

# MGM: Problemset

⌚ Category	
📎 Files	
🕒 Created	@May 3, 2023 10:44 AM
📅 Reminder	
⌚ Status	<button>Open</button>
🔗 URL	
🕒 Updated	@May 5, 2023 2:28 PM

## Exercise 0.8

### 1.

A fundamental problem in statistics and machine learning is to come up with useful measures of “distance” between pairs of probability distributions.

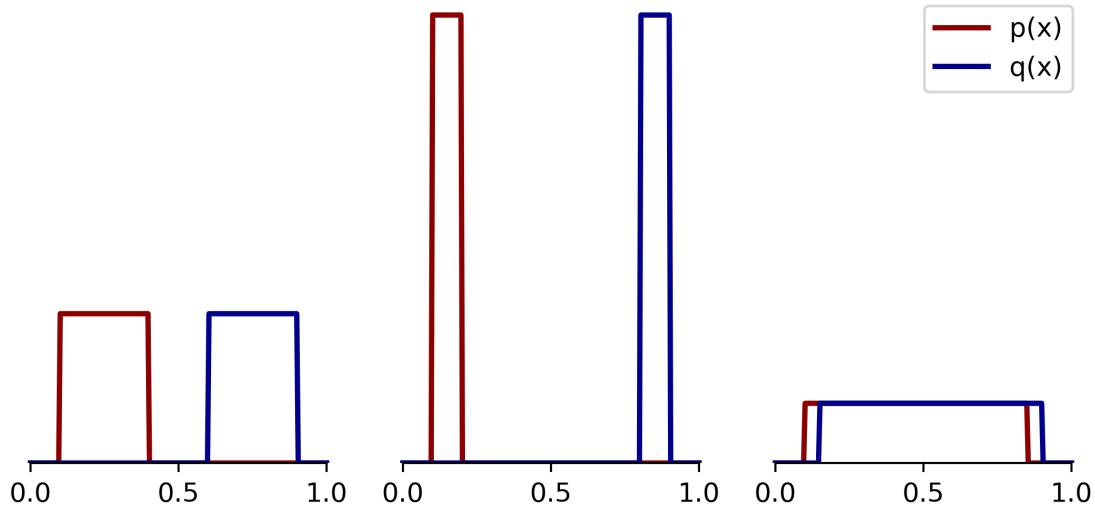
Two desirable properties of a distance function are symmetry and the triangle inequality. Unfortunately, many notions of “distance” between probability distributions do not satisfy these properties. These weaker notions of distance are often called *divergences*. Perhaps the most well-known divergence is the **Kullback-Liebler (KL) divergence**:

$$D_{KL}(P||Q) = \int p(x) \cdot \log \left( \frac{p(x)}{q(x)} \right)$$

While the KL divergence is incredibly useful and fundamental in information theory, it also has its shortcomings:

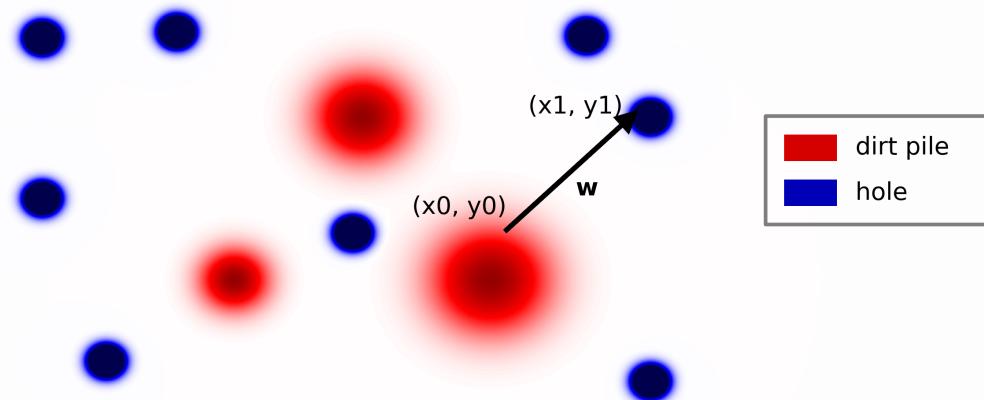
- For instance, one of the first things we learn about the KL divergence is that it is not symmetric,  $D_{KL}(P||Q) \neq D_{KL}(Q||P)$ . This is arguably not a huge problem, since various symmetrized analogues to the KL divergence exist.
- A bigger problem is that the divergence may be infinite if the support of  $P$  and  $Q$  are not equal.  
One way of circumventing this is to smooth (i.e. add blur to) the distributions

before computing the KL divergence, so that the support of  $P$  and  $Q$  matches. However, choosing the bandwidth parameter of the smoothing kernel is not always straightforward.



Optimal transport theory is one way to construct an alternative notion of distance between probability distributions.

**Optimal transport** (also known as Monge-Kantorovich or transportation) problem is a mathematical optimization problem that aims to find the optimal way of transporting a set of objects from one location to another, given some cost or distance measure between the objects and the transportation cost.



More formally, let  $X$  and  $Y$  be two metric spaces, and let  $\mu$  and  $\nu$  be two probability measures on  $X$  and  $Y$ , respectively. The optimal transport problem seeks to find a transport plan  $\pi$  that minimizes the cost functional:

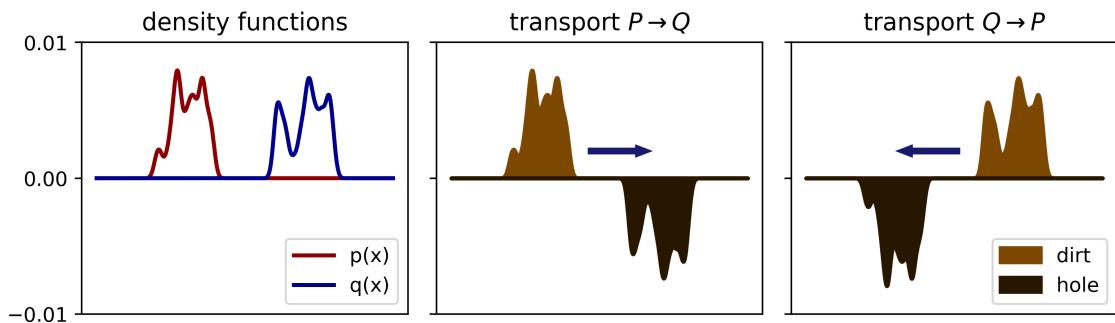
$$\inf_{\pi \in \Pi(\mu, \nu)} \int_{X \times Y} c(x, y) d\pi(x, y)$$

subject to the marginal constraints:

$$\begin{aligned} \pi(A \times Y) &= \mu(A) \\ \text{and} \\ \pi(X \times B) &= \nu(B) \end{aligned}$$

where,

- $c(x, y)$  is the cost or distance function between  $x$  and  $y$
- $\Pi(\mu, \nu)$  is the set of all joint probability measures on  $X \times Y$  with marginals  $\mu$  and  $\nu$ .



The total transportation cost overcomes the two weaknesses of KL divergence. First, since the cost function  $C$  is symmetric, the overall cost to transport  $P \rightarrow Q$  is the same as transporting  $Q \rightarrow P$ .

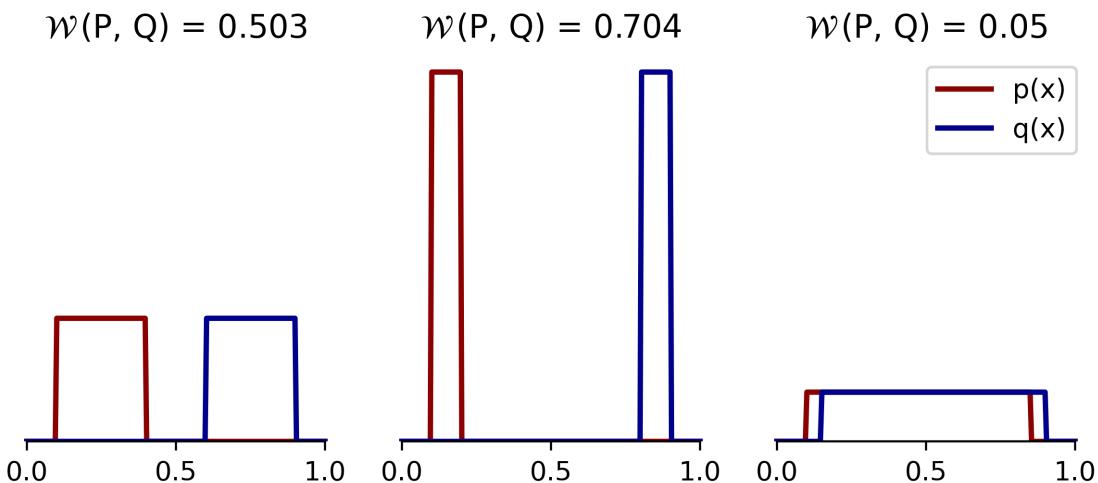
The **Wasserstein distance**, also known as the earthmover's distance, is a special case of the optimal transport problem, where the cost function  $c(x, y)$  is the distance between  $x$  and  $y$  raised to some power  $p$  (usually  $p=1$  or  $p=2$ ).

The Wasserstein distance of order  $p$  between two probability measures  $\mu$  and  $\nu$  is defined as:

$$W_p(\mu, \nu) = \left( \inf_{\pi \in \Pi(\mu, \nu)} \int_{X \times Y} d(x, y)^p d\pi(x, y) \right)^{\frac{1}{p}}$$

where,

- $d(x, y)$  is the distance between  $x$  and  $y$ .



The Wasserstein distance provides a metric on the space of probability measures that captures the similarity between two distributions in terms of their shape and location. This distance is not only symmetric, but also satisfies the triangle inequality.<sup>2</sup>

Recall the second shortcoming of KL divergence — it was infinite for a variety of distributions with unequal support. Not only is the Wasserstein distance finite in all cases, but the distances agree with our natural intuitions: the panel on the right

results in the smallest Wasserstein distance, while the middle panel shows the largest distance.

## 2.

Entropy-regularized optimal transport (EROT) is a modification of the standard optimal transport problem that introduces a strongly convex regularizer to the cost functional.

The strong convex regularizer encourages the optimal transport plan to be more diffuse or spread out, leading to smoother and more stable solutions.

$$d_c^\lambda(p, q) = \min_{\pi} \langle c, \pi \rangle - \lambda \cdot H(\pi)$$

Where,

- $H(\pi) = - \sum \pi(x, y) \log(\pi(x, y))$
- Negative of  $H$  / Shannon Entropy is strongly convex
- Thus,  $d = \text{linear} + \text{strongly convex} = \text{strongly convex}$
- Thus, this is guaranteed to have unique minimizer

The role of entropy regularization in EROT is to overcome some of the limitations of the standard optimal transport problem, such as the difficulty of computing the optimal transport plan and the sensitivity to noise and outliers in the data.

Entropy regularization allows for efficient and stable computation of the optimal transport plan using convex optimization techniques, and it also provides a smooth and continuous approximation to the optimal transport problem, which can be easier to work with in practice

## 3.

(3)

Let ' $f$ ' be linear.

Let ' $g$ ' be strongly convex.

$$\Rightarrow h = f + g$$

$\because f$  is also convex, we can write:

$$f(y) \geq f(u) + \nabla f(u)^T(y-u) \quad -\textcircled{1}$$

[1<sup>st</sup> order cond]

$\because g$  is strongly convex, we can write:

[strong conv<sup>n</sup> cond]

$$g(y) \geq g(x) + \nabla g(x)^T(y-x) + \frac{\lambda}{2} \|y-x\|^2 \quad -\textcircled{2}$$

(for some true const.  $\lambda$ )

$\textcircled{1} + \textcircled{2}$

$$[f(y) + g(y)] \geq [f(x) + g(x)] + (\nabla f(x)^T + \nabla g(x)^T)(y-x) + \frac{\lambda}{2} \|y-x\|^2$$

$$\Rightarrow h(y) \geq h(u) + h(u)^T(y-u) + \frac{\lambda}{2} \|y-u\|^2$$

Thus, ' $h$ ' is strongly convex

#### 4.

The entropy-regularized optimal transport problem can be solved using the alternating projection algorithm, which alternates between finding the optimal

transport plan with respect to the cost function and the entropy term.

(4)

$$d_c^\lambda(p, q) = \min_{\pi} \langle c, \pi \rangle - \lambda \cdot H(\pi)$$

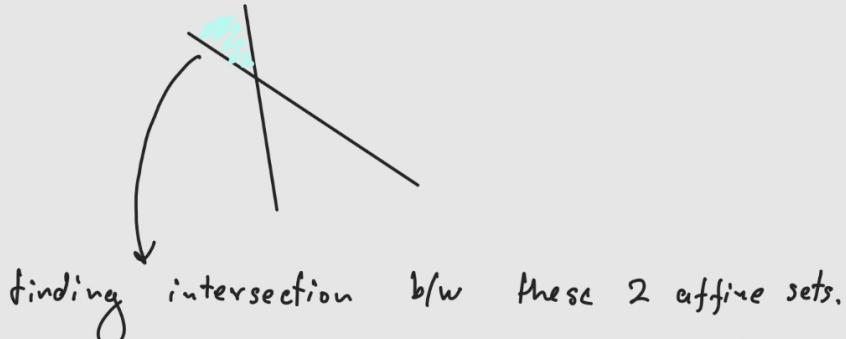
$\therefore d_c^\lambda(p, q)$  is strongly convex

↳ has existence of a unique sol<sup>n</sup>:

$$\underset{\pi \in \Pi(p, q)}{\operatorname{argmin}} D_{KL}(\pi \| p_k^\lambda)$$

Let  $\Pi(p)$ ,  $\Pi(q)$  denote row and column marginal constraints.

$$\Pi(p, q) = \Pi(p) \cap \Pi(q)$$



To do this we use Alternating Projection which finds intersection b/w 2 convex sets  $C_1, C_2$ .

Initialize

$\Pi_\lambda^0 = p_k^\lambda$  & define alternating projections.

$$\Pi_\lambda^{(l+1)} = \begin{cases} \underset{\pi \in \Pi(p)}{\operatorname{argmin}} D_{KL}(\pi \| \Pi_\lambda^l), & l = \text{even} \\ \underset{\pi \in \Pi(q)}{\operatorname{argmin}} D_{KL}(\pi \| \Pi_\lambda^{(l+1)}), & l = \text{odd} \end{cases}$$



Note

$$\Pi(p) = \{ \pi \mid \pi \mathbb{1}_m = p \}$$

$$\Pi(q) = \{ \pi \mid I_n^T \pi = q \}$$

↳ affine sets

↳ convex, can use alternating projection

For  $\lambda$  even: Need to find

$$\arg \min_{\pi^{\text{2m+1}}} D_{\text{KL}}(\pi \| \pi_\lambda^{(t)})$$

Since constrained problem  $\rightarrow$  introduce Lagrange mul.

$$\Rightarrow \arg \min_{\substack{\pi \\ f}} \underset{\lambda}{\text{man}} D(\pi \| \pi_\lambda^{(t)}) + \langle f, p - \pi 1_m \rangle$$

↓  
Slater's condn  
Lagrange multiplier

Assume strong duality interchange min & man, and then look at the min prob.

From First order optimality:

$$\frac{\delta}{\delta \pi} [D(\pi \| \pi_\lambda^{(t)}) + \langle f, p - \pi 1_m \rangle] = 0$$

But current  $\pi$  is  $\pi_{\lambda, f}^{(t+1)}$ :

$$\begin{aligned} & \frac{\delta}{\delta \pi_{\lambda}^{t+1}} [D_{\text{KL}}(\pi_\lambda^{t+1} \| \pi_\lambda^t)] \\ \Rightarrow & \frac{\delta}{\delta \pi_\lambda^{t+1}} [\pi_\lambda^{t+1}(u, y) \log \pi_\lambda^{t+1}(u, y) - \pi_\lambda^{t+1}(u, y) \log \pi_\lambda^t(u, y)] \\ \Rightarrow & 1 + \log \pi_\lambda^{t+1}(u, y) - \log \pi_\lambda^t(u, y) - ① \end{aligned}$$

↳ acts like a const.

& we have

$$\left. \frac{\delta}{\delta \pi_{\lambda}^{t+1}} \langle f, p - \pi_\lambda^{t+1} 1_m \rangle \right|_{\substack{\text{for fixed} \\ (u, y)}} = -f_u - ②$$

Combining ① & ② (1<sup>st</sup> order optimality)

$$1 + \log \pi_{\lambda}^{(l+1)}(n, y) - \log \pi_{\lambda}^l(n, y) - t_n = 0$$

Solve for  $\pi_{\lambda}^{(l+1)}(n, y)$ ,

$$\Rightarrow e^{t_n-1} = \frac{\pi_{\lambda}^{l+1}(n, y)}{\pi_{\lambda}^l(n, y)}$$

$\therefore \pi_{\lambda}^{l+1} \in \Pi(\rho)$ , we have:

$$\rho(n) = \sum \pi_{\lambda}^{l+1}(n, y) = e^{t_n-1} \sum \pi_{\lambda}^l(n, y) \quad \rightarrow ③$$

$$\pi_{\lambda}^{(l+1)}(n, y) = \frac{\rho(n)}{\sum \pi_{\lambda}^l(n, y)} \pi_{\lambda}^l(n, y)$$

Changing  $l+1$  to  $2l$  (even)  
 ↓       $l$  to  $2l-1$

$$\pi_{\lambda}^{(2l)}(n, y) = \frac{\rho(n)}{\sum \pi_{\lambda}^{(2l-1)}(n, y)} \pi_{\lambda}^{(2l-1)}(n, y)$$

vectorize for  $\mathbf{A}(n, y)$ :

$$\pi_{\lambda}^{(2l)} = \text{diag}\left(\frac{\rho}{\pi_{\lambda}^{(2l-1)} \mathbf{1}_m}\right) \pi_{\lambda}^{(2l-1)} \quad \rightarrow ④$$

Similarly for odd:

$$\pi_{\lambda}^{(2l+1)} = \text{diag}\left(\frac{\rho}{\mathbf{I}_n^\top \pi_{\lambda}^{(2l)}}\right) \pi_{\lambda}^{(2l)} \quad \rightarrow ⑤$$

## 5.

### Sinkhorn's Theorem

The Sinkhorn's theorem states that every square matrix with positive elements can be transformed into a doubly stochastic matrix  $D_1 A D_2$ . Where  $D_1$  and  $D_2$  are diagonal matrices with all positive main diagonals. The matrices  $D_1$  and  $D_2$  are themselves unique up to a constant factor.

A doubly stochastic matrix is an all non-negative square matrix that is both row-normalized and column-normalized.

### Sinkhorn–Knopp algorithm

A simple iterative method to approach the double stochastic matrix is to alternately rescale all rows and all columns of  $A$  to sum to 1.

```
def sinkhorn(A, N, L):
    # Pseudo-Code for calculating the doubly stochastic matrix
    # using Sinkhorn-Knopp algorithm.
    #
    # -----
    # Input: positive matrix A[N x N], max iteration L

    for i in range(L):
        A = A / np.matmul(A, np.ones(N, 1))
        A = A / np.matmul(np.ones(1, N), A)

        # Test for convergence and early stop.
        if _converge:
            break

    return A
```

### Optimal Transport

The Sinkhorn algorithm can be used for computing the optimal transport plan between two probability measures. The algorithm is based on the idea of alternating projections, and uses a scaling parameter to regularize the computations.

#### Proposition:

Let  $K \in \mathbb{R}^{n \times n}$  with  $K_{x,y} = K(x,y)$ . For some  $u \in \mathbb{R}^n, v \in \mathbb{R}^m$  :

$$\pi_\lambda = \text{diag}(u) \cdot K \cdot \text{diag}(v)$$

We Approximate the solution by initializing:

- $u^{(1)} = \mathbf{1}_n, v^{(1)} = \mathbf{1}_m$
- $u^{(l+1)} = \frac{p}{k \cdot v^{(l)}}, v^{(l+1)} = \frac{q}{k^T \cdot u^{(l)}}$

$$\tilde{\pi} = \text{diag}(u) \cdot K \cdot \text{diag}(v)$$

$$\tilde{\pi}^{(2l)} = \text{diag}(u^{(l+1)}) \cdot K \cdot \text{diag}(v^{(l)})$$

$$\tilde{\pi}^{(2l+1)} = \text{diag}(u^{(l+1)}) \cdot K \cdot \text{diag}(v^{(l)})$$

## 6.

### Matrix Scaling

Matrix Scaling problem is about finding  $u$  and  $v$  that scales rows and columns of  $K$ , such that it matches the row and column sum of the given distribution  $p$  and  $q$ .

In our case,

$$p = \Pi \cdot \mathbf{1}_m = \text{diag}(u) K v$$

$$q = \Pi \cdot \mathbf{1}_m = \text{diag}(v) K^T u$$

## 7.

$$\left. \begin{array}{l} p = \pi^T \mathbf{1}_m = \text{diag}(u) K v \\ q = \pi^T \mathbf{1}_n = \text{diag}(v) K^T u \end{array} \right\} - \textcircled{A}$$

We Approximate the sol<sup>n</sup> to  $\textcircled{A}$  by initializing

$$u^{(1)} = \mathbf{1}_m, \quad v^{(1)} = \mathbf{1}_m$$

$$u^{(2\ell+1)} = \frac{p}{K v^{(2\ell)}}, \quad v^{(2\ell+1)} = \frac{q}{K^T u^{(2\ell)}}$$

With these choices we get primal iterates:

Proof:

$$\tilde{\pi} = \text{diag}(u) K \text{diag}(v)$$

$$\tilde{\pi}^{(2\ell)} = \text{diag}(u^{(2\ell)}) K \text{diag}(v^{(2\ell)}) - \textcircled{B}$$

$$\tilde{\pi}^{(2\ell+1)} = \text{diag}(u^{(2\ell+1)}) K \text{diag}(v^{(2\ell+1)}) - \textcircled{D}$$

Rearranging terms from  $\textcircled{B}$ , ( $\ell = \ell-1$ )

$$K \text{diag}(v^{(2\ell)}) = \frac{\tilde{\pi}^{(2\ell-1)}}{\text{diag}(u^{(2\ell)})} - \textcircled{C}$$

Using in  $\textcircled{A}$ ,

$$\tilde{\pi}^{(2\ell)} = \text{diag}(u^{(2\ell+1)}) \frac{\tilde{\pi}^{(2\ell-1)}}{\text{diag}(u^{(2\ell)})}$$

$$= \text{diag}\left(\underbrace{\frac{p}{K v^{(2\ell)}}}_{\text{from } \textcircled{B}}\right) \frac{\tilde{\pi}^{(2\ell-1)}}{\text{diag}(u^{(2\ell)})}$$

$$= \text{diag}\left(\underbrace{\frac{p}{\text{diag}(u^{(2\ell)}) K v^{(2\ell)}}}_{\text{Recall } \textcircled{B} (\ell = \ell-1)}\right) \tilde{\pi}^{(2\ell-1)}$$

from  $\textcircled{B}$

$$= \text{diag}\left(\frac{p}{\tilde{\pi}^{(2\ell-1)} \mathbf{1}_m}\right) \tilde{\pi}^{(2\ell-1)}$$

Similarly

$$\tilde{\pi}^{(2\ell)} = \text{diag}\left(\frac{q}{\mathbf{1}_n^T \tilde{\pi}^{(2\ell)}}\right) \tilde{\pi}^{(2\ell)}$$

## 8.

### Sinkhorn Generative Modelling

Paper: Learning generative modelling with Sinkhorn Divergence, AISTATS, 2018

Algorithm:

- Draw Mini-batches:

$$x_1, x_2, \dots, x_B \sim p$$

$$y_1, y_2, \dots, y_B \sim p_\theta$$

- Approximate  $W_1(p, p_\theta) \approx W_1(\hat{p}, \hat{p}_\theta)$

$$\hat{p}(x) = \frac{1}{B} \sum_{i=1}^B \mathbf{1}_{x_i=x}$$

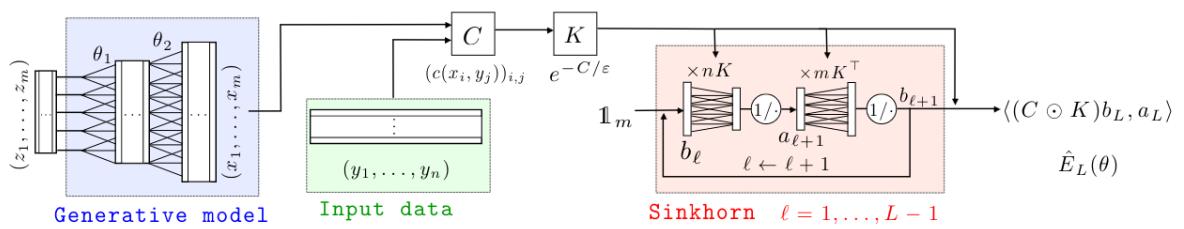
$$\hat{p}_\theta(x) = \frac{1}{B} \sum_{i=1}^B \mathbf{1}_{y_i=x}$$

- Estimate  $W_1(\hat{p}, \hat{p}_\theta)$  using Sinkhorn Algorithm

- Compute Gradient Update:

$$\theta' = \theta - \eta \nabla_\theta W_1(\hat{p}, \hat{p}_\theta)$$

## 9.



## 10.

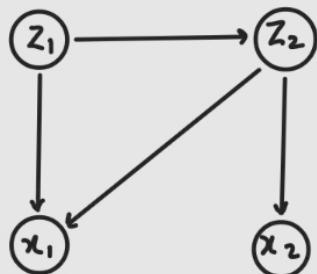
Probabilistic Graphical Models are powerful models that use graphs to describe dependency.

## 11.

(11)



$$p(x, z) = p(x|z) p(z)$$



$$p(x_1, x_2, z_1, z_2) = p(x_1|z_1, z_2) p(x_2|z_2) p(z_1) p(z_2|z_1)$$

## 12.

Variational Autoencoder (VAE) is called a latent variable model because it learns a compact and continuous representation of the input data, called the latent space, which captures the underlying structure and important features of the data.

## 13.

**ELBO**

Data likelihood:  $p_\theta(z) = \int p_\theta(z) p_\theta(z|z) dz$

we want                              ↓  
to minimize                            intractable  
this

$\Rightarrow$  Using variational inference to approximate the unknown posterior distribution from only the observed data.

VAE objective is to maximize

$$\log(p_\theta(z)) = \underset{z \sim q_\phi(z|x)}{\mathbb{E}} [\log p_\theta(z)] \quad (p_\theta(z) \text{ doesn't depend on } z)$$

$$= \underset{z}{\mathbb{E}} \left[ \log \left( \frac{p_\theta(z|x) p_\theta(z)}{q_\phi(z|x)} \right) \right] \quad (\text{Bayes})$$

$$= \underset{z}{\mathbb{E}} \left[ \log \left( \frac{p_\theta(z|x) p_\theta(z)}{q_\phi(z|x)} \cdot \frac{q_\phi(z|x)}{q_\phi(z|x)} \right) \right] \quad (\text{importance sampling})$$

$$\begin{aligned} &= \underset{z}{\mathbb{E}} [\log(p_\theta(z|x))] \\ &\quad - \underset{z}{\mathbb{E}} \left[ \log \left( \frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \\ &\quad + \underset{z}{\mathbb{E}} \left[ \log \left( \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right] \end{aligned}$$

↑ intractable

$$= \underset{z}{\mathbb{E}} [\log p_\theta(z|x)] - D_{KL}[q_\phi(z|x) \| p_\theta(z)] + \underbrace{D_{KL}[q_\phi(z|x) \| p_\theta(z|x)]}_{\geq 0}$$

$$\geq \underbrace{\underset{z}{\mathbb{E}} [\log p_\theta(z|x)] - D_{KL}[q_\phi(z|x) \| p_\theta(z)]}_{\text{Tractable lower Bound (ELBO)}}$$

↓  
Tractable lower Bound  
(ELBO)

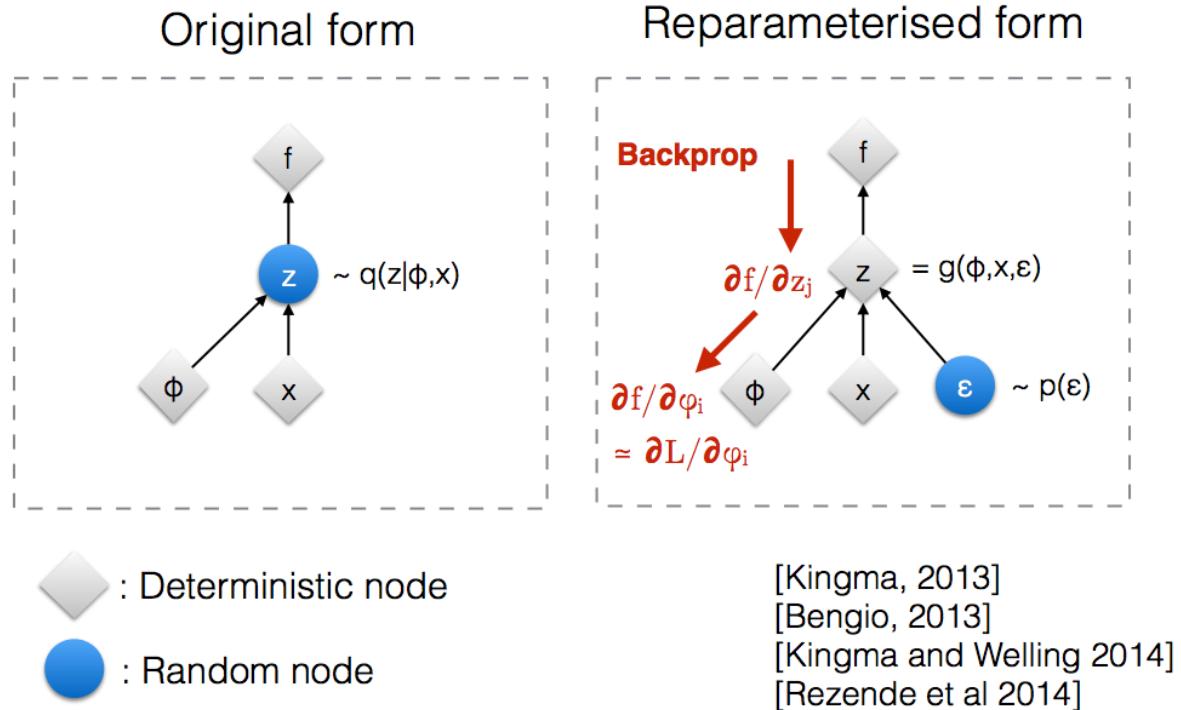
↳ VAE objective function

VAE Loss = - VAE Objective

$$= -\text{ELBO}$$

$$= D_{KL}[q_\phi(z|x) \| p_\theta(z)] - \underset{z}{\mathbb{E}} [\log(p_\theta(z|x))]$$

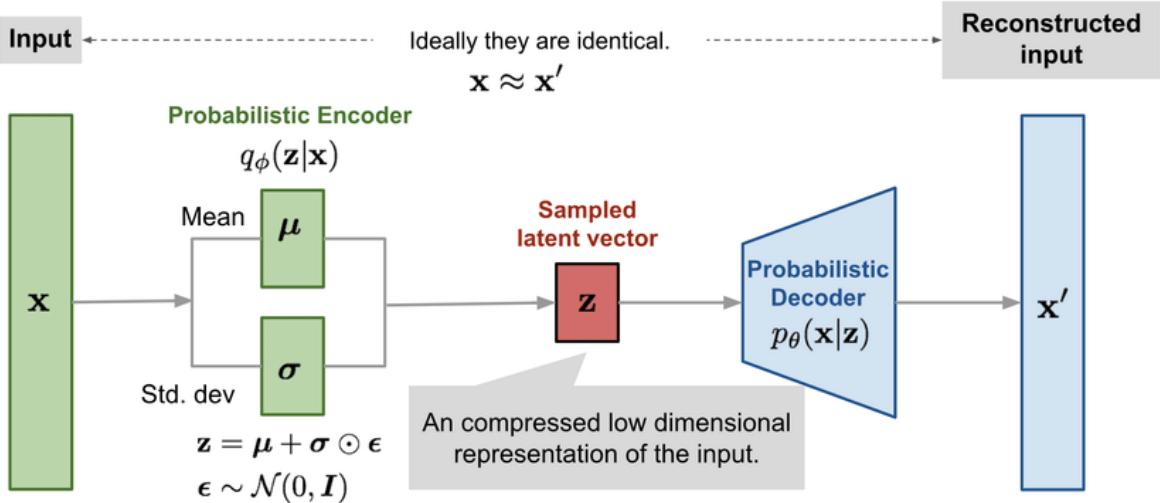
## Re-parameterization



The reparametrization trick is a technique used in Variational Autoencoder (VAE) models to enable backpropagation through the stochastic layer, which is typically represented by a random sample from a normal distribution.

The reparametrization trick separates the random sampling operation from the model parameters, allowing gradients to flow through the stochastic layer and enabling efficient training via stochastic gradient descent.

## Schematic Diagram



## Training Procedure

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

```

 $\theta, \phi \leftarrow$  Initialize parameters
repeat
     $\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)
     $\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$ 
     $\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))
     $\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])
until convergence of parameters  $(\theta, \phi)$ 
return  $\theta, \phi$ 

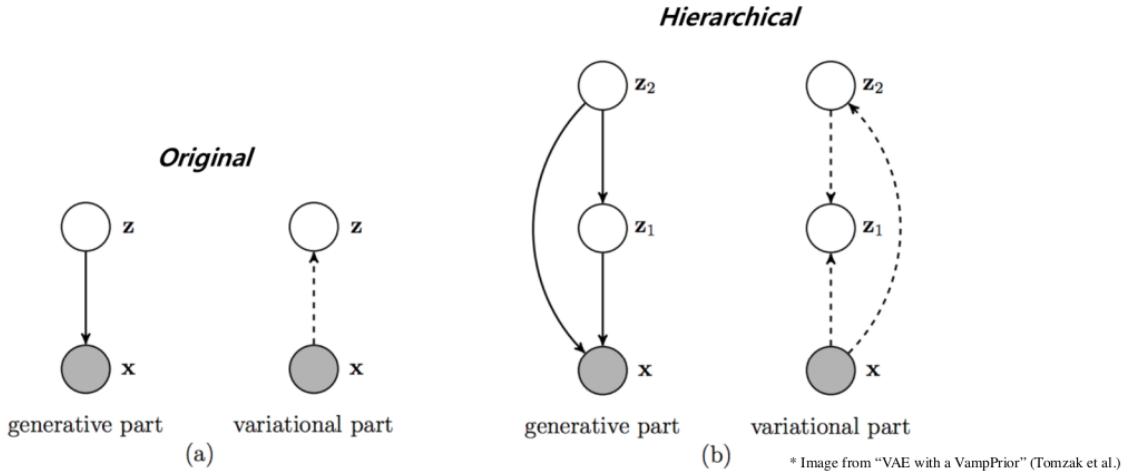
```

---

## 14.

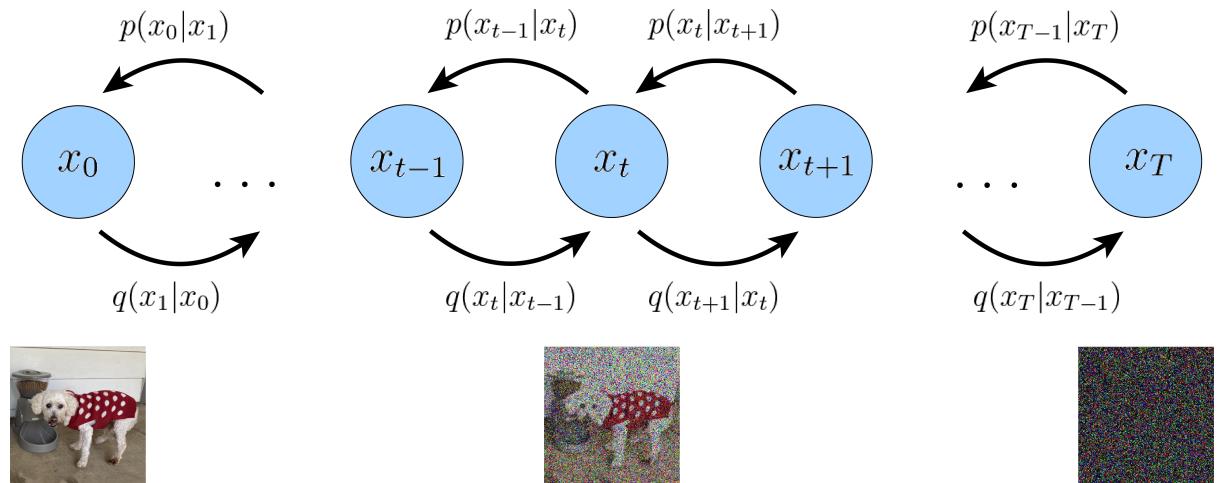
### Hierarchical VAE (HVAE)

- Generalize VAE to multiple hierarchies over latent variables
- Latent variables themselves are generated from other higher level, more abstract latents.
- General HVAE is with  $T$  hierarchical levels, each latent is allowed to condition on all previous labels.



## Markovian Hierarchical VAE (MHVAE)

- Assume the hierarchy to be Markovian.
- Decoding each latent  $Z_t$  only conditions on previous latent  $Z_{t+1}$



## 15.

Papers of Variational Diffusion Models:

1. “Deep Unsupervised Learning using Non-equilibrium Thermodynamics”, ICML 2015
2. “Denoising Diffusion Probabilistic Models”, NeurIPS 2020

- “Variational Diffusion Models”, Kingma, NeurIPS 2021

## 16.

### Modelling Assumptions in VDM:

- Latent Dimension is exactly equal to Data Dimension.
- Structure of the latent encoder is not learnt at each time step.  
It is pre-defined as linear Gaussian model.
- Gaussian parameters of the latent encoder vary over time in such a way that the distribution of the latent variable at the final time step  $T$  is a standard Gaussian.

## 17.

### Implications of Modelling Assumptions in VDM:

- Can represent both true data samples and latent variables as  $x_t$   
where,  
 $x_0$  = true data samples  
 $x_t$  = corresponding latent with hierarchy indexed by  $t \in [1, T]$

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

- The distribution of each latent variable in the encoder is a Gaussian centred around its previous hierarchical latent.

#### Mean

$$\mu_t(x_t) = \sqrt{\alpha_t} \cdot x_{t-1}$$

#### Variance

$$\Sigma_t(x_t) = (1 - \alpha_t) I$$

Where, the form of coefficient  $\alpha_t$  are chosen such that the variance of the latent variables stay at a similar scale.

- Encoding process is variance preserving

b. Alternate parameterizations exist

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t) I)$$

3.  $\alpha_t$  evolves over time according to a fixed or learnable schedule, such that the distribution of the final latent  $p(x_T)$  is a standard Gaussian.

Hence, HVAE becomes:

$$\begin{aligned} p(x_{0:T}) &= p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \\ p(x_T) &= \mathcal{N}(x_T; 0, I) \end{aligned}$$

All assumptions together:

Steady notification of an input image over time, we progressively corrupt an image by adding Gaussian noise until it eventually becomes completely identical to pure Gaussian noise.

## 18.

- $p(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$
- $p(x_T) = \mathcal{N}(x_T; 0, I)$

## 19.

Diffusion models are a type of probabilistic model used for learning patterns in data. They work by simulating a process of gradual diffusion, where information is spread slowly and smoothly across the data space. This diffusion process is controlled by a set of learned parameters that determine the rate and direction of diffusion. By iteratively diffusing the data over many steps, the model is able to capture complex patterns and dependencies in the data, and can be used for tasks such as image generation, denoising, and inpainting. Overall, diffusion models provide a powerful tool for understanding and modeling complex data distributions.

**20.**

like HVAE, VDM can be optimized by maximizing ELBO:

$$\log(p(\mathbf{u})) = \log \int p(\mathbf{u}_{0:T}) d\mathbf{u}_{0:T}$$

$$= \log \int \frac{p(\mathbf{u}_{0:T}) q(\mathbf{u}_{1:T} | \mathbf{u}_0)}{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} d\mathbf{u}_{1:T} \quad \begin{pmatrix} \text{Importance} \\ \text{Sampling} \end{pmatrix}$$

$$= \log \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \frac{p(\mathbf{u}_{0:T})}{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \right]$$

$$\geq \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \log \frac{p(\mathbf{u}_{0:T})}{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \right] \quad (\text{Jensen's inequality})$$

$$= \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \log \frac{p(\mathbf{u}_T) \prod_{t=1}^T p_\theta(\mathbf{u}_{t+1} | \mathbf{u}_t)}{\prod_{t=1}^T q(\mathbf{u}_t | \mathbf{u}_{t-1})} \right]$$

make the indexing same

$$= \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \log \frac{p(\mathbf{u}_T) p_\theta(\mathbf{u}_0 | \mathbf{u}_1) \prod_{t=2}^T p_\theta(\mathbf{u}_{t+1} | \mathbf{u}_t)}{q(\mathbf{u}_T | \mathbf{u}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{u}_t | \mathbf{u}_{t-1})} \right]$$

$$= \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \log \frac{p(\mathbf{u}_T) p_\theta(\mathbf{u}_0 | \mathbf{u}_1) \prod_{t=1}^{T-1} p_\theta(\mathbf{u}_t | \mathbf{u}_{t+1})}{q(\mathbf{u}_T | \mathbf{u}_{T-1}) \prod_{t=1}^{T-1} q(\mathbf{u}_t | \mathbf{u}_{t-1})} \right]$$

$$= \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \log \frac{p(\mathbf{u}_T) p_\theta(\mathbf{u}_0 | \mathbf{u}_1)}{q(\mathbf{u}_T | \mathbf{u}_{T-1})} \right]$$

$$+ \mathbb{E}_{q(\mathbf{u}_{1:T} | \mathbf{u}_0)} \left[ \log \prod_{t=1}^{T-1} \frac{p_\theta(\mathbf{u}_t | \mathbf{u}_{t+1})}{q(\mathbf{u}_t | \mathbf{u}_{t-1})} \right]$$

$$\begin{aligned}
&= \underset{q(n_1:n_T|n_0)}{\mathbb{E}} \left[ \log p_\theta(n_0|n_1) \right] + \underset{q(n_1:n_T|n_0)}{\mathbb{E}} \left[ \log \frac{p(n_T)}{q(n_T|n_{T-1})} \right] \\
&\quad + \sum_{t=1}^{T-1} \underset{q(n_1:n_T|n_0)}{\mathbb{E}} \left[ \log \frac{p(n_T)}{q(n_T|n_{T-1})} \right] \\
&\quad \xrightarrow{\text{simplify these to relevant vars}} \\
&= \underset{q(n_1|n_0)}{\mathbb{E}} \left[ \log p_\theta(n_0|n_1) \right] + \underset{q(x_{T-1}, n_T|n_0)}{\mathbb{E}} \left[ \log \frac{p(n_T)}{q(n_T|x_{T-1})} \right] \\
&\quad + \sum_{t=1}^{T-1} \underset{q(x_t, n_t, n_{t+1}|n_0)}{\mathbb{E}} \left[ \log \frac{p_\theta(n_t|x_{t+1})}{q(n_t|x_{t-1})} \right] \quad - \textcircled{1} \\
&\quad \xrightarrow{\text{ELBO}}
\end{aligned}$$

Further Simplification

$$\begin{aligned}
&= \underset{\textcircled{1}}{\mathbb{E}_{q(n_1|n_0)}} \left[ \log p_\theta(n_0|n_1) \right] - \underset{\textcircled{2}}{\mathbb{E}_{q(n_{T-1}|n_0)}} \left[ D_{KL} \left( q(n_T|n_{T-1}) \parallel p(n_T) \right) \right] \\
&\quad - \underset{\textcircled{3}}{\sum_{t=1}^{T-1} \mathbb{E}_{q(x_{t-1}, n_t, n_{t+1}|n_0)} \left[ D_{KL} \left( q(n_t|x_{t-1}) \parallel p_\theta(n_t|x_{t+1}) \right) \right]}
\end{aligned}$$

(1)  $\rightarrow$  Reconstruction term

(2)  $\rightarrow$  Prior Matching term

(3)  $\rightarrow$  Consistency term

denoising or matching

21.

(21)

$$q(u_{t-1} | u_t, u_0) = \frac{q(u_t | u_{t-1}) q(u_{t-1} | u_0)}{q(u_t | u_0)}$$

$$\Rightarrow q(u_{t-1} | u_t, u_0) = \frac{N(u_t; \sqrt{\alpha_t} u_{t-1}, (1-\alpha_t)I) N(u_{t-1}; \sqrt{\bar{\alpha}_{t-1}} u_0, (1-\bar{\alpha}_{t-1})I)}{N(u_t, \sqrt{\bar{\alpha}_t} u_0, (1-\bar{\alpha}_t)I)}$$

$$\propto \exp \left\{ - \left[ \frac{(u_t - \sqrt{\alpha_t} u_{t-1})^2}{2(1-\alpha_t)} + \frac{(u_{t-1} - \sqrt{\bar{\alpha}_{t-1}} u_0)^2}{2(1-\bar{\alpha}_{t-1})} - \frac{(u_t - \sqrt{\bar{\alpha}_t} u_0)^2}{2(1-\bar{\alpha}_t)} \right] \right\}$$

$$= \exp \left\{ -\frac{1}{2} \frac{1}{\frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)}} \left[ \frac{u_{t-1}^2 - 2\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})u_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)u_0^2}{1-\bar{\alpha}_t} \right] \right\}$$

$$\propto N(u_{t-1}; \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})u_t + \sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)u_0}{1-\bar{\alpha}_t}, \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} I)$$

$\therefore q(u_{t-1} | u_t, u_0) \rightarrow \text{Normal Distribution}$

22.

(22)

$$\log p(w) \geq \underbrace{\mathbb{E}_{q(n_t|n_0)} [\log p_\theta(n_0|n_t)]}_{\text{reconstruction error}} - \underbrace{D_{KL}(q(n_t|n_0) \| p(n_t))}_{\text{prior matching term}} \\ - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(n_t|n_0)} [D_{KL}(q(n_{t-1}|n_t, n_0) \| p_\theta(n_{t-1}|n_t))]}_{\text{denoising matching term}}$$

$$\begin{aligned} & \arg \min_{\theta} D_{KL}(q(n_{t-1}|n_t, n_0) \| p_\theta(n_{t-1}|n_t)) \\ &= \arg \min_{\theta} D_{KL}\left(N(n_{t-1}; \mu_q, \Sigma_q(t)) \| N(n_{t-1}; \mu_\theta, \Sigma_\theta(t))\right) \\ &\quad \xrightarrow{\text{we set the variances of these terms to be exactly same.}} \\ &= \arg \min_{\theta} \frac{1}{2} \left[ \log \frac{|\Sigma_q(t)|}{|\Sigma_\theta(t)|} - d + \text{tr}(\underbrace{\Sigma_q(t)^{-1} \Sigma_\theta(t)}_{I \in \mathbb{R}^{d \times d}}) \right. \\ &\quad \left. + (\mu_\theta - \mu_q)^\top \Sigma_q(t)^{-1} (\mu_\theta - \mu_q) \right] \\ &= \arg \min_{\theta} \frac{1}{2} \left[ (\mu_\theta - \mu_q)^\top \underbrace{\Sigma_q(t)^{-1}}_{(\Sigma_q^2(t)I)^{-1}} (\mu_\theta - \mu_q) \right] \\ &\approx \arg \min_{\theta} \frac{1}{2 \Sigma_q^2(t)} \left[ \|\mu_\theta - \mu_q\|_2^2 \right] \quad -\textcircled{1} \end{aligned}$$

where,

$$\mu_\theta = \mu_\theta(n_t, t) , \quad \mu_q = \mu_q(n_t, n_0)$$

Recall

$$d_{\eta}(n_t, n_0) = \frac{\sqrt{\alpha_t(1-\bar{\alpha}_{t-1})} n_t + \sqrt{\bar{\alpha}_{t-1}(1-\alpha_t)} n_0}{1-\bar{\alpha}_t} \quad - \textcircled{2}$$

As  $d_{\eta}(n_t)$  also cond's on  $n_t$ , we can match  $d_{\eta}(n_t, n_0)$  closely by setting it as:

$$d_{\eta}(n_t, t) = \frac{\sqrt{\alpha_t(1-\bar{\alpha}_{t-1})} n_t + \sqrt{\bar{\alpha}_{t-1}(1-\alpha_t)} \hat{n}_0(n_t, t)}{1-\bar{\alpha}_t} \quad - \textcircled{3}$$

where,

$\hat{n}_0(n_t, t) \rightarrow$  parameterized by NN  
to predict  $n_0$  from noisy image  
 $n_t$  at time index 't'.

Using  $\textcircled{2}, \textcircled{3}$  in  $\textcircled{1}$

$$\Rightarrow \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma_{\eta}^2(t)} \frac{\bar{\alpha}_{t-1}(1-\alpha_t)^2}{(1-\bar{\alpha}_t)^2} \left[ \left\| \hat{n}_0(n_t, t) - n_0 \right\|_2^2 \right]$$

23.

(23)

Model  $\alpha_t$  using a neural network  $\hat{\alpha}_n(t)$ ,  $\eta$  parameters

$$\frac{1}{2 \sigma_y^2(t)} \frac{\bar{\alpha}_{t-1}(1-\bar{\alpha}_t)^2}{(1-\bar{\alpha}_t)^2} \left[ \|\vec{n}_0(n_t, t) - n_0\|_2^2 \right]$$

$$\xrightarrow{\frac{(1-\bar{\alpha}_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}}$$

$$= \frac{1}{2} \left( \frac{\bar{\alpha}_{t-1}}{1-\bar{\alpha}_{t-1}} - \frac{\bar{\alpha}_t}{1-\bar{\alpha}_t} \right) \left[ \|\vec{n}_0(n_t, t) - n_0\|_2^2 \right] - \textcircled{1}$$

Recall

$$q(n_t | n_0) \sim N(n_t; \sqrt{\alpha_t} n_0, (1-\bar{\alpha}_t) I)$$

Define signal-to-noise ratio (SNR) =  $\frac{\alpha_t^2}{\sigma^2}$

$$SNR(t) = \frac{\bar{\alpha}_t}{1-\bar{\alpha}_t} \quad - \textcircled{2}$$

Using  $\textcircled{2}$  in  $\textcircled{1}$

$$= \frac{1}{2} (SNR(t-1) - SNR(t)) \left[ \|\vec{n}_0(n_t, t) - n_0\|_2^2 \right]$$

24.

(24.)

Solving for  $n_0$ :

$$n_0 = \frac{n_t - \sqrt{1-\bar{\alpha}_t} \epsilon_0}{\sqrt{\bar{\alpha}_t}} \quad -①$$

$$\log(n_t, n_0) = \frac{\sqrt{\alpha_t} (1-\bar{\alpha}_{t-1}) n_t + \sqrt{\bar{\alpha}_t} (1-\alpha_t) n_0}{1 - \bar{\alpha}_t}$$

Using ①.

$$= \frac{1}{\sqrt{\alpha_t}} n_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0 \quad -②$$

Using ② we can approximate  $\log(n_t, t)$ :

$$\log(n_t, t) = \frac{1}{\sqrt{\alpha_t}} n_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_0(n_t, t) \quad -③$$

Hence corresponding optimization problem becomes:

$$\arg \min_{\theta} D_{KL}(q(n_{t-1} | n_t, n_0) || p_\theta(n_{t-1} | n_t))$$

$$= \arg \min_{\theta} D_{KL}\left(N(n_{t-1}; \log, \Sigma_q(t)) || N(n_{t-1}; \log, \Sigma_\theta(t))\right)$$

Using ②, ③

$$= \begin{bmatrix} \frac{1}{2 \sigma_q^2(t)} & \frac{(1-\alpha_t)^2}{(1-\bar{\alpha}_t)\alpha_t} \left[ \| \epsilon_0 - \hat{\epsilon}_0(n_t - t) \|_2^2 \right] \end{bmatrix}$$

25.

(23.)

Tweedie's formula:

$$E[\mu_z | z] = z + \sum_z \nabla_z \log p(z) \quad \textcircled{1}$$

$$q(n_t | n_0) = N(n_t; \sqrt{\alpha_t} n_0, (1-\alpha_t) I)$$

using ①

$$E[\mu_n | n_t] = n_t + (1-\alpha_t) \nabla_{n_t} \log p(n_t)$$

$$\Rightarrow \sqrt{\alpha_t} n_0 = n_t + (1-\alpha_t) \nabla \log p(n_t) \quad \textcircled{2}$$

Using  $n_0$  in ground truth denoising transition  $\log(n_t, n_0)$ :

$$\log(n_t, n_0) = \frac{\sqrt{\alpha_t} (1-\alpha_{t-1}) n_t + \sqrt{\alpha_{t-1}} (1-\alpha_t) n_0}{1-\alpha_t}$$

using ②

$$\Rightarrow \log(n_t, n_0) = \frac{1}{\sqrt{\alpha_t}} n_t + \underbrace{\frac{1-\alpha_t}{\sqrt{\alpha_t}} \cdot \nabla \log p(n_t)}_{\text{using ③}} \quad \textcircled{3}$$

using ③

$$\mu_\theta(n_t, t) = \frac{1}{\sqrt{\alpha_t}} n_t + \frac{1-\alpha_t}{\sqrt{\alpha_t}} \underbrace{s_\theta(n_t, t)}_{\text{using ④}} \quad \textcircled{4}$$

Using ②

$$\Rightarrow h_g(n_t, n_0) = \frac{1}{\sqrt{\alpha_t}} n_t + \frac{1-\alpha_t}{\sqrt{\alpha_t}} \cdot \underbrace{\nabla \log p(n_t)}_{\text{using ③}} \quad - ③$$

using ③

$$h_o(n_t, t) = \frac{1}{\sqrt{\alpha_t}} n_t + \frac{1-\alpha_t}{\sqrt{\alpha_t}} \underbrace{s_o(n_t, t)}_{\text{using ④}} \quad - ④$$

Corresponding optm. problem

$$\underset{\theta}{\operatorname{argmin}} D_{KL}(q(n_{t+1} | n_t, n_0) || p_\theta(n_{t+1} | n_t))$$

$$= \underset{\theta}{\operatorname{argmin}} D_{KL}\left(N(n_{t+1}, \mu_q, \Sigma_q(t)) || N(n_{t+1}; \mu_\theta, \Sigma_\theta(t))\right)$$

using ③, ④

$$\underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{\alpha_t} \left[ \|s_o(n_t, t) - \nabla \log p(n_t)\|_2^2 \right]$$

$s_o(n_t, t) \rightarrow NN$  that predicts score func

26.

(26)

Score funcn  $\nabla_{n_t} \log p(n_t)$  looks similar to source noise  $\epsilon_0$ .

From reparametr.

$$n_0 = \frac{n_t - \sqrt{1-\bar{\alpha}_t} \epsilon_0}{\sqrt{\bar{\alpha}_t}} \quad \begin{pmatrix} \text{recal reparam.} \\ \sigma(n_t/n_0) \end{pmatrix} - (1)$$

. From tweedie's formula

$$n_0 = \frac{n_t + (1-\hat{\alpha}_t) \nabla \log p(n_t)}{\sqrt{\hat{\alpha}_t}} - (2)$$

$$\underline{(1) = (2)}$$

$$\nabla \log p(n_t) = \frac{-1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_0 \leftarrow \begin{matrix} \text{original} \\ \text{source noise} \end{matrix}$$

27.

(27)

### Implicit Score Matching (Hyvärinen, 2005)

$$\underset{\theta}{\operatorname{argmin}} \underset{n \sim p}{E} \left[ \frac{1}{2} \| S_{\theta}(n) - \nabla_n \log p(n) \|_2^2 \right] = \underset{\theta}{\operatorname{argmin}} \underset{n \sim p}{E} \left[ \operatorname{tr}(\nabla_n S_{\theta}(n)) + \frac{1}{2} \| S_{\theta}(n) \|_2^2 \right]$$

Proof:

$$\begin{aligned} & \underset{\theta}{\operatorname{argmin}} \underset{n \sim p}{E} \left[ \frac{1}{2} \| S_{\theta}(n) - \nabla_n \log p(n) \|_2^2 \right] \\ &= \underset{\theta}{\operatorname{argmin}} \underset{n \sim p}{E} \left[ \frac{1}{2} \| S_{\theta}(n) \|^2 - 2 \langle S_{\theta}(n), \nabla_n \log p(n) \rangle + \underbrace{\| \nabla_n \log p(n) \|_2^2}_{\text{independent of } \theta} \right] \end{aligned}$$

Just need to show inner product term is equivalent to  $\operatorname{tr}(\nabla_n S_{\theta}(n))$

$$\begin{aligned} & E_{n \sim p} [ S_{\theta}(n)^T \nabla_n \log p(n) ] = \sum_{i=1}^d \int_n S_{\theta}(n)_i \frac{\delta \log p(n)}{\delta n_i} p(n) dn \\ &= - \sum_{i=1}^d \int_n \frac{S_{\theta}(n)_i}{\delta n_i} p(n) dn \\ &= - \int_n \operatorname{tr}(\nabla_n S_{\theta}(n)) p(n) dn \\ &= - E_{n \sim p} [\operatorname{tr}(\nabla_n S_{\theta}(n))] \end{aligned}$$

M.P

28.

(28)

### Conditional Diffusion Models

$$p(n_{0:T} | \underline{y}) = p(n_T) \prod_{i=1}^T p_\theta(n_{t-1} | n_t, \underline{y}_t)$$

### Classifier Guidance

explicitly control amount of weight to cond<sup>n</sup>.

### Baye's rule

$$\begin{aligned}\nabla_{n_t} \log p(n_t | y) &= \nabla_{n_t} \log \left( \frac{p(n_t) p(y | n_t)}{p(y)} \right) \\ &= \underbrace{\nabla_{n_t} \log p(n_t)}_{\text{unconditional score}} + \underbrace{\nabla_{n_t} \log p(y | n_t)}_{\text{adversarial gradient}} - \nabla_{n_t} \log p(y)\end{aligned}$$

→ score of uncond<sup>n</sup> diff mode ( $\nabla_{n_t} \log p(n_t)$ ) is learned.

→ Alongside, classifier that takes arbitrary noise ' $n_t$ ' attempts to predict ' $y$ '.

Introduce fine-grained control:

$$\nabla \log p(n_t | y) = \nabla \log p(n_t) + \boxed{Y} \nabla \log p(y | n_t)$$

### Drawbacks

- Requires separately learnt Classifier
- Classifier must handle arbitrary noise inputs, which most pre-trained models are unable to do.

## Classifier-free Guidance

Conditioning without explicit classification

$$\nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

$$\Rightarrow \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

getting rid  
of classifier

$$\Rightarrow \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \underbrace{\gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})}_{\text{Conditional score}} + \underbrace{(1-\gamma) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{Unconditional score}}$$

- $\gamma \rightarrow$  controls learned cond<sup>n</sup> model
- $\gamma=1 \rightarrow$  explicitly learns vanill cond<sup>n</sup> distr.
- $\gamma > 1 \rightarrow$  prioritizes canon score, & also moves away from uncond<sup>n</sup> score.

## 29.

Autoregressive Flows	Normalizing Flows
1) AR model the joint distribution using chain of conditional distribution	NF model probability distribution by applying a series of invertible transforms.
2) Are computationally efficient, but may suffer from slow convergence and are limited to modelling low-dim data.	NF are flexible in modelling high-dim data, but can be computationally expensive due to the need to evaluate the determinant of the Jacobian in each transformation.
3) Conditional distributions are typically modelled as Gaussians or discrete distributions.	In NF, transformations can be any invertible functions.

## 30.

### Linear Autoregressive Models

$\epsilon_t \sim \mathcal{N}(0, \sigma^2)$  (i.i.d), consider a scalar separation of observations  $x \in \mathbb{R}^{T \times 1}$  defined by the linear dynamics:

$$x_t = \rho(x_{t-1} - \mu) + \mu + \epsilon_t$$

$$x_t = \mu + \epsilon_o$$

These dynamics are referred to as autoregressive process of order 1, or AR(1)

[Ref: "A course in time series analysis"]

- Order - 1, refers to the fact that  $x_t$  is a function of only one previous observation  $x_{t-1}$
- Also known as Markov property

## 31.

### Wide-Sense Stationary Process

A time series  $X$  is called wide-sense stationary if it has stationary mean and

$$Cov(x_t, x_{t-k}) = Cov(x_s, x_{s-l}), \forall s, t, k \in N$$

Define covariance function,  $\gamma : N \rightarrow \mathbb{R}$

$$\gamma(k) = Cov(x_{k+1}, x_1)$$

## 32.

### Fully Visible Sigmoid Belief Network

[Neal, "Connectionist learning of belief networks", 1992]

Def.:

Let  $P \in \mathbb{R}^T$ , and define  $f_t : \mathbb{R}^{(t-1)} \rightarrow \mathbb{R}$  by:

$$f_t(x_{ct}) = \langle p_{ct}, x_{ct} \rangle$$

Then fully visible sigmoid belief network is given by:

$$x_t \sim \text{Bernoulli}(\text{sigmoid}(f_t(x_{ct})))$$

$$x_1 \sim \text{Bernoulli}(P_o)$$

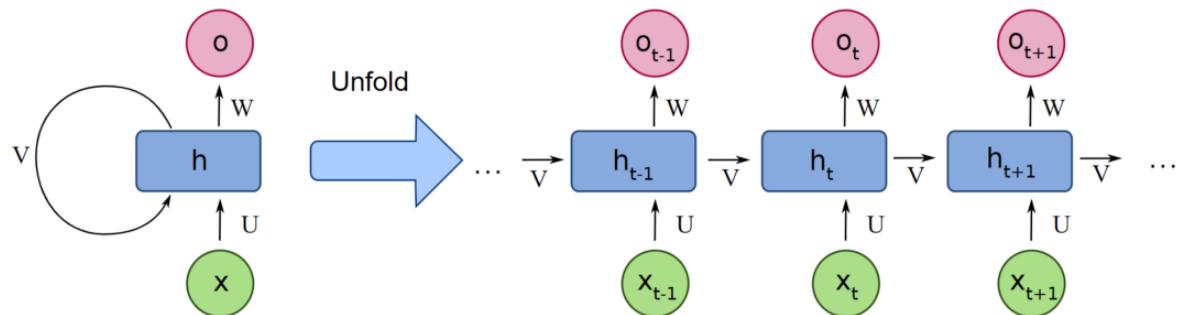
The goal is to predict  $x_t \in \{0, 1\}$  using a log-linear combination of  $x_{ct}$

### Neural Autoregressive Distribution Estimation

Neural Autoregressive Distribution Estimation, on the other hand, is a type of autoregressive model where the joint probability distribution of the data is modeled as a product of conditional probabilities, one for each dimension of the data. In NADE, the conditional probability of each dimension is modeled using a feedforward neural network that takes as input the previous dimensions and outputs the conditional probability of the current dimension.

## 33.

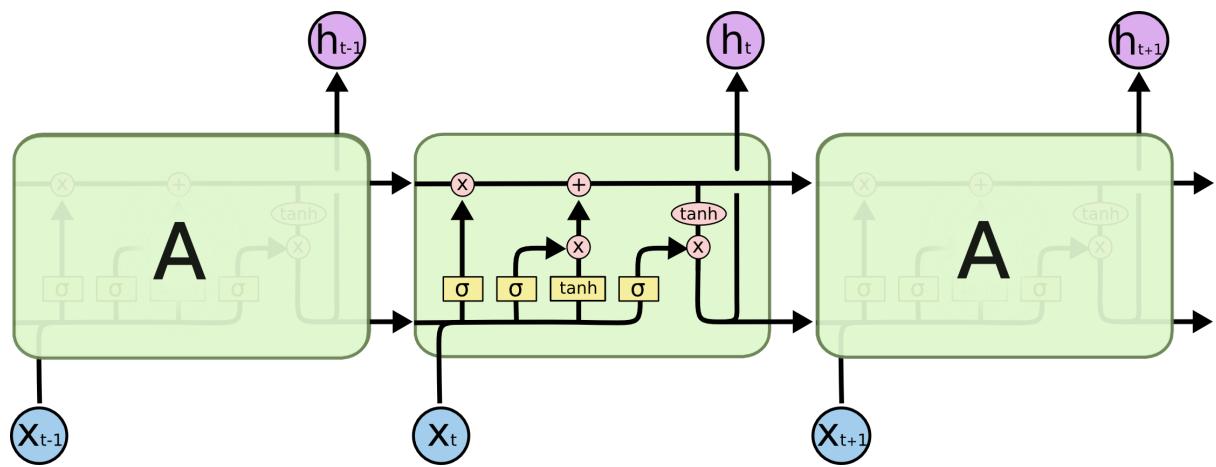
### RNN



- RNNs are a type of neural network that can model sequential data by maintaining an internal state
- That is updated at each time step based on the current input and the previous hidden state.

- The output of an RNN at each time step is a function of the current hidden state and the input at that time step.
- The RNN is trained using backpropagation through time
- RNNs can suffer from the vanishing gradient problem, which occurs when the gradients become very small as they propagate backwards through time. This can lead to difficulties in learning long-term dependencies and can limit the performance of the network.

## LSTM



- LSTM networks are a type of RNN that use memory cells and gates to control the flow of information in the network.
- The memory cell is a long-term memory unit that can store information over many time steps, and the gates control the information flow into and out of the memory cell.
- The gates consist of sigmoidal functions that output values between 0 and 1, which are used to decide which information to keep and which to discard.
- LSTM networks have been shown to be effective in modeling long-term dependencies in sequential data and have been used in various applications such as speech recognition, natural language processing, and image captioning.
- LSTM networks are trained using backpropagation through time with gradient clipping to prevent exploding gradients.

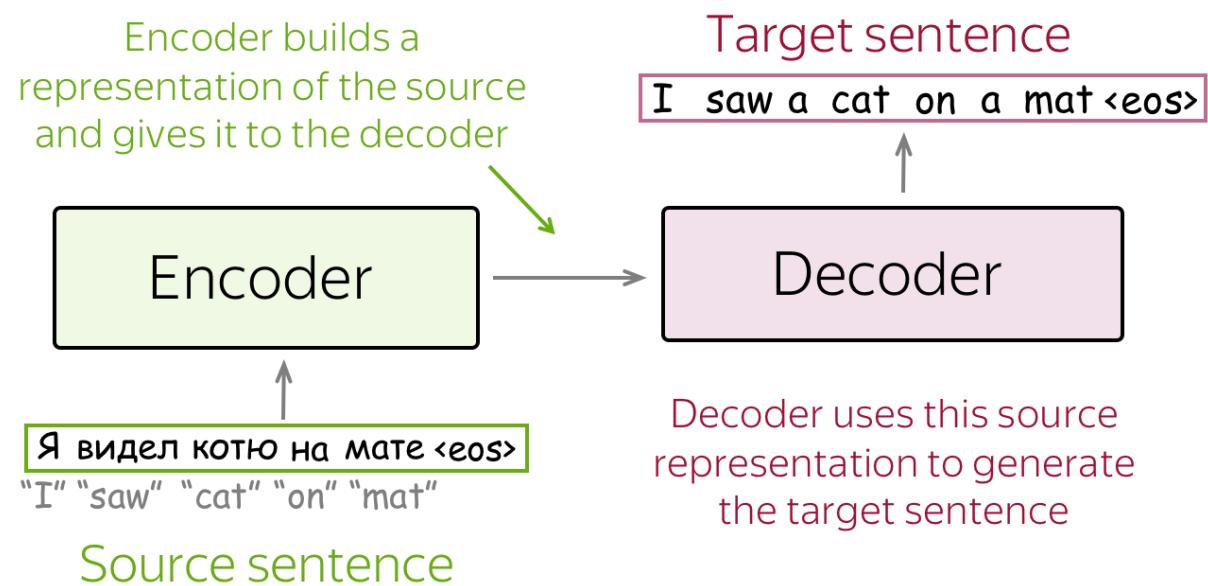
## 34.

### Advantage of LSTM over RNN

- The main advantage of LSTM (Long Short-Term Memory) networks over traditional RNN (Recurrent Neural Network) models is their ability to effectively capture long-term dependencies in sequential data.
- In traditional RNNs, the hidden state at each time step is a function of the current input and the previous hidden state. However, as the sequence length increases, the gradients can become very small or even vanish
- LSTM networks address this issue by introducing a memory cell and several gates that control the flow of information into and out of the cell. The gates are designed to selectively forget or remember information based on the current input and the previous hidden state

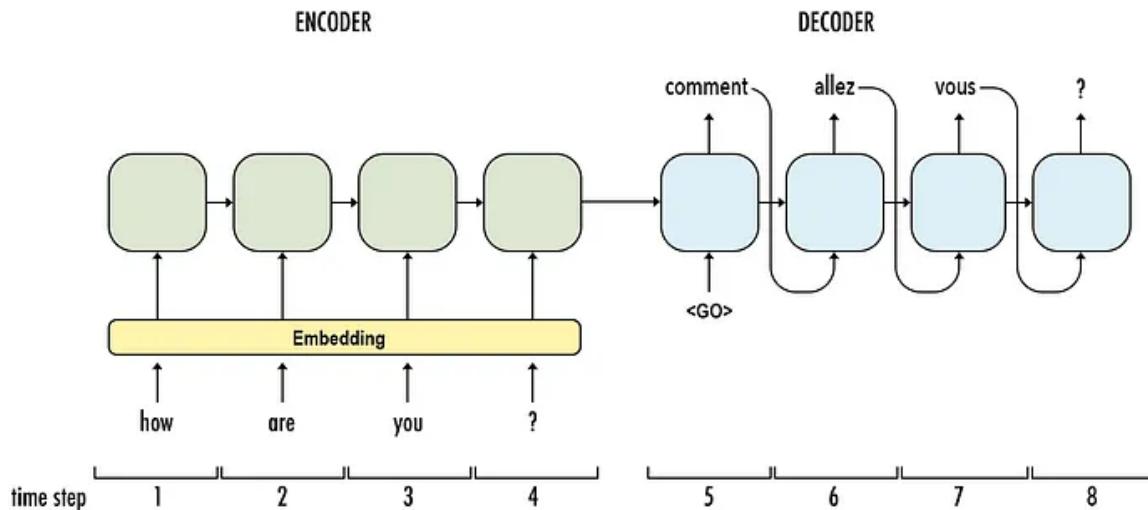
## 35.

### Seq2Seq Model



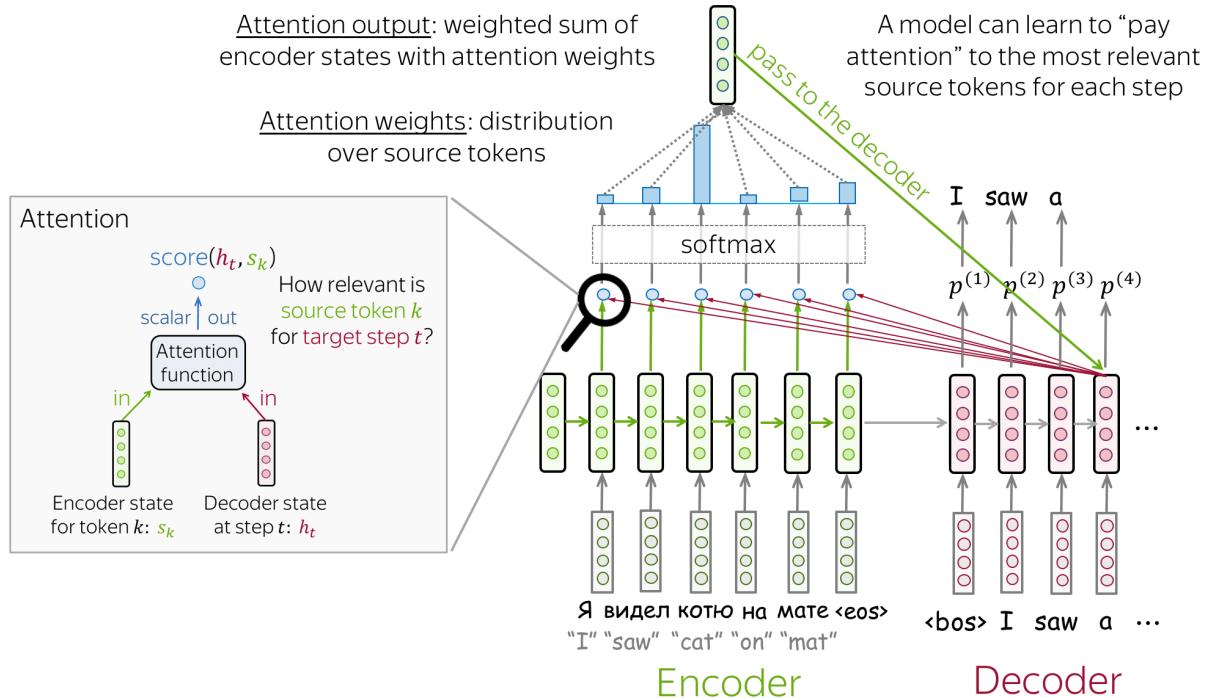
- Sequence-to-sequence (seq2seq) models are a type of neural network architecture used for tasks that involve transforming an input sequence into an output sequence, such as machine translation, text summarization, or speech recognition.

- Seq2seq models consist of an encoder and a decoder
- The encoder maps the input sequence to a fixed-length vector representation
- The decoder generates the output sequence from this representation.



## Seq2Seq Model + Attention

- Seq2seq models with attention, also known as attentional seq2seq models, extend the basic seq2seq architecture by adding an attention mechanism to the decoder.
- The attention mechanism allows the decoder to selectively focus on different parts of the input sequence, depending on the current output generated by the decoder.
- The attention mechanism works by computing a set of attention weights for each input sequence element, which indicate how much attention the decoder should pay to that element when generating the current output.
- The attention weights are computed as a function of the decoder's current hidden state and the encoder's output at each time step.



- The computation of attention weights can be done in different ways, but a common approach is to use a dot product between the decoder's current hidden state and the encoder's output at each time step
- Followed by a softmax operation to obtain a probability distribution over the input sequence.
- The weighted sum of the encoder's output, using the attention weights as weights, is then concatenated with the decoder's input at the current time step, and fed into the decoder.

Dot-product	Bilinear	Multi-Layer Perceptron
$\begin{bmatrix} h_t^T \\ \vdots \\ h_t^T \end{bmatrix} \times \begin{bmatrix} s_k \\ \vdots \\ s_k \end{bmatrix}$	$\begin{bmatrix} h_t^T \\ \vdots \\ h_t^T \end{bmatrix} \times \begin{bmatrix} W \\ \vdots \\ W \end{bmatrix} \times \begin{bmatrix} s_k \\ \vdots \\ s_k \end{bmatrix}$	$\begin{bmatrix} w_2^T \\ \vdots \\ w_2^T \end{bmatrix} \times \tanh \left[ \begin{bmatrix} W_1 \\ \vdots \\ W_1 \end{bmatrix} \times \begin{bmatrix} h_t \\ \vdots \\ h_t \end{bmatrix} \right] \times \begin{bmatrix} s_k \\ \vdots \\ s_k \end{bmatrix}$
$\text{score}(h_t, s_k) = h_t^T s_k$	$\text{score}(h_t, s_k) = h_t^T W s_k$	$\text{score}(h_t, s_k) = w_2^T \cdot \tanh(W_1[h_t, s_k])$

36.

**Tokenisation** is the task of splitting the text into tokens which are then converted to numbers. These numbers are in turn used by the machine learning models for further processing and training.

**Subword tokenization** algorithms rely on the principle that frequently used words should not be split into smaller subwords, but rare words should be decomposed into meaningful subwords.

For instance "annoyingly" might be considered a rare word and could be decomposed into "annoying" and "ly". Both "annoying" and "ly" as stand-alone subwords would appear more frequently while at the same time the meaning of "annoyingly" is kept by the composite meaning of "annoying" and "ly".

This is especially useful in agglutinative languages such as Turkish, where you can form (almost) arbitrarily long complex words by stringing together subwords.

Subword tokenization allows the model to have a reasonable vocabulary size while being able to learn meaningful context-independent representations. In addition, subword tokenization enables the model to process words it has never seen before, by decomposing them into known subwords.

## 37.

# Transformer Architecture

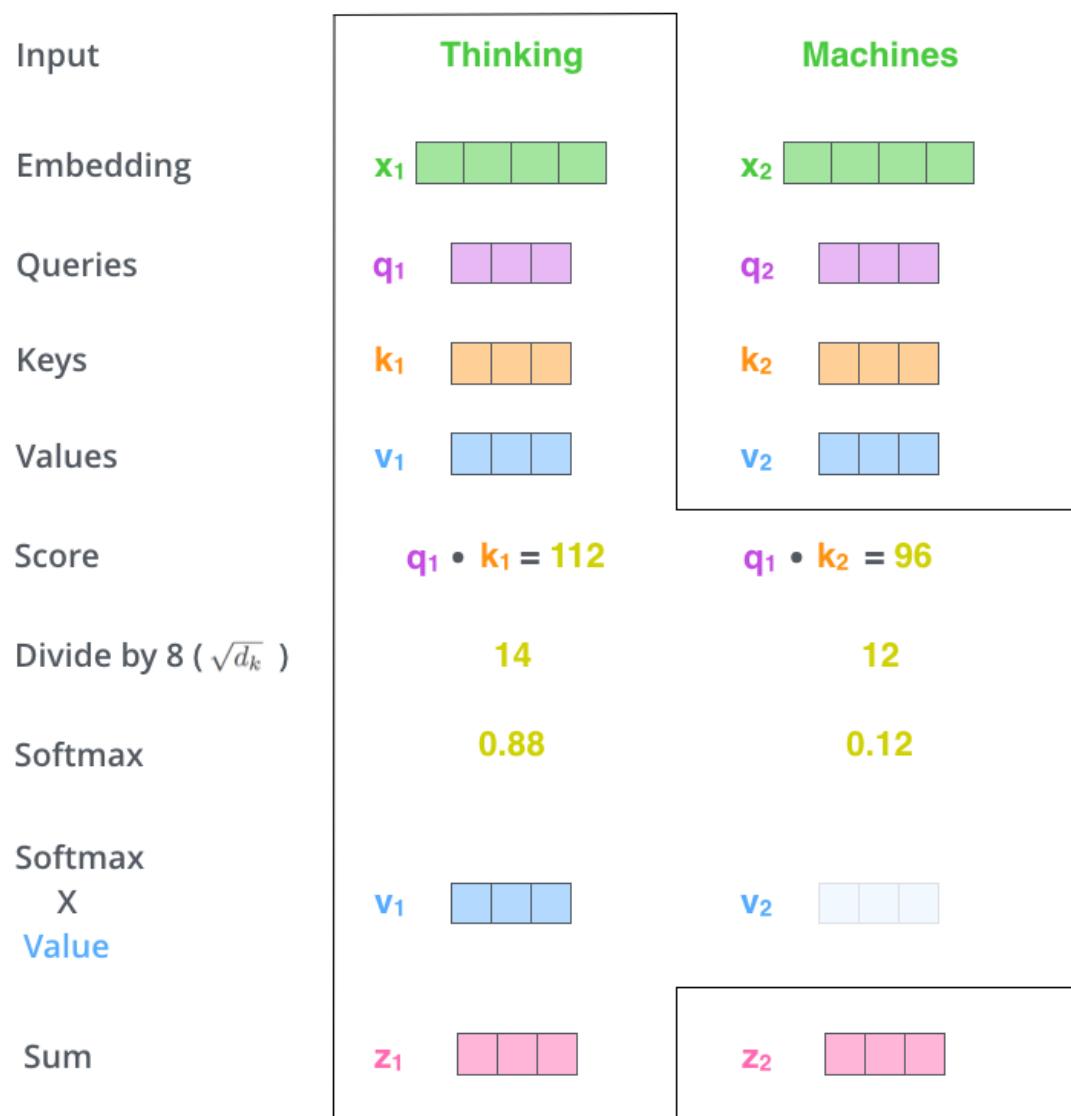
## Scaled Self-Attention

- The **first step** in calculating self-attention is to create three vectors from each of the encoder's input vectors (in this case, the embedding of each word). So for each word, we create a Query vector, a Key vector, and a Value vector. These vectors are created by multiplying the embedding by three matrices that we trained during the training process.
- The **second step** in calculating self-attention is to calculate a score. Say we're calculating the self-attention for the first word in this example, "Thinking". We need to score each word of the input sentence against this word. The score determines how much focus to place on other parts of the input sentence as we

encode a word at a certain position.

The score is calculated by taking the dot product of the query vector with the key vector of the respective word we're scoring. So if we're processing the self-attention for the word in position #1, the first score would be the dot product of  $q_1$  and  $k_1$ . The second score would be the dot product of  $q_1$  and  $k_2$ .

- The **third and fourth steps** are to divide the scores by 8 (the square root of the dimension of the key vectors used in the paper – 64. This leads to having more stable gradients. There could be other possible values here, but this is the default), then pass the result through a softmax operation. Softmax normalizes the scores so they're all positive and add up to 1.
- The **fifth step** is to multiply each value vector by the softmax score (in preparation to sum them up). The intuition here is to keep intact the values of the word(s) we want to focus on, and drown-out irrelevant words (by multiplying them by tiny numbers like 0.001, for example).
- The **sixth step** is to sum up the weighted value vectors. This produces the output of the self-attention layer at this position (for the first word).



## Matrix Form

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

The diagram illustrates the multiplication of matrix  $\mathbf{X}$  by matrix  $\mathbf{W}^Q$  to produce matrix  $\mathbf{Q}$ . Matrix  $\mathbf{X}$  is a 2x4 green grid. Matrix  $\mathbf{W}^Q$  is a 4x4 purple grid. The result, matrix  $\mathbf{Q}$ , is a 2x2 purple grid.

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

The diagram illustrates the multiplication of matrix  $\mathbf{X}$  by matrix  $\mathbf{W}^K$  to produce matrix  $\mathbf{K}$ . Matrix  $\mathbf{X}$  is a 2x4 green grid. Matrix  $\mathbf{W}^K$  is a 4x4 orange grid. The result, matrix  $\mathbf{K}$ , is a 2x2 orange grid.

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

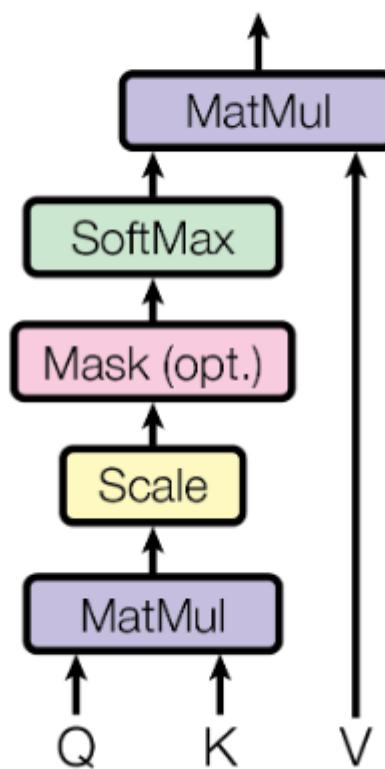
The diagram illustrates the multiplication of matrix  $\mathbf{X}$  by matrix  $\mathbf{W}^V$  to produce matrix  $\mathbf{V}$ . Matrix  $\mathbf{X}$  is a 2x4 green grid. Matrix  $\mathbf{W}^V$  is a 4x4 blue grid. The result, matrix  $\mathbf{V}$ , is a 2x2 blue grid.

$$\text{softmax} \left( \frac{\begin{matrix} \mathbf{Q} & \times & \mathbf{K}^T \\ \begin{matrix} \text{pink grid} \end{matrix} & \times & \begin{matrix} \text{orange grid} \end{matrix} \end{matrix}}{\sqrt{d_k}} \right) \mathbf{V}$$

$$= \begin{matrix} \mathbf{Z} \\ \text{pink grid} \end{matrix}$$

## Architecture

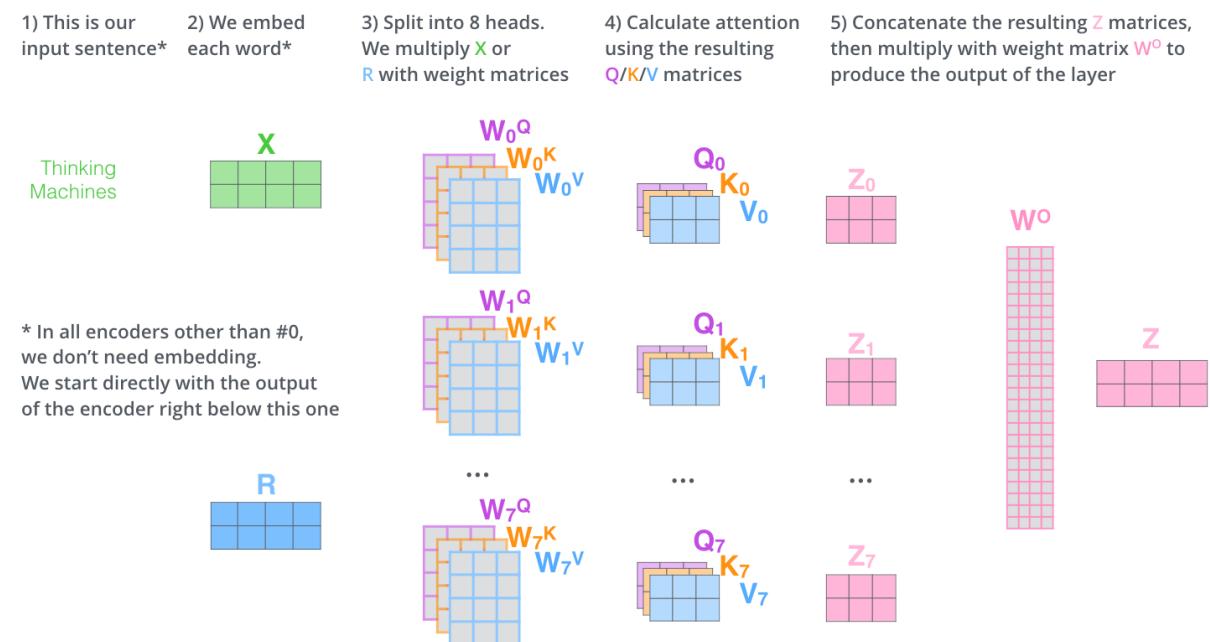
### Scaled Dot-Product Attention



## Multi-Head Attention

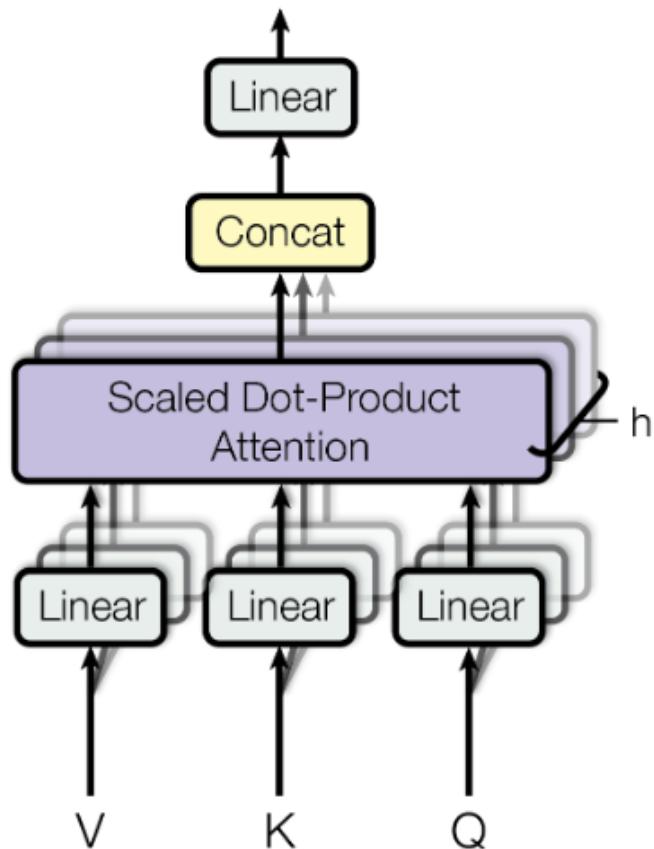
The paper further refined the self-attention layer by adding a mechanism called “multi-headed” attention. This improves the performance of the attention layer in two ways:

1. It expands the model’s ability to focus on different positions. Yes, in the example above,  $z_1$  contains a little bit of every other encoding, but it could be dominated by the actual word itself. If we’re translating a sentence like “The animal didn’t cross the street because it was too tired”, it would be useful to know which word “it” refers to.
2. It gives the attention layer multiple “representation subspaces”. As we’ll see next, with multi-headed attention we have not only one, but multiple sets of Query/Key/Value weight matrices (the Transformer uses eight attention heads, so we end up with eight sets for each encoder/decoder). Each of these sets is randomly initialized. Then, after training, each set is used to project the input embeddings (or vectors from lower encoders/decoders) into a different representation subspace.



## Architecture

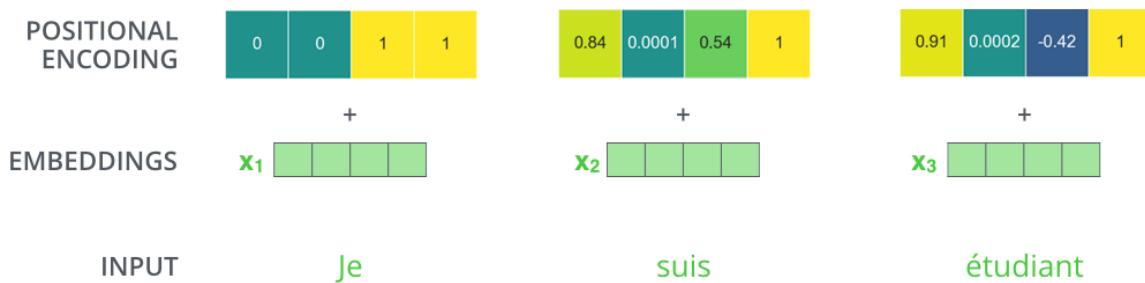
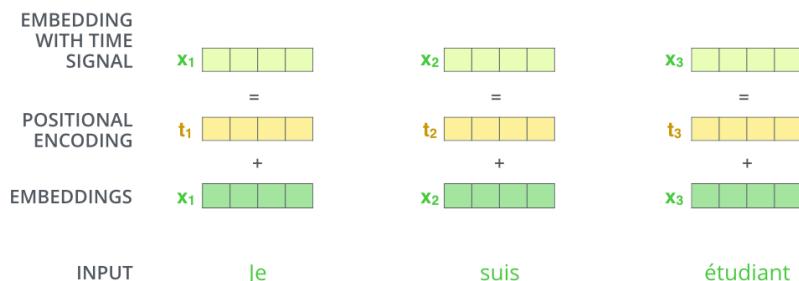
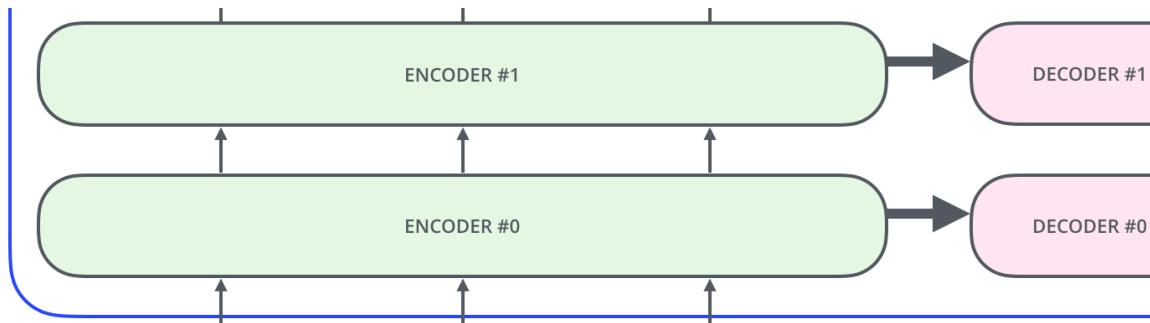
## Multi-Head Attention



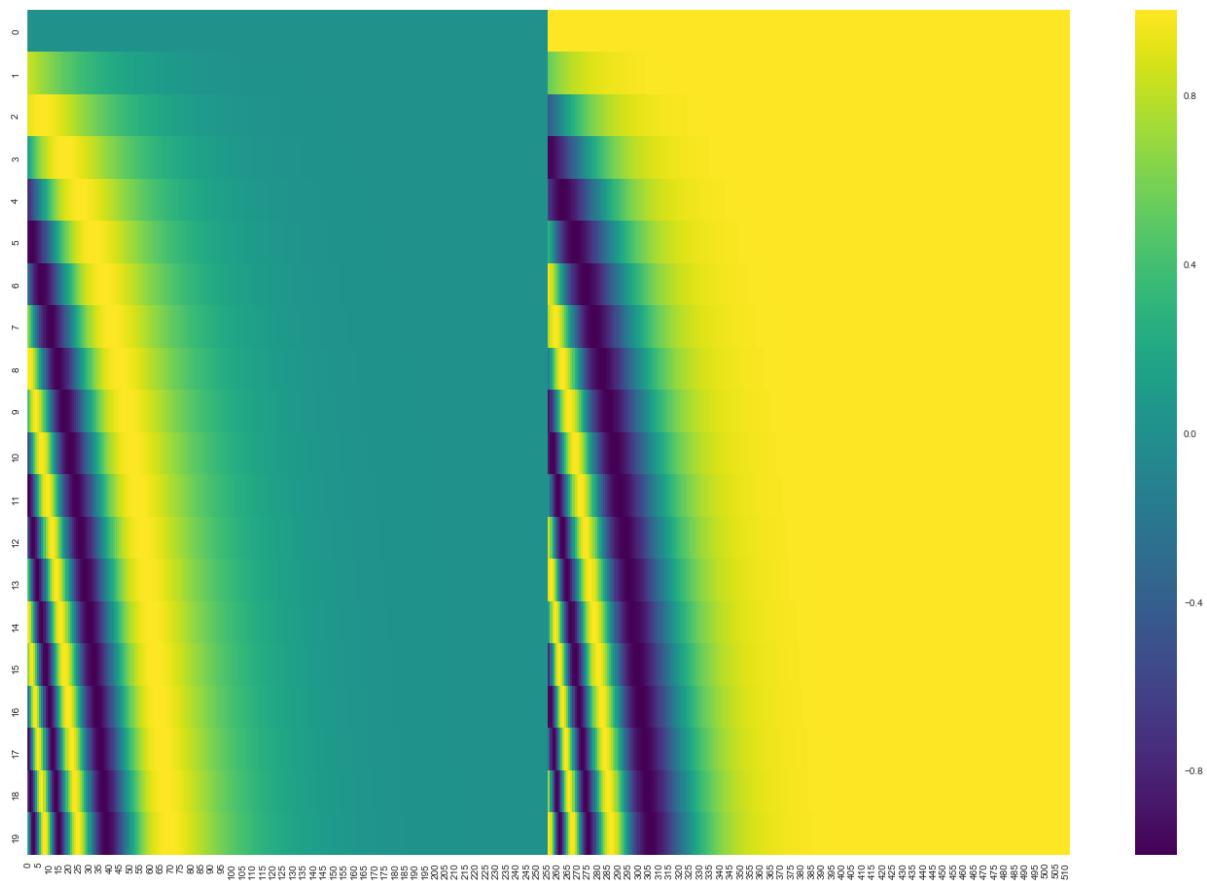
## Positional Encoding

One thing that's missing from the model as we have described it so far is a way to account for the order of the words in the input sequence.

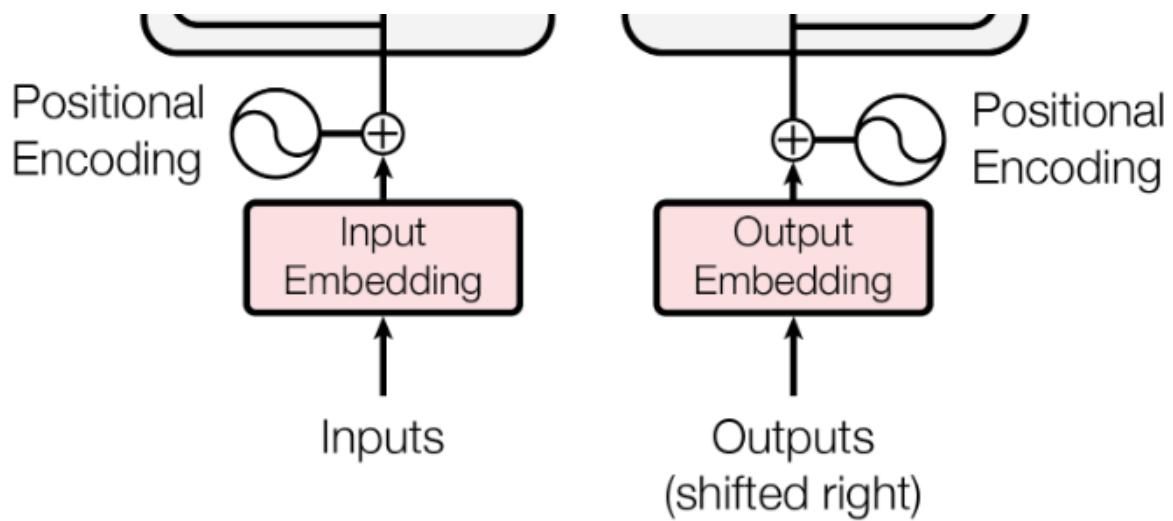
To address this, the transformer adds a vector to each input embedding. These vectors follow a specific pattern that the model learns, which helps it determine the position of each word, or the distance between different words in the sequence. The intuition here is that adding these values to the embeddings provides meaningful distances between the embedding vectors once they're projected into Q/K/V vectors and during dot-product attention.



In the following figure, each row corresponds to a positional encoding of a vector. So the first row would be the vector we'd add to the embedding of the first word in an input sequence. Each row contains 512 values – each with a value between 1 and -1. We've color-coded them so the pattern is visible.



## In Architecture



## Layer-Normalization and Skip Connections

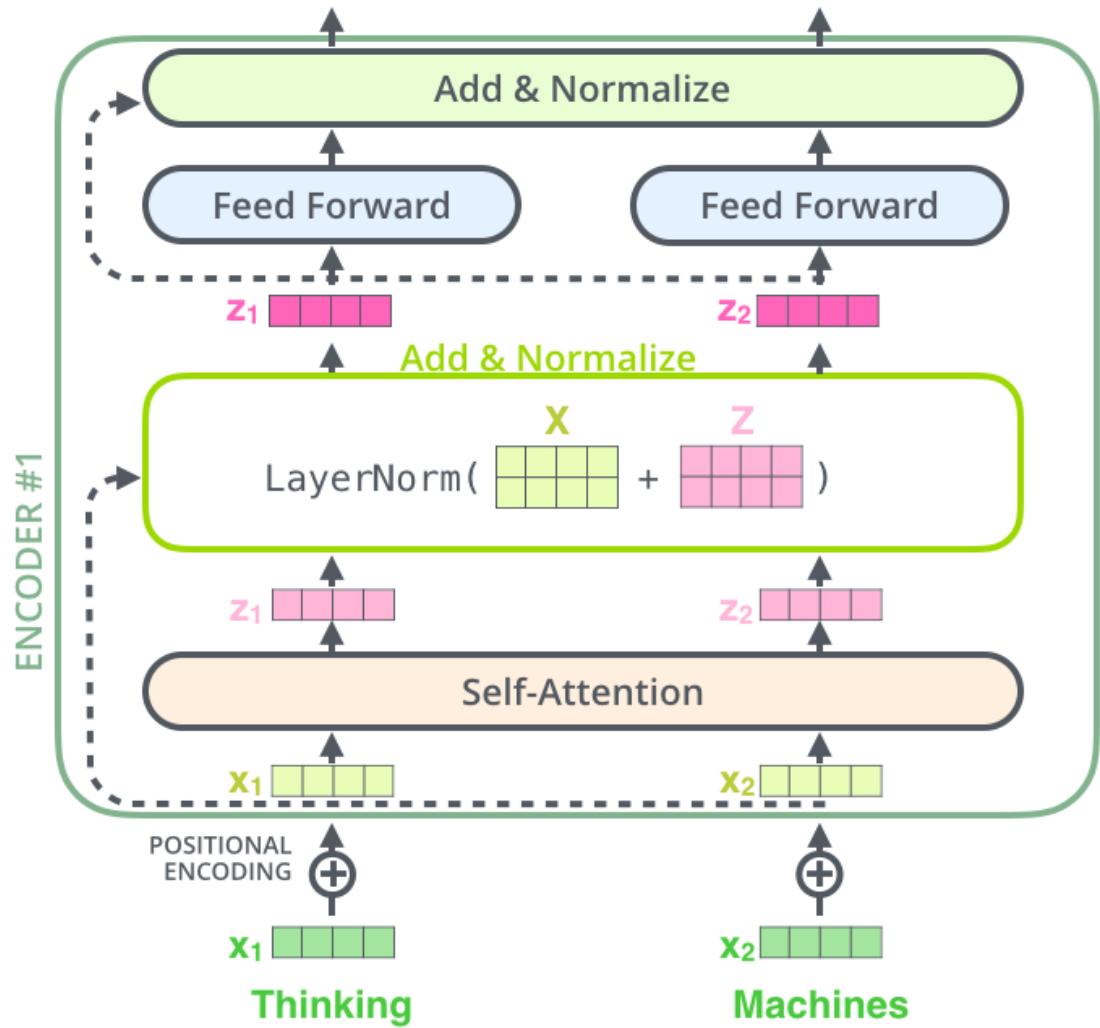
One detail in the architecture of the encoder that we need to mention before moving on, is that each sub-layer (self-attention, ffnn) in each encoder has a residual connection around it, and is followed by a layer-normalization step.

Residual Connection → Help regain context

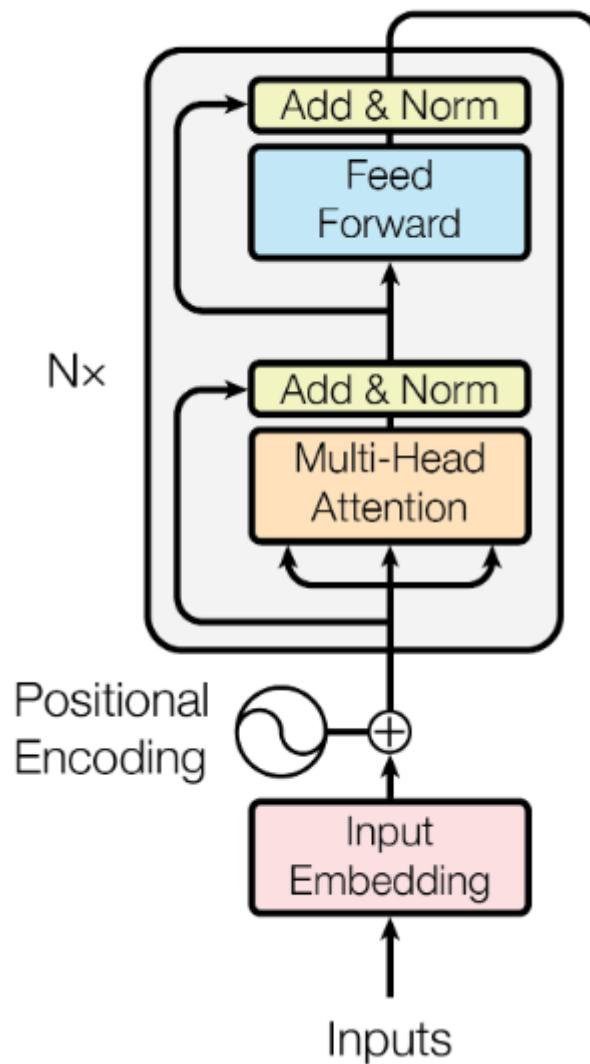
A skip connection is added to allow the output of the block to bypass the feedforward layer and be added to the input of the next block, which can help to mitigate the vanishing gradient problem and improve the training of the model.

Layer-Normalization → Form of Regularization

Layer normalization is a normalization technique used to reduce the internal covariate shift in neural networks, which can help to improve the stability and performance of the model.



## Encoder Architecture



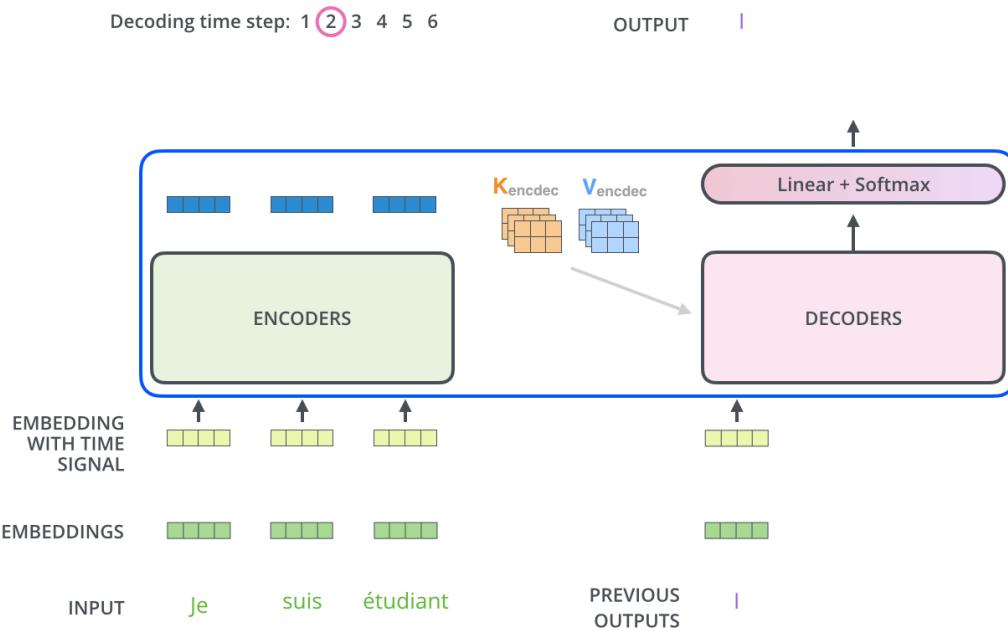
## Decoder

The encoder starts by processing the input sequence. The output of the top encoder is then transformed into a set of attention vectors  $K$  and  $V$ .

These are to be used by each decoder in its “encoder-decoder attention” layer which helps the decoder focus on appropriate places in the input sequence:

The following steps repeat the process until a special symbol is reached indicating the transformer decoder has completed its output. The output of each step is fed to the bottom decoder in the next time step, and the decoders bubble up their decoding results just like the encoders did. And just like we did with the encoder inputs, we

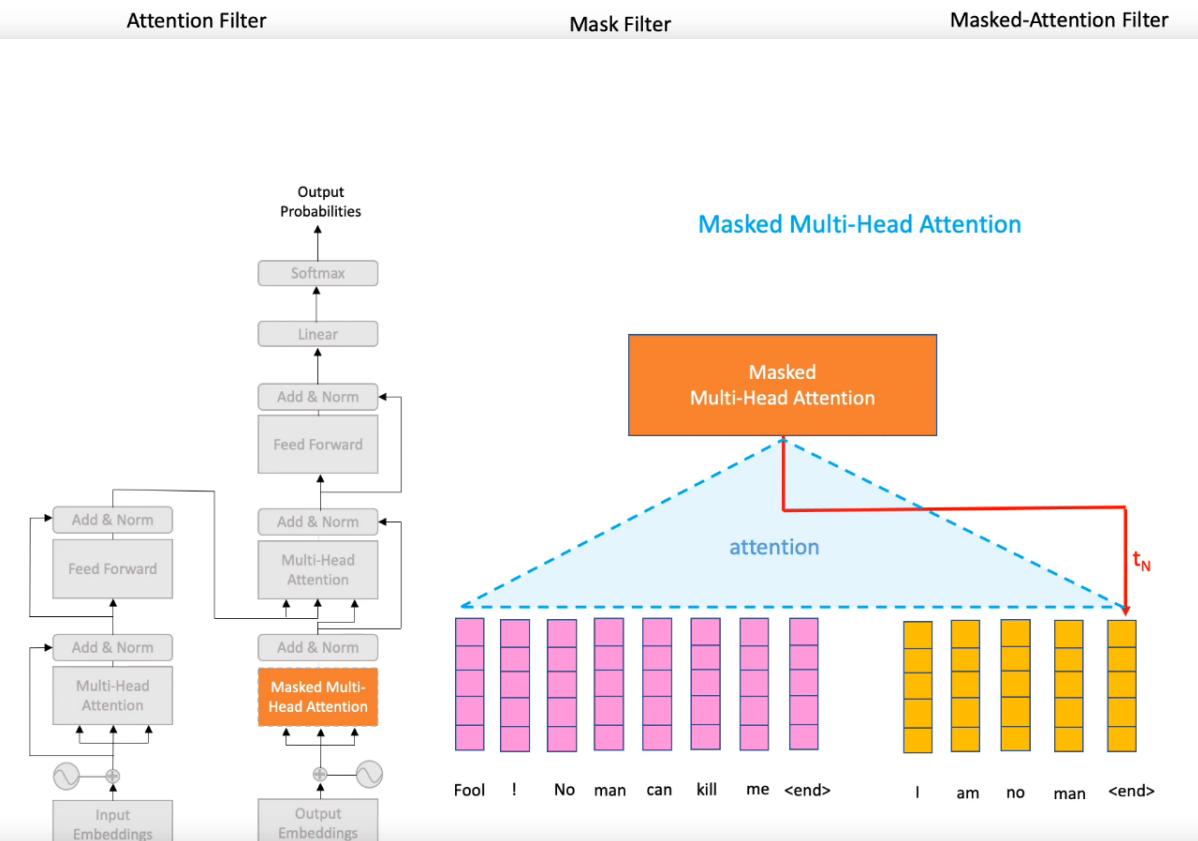
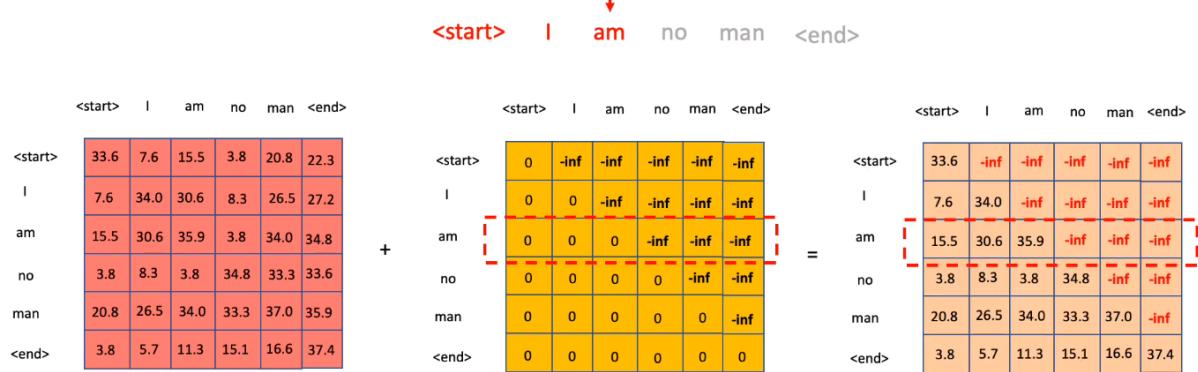
embed and add positional encoding to those decoder inputs to indicate the position of each word.



### Masked-self Attention

The input to the decoder block is the output of the previous decoder block or the encoded input sequence, depending on the specific architecture used. The masked self-attention mechanism allows the decoder to only attend to previous tokens in the output sequence, while masking out future tokens that have not yet been generated.

# Masking



## Cross-multihead Attention

The output of the masked self-attention mechanism is then passed through a cross multihead attention layer, which allows the decoder to attend to the encoded input sequence. This layer enables the model to generate contextually appropriate outputs based on the input sequence.

## Architecture

The decoder block also includes positional encoding, layer normalization, and skip connections, similar to the encoder block.

