used to analyze bias in CLIP?

A: FairFace – diverse face dataset categorized by race, gender, and age

Q17: What are CLIP's bias-related issues?

Zero-shot CLIP underperforms on non-white races compared to white.

Class design bias: If the user defines biased class labels (e.g., "criminal"), CLIP may reinforce harmful stereotypes.

Representation harms: Model classifies people of color more often into "animal" or "crime" categories in label sets containing such terms

Q18: How was intersectional bias studied?

A: Accuracy measured across race x gender categories. CLIP performs above 95% across most gender categories but shows disparity in age and race

Q19: What other methods were used to detect duplicates or data contamination?

A: A custom near-duplicate detector using synthetic augmentations (zoom, rotation, crop, JPEG compression, etc.) trained using a modified ResNet-50 and InfoNCE loss

8 BROADER IMPACTS

Q20: What are potential misuses of CLIP?

A: Tasks like:

Surveillance (e.g., classifying "shoplifters")

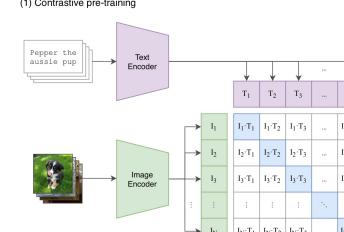
Biased classification based on race/gender

Unethical applications in facial recognition or social profiling

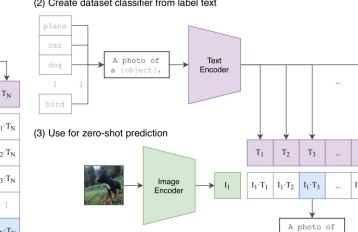
Q21: Why is CLIP's interface considered "flexible but dangerous"?

A: Users can easily define arbitrary categories using natural language. This makes CLIP powerful but vulnerable to misuse if it prompts encode bias or offensive semantics

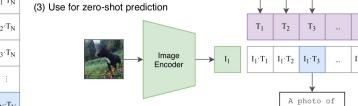
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



Q22: Why does CLIP fine-tune on ImageNet?

A: Accuracy increases by 9.2% on ImageNet, but robustness decreases on other datasets, suggesting overfitting to ImageNet distribution.

ishma haliez notes
repst heet

13. TVLT

Abstract & Introduction

TVLT is a transformer-based model that learns visual-linguistic representations from raw visual and audio inputs—without relying on text, tokenization, or ASR (Automatic Speech Recognition). Instead of converting speech to text, TVLT directly aligns video frames with audio spectrograms using two core training strategies: masked autoencoding and vision-audio matching. The model is efficient (1/3 the parameters, 28x faster) and achieves comparable or better results than text-based counterparts on tasks like VQA, sentiment analysis, and retrieval.

Related Work

The paper contrasts TVLT with:

- Text-based models like BERT, which use written language.
- Audio-based models that use spectrograms and modality-specific architectures.

VL (vision-language) models that rely on ASR for aligning text and vision. TVLT fills the gap as the first modality-agnostic, homogeneous transformer handling raw video + audio without text.

Architecture: TVLT Overview

Input Embeddings:

Video: 224x224 images \rightarrow 16x16 patches \rightarrow Linear projection \rightarrow 768-dim embeddings.

Audio: Log-mel spectrograms (128 freq bins) \rightarrow patches \rightarrow 768-dim embeddings.

Added: modality embeddings + temporal/spatial (for video) or temporal/frequency (for audio) embeddings

Transformer Encoder: 12-layer, modality-agnostic.

Transformer Decoder: 8-layer, used only for masked autoencoding reconstruction.

Training Objectives

Two losses:

Vision-Audio Matching (VAM): Predict if paired audio + video belong together. Uses [CLS] token with a sigmoid classifier. Binary cross-entropy loss

Masked Autoencoding (MAE): Randomly mask 75% of patches in both video and audio inputs. Decoder reconstructs original data. Uses MSE loss on only masked patches

ARCHITECTURE & ENCODING

Q1: What are the input types accepted by TVLT?

A1: TVLT accepts raw video frames and audio spectrograms, not text. These are processed into patch embeddings with added modality, spatial/temporal, and frequency encodings.

Q2: What is the size of each visual patch, and how are visual embeddings created?

A: Each 224x224 frame is divided into 16x16 patches. Each patch is linearly projected into a 768-dimensional embedding.

Q3: How are audio signals transformed for use in TVLT?

A: 1D waveforms are converted into 128-dimensional log mel-spectrograms (Tx128), divided into patches (16x16 or 2x128), and projected to 768-dimensional embeddings.

Q4: What is the purpose of modality embeddings in TVLT?

A: They are trainable vectors that inform the model whether the input patch is from vision or audio.

Q5: What is the architecture of the encoder and decoder in TVLT?

A: The encoder is a 12-layer transformer (hidden size 768). The decoder is an 8-layer transformer (hidden size 512), used only during pretraining for masked autoencoding.

PRETRAINING OBJECTIVES

Q6: Why does TVLT use a dual-objective setup for pretraining?

A: To learn strong unimodal and cross-modal representations: Vision-Audio Matching (VAM) captures cross-modal alignment, while Masked Autoencoding (MAE) strengthens unimodal reconstructions.

Q7: How are positive and negative examples generated for the VAM task?

A: Positive: matched video and audio from the same clip. Negative: video from one clip paired with audio from a different one.

Q8: What is the exact loss used for VAM?

A: Binary cross-entropy loss applied to the [CLS] token output of the encoder after passing through a sigmoid classifier.

Q9: How much of the input is masked during MAE?

A: 75% of both visual and audio patches are masked independently.

Q10: What is the final pretraining loss formula?

A: loss = $\lambda VAM \times loss_{VAM} + \lambda MAE \times loss_{MAE}$ (where $\lambda VAM = 1.0$ and $\lambda MAE = 0.3$)

Efficiency and Design

ASR modules in traditional pipelines dominate inference time (e.g., 2890ms just for ASR).

TVLT runs in 103ms for same input due to bypassing ASR. 1/3 the parameters (88M vs. 283M)

Results on Benchmark Tasks

TVLT outperforms text-based counterparts in audio-to-video retrieval and sentiment analysis (e.g., on CMU-MOSEI).

Slightly underperforms on audio-to-image retrieval and VQA but still competitive

Handles emotion categories better than text-based models except for "disgust"

Ablation Studies

Joint encoder > separate encoders (better fusion).

Separate decoders > joint decoders (better reconstruction).

VAM + MAE > each alone

Diagram Explanations

Figure 1 – Efficiency Diagram

Shows traditional VL models (with ASR) vs. TVLT:

ASR alone takes ~2.6s.

TVLT Pipeline takes ~100ms.

TVLT is 27-28x faster for inference.

Figure 2 – TVLT Architecture

(a) Vision-Audio Matching: Both inputs go through the encoder. A [CLS] token predicts if they are matched (label 0 or 1).

(b) Masked Autoencoding: Masked input \rightarrow encoder \rightarrow reconstructed output. Applies to both spectrogram and video.

Figure 3 & 4 – Reconstruction Visualization

Left: masked input. Middle: reconstruction. Right: original.

Shows TVLT can successfully reconstruct masked video and audio using its learned embeddings

Exam Practice Questions & Answers

Q1: What is the main novelty of the TVLT model compared to traditional VL models?

A: TVLT removes reliance on text and ASR by using a modality-agnostic transformer to learn directly from video frames and audio spectrograms.

Q2: What are the two main training objectives used in TVLT?

A: Vision-Audio Matching (VAM) and Masked Autoencoding (MAE).

Q3: How does Vision-Audio Matching work?

A: It creates positive (matched) and negative (random) video-audio pairs and uses binary cross-entropy loss to train a [CLS] token classifier.

Q4: How is Masked Autoencoding applied in TVLT?

A: Random patches from both audio and video are masked; the decoder reconstructs the original input using MSE loss only on masked patches.

Q5: What does the efficiency comparison show about TVLT?

A: TVLT is over 27x faster in inference than text-based VL models due to removing ASR, and it uses 1/3 the parameters.

Q6: Why is TVLT considered more "green" than other models?

A: It reduces computation by eliminating ASR modules and uses a simpler, unified architecture, although large-scale pretraining is still required.

Q7: What happens when you use joint vs. separate encoders/decoders?

A: Joint encoders perform better than separate ones, while separate decoders perform better than joint ones in masked autoencoding.

Q8: How does TVLT perform on CMU-MOSEI emotion analysis?

A: TVLT outperforms text-based models in recognizing emotions like happy, sad, angry, and fear, thanks to its ability to learn tone and loudness from raw audio.

Q9: What kind of input preprocessing is used for audio in TVLT?

A: Raw audio is converted to 128-bin log-mel spectrograms; patches are created and embedded with temporal/frequency encodings.

Q10: Why does masking speech spans help in audio MAE?

A: Because speech spans contain semantic information; masking them improves the model's ability to learn useful acoustic features.

ABLATION STUDIES

Q11: What did the decoder architecture ablation show?

A: Using separate decoders (one for audio, one for vision) outperformed using a joint decoder.

Q12: What did the encoder ablation study reveal?

A: A joint encoder (shared for vision and audio) outperformed separate encoders followed by a fusion layer.

Q13: How did speech-span masking affect performance?

A: Masking only the speech-active regions in the spectrogram improved performance on both MSR-VTT and VQAv2.

Q14: Between 16x16 and 2x128 patch sizes for audio, which is better?

A: Mixed results. 2x128 was better on MSR-VTT; 16x16 was better on VQAv2. Default chosen is 16x16 for modality alignment.

PERFORMANCE & RESULTS

Q15: Which benchmark tasks were used for evaluation?

A: Audio-to-video retrieval (MSR-VTT, YouCook2, CrossTask), multimodal sentiment (CMU-MOSEI), audio-to-image (Places-400k), and visual QA (VQA1/VQA2).

Q16: How does TVLT compare to its text-based version on video tasks?

A: TVLT consistently outperforms the text-based version on video tasks like retrieval and sentiment/emotion analysis.

Q17: Why is TVLT more robust in emotion classification?

A: It directly captures acoustic features like tone and loudness from raw audio, unlike text-based models limited to transcript content.

Q18: In which areas does the text-based model slightly outperform TVLT?

A: On tasks with clean, written text like VQA and audio-to-image retrieval, text-based models have a slight edge.

Q19: What efficiency gains does TVLT offer over traditional ASR pipelines?

A: TVLT is up to 28x faster in inference and uses only 88M parameters vs. 283M+ in ASR-based pipelines.

Q20: What task head is used during finetuning on retrieval tasks?

A: A two-layer MLP maps the encoder's [CLS] token to a score $\in [0, 1]$ used for match classification.

DATASETS & TRAINING

Q21: Which datasets were used for pretraining?

A: HowTo100M and a 20% subset of YTtemporal180M (called YTT-S).

Q22: How long does pretraining TVLT take?

A: 200k steps over 2 weeks using 4 NVIDIA RTX A6000 GPUs (49GB each).

Q23: How is ASR used in the text-based version?

A: ASR (like SpeechBrain) transcribes speech audio to text, which is then tokenized and embedded.

Q24: How is TTS used in the paper?

A: Text questions from VQA1/VQA2 are converted to synthetic speech via WaveNet and used as audio queries.

Q25: What are the hyperparameters used during pretraining?

A: LR = 1e-5, batch size = 4096, weight decay = 0.001, cosine learning rate schedule.

BROADER INSIGHTS

Q26: What key insight about vision-language learning does TVLT challenge?

A: That high-quality vision-language representations require written language. TVLT shows this can be achieved with raw audio instead.

Q27: Why is TVLT considered more compact?

A: Because it avoids heavy modules like ASR and text tokenizers, reducing overall architecture complexity and compute.

Q28: What advantage does TVLT have in real-time systems like smart assistants?

A: Faster inference and direct audio input makes it ideal for deployment in systems with limited latency budgets.

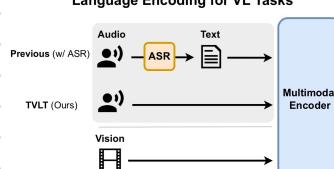
Q29: How does TVLT handle image-only tasks like VQA?

A: By converting the text questions into synthetic speech via TTS, enabling consistent input modality (audio + vision).

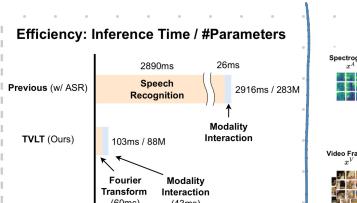
Q30: What is a practical use case where TVLT would outperform text-based models?

A: Emotion recognition from speech and video in real-time, such as in therapeutic AI or expressive virtual agents.

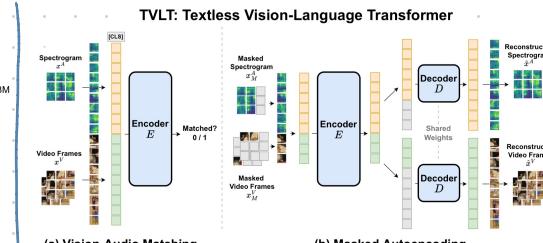
Language Encoding for VL Tasks



Efficiency: Inference Time / #Parameters



TVLT: Textless Vision-Language Transformer



Q2: Attention Masking

You are training GPT on this sequence:

"The dog barked loudly."

Which tokens can the word "barked" attend to during training?

We are predicting the word "barked".

It can only attend to:

"The"

"dog"

It cannot attend to:

"loudly" (it comes after)

Q3: Token Flow

Draw or describe the token flow in a BERT-based classification task using this input:

[CLS] The movie was great [SEP]

What is the role of the [CLS] token?

What happens to the token embeddings after self-attention?

Tokens: [CLS], The, movie, was, great, [SEP]

All tokens go into self-attention.

Each token attends to all other tokens (because it's bidirectional).

The final representation of the [CLS] token becomes the sentence embedding.

For classification tasks (like sentiment analysis), the [CLS] token is passed to a classifier head.

Feature	BERT	GPT	T5
Architecture	Encoder-only	Decoder-only	Encoder-Decoder
Direction	Bidirectional	Left-to-right	Bidirectional (enc), causal (dec)
Pretraining Task	MLM + NSP	Next token prediction	Span corruption
Input Format	Tokenized text + [CLS]	Plain text	Text-to-text
Token Masking?	Yes	Yes (causal)	Yes (with sentinel tokens)
Can generate text?	No (not directly)	Yes	Yes

what's projection?

Suppose you have an input word embedding:

$x \in \mathbb{R}^{512}$ (i.e., 512-dimensional vector)

You apply a projection matrix $W_Q \in \mathbb{R}^{64 \times 512}$

$q = W_Q x$ Then $q \in \mathbb{R}^{64}$ is your projected query vector projection from 512-D into a 64-D

All Model Compositions

Model	Architecture	Training Objective	Strengths	Limitations	Typical Use Cases
Diffusion Models	Markov chain noise addition + iterative denoising	Denoising score matching	Stable training, diverse & high-fidelity samples	Slow sampling, computationally heavy	High-quality image generation, inpainting
GANs	Generator + Discriminator (adversarial)	Adversarial loss	Fast generation, sharp images	Training instability, mode collapse	Real-time image/video synthesis, deepfakes
VAEs	Encoder-decoder with probabilistic latent space	Variational lower bound	Stable, diverse outputs, latent space	Blurry outputs, posterior collapse	Representation learning, anomaly detection
AEs	Encoder-decoder (deterministic)	Reconstruction loss	Simple, fast, unsupervised learning	No explicit generative modeling	Dimensionality reduction, denoising
BERT	Transformer encoder (bidirectional)	Masked Language Modeling + NSP	Deep understanding, strong NLP performance	Not generative, heavy compute	NLP understanding, classification
GPT	Transformer decoder (unidirectional)	Next-token prediction	Strong text generation	Limited context, large compute	Text generation, chatbots, summarization

Formula Sheet

Convolutional Layers (CNN)

• Conv2D Output Size:

$$\text{Output size} = \left\lfloor \frac{\text{Input size} - 2 \cdot \text{padding} - (\text{kernel size} - 1)}{\text{stride}} \right\rfloor + 1$$

• Pooling Output:

$$\text{Output size} = \left\lfloor \frac{\text{Input size} - \text{kernel size}}{\text{stride}} \right\rfloor + 1$$

• Number of Weights (Conv Layer):

$$\text{Weights} = \text{Out channels} \times \text{In channels} \times \text{Kernel height} \times \text{Kernel width}$$

• Number of Biases (Conv Layer):

$$\text{Biases} = \text{Out channels}$$

Positional Encoding (Transformers)

$$\text{PE}(\text{pos}, i) = \begin{cases} \sin\left(\frac{\text{pos} \cdot \text{pos}_i}{10000^{(i-1)/d}}\right) & \text{if } i \text{ is even} \\ \cos\left(\frac{\text{pos} \cdot \text{pos}_i}{10000^{(i-1)/d}}\right) & \text{if } i \text{ is odd} \end{cases}$$

Scaled Dot Product Attention

• Attention Score:

$$\text{score} = \frac{\vec{q} \vec{v}^T}{\sqrt{d_k}}$$

• Softmax:

$$\text{Softmax}(\vec{x}_i) = \frac{e^{\vec{x}_i}}{\sum_j e^{\vec{x}_j}}$$

• Weighted Sum (Attention output):

$$\text{Output} = \sum_i \text{softmax}(\text{score}_i) \cdot v_i$$

Diffusion/VAEs

KL Divergence (VAE Loss)

$$D_{KL}(q(z|x)||p(z)) = \frac{1}{2} \sum_i (\mu^2 + \sigma^2 - \log(\sigma^2) - 1)$$

$$D_{KL}(P||Q) = \sum_i P(i) \log \left(\frac{P(i)}{Q(i)} \right)$$

Note: Use ln (natural log)

Harris Corner Detection

Used for detecting corners in an image.

Step 1: Compute image gradients

$$I_x = \frac{\partial I}{\partial x}, \quad I_y = \frac{\partial I}{\partial y}$$

Step 2: Compute structure matrix MMM

$$M = \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \quad (\text{smoothed over a window})$$

Step 3: Harris Response

$$R = \det(M) - k \cdot (\text{trace}(M))^2$$

$$\bullet \det(M) = I_x^2 I_y^2 - (I_x I_y)^2$$

$$\bullet \text{trace}(M) = I_x^2 + I_y^2$$

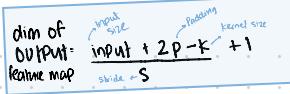
$$\bullet k \in [0.04, 0.06]$$

If R

• Large positive \rightarrow corner

• Large negative \rightarrow edge

• Close to zero \rightarrow flat region



grey scale image - 1 channel

RGB image - 3 channels

what's projection?

Suppose you have an input word embedding:

$x \in \mathbb{R}^{512}$ (i.e., 512-dimensional vector)

You apply a projection matrix $W_Q \in \mathbb{R}^{64 \times 512}$

$q = W_Q x$ Then $q \in \mathbb{R}^{64}$ is your projected query vector projection from 512-D into a 64-D

Positional Encoding

position $z_i = x_i + \text{Positional encoding}_i$

Step 1: Get token embeddings

Imagine the embedding layer gives us fixed vectors for tokens:

Token 1: $[0, 1, 0, 2, 3, 0, 4, 0, 5, 0, 6]$

love: $[0, 5, 0, 4, 0, 3, 2, 0, 1, 0, 0]$

AI: $[0, 9, 0, 8, 0, 7, 0, 6, 0, 5, 0, 4]$

Step 2: Calculate positional encoding vectors

$$\text{formula: } \text{PE}(\text{pos}, i) = \sin\left(\frac{\text{pos} \cdot \text{dim}}{10000^{(i-1)/6}}\right) \text{PE}(\text{pos}, 2i+1) = \cos\left(\frac{\text{pos} \cdot \text{dim}}{10000^{(i-1)/6}}\right) \text{PE}(\text{pos}, 2i)$$

For $i=6$, dimensions $2i = 0, 1, 2$:

Calculate denominator $10000^0 = 1$ for each i :

$$1: 2i/d = \text{Denominator} = 10000^{0/6} = 1$$

$$0: 10000^1 = 1$$

$$1: 2/6 = 0.333 \quad 10000^{0.333} \approx 21.54$$

$$2: 4/6 = 0.666 \quad 10000^{0.666} \approx 46.16$$

Calculate each position:

$$\text{For pos = 0: } \sin(0/1) = 0, \quad \cos(0/1) = 1$$

$$\sin(1/21.54) \approx 0.0464, \cos(1/21.54) \approx 0.9988$$

$$\sin(1/46.16) \approx 0.0225, \cos(1/46.16) \approx 0.999977$$

$$\text{PE}_1 = [0.8415, 0.5403, 0.0464, 0.9989, 0.0215, 0.999997]$$

For pos = 1:

$$\sin(2/21.54) \approx 0.0927, \cos(2/21.54) \approx 0.9957$$

$$\sin(2/46.16) \approx 0.00431, \cos(2/46.16) \approx 0.999997$$

$$\text{PE}_2 = [0.9983, -0.4161, 0.0927, 0.6957, 0.00431, 0.999997]$$

Step 3: Add positional encoding to token embeddings

Token: Read input to transform

$$1: [0, 1, 0, 2, 0, 3, 0, 4, 0, 5, 0]$$

$$0: 10-10 = [0, 1, 2, 3, 4, 5, 6, 1, 0, 5, 1, 6]$$

$$\text{love: } [1.3415, 0.9403, 0.3464, 1.1989, 0.10215, 0.999997]$$

$$\text{AI: } [1.8093, 0.3839, 0.7927, 1.5957, 0.50431, 1.3999997]$$

Additional Attention/Rishanau Attention:

You are doing machine translation using seq2seq + attention. We are at decoder timestep 1.

We have 3 encoder hidden states: h_1, h_2, h_3

We have 1 decoder hidden state from previous timestep: s_0

$$h_1 = [1, 0], \quad h_2 = [0, 1], \quad h_3 = [1, 1], \quad s_0 = [1, 2]$$

We'll use a score function of the form:

$$s_{ij} = v^T \tanh(W_s s_i + W_h h_j)$$

$$W_s = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad W_h = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad v = [1, 1]^T$$

For $h_1 = [1, 0]$

$$W_s h_1 = [1 \times 1 + 1 \times 0, 1 \times 1 + (-1) \times 0] = [1, 1]$$

$$W_s s_0 + W_h h_1 = [1 + 1, 2 + (-1)] = [2, 1]$$

$$\tanh([2, 1]) \approx [0.964, 0.995]$$

$$e_{11} = v^T \cdot [0.964, 0.995] = 0.964 + 0.995 = ** 1.959 **$$

For $h_2 = [0, 1]$

$$W_s h_2 = [0 \times 1 + 1 \times 1, 0 \times 1 + (-1) \times 1] = [-1, 1]$$

$$W_s s_0 + W_h h_2 = [1 + 1, 2 + (-1)] = [2, 1]$$

$$\tanh([2, 1]) \approx [0.964, 0.994]$$

$$e_{12} = 0.964 + 0.994 = ** 0.964 **$$

For $h_3 = [1, 1]$

$$W_s h_3 = [1 \times 1 + 1 \times 1, 1 \times 1 + (-1) \times 1] = [2, 0]$$

$$W_s s_0 + W_h h_3 = [1 + 2, 2 + (-1)] = [3, 2]$$

$$\tanh([3, 2]) \approx [0.995, 0.964]$$

$$e_{13} = 0.995 + 0.964 = ** 1.959 **$$

Apply softmax to get attention weights a_{ij}

$$\text{softmax}(e_{11}, e_{12}, e_{13}) = \frac{e^{e_{11}}}{\sum e^{e_{ij}}}$$

$$e^{1.959} \approx 7.09, \quad e^{1.725} \approx 5.61$$

$$\text{Denominator} = 7.09 + 5.61 + 7.09 = 19.79$$

$$a_{11} = \frac{7.09}{19.79} \approx 0.358 \quad a_{12} = \frac{5.61}{19.79} \approx 0.284 \quad a_{13} = \frac{7.09}{19.79} \approx 0.358$$

Compute context vector c_i

$$h_1 = [1, 0], \quad h_2 = [0, 1], \quad h_3 = [1, 1]$$

$$c_i = [0.358, 1, 0] + [0.283, 0, 1] + [0.358, 1, 1] = [0.358, 0] + [0, 283] + [0.358, 0.358] = [0.716, 0.641]$$

Attention weights $[0.358, 0.283, 0.358]$

$$\text{Context vector } c_i = [0.716, 0.641]$$

$$v_1 \quad v_2 \quad v_3 \quad v_4 \quad [2] \quad \text{Ans}$$

Transformers

$$8) \quad x = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w_i = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w_k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, w_v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

1. Compute θ, k, v

$$\theta = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^T = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

2. Compute Attention Scores

$$\theta = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad k^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$k^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

scaled by $\sqrt{2}$ (1/1)

$$\frac{\theta k^T}{\sqrt{2}} = \begin{bmatrix} 0.71 & 0.71 \\ 0.71 & 0.71 \end{bmatrix}$$

3. Apply softmax

$$\text{softmax}(0.71, 0.71) = (0.5, 0.5)$$

$$\text{softmax}\left(\frac{0.71}{\sqrt{2}}\right) = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

4. Multiply V

2. Softmax.V

$$\begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} = \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix}$$

5. Residual Connection

= $x + z$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0.25 & 0.25 \\ 0.25 & 0.25 \end{bmatrix}$$

$$= \begin{bmatrix} 1.25 & 0.25 \\ 0.25 & 1.25 \end{bmatrix}$$

Hello [1,2] world [3,M]

$$8) \quad x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, w_i = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w_k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, w_v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

1. Compute θ, k, v

$$\theta = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$k = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$v = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^T = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

2. Compute Attention Scores

$$\theta = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad k^T = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$k^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\frac{\theta k^T}{\sqrt{2}} = \begin{bmatrix} 1.636 & 14.887 \\ 14.887 & 34.649 \end{bmatrix}$$

3. Apply softmax

$$\text{softmax}\left(\frac{\theta k^T}{\sqrt{2}}\right) = \begin{bmatrix} 0.00021 & 0.99979 \\ 0 & 1 \end{bmatrix}$$

4. Multiply V

$$2. \text{ Softmax.V} = \begin{bmatrix} 0.5 & 1 \\ 1.5 & 2 \end{bmatrix}$$

$$2. \begin{bmatrix} 0.00021 & 0.99979 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0.5 & 1 \\ 1.5 & 2 \end{bmatrix} = \begin{bmatrix} 1.49979 & 1.99979 \\ 1.5 & 2 \end{bmatrix}$$

5. Residual Connection

= $x + z$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + \begin{bmatrix} 1.49979 & 1.99979 \\ 1.5 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 2.5 & 4 \\ 4.5 & 6 \end{bmatrix}$$

18/07/2023
Audio is aligned with the video

Do frame by frame but it has less accuracy than video based.

3D-CNN for videos but has more parameters.
2D web-based has limited parameters.

Department of Education

- It's very slow

(b) The text is not bringing very big impact into the picture, remove it.

title

By applying Fourier transformation you can map convert 1D into 2D signal

As Image encoder is 2D & audio encoder is 1D and we end up with 2 encoders slowing us down, using Fourier transform

we end up with only 1 encoder speeds is really good.

It is called spectrogram, when you end up converting 1D signal into 2D signal

ASR (case study)
Automatic Speech recognition 02/05/2023

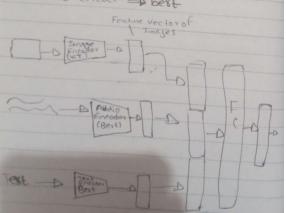
ASO

Kids with read autism

Web-based Emotion recognition
It takes video as input

End goal off

Text encoder -> best
Image Encoder -> VIT
Audio encoder -> best



Paper is called CLIP
Contrastive Language Pre-training

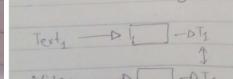
Learning Transferable Visual Models from Natural Language Supervision

Used to extract features from Text.

From features we would use Vision encoder's.

CLIP (Research paper reading assistant)

Contrastive Learning Image Pre-training



If those 2 belong to the same sample we put the class in the feature space.

ishma halfeez notes

represent

CLIP
Contrastive learning, very important used in computer vision
Masked auto encoding