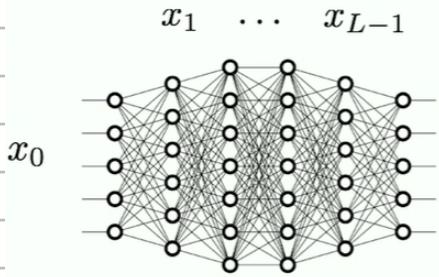


Resnet



$$x_{k+1} = x_k + \underbrace{g(W_k x_k + b_k)}_{\text{Nonlinear function}} \quad k = 1, 2, \dots, L-1$$

Affine function
 $f(x) = ax + b$

continuous time.
 $\dot{x}(t) = g(W_t x(t) + b_t)$
 direction change

$NN = \text{Map } x_0 \rightarrow x_L$
 $NN = \text{flow map } x(0) \rightarrow x(T)$

Transformer

transformer input $(x_1, x_2, \dots, x_n) \in \mathbb{R}^d$: tokens have position encoding.

$$(x_1, x_2, \dots, x_n) \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \delta x_i$$

prompt = input = probability measure on tokens \Rightarrow likelihood of token being next

transformer output: predict the next token

output = probability measure on tokens \Rightarrow empirical dist. of tokens in prompt

transformer = flow map $f_0: \mu(0) \rightarrow \mu(T)$

$$\mu(\mathbb{R}^d) \rightarrow \mu(\mathbb{R}^d)$$

$\mu(0) = \frac{1}{n} \sum_{i=1}^n \delta x_i \Rightarrow \delta: \text{dot indicator, indicates that there is a token at this coordinate}$

Transformer is a "Mean-Field interacting particle system"

$$\dot{x}_i(t) = X_t[\mu(t)](x_i(t)) \quad i=1, 2, \dots, n$$

↑ taken
 ↓ rate of change mean-field: depends on all of the other tokens' aggregate distribution

Continuity Equation

Divergence: $\begin{cases} \text{div} < 0: \text{gather} \\ \text{div} > 0: \text{diffusion} \end{cases}$

$$\partial_t \mu(t) + \text{div}(\mu(t) X_t[\mu(t)]) = 0$$

density. movement speed, depends on $\mu(t)$

self-attention dynamics

$$x_i(t) = X_t[\mu(t)](x_i(t))$$

Self-attention dynamics = special choice of $X_t[\mu(t)](\cdot)$

$$X_t[\mu(t)](x) = V_t \frac{\int e^{\langle Q_t x, K_t y \rangle} \mu(t)(dy)}{\int e^{\langle Q_t x, K_t y \rangle} \mu(t)(dy)} = V_t \mathbb{E}_{y|x} [Y_t]$$

V_t, Q_t, K_t : $d \times d$ matrices learned during training.

$Q_t(\text{Query})$: What I want eg. "I" wants a verb

$K_t(\text{Key})$: my characteristics eg. "I" is a pronoun

$V_t(\text{Value})$: information package eg. "river": flowing, moist, blue, ...

$\langle \cdot, \cdot \rangle$: inner product, indicating similarity

$\langle Q_t x, K_t y \rangle$ means that the more the features of y (key) match what x wants ($Q_t x$), the larger $\langle Q_t x, K_t y \rangle$ is

$$X_t[\mu(t)](x) = V_t \frac{\int e^{\langle Q_t x, K_t y \rangle} \mu(t)(dy)}{\int e^{\langle Q_t x, K_t y \rangle} \mu(t)(dy)} = V_t \mathbb{E}_{y|x} [Y_t]$$

attention-weighted density

Since $\mu(t) = \frac{1}{n} \sum_{i=1}^n \delta x_i(t)$

$$X_t[\mu(t)](x_i(t)) = V_t \left(\sum_{j=1}^n p_{ij}(t) x_j(t) \right)$$

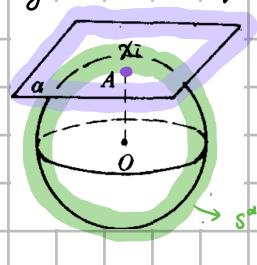
$$p_{ij}(t) = \frac{e^{\langle Q_t x_i(t), K_t x_j(t) \rangle}}{\sum_{k=1}^n e^{\langle Q_t x_i(t), K_t x_k(t) \rangle}}$$

Layer Normalization

Layer Norm (prevents explosion)

Each self-attention involves a weighted sum of a large number of vectors, causing the magnitude of the embedding vectors to increase (gradient exploding) or decrease (gradient vanishing) sharply between layers.

Dynamics on the sphere



\rightarrow tangent plane of x_i

$$x_2(t) = \text{Proj}_{T_{x_i}(t)S^{d-1}} V_t \left(\sum_{j=1}^n p_{ij}(t) x_j(t) \right)$$

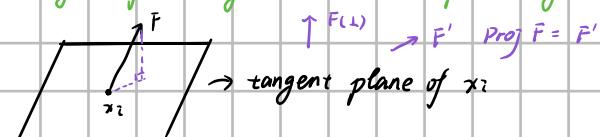
correction operation

original self-attention speed (F)

tangent plane: $T_{x_i} S^{d-1}$

S^{d-1} : region representing vectors with a fixed magnitude

Proj:



Bells And Whistles

MLP layer in-between two self-attention layers

$$x_2(t) = \text{Proj}_{T_{x_i}(t)S^{d-1}} \text{MLP}_t V_t \left(\sum_{j=1}^n p_{ij}(t) x_j(t) \right)$$

Multiple heads

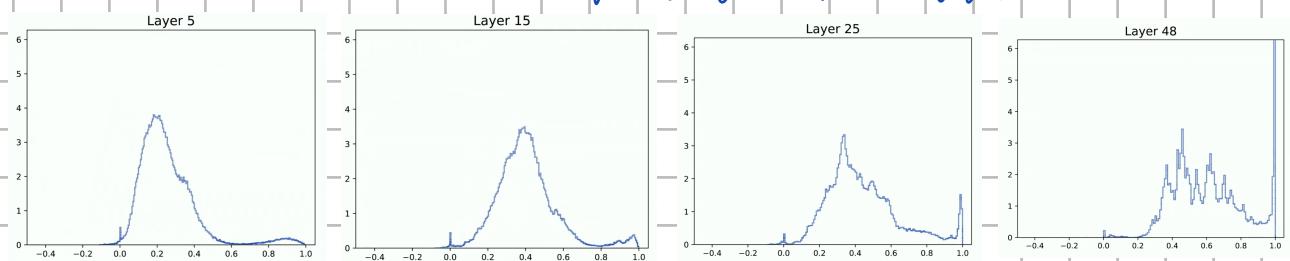
$$x_2(t) = \frac{1}{K} \sum_{k=1}^K \text{Proj}_{T_{x_i}(t)S^{d-1}} \text{MLP}_t^{[k]} V_t^{[k]} \left(\sum_{j=1}^n p_{ij}^{[k]}(t) x_j(t) \right)$$

Long - Time Asymptotics

$$x_2(t) = \text{Proj}_{T_{x_i}(t)S^{d-1}} V_t \left(\sum_{j=1}^n p_{ij}(t) x_j(t) \right) \quad \text{when } t \rightarrow \infty \quad \mu(t) \rightarrow \text{clustering}$$

e.g. $\mu(t)$ can be confusing: "bank" may be in the financial and geographical zone.

After multiple layers ($t+T$), all tokens representing 'finance' converge into a financial cluster, while all tokens representing 'geography' converge into a geographical cluster



A Simple Attention Model

$$x_2(t) = \text{Proj}_{T_{x_i}(t)S^{d-1}} V_t \left(\sum_{j=1}^n p_{ij}(t) x_j(t) \right)$$

$$p_{ij}(t) = \frac{e^{x_i(t)^T Q_t K_t x_j(t)}}{\sum_{k=1}^n e^{x_i(t)^T Q_t K_t x_k(t)}}$$

\Rightarrow simplify: $V_t = \text{Id}$ $Q_t^T K_t = \beta \text{Id}$

$$x_2(t) = \text{Proj}_{T_{x_i}(t)S^{d-1}} \frac{\sum_{j=1}^n e^{\beta \langle x_i(t), x_j(t) \rangle} x_j(t)}{\sum_{k=1}^n e^{\beta \langle x_i(t), x_k(t) \rangle}}$$

$$x_2(t) = \text{Proj}_{T_{x_i}(t)S^{d-1}} \frac{\sum_{j=1}^n e^{\beta \langle x_i(t), x_j(t) \rangle} x_j(t)}{\sum_{k=1}^n e^{\beta \langle x_i(t), x_k(t) \rangle}}$$

$$F(\mu) = \frac{1}{2} \sum_{i,j} e^{-\frac{\beta}{2} \|x_i(t) - x_j(t)\|^2}$$

it is a gradient flow of the repulsive interaction energy

Gradient Flow

self - attention : semantic appeal

"apple" must be close to "fruit"

energy F : geometric repulsion

"apple" and "fruit" cannot be too close, otherwise the energy of the whole system will explore.

if the system is only affected by attraction, all tokens will gather to a point, causing all semantic information to overlap and the model to lose distinguishability.

Due to the presence of the repulsion energy F , transformers finally reach a state

of equilibrium (clustering) not a chaotic state, but rather a layout that maximizes geometric spacing on the sphere

What Happened To The Data



$\frac{d}{dt} \ln \log M_t = \Theta(t)$ → rate of linear change: representing a quantity that is linearly related to the rate of change
 free energy: representing the extent from chaos to order

Each layer of the transformer is transforming chaotic input data into an organized semantic structure at a linear, controllable pace

Schematic Diagram - from "Attention is all you need"

