

Deep Learning: advanced models

Generative models, Transformers, GNNs...

Nicolas Keriven
CNRS, IRISA, Rennes

ENSTA 2025

Pros and cons of deep learning

Deep Neural Nets are :

- ▶ **Difficult to train** : they require tips and tricks (Batch norm, dropout...)

Pros and cons of deep learning

Deep Neural Nets are :

- ▶ **Difficult to train** : they require tips and tricks (Batch norm, dropout...)
- ▶ **Power-hungry** : computationally expensive to train

Pros and cons of deep learning

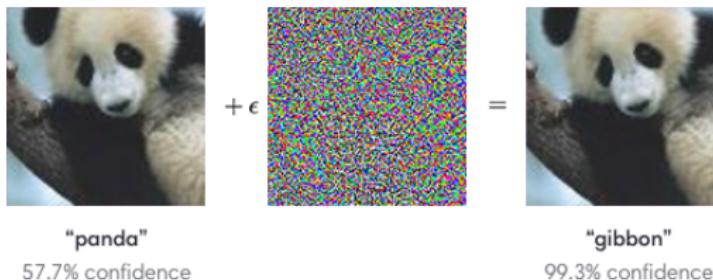
Deep Neural Nets are :

- ▶ **Difficult to train** : they require tips and tricks (Batch norm, dropout...)
- ▶ **Power-hungry** : computationally expensive to train
- ▶ **Difficult to interpret** : hard to get confidence intervals, generalization guarantees...

Pros and cons of deep learning

Deep Neural Nets are :

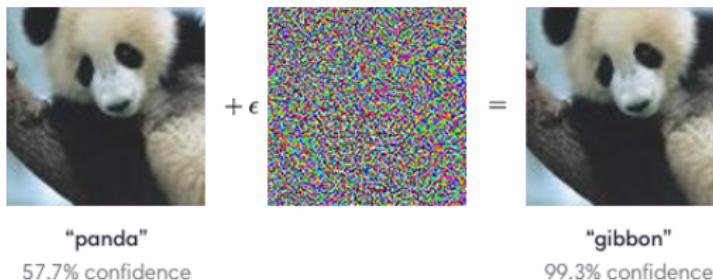
- ▶ **Difficult to train** : they require tips and tricks (Batch norm, dropout...)
- ▶ **Power-hungry** : computationally expensive to train
- ▶ **Difficult to interpret** : hard to get confidence intervals, generalization guarantees...
- ▶ **Unstable** : well-engineered micro-change in the data can fool them. Such change almost never happens in practice, but can be engineered by a malicious party !



Pros and cons of deep learning

Deep Neural Nets are :

- ▶ **Difficult to train** : they require tips and tricks (Batch norm, dropout...)
- ▶ **Power-hungry** : computationally expensive to train
- ▶ **Difficult to interpret** : hard to get confidence intervals, generalization guarantees...
- ▶ **Unstable** : well-engineered micro-change in the data can fool them. Such change almost never happens in practice, but can be engineered by a malicious party !



- ▶ A **bazooka**, when sometimes all you need to do is kill a fly ! In many real-life scenarios, linear models, PCA, k-means... work amazingly well

Pros and cons of deep learning

But ! Deep Neural Nets :

- ▶ are easy to implement and optimize with modern libraries (Pytorch, JAX...),
with **automatic differentiation and GPU compatibility completely transparent**
but there are many tricks to train them ! Can be frustrating...

Pros and cons of deep learning

But ! Deep Neural Nets :

- ▶ are easy to implement and optimize with modern libraries (Pytorch, JAX...),
with **automatic differentiation and GPU compatibility completely transparent**
but there are many tricks to train them ! Can be frustrating...

- ▶ work **mysteriously well** for problems we thought out-of-reach for decades !

Pros and cons of deep learning

But ! Deep Neural Nets :

- ▶ are easy to implement and optimize with modern libraries (Pytorch, JAX...),
with **automatic differentiation and GPU compatibility completely transparent**
but there are many tricks to train them ! Can be frustrating...
- ▶ work **mysteriously well** for problems we thought out-of-reach for decades !
- ▶ are **amazingly flexible** ! You can plug “anything anywhere”, and it will almost
always “work”. Huge playground ! This can be dangerous ! Always be aware of
what you do, or bugs will become impossible to fix.

Pros and cons of deep learning

But ! Deep Neural Nets :

- ▶ are easy to implement and optimize with modern libraries (Pytorch, JAX...), with **automatic differentiation and GPU compatibility completely transparent** but there are many tricks to train them ! Can be frustrating...
- ▶ work **mysteriously well** for problems we thought out-of-reach for decades !
- ▶ are **amazingly flexible** ! You can plug “anything anywhere”, and it will almost always “work”. Huge playground ! This can be dangerous ! Always be aware of what you do, or bugs will become impossible to fix.
- ▶ Today, we are going to see a few variants of DNNs, without going into details. Many resources are available online.

Table of Contents

Generative Models

- Generative Adversarial Networks
- Variational Auto-encoder
- Diffusion models

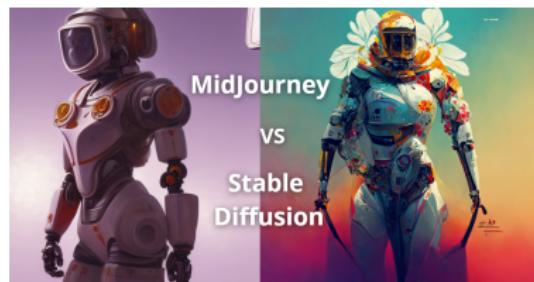
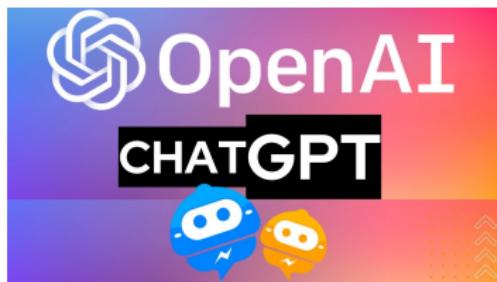
Attention, Transformers

Conclusions and open questions

Some (subjective) non-technical considerations

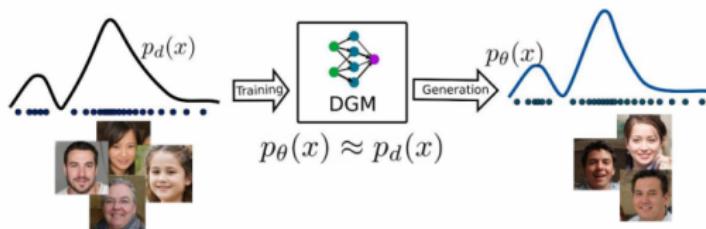
Generative model

- ▶ Goal : generate data resembling other data.
- ▶ Have *de facto* become synonymous with “AI” for the mainstream audience



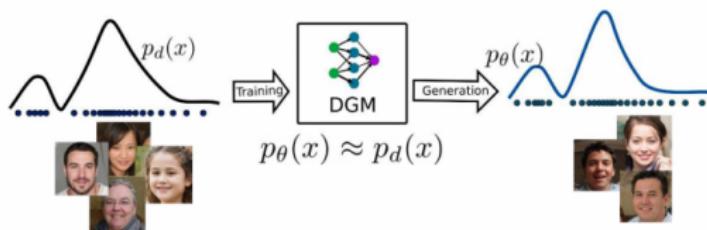
Basic generative model : density estimation !

- ▶ Most basic problem : observe samples from a distribution, generate **new** samples not observed before
- ▶ Just a *density estimation* problem ! But the probability distribution is **far too** complicated to be modelled explicitly (eg images)



Basic generative model : density estimation !

- ▶ Most basic problem : observe samples from a distribution, generate **new** samples not observed before
- ▶ Just a *density estimation* problem ! But the probability distribution is **far too** complicated to be modelled explicitly (eg images)



- ▶ Basic idea : take a very simple distribution in a **latent space**, eg $z \sim \mathcal{N}(0, Id)$, train a **complicated** function f_θ such that $f_\theta(z)$ are good samples.

▶ This is an instantiation of a “transport” problem

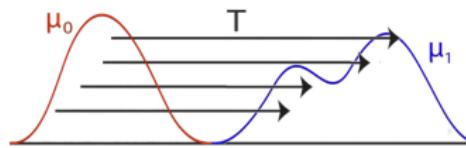


Table of Contents

Generative Models

 Generative Adversarial Networks

 Variational Auto-encoder

 Diffusion models

Attention, Transformers

Conclusions and open questions

Some (subjective) non-technical considerations

Generative Adversarial Networks (GANs)

- ▶ Q : how to train a Neural Network f_θ such that $f_\theta(z)$ produces "realistic images" different from the training set ? Images are already hard to understand !

Generative Adversarial Networks (GANs)

- ▶ Q : how to train a Neural Network f_θ such that $f_\theta(z)$ produces "realistic images" different from the training set ? Images are already hard to understand !
- ▶ Breakthrough idea : only Neural Networks are capable of automatically understanding images... well, let's train **two CNNs** ! One to produce the images, one to criticize them, and make them compete against each other ! **Generative Adversarial Networks** (Goodfellow et al. 2014)

Generative Adversarial Networks (GANs)

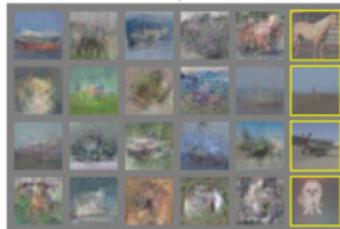
- ▶ Q : how to train a Neural Network f_θ such that $f_\theta(z)$ produces "realistic images" different from the training set ? Images are already hard to understand !
- ▶ Breakthrough idea : only Neural Networks are capable of automatically understanding images... well, let's train **two CNNs**! One to produce the images, one to criticize them, and make them compete against each other! **Generative Adversarial Networks** (Goodfellow et al. 2014)



a)



b)



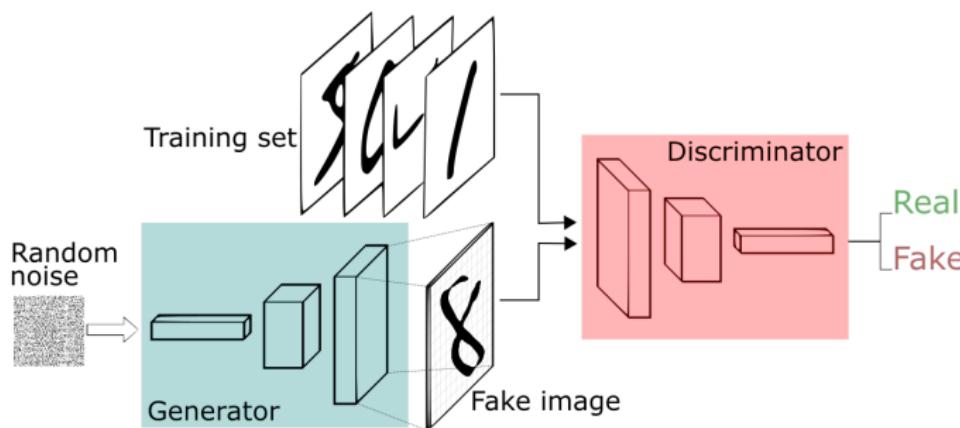
c)



d)

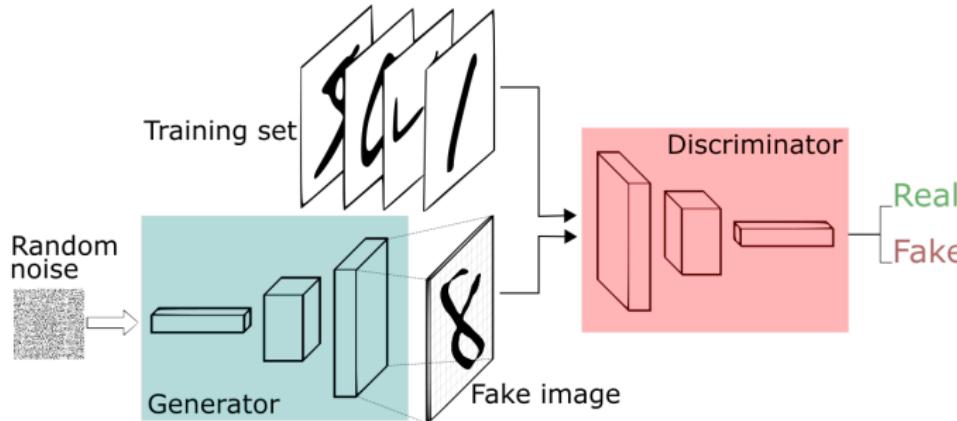
GAN description

- ▶ Train a **generator** network G_{θ_g} such that $G_{\theta_g}(z)$, with z Gaussian, are good samples
- ▶ Train a **discriminator** network D_{θ_d} to classify the samples as being from the training set or from the generator



GAN description

- ▶ Train a **generator** network G_{θ_g} such that $G_{\theta_g}(z)$, with z Gaussian, are good samples
- ▶ Train a **discriminator** network D_{θ_d} to classify the samples as being from the training set or from the generator
- ▶ Make them “compete” : the discriminator will become better at **classifying the fake samples**, the generator will become better at **fooling the discriminator**



GAN training

- ▶ The original GAN paper presents the problem as a **two-players minimax game**, based on the logistic regression loss :

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{x \sim p_{data}} [\log D_{\theta_d}(x)] + \mathbb{E}_{z \sim \mathcal{N}} [\log(1 - D_{\theta_d}(G_{\theta_g}(z)))]$$

GAN training

- ▶ The original GAN paper presents the problem as a **two-players minimax game**, based on the logistic regression loss :

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{x \sim p_{data}} [\log D_{\theta_d}(x)] + \mathbb{E}_{z \sim \mathcal{N}} [\log(1 - D_{\theta_d}(G_{\theta_g}(z)))]$$

- ▶ $D(x)$ represent the probability that x is a true sample. D wants to **maximize the log-likelihood**, ie, outputting 1 for x_i and 0 for $G(z)$.

GAN training

- ▶ The original GAN paper presents the problem as a **two-players minimax game**, based on the logistic regression loss :

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{x \sim p_{data}} [\log D_{\theta_d}(x)] + \mathbb{E}_{z \sim \mathcal{N}} [\log(1 - D_{\theta_d}(G_{\theta_g}(z)))]$$

- ▶ $D(x)$ represent the probability that x is a true sample. D wants to **maximize the log-likelihood**, ie, outputting 1 for x_i and 0 for $G(z)$.
- ▶ G want to **minimize the log-likelihood**, ie, make $G(z)$ ressemble the distribution of the data

GAN training

- ▶ The original GAN paper presents the problem as a **two-players minimax game**, based on the logistic regression loss :

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{x \sim p_{data}} [\log D_{\theta_d}(x)] + \mathbb{E}_{z \sim \mathcal{N}} [\log(1 - D_{\theta_d}(G_{\theta_g}(z)))]$$

- ▶ $D(x)$ represent the probability that x is a true sample. D wants to **maximize the log-likelihood**, ie, outputting 1 for x_i and 0 for $G(z)$.
- ▶ G want to **minimize the log-likelihood**, ie, make $G(z)$ ressemble the distribution of the data

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations do

 for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{data}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log(1 - D(G(z^{(i)}))) \right].$$

 end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Original algo : alternate k steps of gradient ascent for D with one step of gradient descent for G

GAN issues

- ▶ Hard to train : minimax problems are typically hard !

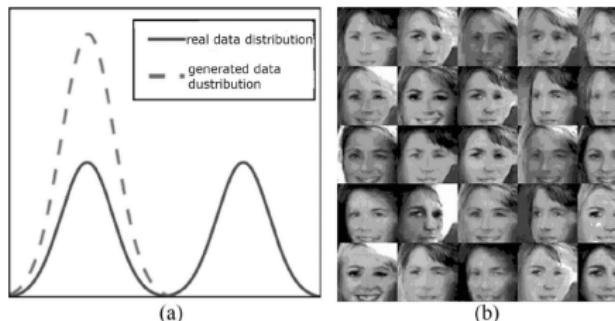
- ▶ they can loop
- ▶ saddle points are unstable
- ▶ convergence is (very) slow, etc.

GAN issues

- ▶ Hard to train : minimax problems are typically hard !
 - ▶ they can loop
 - ▶ saddle points are unstable
 - ▶ convergence is (very) slow, etc.
- ▶ Data-hungry : GANs require a **lot** of data to be trained properly

GAN issues

- ▶ Hard to train : minimax problems are typically hard !
 - ▶ they can loop
 - ▶ saddle points are unstable
 - ▶ convergence is (very) slow, etc.
- ▶ Data-hungry : GANs require a **lot** of data to be trained properly
- ▶ Mode collapse : the generator produces limited diversity (eg, for a dataset with cats and dogs, produces only dogs)



Sota GANs

There are many (many !) tricks and variants to improve GANs. Mostly, bigger models with more data seem to be always better! (BigGAN below).



Sota GANs

There are many (many !) tricks and variants to improve GANs. Mostly, bigger models with more data seem to be always better! (BigGAN below).



But nowadays, they have been mostly replaced with more modern approaches like diffusion models.

Table of Contents

Generative Models

 Generative Adversarial Networks

 Variational Auto-encoder

 Diffusion models

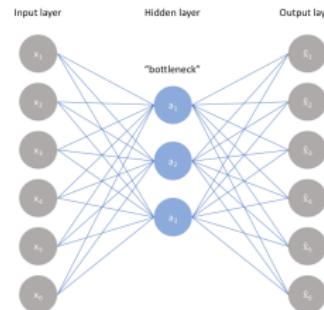
Attention, Transformers

Conclusions and open questions

Some (subjective) non-technical considerations

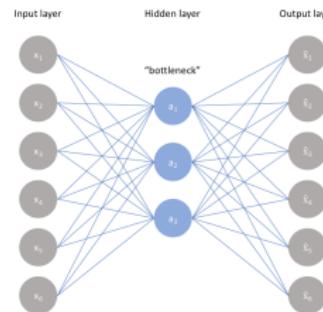
Auto-Encoder

- ▶ We have seen Auto-Encoder :
 - ▶ Encoder : images → compressed **latent space**
 - ▶ Decoder : latent space → image



Auto-Encoder

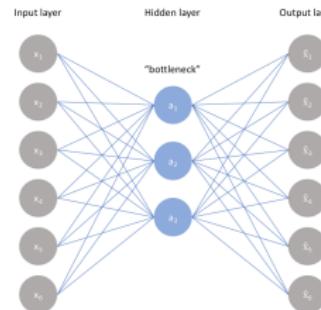
- ▶ We have seen Auto-Encoder :
 - ▶ Encoder : images → compressed **latent space**
 - ▶ Decoder : latent space → image



- ▶ Idea : by working directly in the latent space, the decoder is a generative model !

Auto-Encoder

- ▶ We have seen Auto-Encoder :
 - ▶ Encoder : images → compressed **latent space**
 - ▶ Decoder : latent space → image



- ▶ Idea : by working directly in the latent space, the decoder is a generative model !
- ▶ **Difficulty** : sampling a random vector in the latent space might produce garbage. The latent space is **highly irregular** !

Variation Auto-Encoder (VAE) Kingma and Welling, 2014

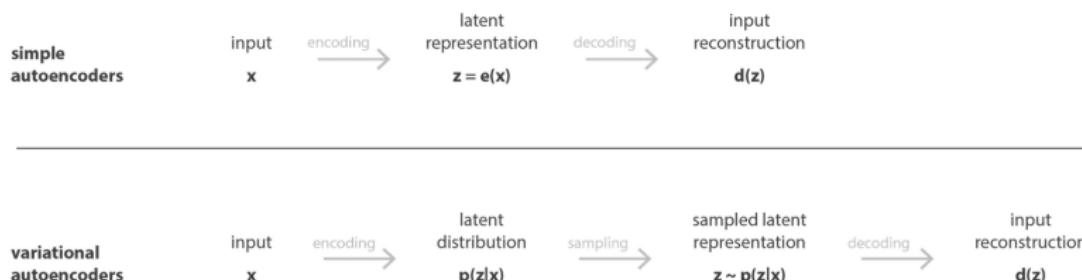
<https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73>

- ▶ Idea : **regularize** the distribution of the encodings, such that the latent space is smoother, and new sample can be produced

Variation Auto-Encoder (VAE) Kingma and Welling, 2014

<https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73>

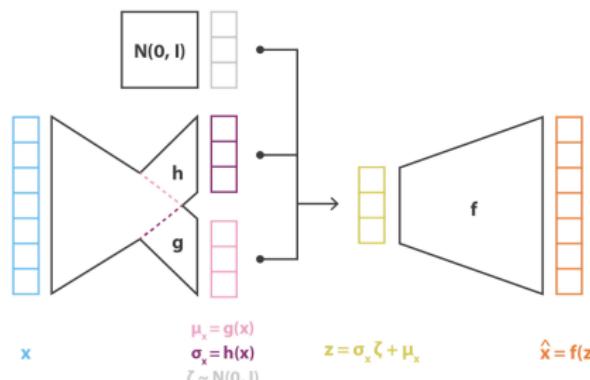
- ▶ Idea : **regularize** the distribution of the encodings, such that the latent space is smoother, and new sample can be produced
- ▶ To do so :
 - ▶ encode inputs as distributions, and feed sample from these distributions to the decoder.



Variational Auto-Encoder (VAE) Kingma and Welling, 2014

<https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73>

- ▶ Idea : **regularize** the distribution of the encodings, such that the latent space is smoother, and new sample can be produced
- ▶ To do so :
 - ▶ encode inputs as distributions, and feed sample from these distributions to the decoder.
 - ▶ Add a regularizing term (Kullback-Leibler divergence) on these distributions.



$$\text{loss} = C \|x - \hat{x}\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)] = C \|x - f(z)\|^2 + \text{KL}[N(g(x), h(x)), N(0, I)]$$

Variation Auto-Encoder (VAE) Kingma and Welling, 2014

- Desired properties : continuity (two close points in the latent space should give close outputs) and completeness (sampling any point should give a reasonable result)



Variation Auto-Encoder (VAE) Kingma and Welling, 2014

- Desired properties : continuity (two close points in the latent space should give close outputs) and completeness (sampling any point should give a reasonable result)
- Connection with Variational Inference (hence the name)



Variation Auto-Encoder (VAE) Kingma and Welling, 2014

- Desired properties : continuity (two close points in the latent space should give close outputs) and completeness (sampling any point should give a reasonable result)
- Connection with Variational Inference (hence the name)
- Easier to train than GANs, but output might have less diversity



Table of Contents

Generative Models

 Generative Adversarial Networks

 Variational Auto-encoder

 Diffusion models

Attention, Transformers

Conclusions and open questions

Some (subjective) non-technical considerations

Diffusion models

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

- ▶ Diffusion models have taken the generative world by storm and have completely replaced GANs (and somewhat VAEs) nowadays : Stable Diffusion, Midjourney, DALL-E, Imagen... Mostly due to their **stability, performance and flexibility**

Diffusion models

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

- ▶ Diffusion models have taken the generative world by storm and have completely replaced GANs (and somewhat VAEs) nowadays : Stable Diffusion, Midjourney, DALL-E, Imagen... Mostly due to their **stability, performance and flexibility**
- ▶ Diffusion models were *not* invented as generative models, and the idea seems a bit ludicrous at first glance !

Diffusion models

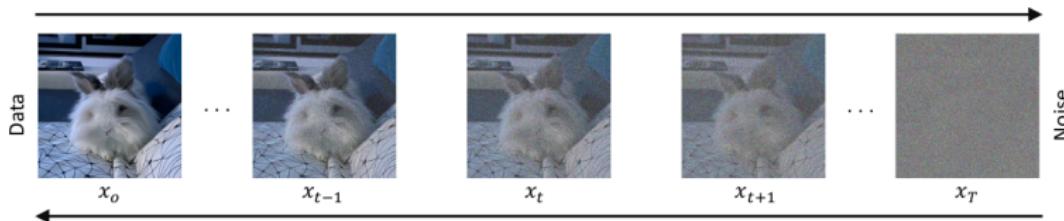
<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

- ▶ Diffusion models have taken the generative world by storm and have completely replaced GANs (and somewhat VAEs) nowadays : Stable Diffusion, Midjourney, DALL-E, Imagen... Mostly due to their **stability, performance and flexibility**
- ▶ Diffusion models were *not* invented as generative models, and the idea seems a bit ludicrous at first glance !
- ▶ Idea (roughly) : train a model to **denoise** an image, at many different noise levels. Then, start with **random noise** (!), and denoise it until an image is produced !

Diffusion models

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

- ▶ Diffusion models have taken the generative world by storm and have completely replaced GANs (and somewhat VAEs) nowadays : Stable Diffusion, Midjourney, DALL-E, Imagen... Mostly due to their **stability, performance and flexibility**
- ▶ Diffusion models were *not* invented as generative models, and the idea seems a bit ludicrous at first glance !
- ▶ Idea (roughly) : train a model to **denoise** an image, at many different noise levels. Then, start with **random noise** (!), and denoise it until an image is produced !
- ▶ Interpretation : we learn to “reverse” the diffusion process that gradually transforms an image into pure noise



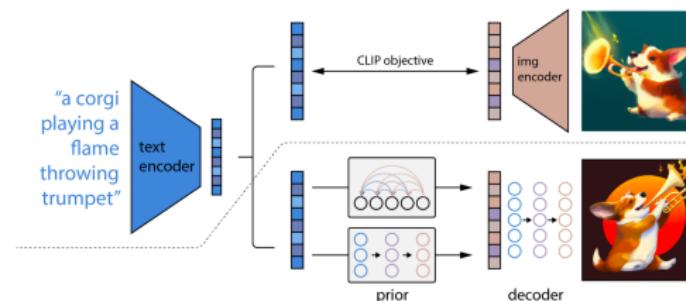
Diffusion models : pros

- ▶ Pro : Flexibility
 - ▶ denoising can happen **in the latent space**; somehow, learn to “project” random noise to a point in the latent space that makes sense

Diffusion models : pros

► Pro : Flexibility

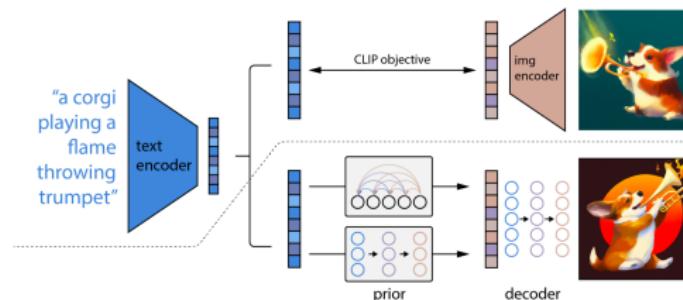
- denoising can happen **in the latent space**; somehow, learn to “project” random noise to a point in the latent space that makes sense
- possible to incorporate additional information/embedding at various stages before **conditional** diffusion happens : basis for all modern **text-to-image** systems ! (DALL-E, Stable Diffusion, Midjourney...)
- see eg <https://www.assemblyai.com/blog/how-dall-e-2-actually-works/>



Diffusion models : pros

► Pro : Flexibility

- ▶ denoising can happen **in the latent space**; somehow, learn to “project” random noise to a point in the latent space that makes sense
- ▶ possible to incorporate additional information/embedding at various stages before **conditional** diffusion happens : basis for all modern **text-to-image** systems ! (DALL-E, Stable Diffusion, Midjourney...)
- ▶ see eg <https://www.assemblyai.com/blog/how-dall-e-2-actually-works/>

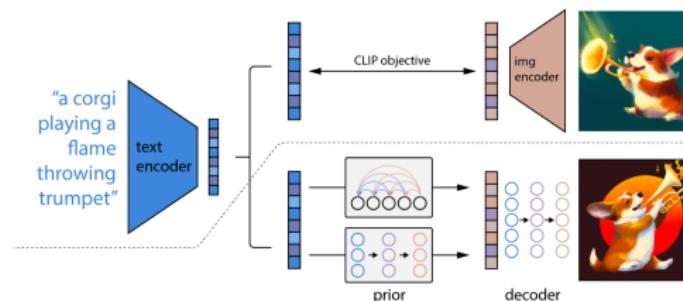


- ▶ Pro : **Tractability and performance.** Every step is analytically tractable, thus improving stability and interpretability. Images are generally more realistic than GANs or VAEs

Diffusion models : pros

► Pro : Flexibility

- denoising can happen **in the latent space**; somehow, learn to “project” random noise to a point in the latent space that makes sense
- possible to incorporate additional information/embedding at various stages before **conditional** diffusion happens : basis for all modern **text-to-image** systems ! (DALL-E, Stable Diffusion, Midjourney...)
- see eg <https://www.assemblyai.com/blog/how-dall-e-2-actually-works/>



- Pro : **Tractability and performance.** Every step is analytically tractable, thus improving stability and interpretability. Images are generally more realistic than GANs or VAEs
- Con : **slow sampling.** The sampling process is slower, we must perform many steps of denoising.

Diffusion models : the maths ?

- ▶ Diffusion models have renewed the interest in classical ODEs theory (w/ a lot of tool from physics), with an ML flavor
 - ▶ probability path, brownian motion, Fokker-Planck/Langevin equation...
 - ▶ see <https://iclr-blogposts.github.io/2024/blog/diffusion-theory-from-scratch/>

Diffusion : flow

- ▶ Variants/generalization are nowadays referred to as **flows** (matching flows, normalizing flows, conditional flow...)
 - ▶ see <https://dl.heeere.com/conditional-flow-matching/blog/conditional-flow-matching/>

Table of Contents

Generative Models

Generative Adversarial Networks

Variational Auto-encoder

Diffusion models

Attention, Transformers

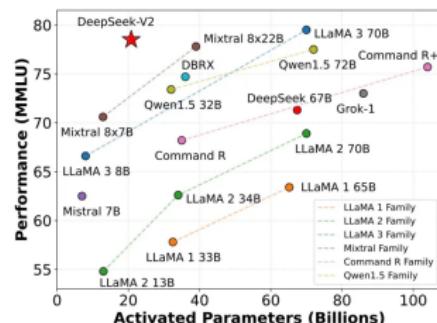
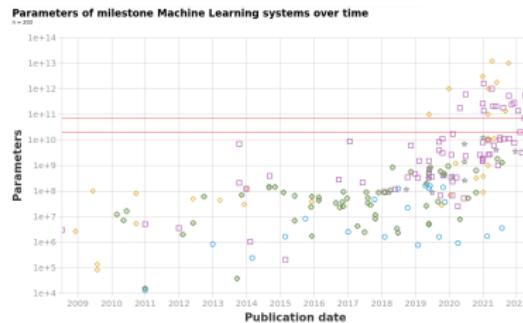
Conclusions and open questions

Some (subjective) non-technical considerations

Transformers : what

<https://towardsdatascience.com/transformers-explained-visually-part-1-overview-of-functionality-95a6dd460452>

- ▶ **Transformers** (Vaswani et al.), introduced by Google in 2017, have revolutionized first Natural Language Processing (NLP), then many other domains
- ▶ Documents classification, translation, text generation, text-to-image, image-to-text, chatbots...
- ▶ Well-known models, historically from Google/Meta/OpenAI, now from everyone and their neighbors (go explore HuggingFace !)

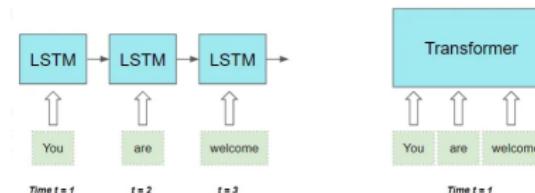


<https://epoch.ai/blog/machine-learning-model-sizes-and-the-parameter-gap>

<https://gradientflow.com/deepseek-v2-unpacked/>

Transformers : how ?

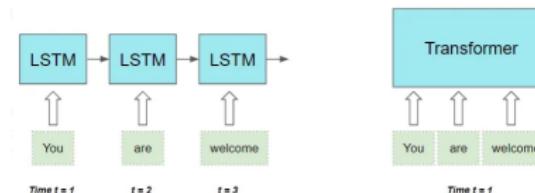
- ▶ Unlike RNNs, Transformers almost totally discard the “sequential” nature of texts : everything depends on everything !



- ▶ far, **far** more parameters ! (usually in the billions) : **huge training data, pre-trained components**. Modern Transformers can almost only be trained by companies with massive funding and resources. Ethical problems...

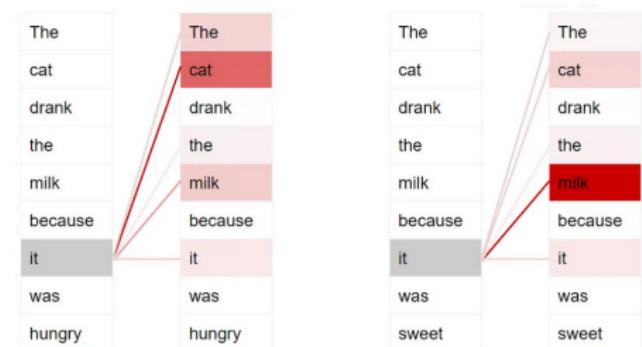
Transformers : how ?

- ▶ Unlike RNNs, Transformers almost totally discard the “sequential” nature of texts : everything depends on everything !



- ▶ far, far more parameters ! (usually in the billions) : huge training data, pre-trained components. Modern Transformers can almost only be trained by companies with massive funding and resources. Ethical problems...

The key to Transformers is **attention** : a (deceptively) simple mechanism to **model dependencies** between words/objects



Attention : equations

- ▶ In NLP, words (or subwords) are *token*, embedded using a [tokenizer](#).

Algorithm 3: Basic single-query attention.

Input: $e \in \mathbb{R}^{d_{\text{in}}}$, vector representation of the current token

Input: $e_t \in \mathbb{R}^{d_{\text{in}}}$, vector representations of context tokens $t \in [T]$.

Output: $\tilde{v} \in \mathbb{R}^{d_{\text{out}}}$, vector representation of the token and context combined.

Parameters: $W_q, W_k \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{in}}}$,
 $b_q, b_k \in \mathbb{R}^{d_{\text{attn}}}$, the query and key linear projections.

Parameters: $W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $b_v \in \mathbb{R}^{d_{\text{out}}}$, the value linear projection.

- 1 $q \leftarrow W_q e + b_q$
 - 2 $\forall t : k_t \leftarrow W_k e_t + b_k$
 - 3 $\forall t : v_t \leftarrow W_v e_t + b_v$
 - 4 $\forall t : \alpha_t = \frac{\exp(q^\top k_t / \sqrt{d_{\text{attn}}})}{\sum_u \exp(q^\top k_u / \sqrt{d_{\text{attn}}})}$
 - 5 **return** $\tilde{v} = \sum_{t=1}^T \alpha_t v_t$
-

Attention : equations

- ▶ In NLP, words (or subwords) are *token*, embedded using a [tokenizer](#).
- ▶ Given a token e , and **context tokens** e_t , **attention coefficients** are computed to relate e to the context

Algorithm 3: Basic single-query attention.

Input: $e \in \mathbb{R}^{d_{\text{in}}}$, vector representation of the current token

Input: $e_t \in \mathbb{R}^{d_{\text{in}}}$, vector representations of context tokens $t \in [T]$.

Output: $\tilde{v} \in \mathbb{R}^{d_{\text{out}}}$, vector representation of the token and context combined.

Parameters: $W_q, W_k \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{in}}}$,
 $b_q, b_k \in \mathbb{R}^{d_{\text{attn}}}$, the query and key linear projections.

Parameters: $W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $b_v \in \mathbb{R}^{d_{\text{out}}}$, the value linear projection.

- 1 $q \leftarrow W_q e + b_q$
 - 2 $\forall t : k_t \leftarrow W_k e_t + b_k$
 - 3 $\forall t : v_t \leftarrow W_v e_t + b_v$
 - 4 $\forall t : \alpha_t = \frac{\exp(q^\top k_t / \sqrt{d_{\text{attn}}})}{\sum_u \exp(q^\top k_u / \sqrt{d_{\text{attn}}})}$
 - 5 **return** $\tilde{v} = \sum_{t=1}^T \alpha_t v_t$
-

Attention : equations

- ▶ In NLP, words (or subwords) are *token*, embedded using a [tokenizer](#).
- ▶ Given a token e , and [context tokens](#) e_t , [attention coefficients](#) are computed to relate e to the context
- ▶ It can be [self-attention](#) (a word within its sentence), or [attention to a target context](#)

Algorithm 3: Basic single-query attention.

Input: $e \in \mathbb{R}^{d_{\text{in}}}$, vector representation of the current token

Input: $e_t \in \mathbb{R}^{d_{\text{in}}}$, vector representations of context tokens $t \in [T]$.

Output: $\tilde{v} \in \mathbb{R}^{d_{\text{out}}}$, vector representation of the token and context combined.

Parameters: $W_q, W_k \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{in}}}$,
 $b_q, b_k \in \mathbb{R}^{d_{\text{attn}}}$, the query and key linear projections.

Parameters: $W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $b_v \in \mathbb{R}^{d_{\text{out}}}$, the value linear projection.

- 1 $q \leftarrow W_q e + b_q$
 - 2 $\forall t : k_t \leftarrow W_k e_t + b_k$
 - 3 $\forall t : v_t \leftarrow W_v e_t + b_v$
 - 4 $\forall t : \alpha_t = \frac{\exp(q^\top k_t / \sqrt{d_{\text{attn}}})}{\sum_u \exp(q^\top k_u / \sqrt{d_{\text{attn}}})}$
 - 5 **return** $\tilde{v} = \sum_{t=1}^T \alpha_t v_t$
-

Attention : equations

- ▶ In NLP, words (or subwords) are *token*, embedded using a [tokenizer](#).
- ▶ Given a token e , and [context tokens](#) e_t , [attention coefficients](#) are computed to relate e to the context
- ▶ It can be [self-attention](#) (a word within its sentence), or [attention to a target context](#)
- ▶ The key operation is [softmax](#)
 $a_k : x \mapsto \frac{e^{x_k}}{\sum_i e^{x_i}}$. Attention is nothing else than common linear operations with learned parameters, [with proper normalization](#).

Algorithm 3: Basic single-query attention.

Input: $e \in \mathbb{R}^{d_{\text{in}}}$, vector representation of the current token

Input: $e_t \in \mathbb{R}^{d_{\text{in}}}$, vector representations of context tokens $t \in [T]$.

Output: $\tilde{v} \in \mathbb{R}^{d_{\text{out}}}$, vector representation of the token and context combined.

Parameters: $W_q, W_k \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{in}}}$, $b_q, b_k \in \mathbb{R}^{d_{\text{attn}}}$, the query and key linear projections.

Parameters: $W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $b_v \in \mathbb{R}^{d_{\text{out}}}$, the value linear projection.

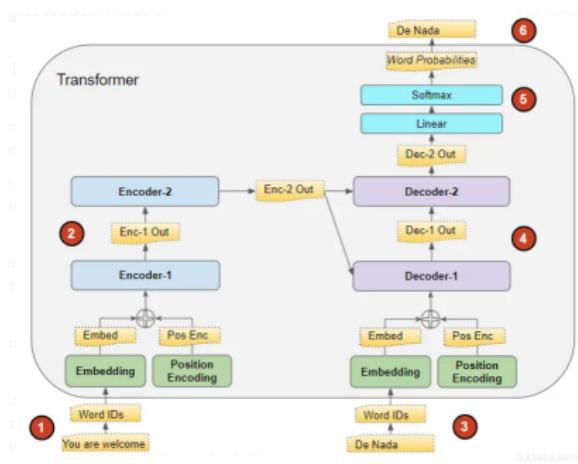
- 1 $q \leftarrow W_q e + b_q$
 - 2 $\forall t : k_t \leftarrow W_k e_t + b_k$
 - 3 $\forall t : v_t \leftarrow W_v e_t + b_v$
 - 4 $\forall t : \alpha_t = \frac{\exp(q^\top k_t / \sqrt{d_{\text{attn}}})}{\sum_u \exp(q^\top k_u / \sqrt{d_{\text{attn}}})}$
 - 5 **return** $\tilde{v} = \sum_{t=1}^T \alpha_t v_t$
-

Transformer : example of a (high-level) seq2seq architecture

https:

//towardsdatascience.com transformers-explained-visually-part-1-overview-of-functionality-95a6dd460452

- ▶ Training : takes a query x , a target y , produces an output

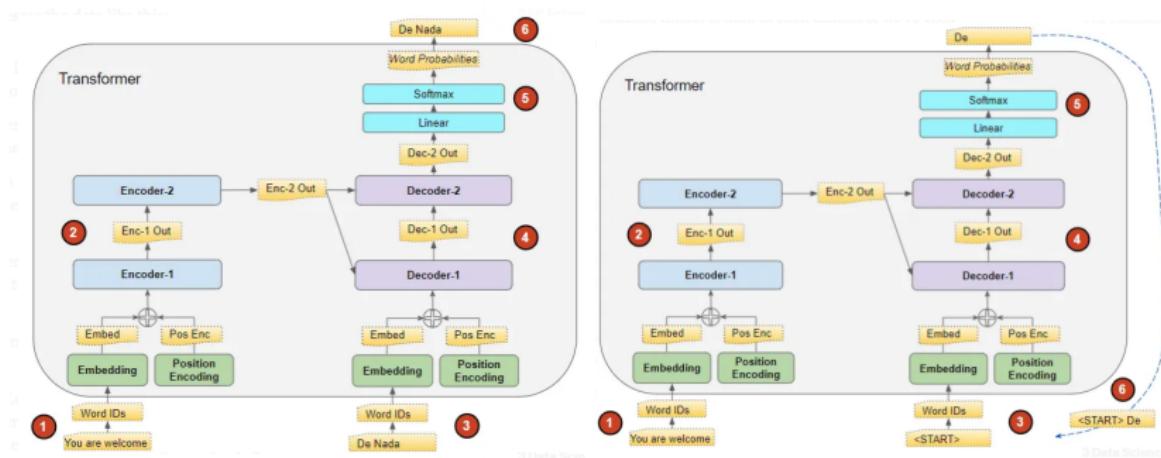


Transformer : example of a (high-level) seq2seq architecture

[https:](https://towardsdatascience.com/transf...)

//towardsdatascience.com/transf...mers-explained-visually-part-1-overview-of-functionality-95a6dd460452

- ▶ Training : takes a query x , a target y , produces an output
- ▶ Inference : takes a query x , start with an empty target, generate words one by one, by feeding the sequence generated until now as the target ("next word prediction")

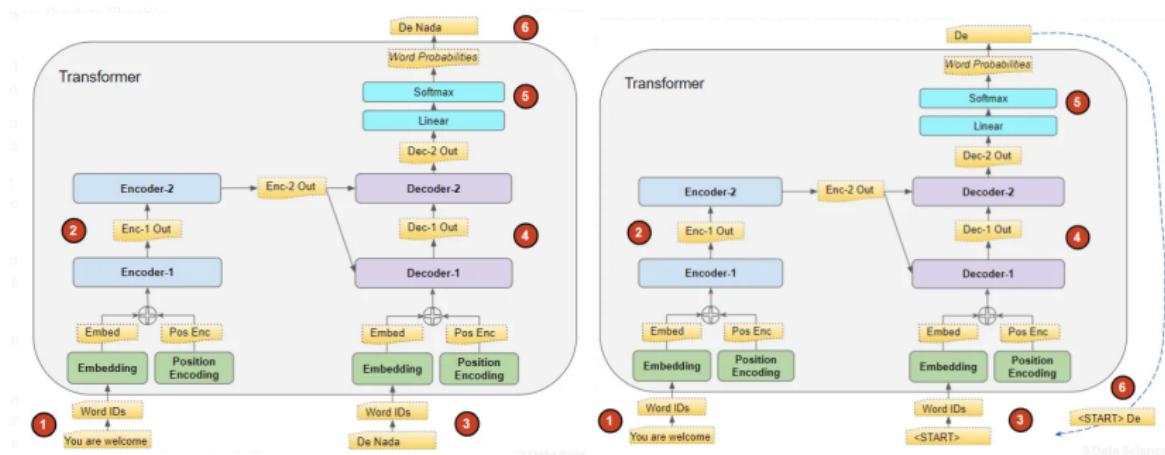


Transformer : example of a (high-level) seq2seq architecture

[https:](https://towardsdatascience.com/transf...)

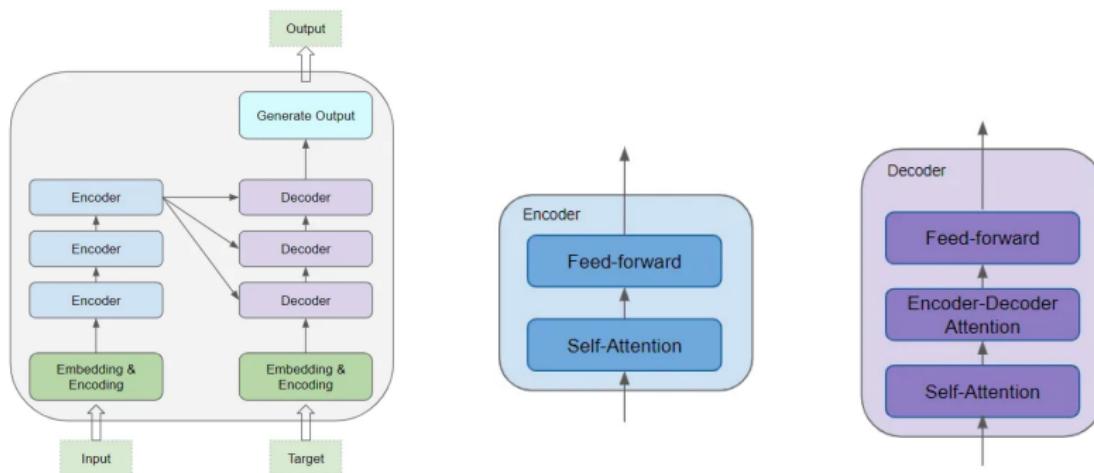
//towardsdatascience.com/transf...mers-explained-visually-part-1-overview-of-functionality-95a6dd460452

- ▶ Training : takes a query x , a target y , produces an output
- ▶ Inference : takes a query x , start with an empty target, generate words one by one, by feeding the sequence generated until now as the target ("next word prediction")
- ▶ Unlike RNNs/LSTMs, the next word generated directly depends on everything else (through attention and self-attention), not just by recurrent mechanisms



Transformer : example of a (high-level) seq2seq architecture

- ▶ Familiar concepts : encoding/decoding, latent space... but with attention layers !

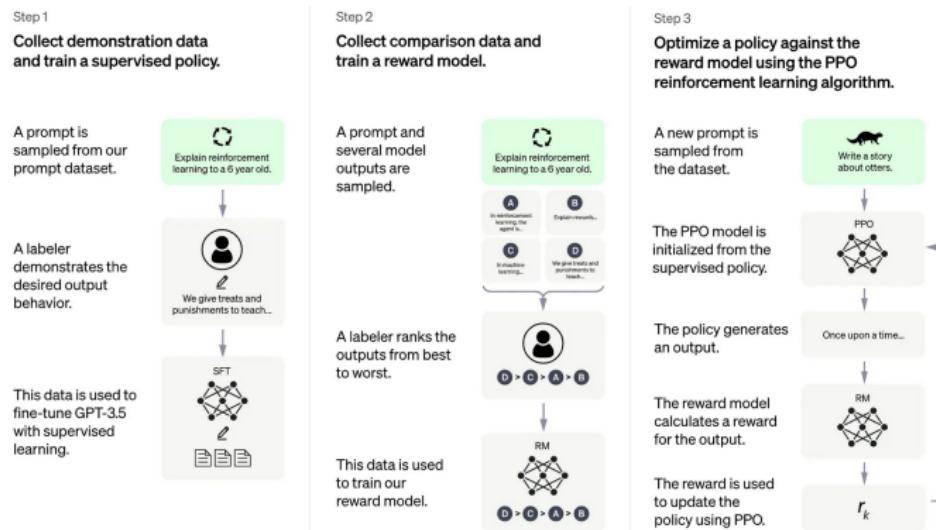


Transformers + RL

- ▶ Combined with RL strategies (action/rewards), Transformers become incredible chatbots.

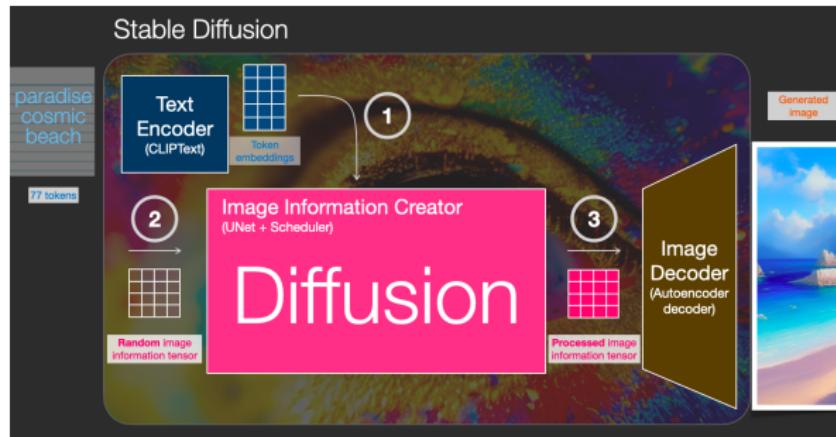
Transformers + RL

- ▶ Combined with RL strategies (action/rewards), Transformers become incredible chatbots.
- ▶ ChatGPT is based on the GPT3 model (Generative Pre-trained Transformer), trained with complex RL strategies and **HUGE** amount of data.



Transformers + diffusion models

- ▶ Transformers are great at computing **embeddings** for text, that can be used for other tasks. The BERT model is only an embedding model !
- ▶ These embeddings can be fed to other generative models. For instance, **all** the modern text-to-images models (Stable Diffusion, Midjourney, Imagen, DALL-E 2...) are based on **Transformers + diffusion models**.
- ▶ All these components are **pretrained** ! And then fine-tuned once put together.
- ▶ Many models/datasets are open-source ! See **HuggingFace** 😊 and eg
<https://www.assemblyai.com/blog/build-a-free-stable-diffusion-app-with-a-gpu-backend/>



Transformers + language models + diffusion models + time series analysis
+ CNN + massive training data + immense computing power + ...

Text-to-video : OpenAI's Sora (from last year)



<https://cdn.openai.com/sora/videos/tokyo-walk.mp4>

Table of Contents

Generative Models

Generative Adversarial Networks

Variational Auto-encoder

Diffusion models

Attention, Transformers

Conclusions and open questions

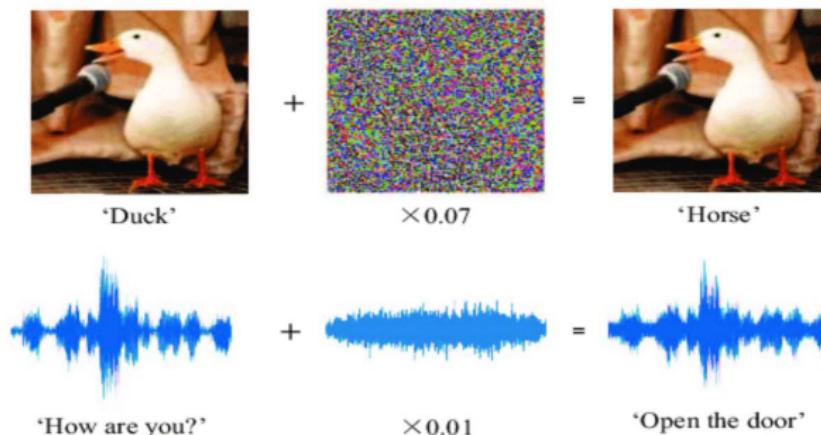
Some (subjective) non-technical considerations

Conclusion

- ▶ Deep learning is still mysterious and frustrating for many reasons
- ▶ They work incredibly well in some situations...
- ▶ ...but don't kill flies with bazookas, and always remember the basic principles ! (which are useful in deep learning too !)
- ▶ Modern libraries make them easy to build, easy to train...
- ▶ ... but they are useful for non-deep stuff too !

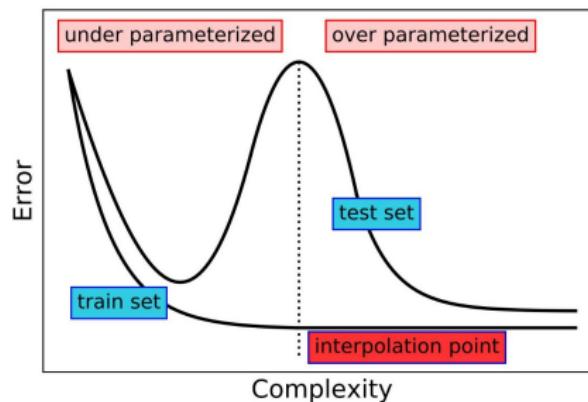
Hot question : adversarial attack

Once trained, a deep NN usually can be fooled easily, by **designing** an attack specifically for this NN. This unstability seems to be an integral part of high-performance NNs !



Hot question : generalization and double descent

- ▶ Modern Deep NN are **wildly overparametrized**, and **do interpolate the training data !!**
- ▶ But they still generatlize well !
- ▶ There is a **double descent phenomenon** : by increasing the complexity of the model past the interpolation point, among *all* interpolating models, it seems DNN+SGD favors models that **generalize well**. Still very much an open topic.



Hot topic : GNN, Geometric Deep Learning

- ▶ Graph Neural Networks (GNNs) work on **irregular data** that lives on network : molecules, proteins, social networks, computer networks, meshes, etc. They use “convolution” in this irregular space
- ▶ Most importantly, they are **permutation-invariant** : their output are left unchanged by a relabelling of nodes of the graph (graph isomorphism)
- ▶ They gave rise to **geometric deep learning**, which is the study of **symmetries and invariances** in DNNs : translation for CNN, permutation for GNN, but also rotation, gauge, etc. Much inspired by physics !

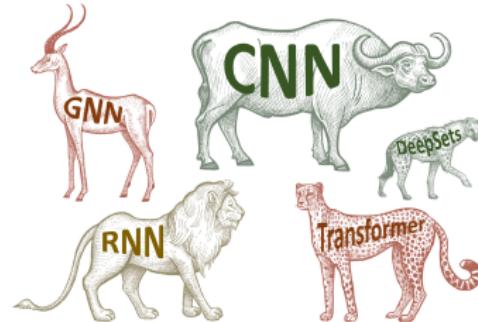
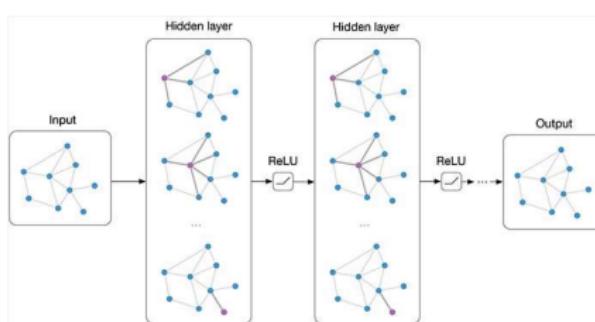


Table of Contents

Generative Models

Generative Adversarial Networks

Variational Auto-encoder

Diffusion models

Attention, Transformers

Conclusions and open questions

Some (subjective) non-technical considerations

Hot topic : the future of AI ?

- ▶ It seems the world's ressource is being poured into "AI". **But what AI?** (AGI?)

Hot topic : the future of AI ?

- ▶ It seems the world's ressource is being poured into "AI". **But what AI?** (AGI?)
- ▶ The **economic model of a chatbot** is discutable 
- ▶ "AI for Science" (healthcare, physics, chemistry, biology...) is promising, but has been for a long time

Hot topic : the future of AI ?

- ▶ It seems the world's ressource is being poured into "AI". **But what AI?** (AGI?)
- ▶ The **economic model of a chatbot** is discutable 
- ▶ "AI for Science" (healthcare, physics, chemistry, biology...) is promising, but has been for a long time
- ▶ "AI" will indeed optimize/automatize processes at many levels (industry, defense...), but this has also been true for a long time

Hot topic : the future of AI ?

- ▶ It seems the world's ressource is being poured into "AI". **But what AI?** (AGI?)
- ▶ The **economic model of a chatbot** is discutable 
- ▶ "AI for Science" (healthcare, physics, chemistry, biology...) is promising, but has been for a long time
- ▶ "AI" will indeed optimize/automatize processes at many levels (industry, defense...), but this has also been true for a long time
- ▶ AGI : a single "super-intelligent" entity is probably
 - ▶ unreachable
 - ▶ essentially useless (?)
 - ▶ a fantasy from the 60s
- ▶ **Multiple systems** interacting for a better (?) society is probably what we should aim at
- ▶ M. Jordan at AI summit : "Beyond the hype, AI might simply be a **new engineering field**, one in which human is in the loop"
 - ▶ <https://www.youtube.com/live/W0QLq4qEmKg?si=umJjs-o3JTZqPDF5&t=3811>

Hot topic : what about the scientific discipline ?

- ▶ Interpretability : Large (Language) Models are increasingly mysterious, people probe them as if they were natural phenomenon, and not human-made
 - ▶ were kernel methods really more interpretable... ?

Hot topic : what about the scientific discipline ?

- ▶ Interpretability : Large (Language) Models are increasingly mysterious, people probe them as if they were natural phenomenon, and not human-made
 - ▶ were kernel methods really more interpretable... ?
- ▶ Already a few years ago, deep learning was accused of being “alchemy” (and things were *incomparable* with today)

Hot topic : what about the scientific discipline ?

- ▶ Interpretability : Large (Language) Models are increasingly mysterious, people probe them as if they were natural phenomenon, and not human-made
 - ▶ were kernel methods really more interpretable... ?
- ▶ Already a few years ago, deep learning was accused of being "alchemy" (and things were *incomparable* with today)
- ▶ Fundamental/academic research *somewhat* lack the ressources (compute, man-power) to compete with large companies, but
 - ▶ it has (or must redefine) its place ! Many fields besides "training the biggest models" : privacy, efficiency, fairness, fundamentally new ideas, etc.
 - ▶ Likeness : companies are building skyscrapers, fundamental research establish the law of physics and the physics of materials.
 - ▶ Right now, the problem is that we overlap too much the two...

Hot topic : privacy

► Data privacy

- ▶ some data are incredibly sensitive (health, defense, critical industries, etc.)
- ▶ mathematical theory exist ("differential privacy")
- ▶ but DNNs interpolate, and thus generally *memorize training data* :(

Hot topic : privacy

► Data privacy

- ▶ some data are incredibly sensitive (health, defense, critical industries, etc.)
- ▶ mathematical theory exist ("differential privacy")
- ▶ but DNNs interpolate, and thus generally *memorize training data* :(

► Model privacy

- ▶ keep sensitive algorithms private/secure
- ▶ but open-sourceness is the only guarantee of transparency



Ethics and Justice, Law, Regulation, and Policy

Privacy in an AI Era: How Do We Protect Our Personal Information?

A new report analyzes the risks of AI and offers potential solutions.

Mar 18, 2024 | Katharine Miller [Twitter](#) [Facebook](#) [YouTube](#) [LinkedIn](#) [Email](#)

Hot topic : fairness

- ▶ By default, a model reproduce the **biases** in its training data (it's not inherently good or bad, it just is)
 - ▶ some sensitive data are not collected, so cannot be used to compensate for bias



Hot topic : fairness

- ▶ By default, a model reproduce the **biases** in its training data (it's not inherently good or bad, it just is)
 - ▶ some sensitive data are not collected, so cannot be used to compensate for bias



- ▶ Mathematical models : unbalanced classes, conditional probabilities...

Hot topic : fairness

- ▶ By default, a model reproduce the **biases** in its training data (it's not inherently good or bad, it just is)
 - ▶ some sensitive data are not collected, so cannot be used to compensate for bias



- ▶ Mathematical models : unbalanced classes, conditional probabilities...
- ▶ Non-open-source models can be *voluntarily* biased !
 - ▶ even open-source models are so complicated that bias can be hard to detect
 - ▶ from ads on the internet to democratic threat... ☹

FEBRUARY 4, 2025 | 5 MIN READ

Inside the NSF's Effort to Scour Research Grants for Violations of Trump's Orders

The U.S. National Science Foundation has unfrozen grant funding, but it continues to scrutinize research projects, sowing turmoil

BY DAN GARISTO, MAX KOZLOV & NATURE MAGAZINE

Hot topic : sovereignty

- ▶ Sovereignty over algorithms/data/models become as critical as other domains

Hot topic : sovereignty

- ▶ Sovereignty over algorithms/data/models become as critical as other domains
- ▶ beyond this, the *entire* ML/data science ecosystem uses critical systems and infrastructure (almost exclusively US-based : AWS...)

Hot topic : sovereignty

- ▶ Sovereignty over algorithms/data/models become as critical as other domains
- ▶ beyond this, the *entire* ML/data science ecosystem uses critical systems and infrastructure (almost exclusively US-based : AWS...)
- ▶ until now, the community was still very “*academic*” : open-source, international, everyone is nice and sharing. Recent events indicate that that may change...

Hot topic : sovereignty

- ▶ Sovereignty over algorithms/data/models become as critical as other domains
- ▶ beyond this, the *entire* ML/data science ecosystem uses critical systems and infrastructure (almost exclusively US-based : AWS...)
- ▶ until now, the community was still very “*academic*” : open-source, international, everyone is nice and sharing. Recent events indicate that that may change...
- ▶ What about France/Europe ? For now, policy-makers have decided to focus on “*ethical*” AI (open-source, fair, frugal, **useful** etc.), while still investing in local infrastructures.

IA : la pépite française Mistral va construire son premier data center en France

ÉDITION

La start-up française Mistral va construire son premier centre de données dédié à l'intelligence artificielle en France, a annoncé son patron Arthur Mensch sur la chaîne de télévision TF1 dimanche.



**AI Action Summit Conference:
AI, Science, and Society by IP
Paris**



Hot topic : environmental cost

- ▶ AI/ML requires large amount of resources : water, electricity, rare earths and metals...

The image displays two news articles from different sources. On the left is a screenshot of the NPR website, featuring a dark header with the 'npr' logo, a navigation bar with links for NEWS, CULTURE, MUSIC, PODCASTS & SHOWS, and a SEARCH bar. Below this, a headline under the BUSINESS category reads: "AI brings soaring emissions for Google and Microsoft, a major contributor to climate change". On the right is a screenshot of the Dallas Observer website, with a red header containing social media icons and a navigation menu with links for News, Opinion, Food & Drink, Arts & Culture, Music, Cannabis, and Things To Do. Below this, a headline under the ENVIRONMENT and TECHNOLOGY categories reads: "Electricity and Water Are Required To Run Data Centers. Texas Is Running Short on Both."

Hot topic : environmental cost

- ▶ AI/ML requires large amount of resources : water, electricity, rare earths and metals...

The screenshot shows the Dallas Observer website. At the top, there is a navigation bar with links for NEWS, CULTURE, MUSIC, PODCASTS & SHOWS, and SEARCH. On the right side of the header, there are social media icons for Facebook, Twitter, X, LinkedIn, and others. Below the header, there is a main content area with a red banner at the top. The banner contains the Dallas Observer logo and a link to "Score Sweet Dallas Observer Merch From Our Collab With Artist J". The main article title is "Electricity and Water Are Required To Run Data Centers. Texas Is Running Short on Both." The article is categorized under ENVIRONMENT and TECHNOLOGY.

- ▶ AI brings soaring emissions for Google and Microsoft, a major contributor to climate change
- ▶ AI may help “tackle climate change”... but beware the rebound effect !
 - ▶ <https://www.climatechange.ai/>

The screenshot shows the UN Environment Programme website. At the top, there is a navigation bar with links for Who we are, Where we work, What we do, Publications & data, and a search icon. Below the navigation bar, there is a main content area with a large image of a circuit board. Overlaid on the image is a text box containing the following text: "AI has an environmental problem. Here's what the world can do about that." At the bottom left of the image, there is a timestamp: "21 SEP 2024 | STORY | ENVIRONMENT UNDER REVIEW".

Hot topic : creativity and ethics

- ▶ **Copyrights** : are Internet-scrapped data stolen from artists/writers ? Why would a company pay to hire artists ?

Hot topic : creativity and ethics

- ▶ **Copyrights** : are Internet-scraped data stolen from artists/writers ? Why would a company pay to hire artists ?
- ▶ **Art and innovation** : at its core, generative models combine/reproduce the training data. They cannot do otherwise, structurally. Doesn't this encourage stagnation ? Is true creativity gonna be hidden in a sea of AI-generated content ? Moreover, all these are controlled by a few megacorp ?

Hot topic : creativity and ethics

- ▶ **Copyrights** : are Internet-scrapped data stolen from artists/writers ? Why would a company pay to hire artists ?
- ▶ **Art and innovation** : at its core, generative models combine/reproduce the training data. They cannot do otherwise, structurally. Doesn't this encourage stagnation ? Is true creativity gonna be hidden in a sea of AI-generated content ? Moreover, all these are controlled by a few megacorp ?
- ▶ **Misinformation** : fake news, deep fakes, democracies... is the Internet of the future gonna be a collection of not-to-be-trusted fake content ? Through fake news and recommendation algorithms, is the world gonna be more and more polarized ? Are video evidences still gonna be acceptable in court ? Is online bullying gonna take immense proportion ? ...

Hot topic : creativity and ethics

- ▶ **Copyrights** : are Internet-scrapped data stolen from artists/writers ? Why would a company pay to hire artists ?
- ▶ **Art and innovation** : at its core, generative models combine/reproduce the training data. They cannot do otherwise, structurally. Doesn't this encourage stagnation ? Is true creativity gonna be hidden in a sea of AI-generated content ? Moreover, all these are controlled by a few megacorp ?
- ▶ **Misinformation** : fake news, deep fakes, democracies... is the Internet of the future gonna be a collection of not-to-be-trusted fake content ? Through fake news and recommendation algorithms, is the world gonna be more and more polarized ? Are video evidences still gonna be acceptable in court ? Is online bullying gonna take immense proportion ? ...

AI is *not* gonna “destroy” the world, at least not Terminator-style (humans are far more dangerous). But... have we just destroyed / altered fundamentally crucial common goods ? Like the Internet ? Are we the last generation to “mostly” trust what we see online ?