

# **The Mathematics Behind LLM Transformers**

Carles Casacuberta

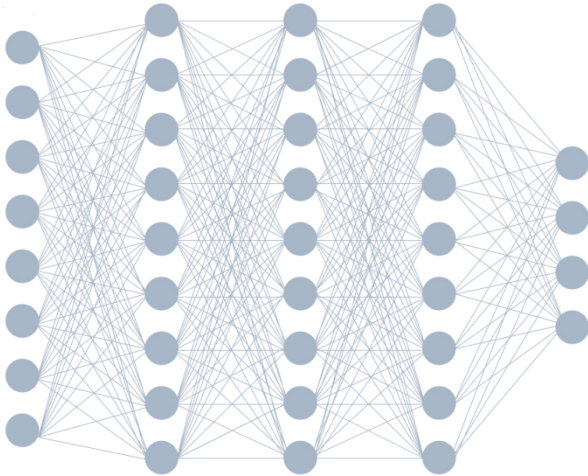
Universitat de Barcelona

**Seminari Cuc**

**Proof Theory and Foundations of Mathematics**

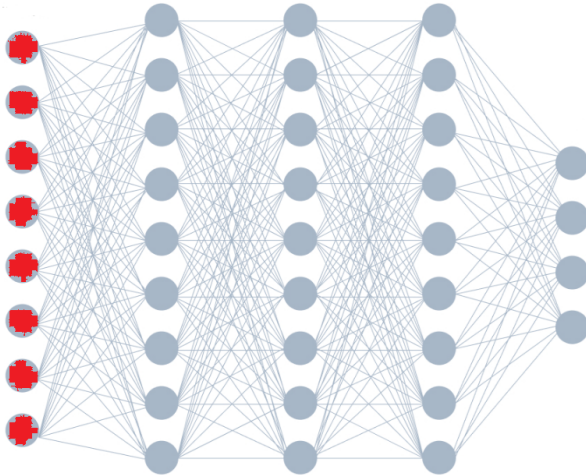
8 May 2025

# Multilayer Perceptron



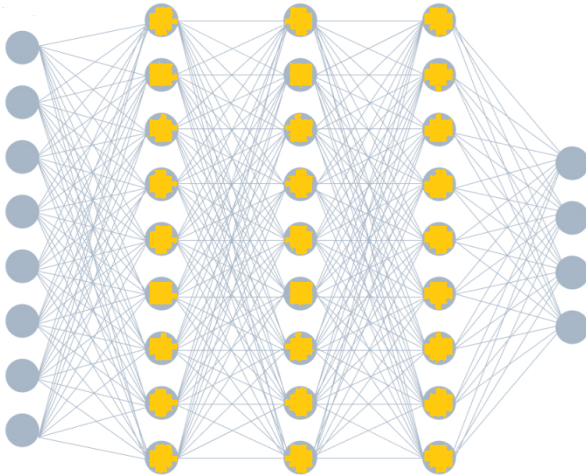
A fully connected **feed-forward** neural network

# Multilayer Perceptron



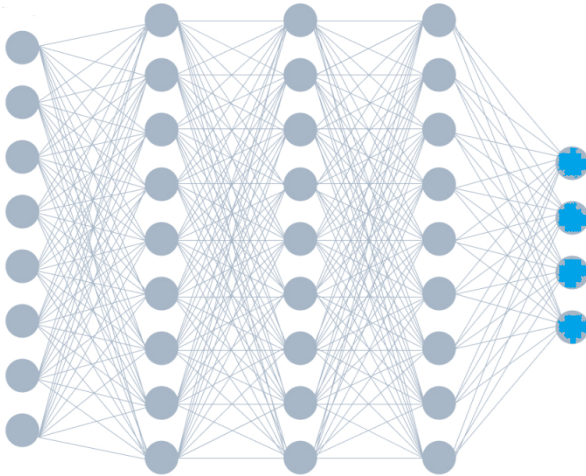
**Input layer**

# Multilayer Perceptron



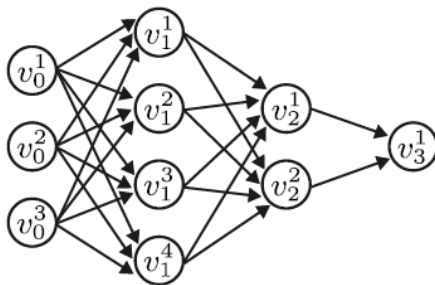
**Hidden layers**

# Multilayer Perceptron



**Output layer**

# Neurons



9

$v_0^1, \dots, v_0^{N_0}$ : Neurons of the input layer.

For  $\ell = 1, \dots, L$  (number of hidden layers plus output layer),

$v_\ell^1, \dots, v_\ell^{N_\ell}$ : Neurons of the  $\ell$ -th layer.

# Activation

For  $x \in \mathbb{R}^{N_0}$  and  $\ell = 0, \dots, L$ , compute **hidden states**  $h_\ell(x) \in \mathbb{R}^{N_\ell}$  recursively as  $h_0(x) = x$  and

$$h_\ell(x) = \sigma(W_\ell h_{\ell-1}(x) + b_\ell)$$

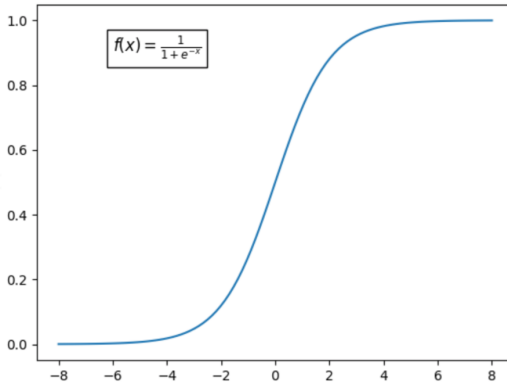
where  $\sigma$  is a nonlinear, almost everywhere smooth **activation function** applied entrywise,  $W_\ell$  is a matrix of **weights**, and  $b_\ell$  is a vector of **biases**. The set  $\theta$  of entries of  $W_\ell$  and  $b_\ell$  is the set of **learnable parameters** of the network.

For an ordered collection  $x_1, \dots, x_n$  of **input data**, the vector

$$(h_\ell(x_1)^k, \dots, h_\ell(x_n)^k)$$

is the **activation vector** of the neuron  $v_\ell^k$  at  $x_1, \dots, x_n$ .

# Activation



The **sigmoid** is a common activation function



# Loss Function

In a **supervised task**, one considers a function  $\phi: X \rightarrow Y$ , where  $X$  is a set of **inputs** and  $Y$  is a set of **labels**, e.g.,  $Y = \mathbb{R}$  for a **regression** task, or  $Y$  finite for a **classification** task.

In practice,

- ▶ a **training set**  $S_{\text{train}} \subseteq X \times Y$  and
- ▶ a **test set**  $S_{\text{test}} \subseteq X \times Y$  are chosen.

If  $\mathcal{L}: Y \times Y \rightarrow \mathbb{R}$  is a suitable **loss function** and a subset  $S = \{(x_i, y_i)\}_{i=1}^n \subseteq X \times Y$  is given, the **empirical risk** is

$$R_S = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(h_L(x_i), y_i, \theta)$$

where  $h_L(x_i)$  are **predicted** values and  $y_i$  are **true** values.

# Generalization Gap

After learning with **backpropagation** the parameter set  $\theta$  by minimizing the empirical risk  $R_{\text{train}}$ , the **generalization gap** is

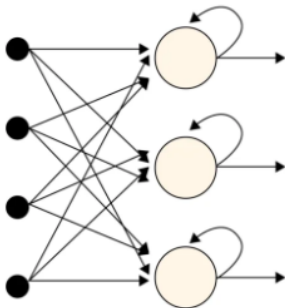
$$G = R_{\text{test}} - R_{\text{train}}$$

To prevent overfitting and reduce the generalization gap, a **regularization function**  $\mathcal{T}$  can be used as follows:

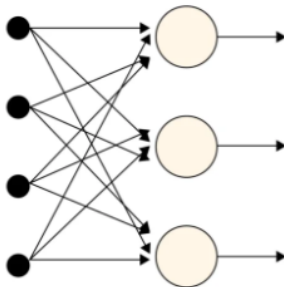
$$\tilde{R}_S = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(h_L(x_i) + \mathcal{T}(x_i), y_i, \theta).$$

# Recurrent Neural Networks

(a) Recurrent Neural Network



(b) Feed-Forward Neural Network



# Recurrent Neural Networks

For a recurrent neural network, input data (**tokens**) are **sequentially** ordered:  $w_1, \dots, w_T$ .

Each token  $w_t$  is converted into a vector  $x_t \in \mathbb{R}^d$  by means of a suitable **embedding** followed by **positional encoding**, which preserves the token order in the sequence  $x_1, \dots, x_T$ .

Hidden states for each input  $x_t$  are computed at layer  $\ell$  as

$$h_\ell(x_t) = \sigma \left( W_\ell^x h_{\ell-1}(x_t) + W_\ell^h h_\ell(x_{t-1}) + b_\ell \right)$$

where  $W_\ell^x$  and  $W_\ell^h$  are matrices of weights to be learned, and  $b_\ell$  is a bias vector, also to be learned.

# Recurrent Neural Networks

The **output** of a RNN is, for each  $t$ ,

$$y_t = W^{\text{out}} h_L(x_t) + b^{\text{out}}.$$

Here  $W^{\text{out}}$  is a matrix of dimensions  $N \times d$ , where  $N$  is the size of the vocabulary. Outputs are converted with the **softmax** function:

$$\hat{y}_t[i] = \text{softmax}(y_t)[i] = \frac{\exp(y_t[i])}{\sum_{j=1}^N \exp(y_t[j])}, \quad i = 1, \dots, N.$$

The model is trained using **cross-entropy** as loss function:

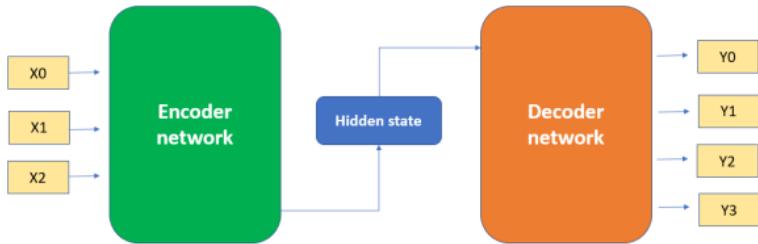
$$\mathcal{L} = -\log \hat{y}_t[y_t],$$

where  $\hat{y}_t[w]$  is the probability assigned by  $\hat{y}_t$  to a token  $w$ .

**Main difficulty:** **Vanishing gradient** during backpropagation.

# Encoders

**Encoder-decoder** architectures are used for sequence learning in **generative** artificial intelligence.



# Self-Attention

**Transformers** rely on **self-attention** in both encoders and decoders, which replaces recurrence by enabling all tokens to interact with each other simultaneously.

**Self-attention** assigns to each state  $h_t \in \mathbb{R}^d$  (either an input or a hidden state from a previous step) three vectors: **queries** ( $q_t$ ), **keys** ( $k_t$ ), and **values** ( $v_t$ ), by means of linear projections:

$$q_t = W^Q h_t, \quad k_t = W^K h_t, \quad v_t = W^V h_t,$$

where  $W^Q$ ,  $W^K$  and  $W^V$  are  $d \times d$  learnable matrices, and  $d$  is called **model length**.

# Self-Attention

Then **alignment scores** are defined for each pair  $h_i, h_j$  as

$$e_{ij} = \frac{k_i \cdot q_j}{\sqrt{d}},$$

where  $d$  is the model length.

Alignment scores are normalized using the **softmax** function:

$$a_{ij} = \text{softmax}(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{t=1}^T \exp(e_{it})}.$$

Thus,  $a_{ij}$  can be interpreted as a probability that state  $h_i$  is relevant to state  $h_j$ .



# Alignment Scores

Normalized alignment scores are used to compute a weighted average of values:

$$C = (c_1, \dots, c_T), \quad c_i = \sum_{t=1}^T a_{it} v_t,$$

which is called the **attention matrix** of  $H = (h_1, \dots, h_T)$ :

$$C = \text{Att}(H) = \text{Att}(Q, K, V) = \text{softmax} \left( \frac{QK^t}{\sqrt{d}} V \right).$$

# Multi-head Attention

There is a **multi-head attention** layer in the transformer architecture. The multi-head attention layer works as **multiple parallel attention mechanisms**, called **heads**.

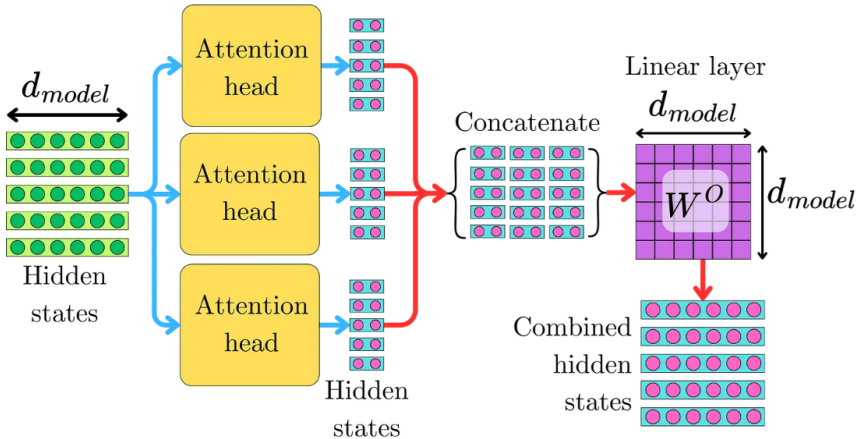
The number  $n$  of heads must divide the model length  $d$ .

Denote by  $H = (h_1, \dots, h_T)$  the incoming hidden states, and let

$$\text{Mult}(H) = \text{concat}(\text{head}_1, \dots, \text{head}_n) W^O, \quad \text{head}_i = \text{Att}_i(H),$$

where  $n$  is the number of heads and  $W^O$  is another learnable matrix. Here the matrices  $W_i^Q$ ,  $W_i^K$  and  $W_i^V$  used to compute  $\text{Att}_i$  for each  $i$  are of size  $d \times (d/n)$ , so  $W^O$  is again a  $d \times d$  matrix.

# Multi-head Attention



# Transformers

In a transformer, each **encoder** has two layers: the self-attention mechanism followed by a two-step feed-forward neural network:

$$h_t^{\text{enc}} = W_2 \sigma(W_1 h_t + b_1) + b_2.$$

Thus, an encoder receives a sequence of tokens  $w_1, \dots, w_T$ . Each token  $w_t$  is converted into a vector  $x_t \in \mathbb{R}^d$  using **positional encoding**. The vectors  $x_1, \dots, x_T$  pass through several **attention blocks**, and the encoder returns a matrix  $H^{\text{enc}}$  of size  $T \times d$ .

# Transformers

The **decoder** receives  $H^{\text{enc}}$  together with  $x_1, \dots, x_T$  during training (**teacher forcing**), or the previous predicted values during generation (**autoregression**). The decoder generates an output matrix  $H^{\text{dec}}$  by means of a **masked** self-attention layer, a **cross-attention** layer, and a feed-forward neural network.

**Masked** self-attention has

$$\text{Att}(\mathbf{q}_t, K, V) = \sum_{j=1}^t a_{tj} v_j;$$

hence, **future** tokens  $j > t$  are not included.

**Cross-attention** recovers information from  $H^{\text{enc}}$  as follows: the query matrix  $Q$  is computed as  $W^Q H^{\text{dec}}$  but the key and value matrices are computed as  $W^K H^{\text{enc}}$  and  $W^V H^{\text{enc}}$ .

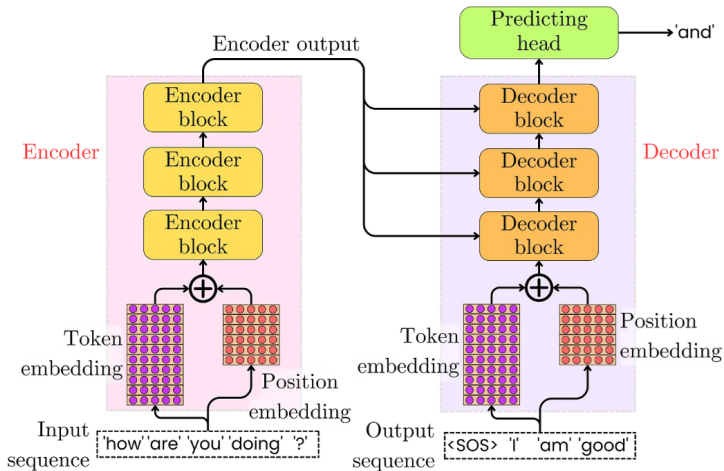
# Transformers

The transformer ends with a **predicting head**: A classifier over the whole token vocabulary made out of a linear layer followed by softmax, converting the decoder's final hidden states into **token probabilities** to predict the next word in a sequence:

$$\hat{y}_t = \text{softmax}(W^{\text{out}}h_t^{\text{dec}} + b).$$

Here  $h^{\text{dec}}$  is the decoder's output and  $\hat{y}_t \in \mathbb{R}^N$  are probabilities for each token, where  $N$  is the length of the vocabulary. The size of the matrix  $W^{\text{out}}$  is  $N \times d$ .

# Transformers



# Transformers

