

# **TTIC 31230, Fundamentals of Deep Learning**

David McAllester, Fall 2023

**Language Modeling**

**Recurrent Neural Networks**

**Machine Translation**

**The Transformer**

# Language Modeling

The recent progress on NLP benchmarks is due to pretraining on language modeling.

Language modeling is based on unconditional cross-entropy minimization.

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim P_{\text{op}}} [-\ln P_{\Phi}(y)]$$

In language modeling  $y$  is a sentence (or fixed length block of text).

## Language Modeling

Let  $V$  be some finite vocabulary of tokens.

Each token is a character sequence that can be used as a part of a rare word such a name in a foreign language. Most English words are a single token.

Tokens typically do not cross word boundaries.

We are interested in probability distributions over  $V^*$  (the finite sequences of tokens).

# Language Modeling

Let  $\text{Pop}$  be a population distribution over sequences of tokens.

We want to train a model  $P_{\Phi}(y)$  for token sequences  $y$

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} E_{y \sim \text{Pop}} [-\ln P_{\Phi}(y)]$$

A structured object, such as a token sequence or an image, has an exponentially small probability.

## Autoregressive Models

An autoregressive model uses the chain rule to represent a distribution on sequences in terms of the conditional probability for each token given the earlier tokens.

$$P_{\Phi}(w_0, w_1, \dots, w_T) = \prod_{t=0}^T P_{\Phi}(w_t \mid w_0, \dots, w_{t-1})$$

Modern language models are actually defining probability distributions on long sequences (thousands of tokens).

## The end-of-sequence token

We want to define a probability distribution over sentences of different length.

For this we require that each sentence is “terminated” with an end of sequence token **<EOS>**.

# Training

$$P_{\Phi}(w_0, w_1, \dots, w_T) = \prod_{t=0}^T P_{\Phi}(w_t \mid w_0, \dots, w_{t-1})$$

For training we need to compute the log loss on a training example.

The log loss on a sequence is the sum of the log losses for each token generation.

## Generation

$$P_{\Phi}(w_0, w_1, \dots, w_T) = \prod_{t=0}^T P_{\Phi}(w_t \mid w_0, \dots, w_{t-1})$$

We can generate from an autoregressive language model by generating one word at a time.

To generate we sample from a probability distribution over the first word. Once this word is generated we compute a probability distribution for the second word and so on.

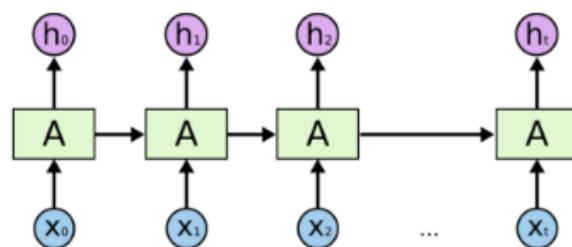


# Word Embeddings

Each word  $w$  is associated with a vector  $e(w)[I]$  called the embedding of word  $w$ .

The matrix  $e$  can be viewed as a dictionary assigning each word  $w$  the vector  $e(w)[I]$ .

# Recurrent Neural Network (RNN) Language Modeling

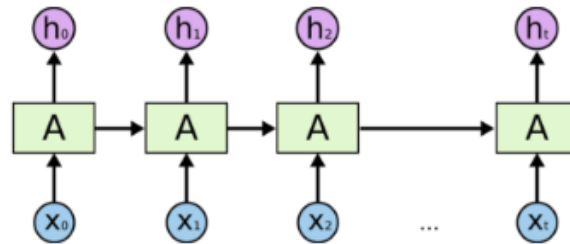


[Christopher Olah]

A typical RNN neural language model has the form

$$P_{\Phi}(w_t \mid w_0, \dots, w_{t-1}) = \underset{w_t}{\text{softmax}} e(w_t)[I]h[t-1, I]$$

# Vanilla RNNs



[Christopher Olah]

A Vanilla RNN uses two-input linear threshold units.

$$h[t, j] = \sigma (W^{h,h}[j, I]h[t - 1, I] + W^{x,h}[j, K]x[t, K] - B[j])$$

## Exploding and Vanishing Gradients

If we avoid saturation of the activation functions then we get exponentially growing or shrinking eigenvectors of the weight matrix.

Note that if the forward values are bounded by sigmoids or tanh then they cannot explode.

However the gradients can still explode.

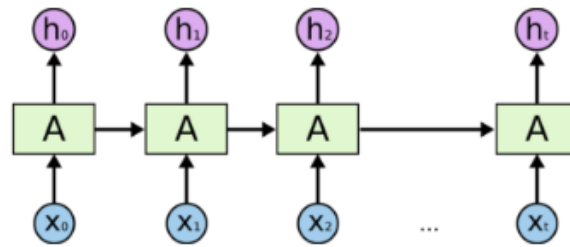
## Exploding Gradients: Gradient Clipping

We can dampen the effect of exploding gradients by clipping them before applying SGD.

$$W.\text{grad}' = \begin{cases} W.\text{grad} & \text{if } \|W.\text{grad}\| \leq n_{\max} \\ n_{\max} W.\text{grad} / \|W.\text{grad}\| & \text{otherwise} \end{cases}$$

See `torch.nn.utils.clip_grad_norm`

## Time as Depth



[Christopher Olah]

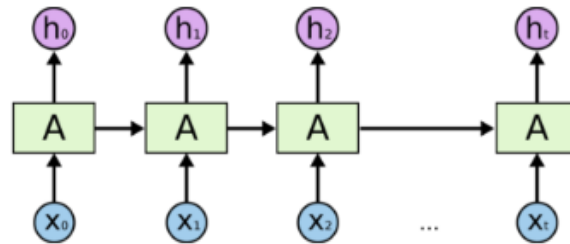
We would like the RNN to **remember and use** information from much earlier inputs.

All the issues with depth now occur through time.

However, for RNNs **at each time step we use the same model parameters**.

In CNNs **at each layer uses its own model parameters**.

## “Residual Connections” Through Time

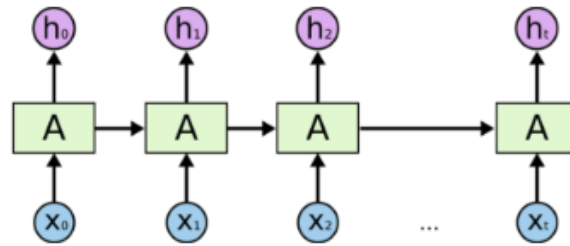


[Christopher Olah]

We would like to have residual connections through time.

However, we have to handle the fact that the same model parameters are used at every time step.

## Gated RNNs



[Christopher Olah]

$$h[t, j] = G_t[t, j]h[t-1, j] + (1 - G[t, j])R[t, j]$$

This is analogous to a residual connection.

Rather than add the “next layer”  $R[t, j]$  to the input  $h[t-1, j]$  as in a residual connection, we take a convex combination determined by a computed “gate”  $G[t, j] \in [0, 1]$ .



## Update Gate RNN (UGRNN)

$$R[t, j] = \tanh (W^{h,R}[j, I]h[t-1, I] + W^{x,R}[j, K]x[t, K] - B^R[j])$$

$$G[t, j] = \sigma (W^{h,G}[j, I]h[t-1, I] + W^{x,G}[j, K]x[t, K] - B^G[j])$$

$$h[t, j] = G[t, j]h[t-1, j] + (1 - G[t, j])R[t, j]$$

$$\Phi = (W^{h,R}, W^{x,R}, B^R, W^{h,G}, W^{x,G}, B^G)$$

$$\tanh(x) \in (-1, 1) \quad \sigma(x) \in (0, 1)$$

## Hadamard product

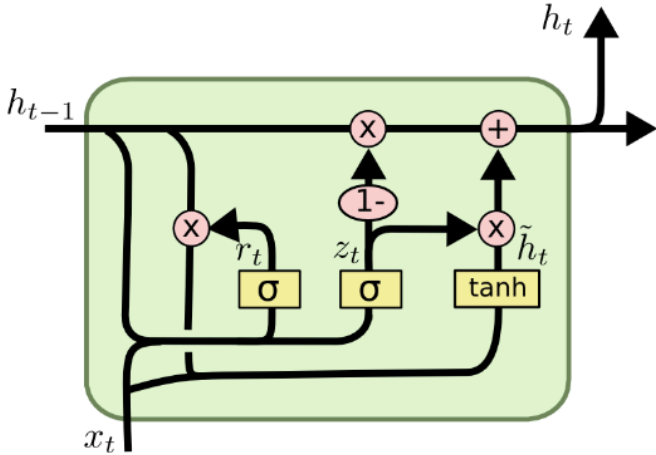
$$h[t, j] = G[t, j]h[t-1, j] + (1 - G[t, j])R[t, j]$$

is sometimes written as

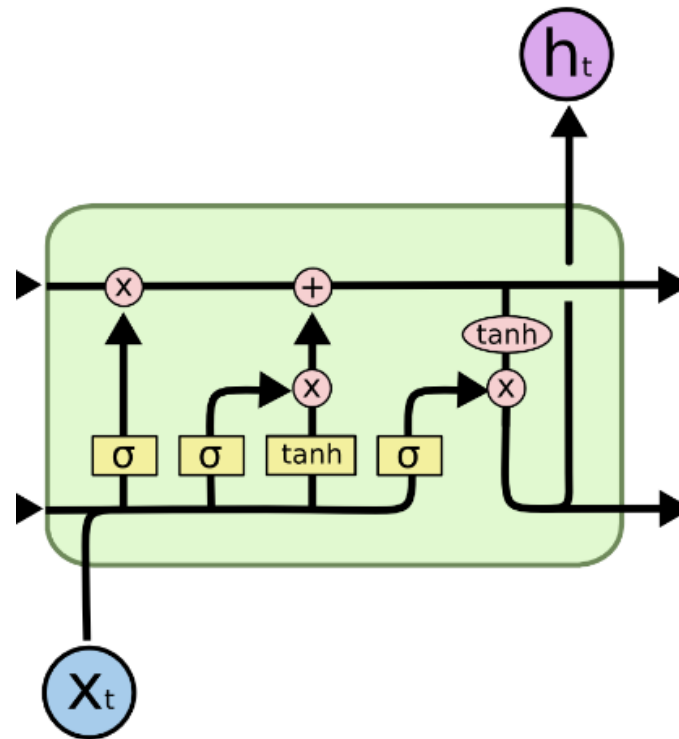
$$h[t, J] = G[t, J] \odot h[t-1, J] + (1 - G[t, J]) \odot R[t, J]$$

$\odot$  is the Hadamard product (componentwise product) on vectors.

# Gated Recurrent Unity (GRU) by Cho et al. 2014



# Long Short Term Memory (LSTM)

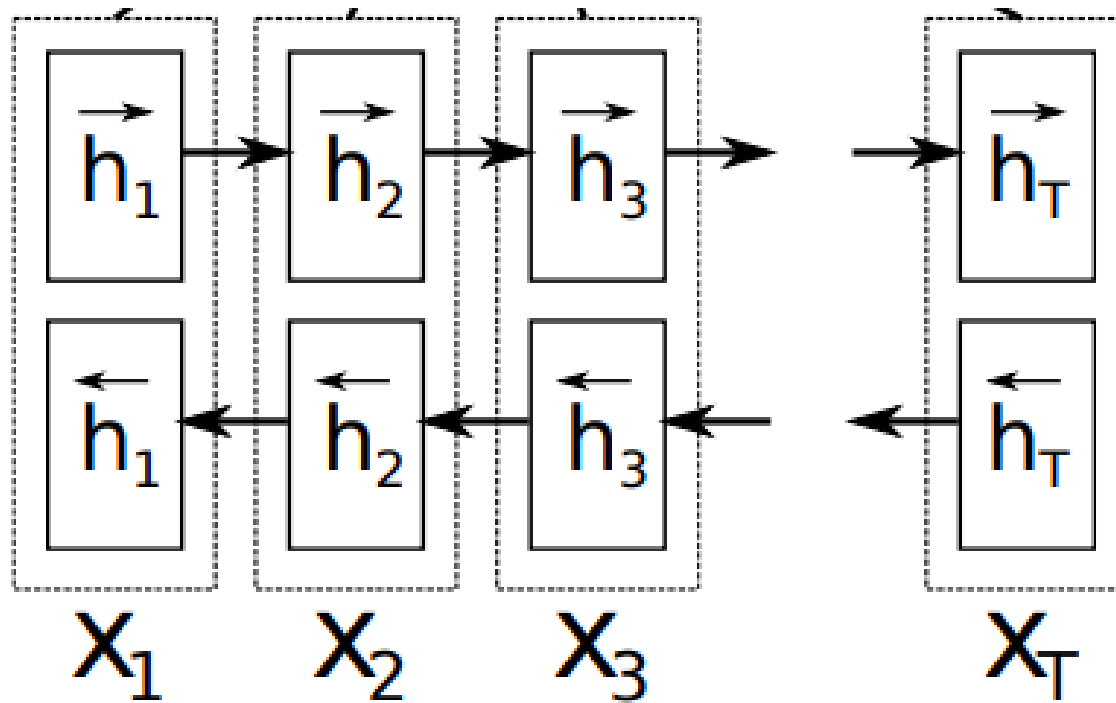


[LSTM: Hochreiter&Schmidhuber, 1997]

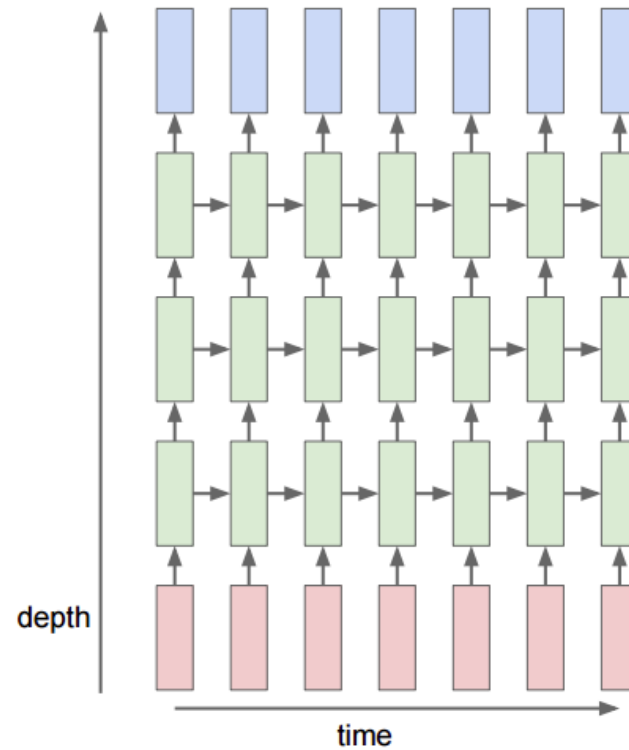
## UGRNN vs. GRUs vs. LSTMs

[Collins, Dickstein and Sussulo 2016] found that the number of parameters is more important than the choice of RNN architectures.

# bidirectional RNNs



# Multi-Layer RNNs



Modern versions would stack layers using residual connections.

# Machine Translation

$$w_0, \dots, w_{T_{\text{in}}} \Rightarrow \tilde{w}_0, \dots, \tilde{w}_{T_{\text{out}}}$$

Translation is a **sequence to sequence** (seq2seq) task.

**Sequence to Sequence Learning with Neural Networks**, Sutskever, Vinyals and Le, NeurIPS 2014, arXiv Sept 10, 2014.



# Machine Translation

We define a model

$$P_{\Phi} (\tilde{w}_0, \dots, \tilde{w}_{T_{\text{out}}} \mid w_0, \dots, w_{T_{\text{in}}})$$

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} E_{\langle x, y \rangle \sim P_{\text{op}}} [-\ln P_{\Phi}(y|x)]$$

## Translation Using Thought Vectors

The final state of a **right-to-left (backward)** RNN is viewed as a “**thought vector**” representation of the input sentence.

We use the thought vector for the input sentence as the initial hidden state of a **left-to-right (forward)** RNN language model generating the output sentence.

Computing the input thought vector backward provides a good start to the forward generation of the output.

# The Introduction of Attention

**Neural Machine Translation by Jointly Learning to **Align** and Translate** Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio, ICLR 2015 (arXiv Sept. 1, 2014)

## Attention

As we generate each word in the output translation we compute an attention over the input.

Intuitively, we want to define an “alignment” between the words in the output and the words in the input.

In modern terminology this is called a “cross attention” — one thing (the output) attending to a different thing (the input).

This is different from the “self attention” used in transformers.

## Encoder-Decoder Models with Cross-Attention

Let  $h_{\text{thought}}$  be the thought vector for the input sentence.

Let  $h_{\text{in}}(t_{\text{in}})$  be a sequence of vectors generated by the encoder for the sequence of input words.

Let  $h_{\text{out}}(t_{\text{out}})$  be a thought vector for the first  $t_{\text{out}}$  words in the output sentence.

$$\text{Attention } \alpha[t_{\text{out}}, t_{\text{in}}] = \text{softmax}_{t_{\text{in}}} e(w_{t_{\text{out}}})[N] h_{\text{in}}(t_{\text{in}})[N]$$

$$\text{Weighted Sum : } \tilde{h}_{\text{out}}(t_{\text{out}}) = \alpha[t_{\text{out}}, T_{\text{in}}] h_{\text{in}}(T_{\text{in}})$$

$$\text{generate : } w_{t_{\text{out}}+1} \sim P_{\Phi}(w_{t_{\text{out}}+1} | h_{\text{thought}}, h_{\text{out}}(t_{\text{out}}), \tilde{h}_{\text{out}}(t_{\text{out}}))$$

## Cross Attention in Image Captioning

We can treat image captioning as translating an image into a caption.

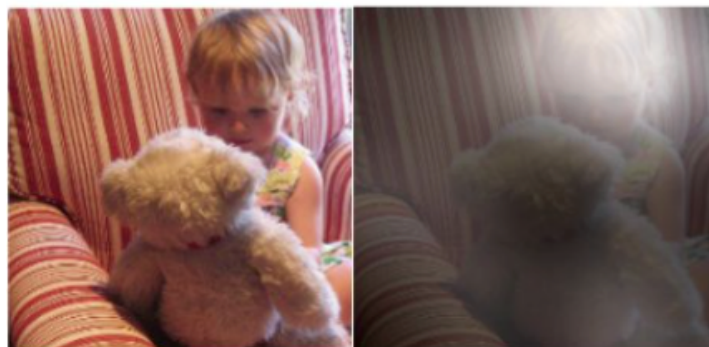
In translation with attention involves an attention over the input aligning output words with positions in the input.

For each output word we get an attention over the image positions.

# Attention in Image Captioning



A woman is throwing a frisbee in a park.



A little girl sitting on a bed with a teddy bear.

Xu et al. ICML 2015

# The Transformer: Self Attention

Attention is All You Need, Vaswani et al., June 2017

We replace the RNNs with self attention.

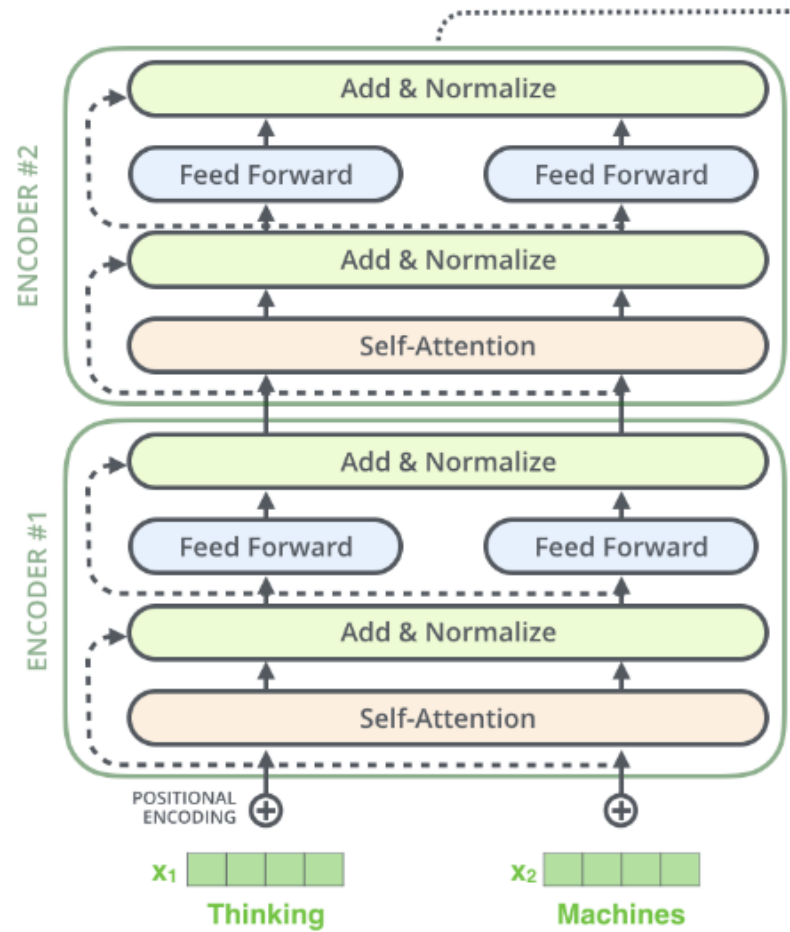
For the encoder we will have a “residual stack”  $h_{0,\text{in}}(t_{\text{in}}), \dots, h_{N,\text{in}}(t_{\text{in}})$ .

Set  $h_{0,\text{in}}(t_{\text{in}}) = e(w_{t_{\text{in}}}); \text{pos}(t_{t_{\text{in}}})$ .

Here semicolon denotes vector concatenation and  $\text{pos}(t_{\text{in}})$  is a “position encoding” for the position  $t$ .



# The Residual Pathway



## Parallel Layer Computation

However, in the transformer we can compute the layer  $L_{\ell+1}[T, J]$  from  $L_{\ell}[T, J]$  in parallel.

This is an important difference from RNNs which compute sequentially over time.

In this respect the transformer is more similar to a CNN than to an RNN.

## Self-Attention

The fundamental innovation of the transformer is the self-attention layer.

For each position  $t$  in the sequence we compute an attention over the other positions in the sequence.

## Transformer Heads

There is an intuitive analogy between the Transformer’s self attention and a dependency parse tree.

In a dependency parse consists of edges between words labeled with grammatical roles such as “subject-of” or “object-of”.

The self attention layers of the transformer we have “heads” which can be viewed as labels for dependency edges.

Self attention constructs a tensor  $\alpha_\ell[k, t_1, t_2]$  — the strength of the attention weight (edge weight) from  $t_1$  to  $t_2$  with head (label)  $k$  at layer  $\ell$  of the network.

## Query-Key Attention

For each head  $k$  and position  $t$  we compute a key vector and a query vector with dimension  $I$  typically smaller than dimension  $J$ .

$$\text{Query}_{\ell+1}[k, t, i] = W_{\ell+1}^{\text{query}}[k, i, J]L_{\ell}[t, J]$$

$$\text{Key}_{\ell+1}[k, t, i] = W_{\ell+1}^{\text{key}}[k, i, J]L_{\ell}[t, J]$$

$$\alpha_{\ell+1}[k, t_1, t_2] = \text{softmax}_{t_2} \frac{1}{\sqrt{I}} \text{Query}_{\ell+1}[k, t_1, I]\text{Key}_{\ell+1}[k, t_2, I]$$

## Computing the Self-Attention Layer

$$\text{Value}_{\ell+1}[k, t, i] = W_{\ell+1}^V[k, i, J]L_{\ell}[t, J]$$

$$\tilde{h}_{\ell+1}[k, t, i] = \alpha[k, t, T]\text{Value}[k, T, i]$$

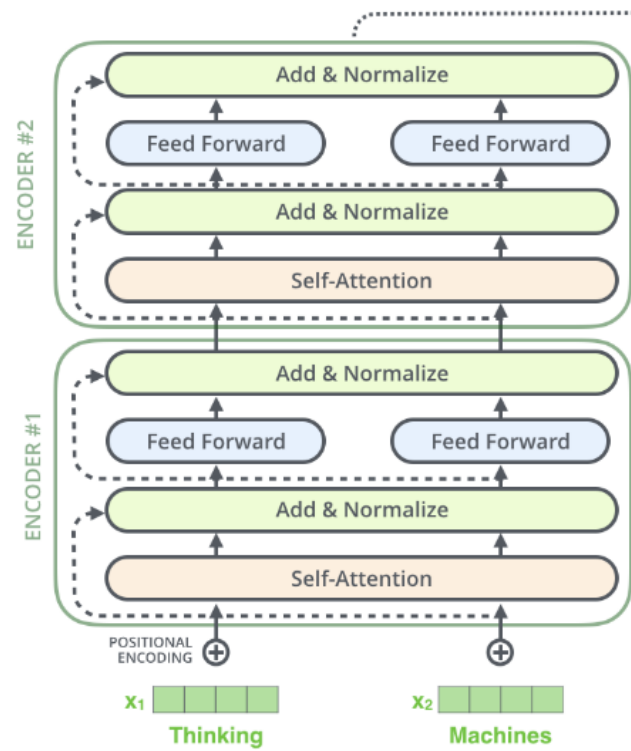
$$\hat{h}_{\ell+1}[t, C] = \tilde{h}_{\ell+1}[0, t, I]; \cdots ; \tilde{h}_{\ell+1}[K - 1, t, I]$$

$$L_{\ell+1}^{SA}[t, j] = W_{\ell+1}^0[j, C]\hat{h}_{\ell+1}[t, C]$$

Here semicolon denotes vector concatenation.

# The Residual Pathway

We have now defined the self-attention layer.



## Feed-Forward Layers

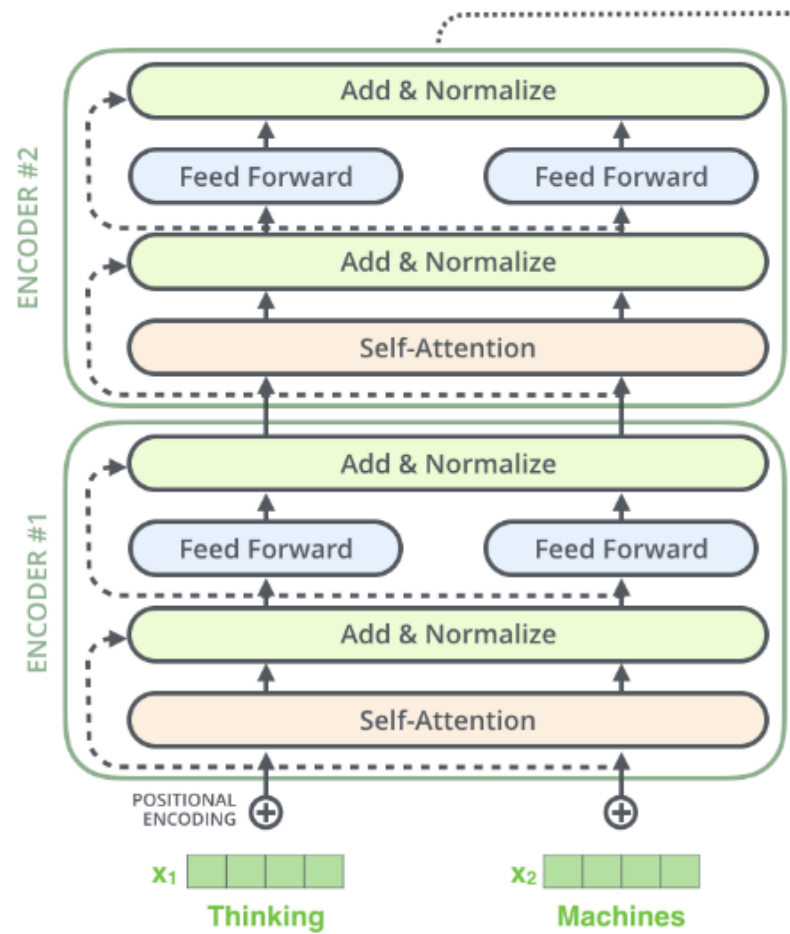
The feed-forward layers apply a two-level multi-layer perceptron (MLP) to the vector at each time position independently.

$$h_{\ell+1}[t, i] = \text{ReLU}(W_{\ell+1}^{\text{FF1}}[i, J] L_{\ell}[t, J] - B_{\ell+1}^{\text{FF1}}[i])$$

$$L_{\ell+1}[t, j] = W_{\ell+1}^{\text{FF2}}[j, I] h_{\ell+1}[t, I] - B_{\ell+1}^{\text{FF2}}[j]$$



# The Transformer



## Recent Trends

LoRA: Low-Rank Adaptation of Large Language Models, Huet al. (June 2021).

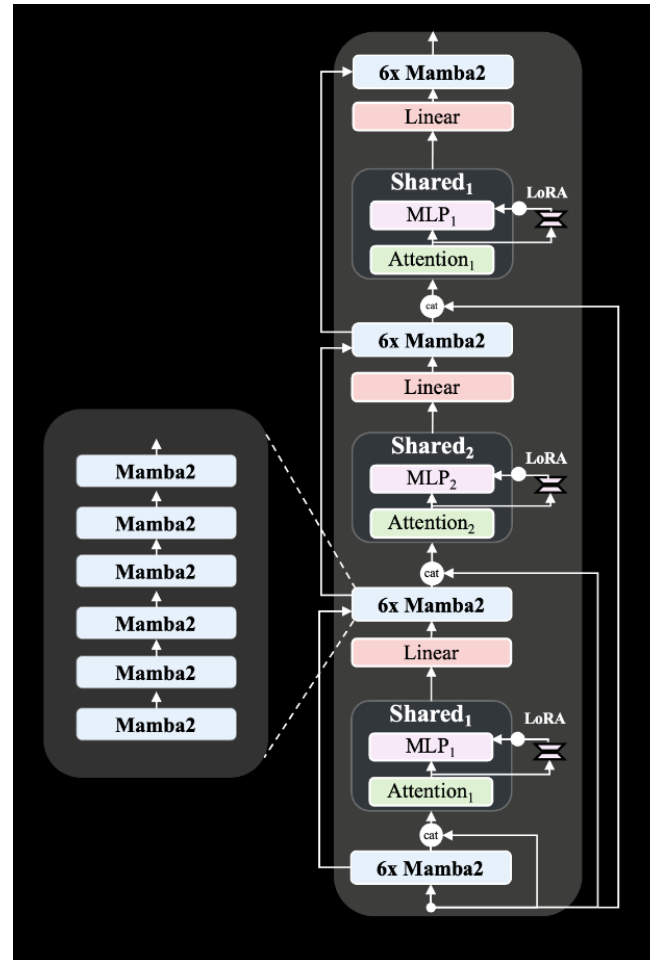
Efficiently Modeling Long Sequences with Structured State Spaces Albert Gu net al. (August 2022)

Mamba: Linear-Time Sequence Modeling with Selective State Spaces Albert Gu and Tri Dao, (December 2023)

## Zambda2-7B (October 14, 2024)

Zyphra is excited to release Zamba2-7B, a state-of-the-art small language model. **At the 7B scale**, we outperform the leading models of Mistral, Google's Gemma and Meta's Llama3 series in both quality and performance ...

# Zamba2



**END**