



aivancity

SCHOOL FOR

TECHNOLOGY, BUSINESS & SOCIETY

PARIS-CACHAN

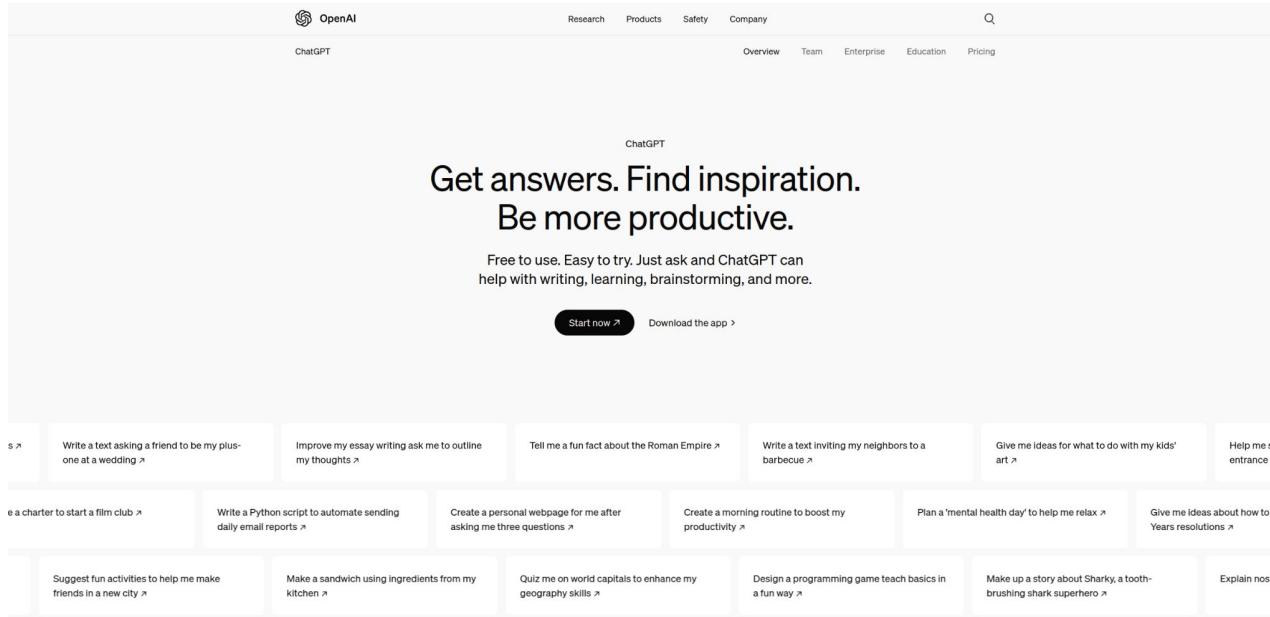


05/02/2025

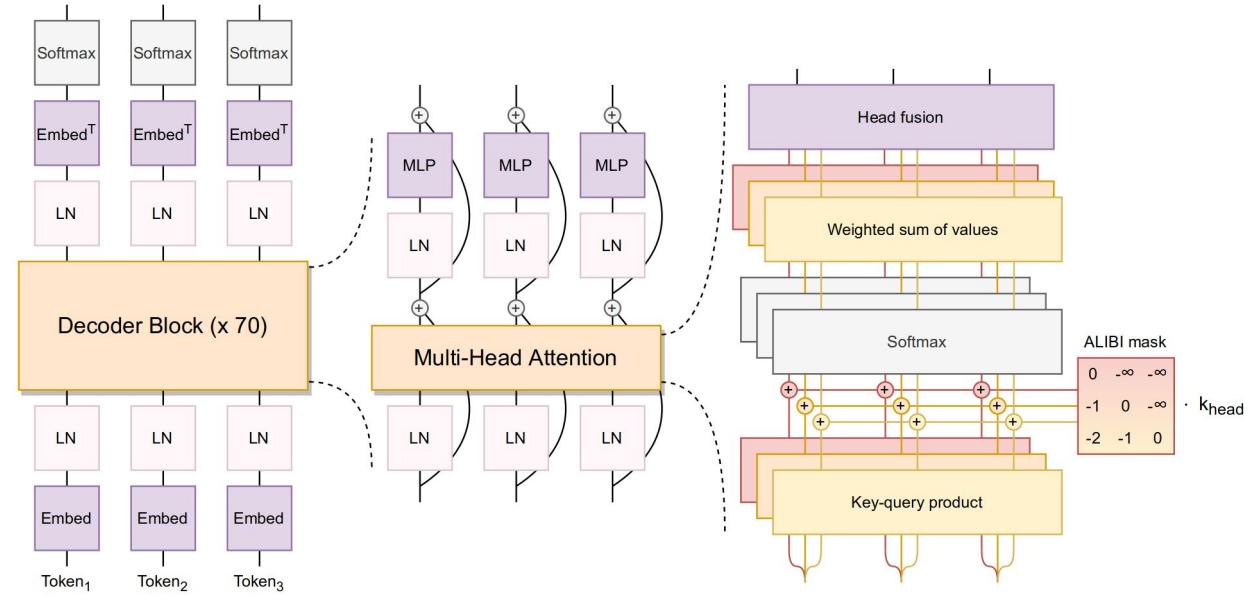
# Natural Language Processing (NLP)

*LLM Architectures: Attention Mechanism and Transformers*

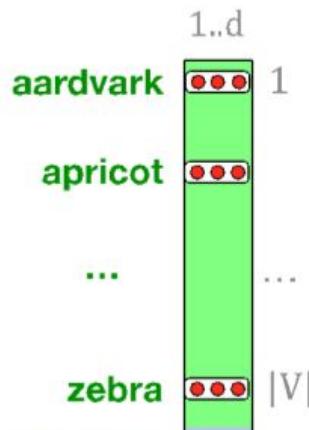
# Chatbots like ChatGPT rely on LLM



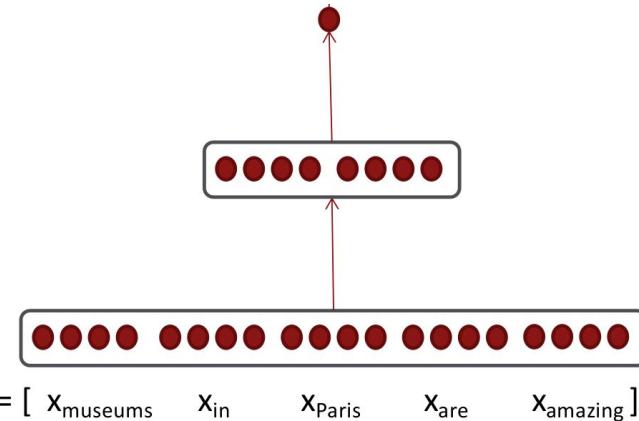
# LLMs rely on Attention and Transformer



# What can we do so far?

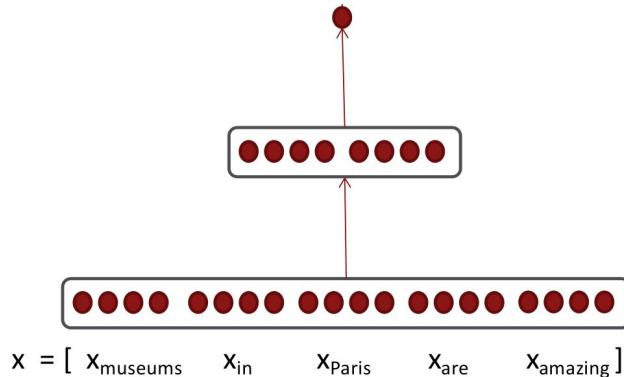


- Embedding matrix: a vector for each word
- Fine for **classification** (e.g. Named Entity Recognition)



# Continuous bag of words (CBOW)

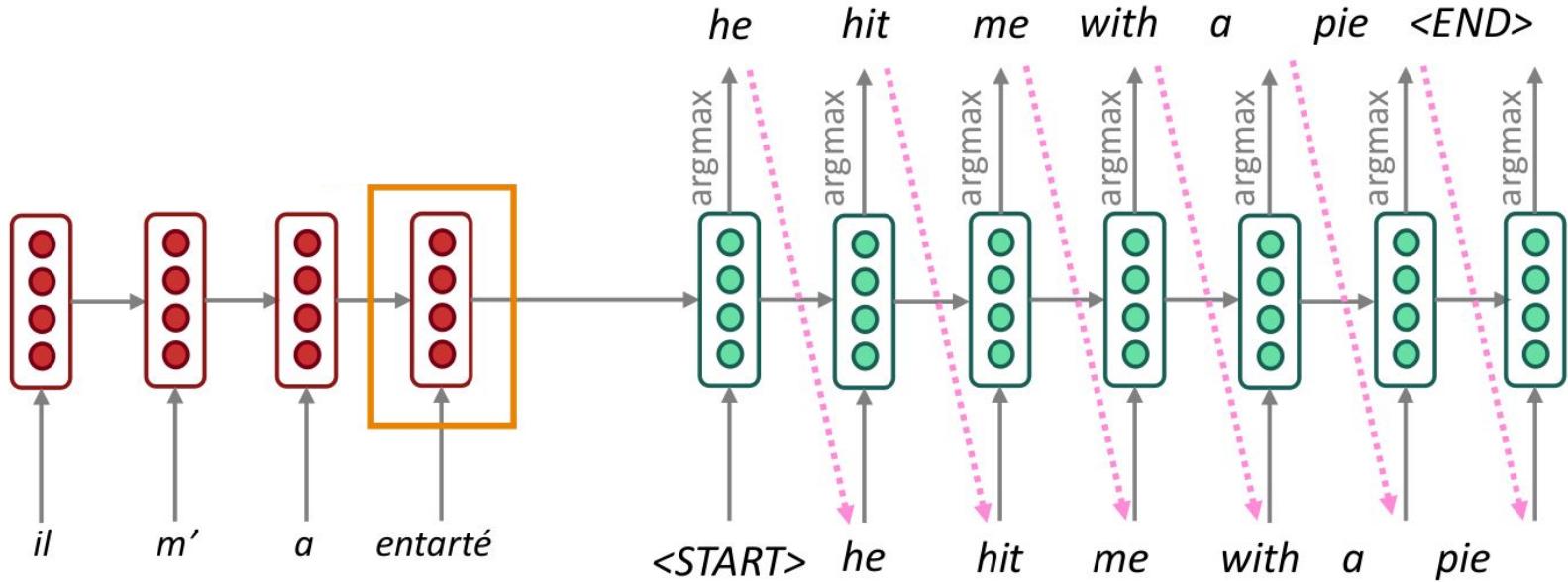
"bag of words" because does not model word order, puts all words in the same "bag"



- $x_1$ : yes , we have no bananas
- $x_2$ : say yes for bananas
- $x_3$ : no bananas , we say

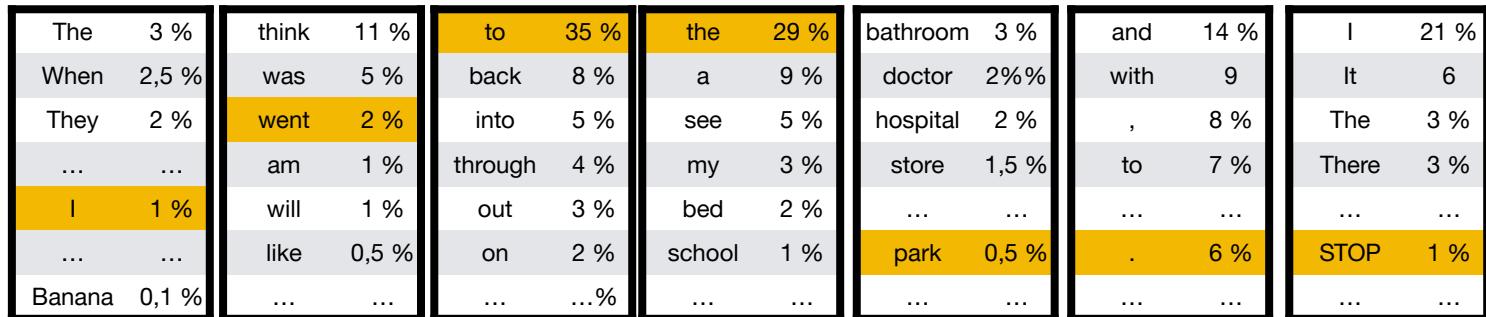
	1	2	3
,	1	0	1
bananas	1	1	1
for	0	1	0
have	1	0	0
no	1	0	1
say	0	1	1
we	1	0	1
yes	1	1	0

# Sequence-to-Sequence (Translation)

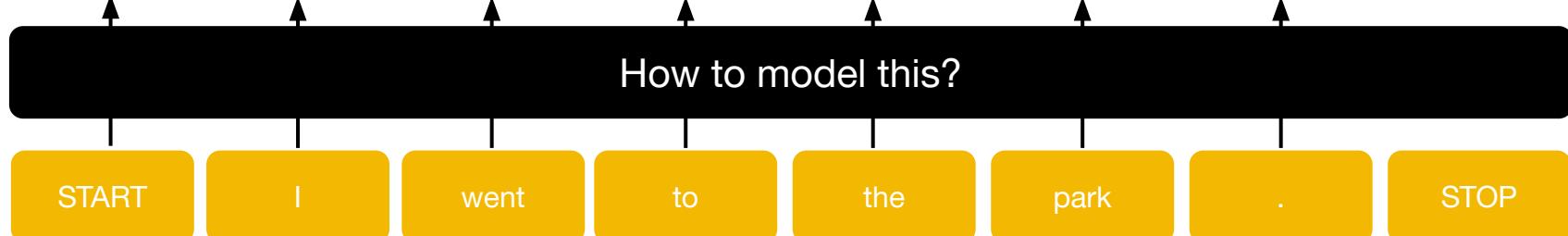


# Language Modeling

$p(x|\text{START})p(x|\text{START I})p(x|\dots\text{went}) \ p(x|\dots\text{to}) \ p(x|\dots\text{the}) \ p(x|\dots\text{park}) \ p(x|\text{START I went to the park.})$

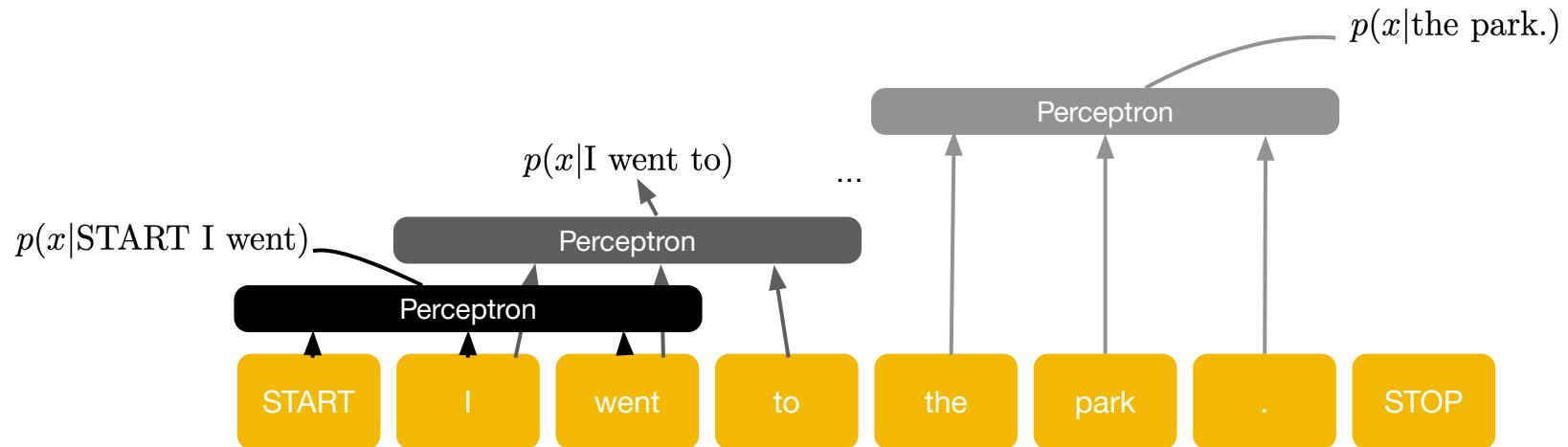


How to model this?



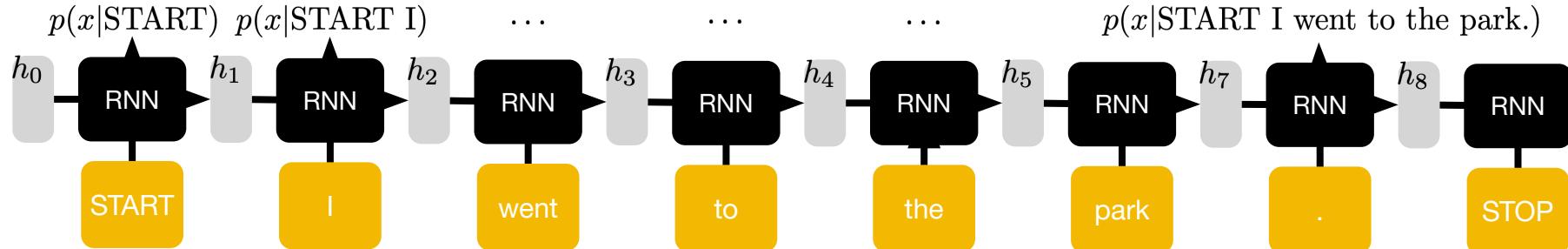
# v0: sliding window

- Sliding window of size N, like CBOW or n-grams (next class):  
(N – 1)th-order Markov assumption (prediction only depends on N-1 last)

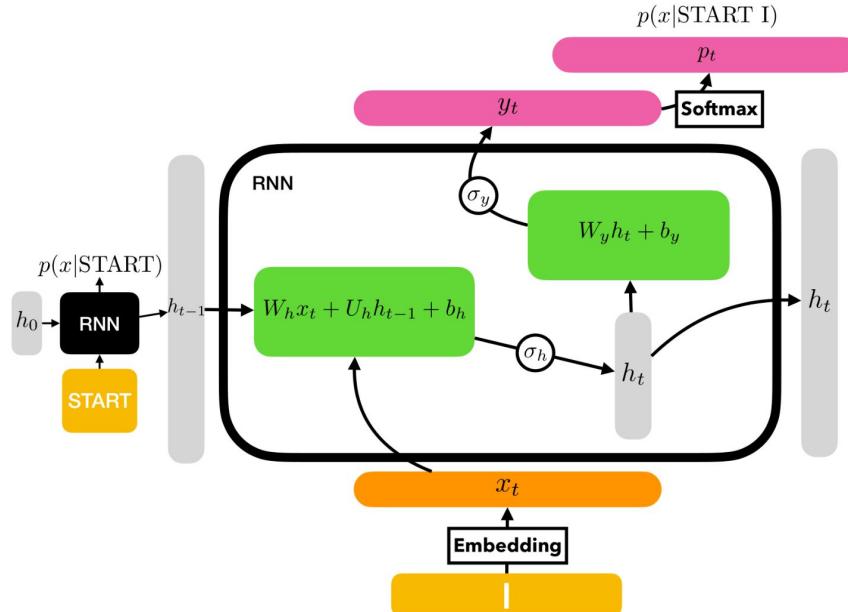


# v1: Recurrent Neural Networks (RNN)

- One RNN is applied recursively to the sequence
- Inputs: previous hidden state  $h_{t-1}$ , observation  $x_t$
- Outputs: next hidden state  $h_t$ , (optionally) output  $y_t$
- Memory about history is passed through hidden states



# v1: Recurrent Neural Networks (RNN)



## Variables:

- $x_t$ : input (embedding) vector
- $y_t$ : output vector (logits)
- $p_t$ : probability over tokens
- $h_{t-1}$ : previous hidden vector
- $h_t$ : next hidden vector
- $\sigma_h$ : activation function for hidden state
- $\sigma_y$ : output activation function

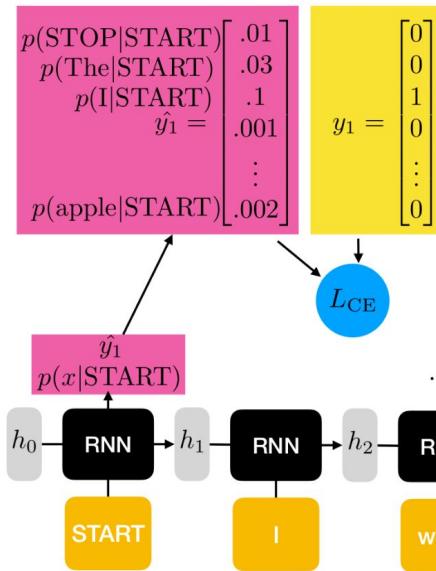
## Equations:

$$h_t := \sigma_h(W_h x_t + U_h h_{t-1} + b_h)$$

$$y_t := \sigma_y(W_y h_t + b_y)$$

$$p_{t_i} = \frac{\exp(y_{t_i})}{\sum_{i=j}^d \exp(y_{t_j})}$$

# Generation as a Sequence of Classifications



2. Run model on (batch of)  $x$ 's from data  $\mathcal{D}$  to get probability distributions  $\hat{y}$
3. Calculate loss compared to true  $y$ 's (Cross Entropy Loss)

$$\hat{y}_7 = \begin{bmatrix} .2 \\ .03 \\ .12 \\ .01 \\ \vdots \\ .001 \end{bmatrix} \quad y_7 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$L_{\text{CE}}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

$\hat{y}_7$   
 $p(x|\text{START I went to the park.})$

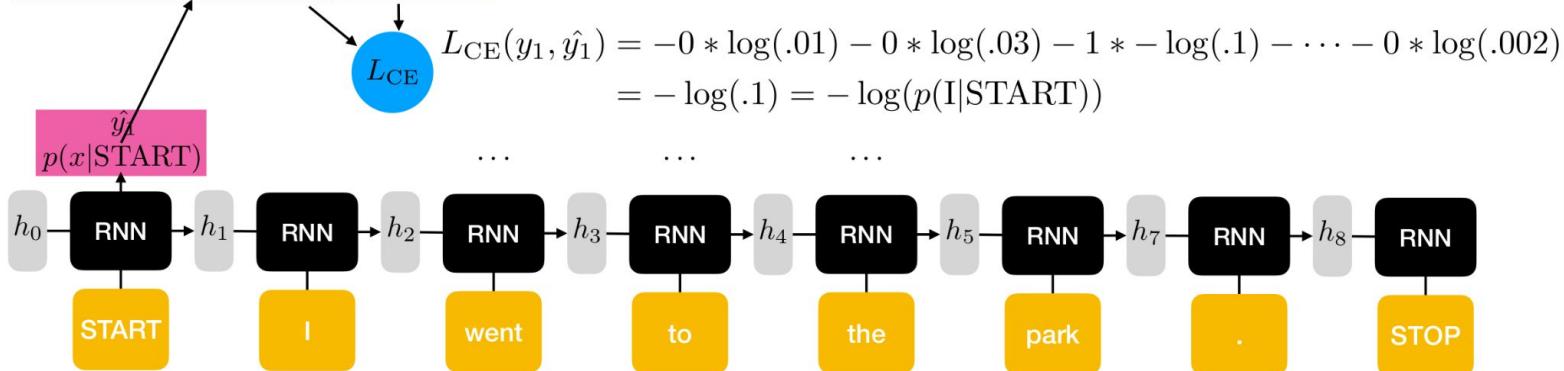
# Cross-Entropy

$$L_{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

$$\hat{y}_1 = \begin{bmatrix} .01 \\ .03 \\ .1 \\ .001 \\ \vdots \\ .002 \end{bmatrix}$$

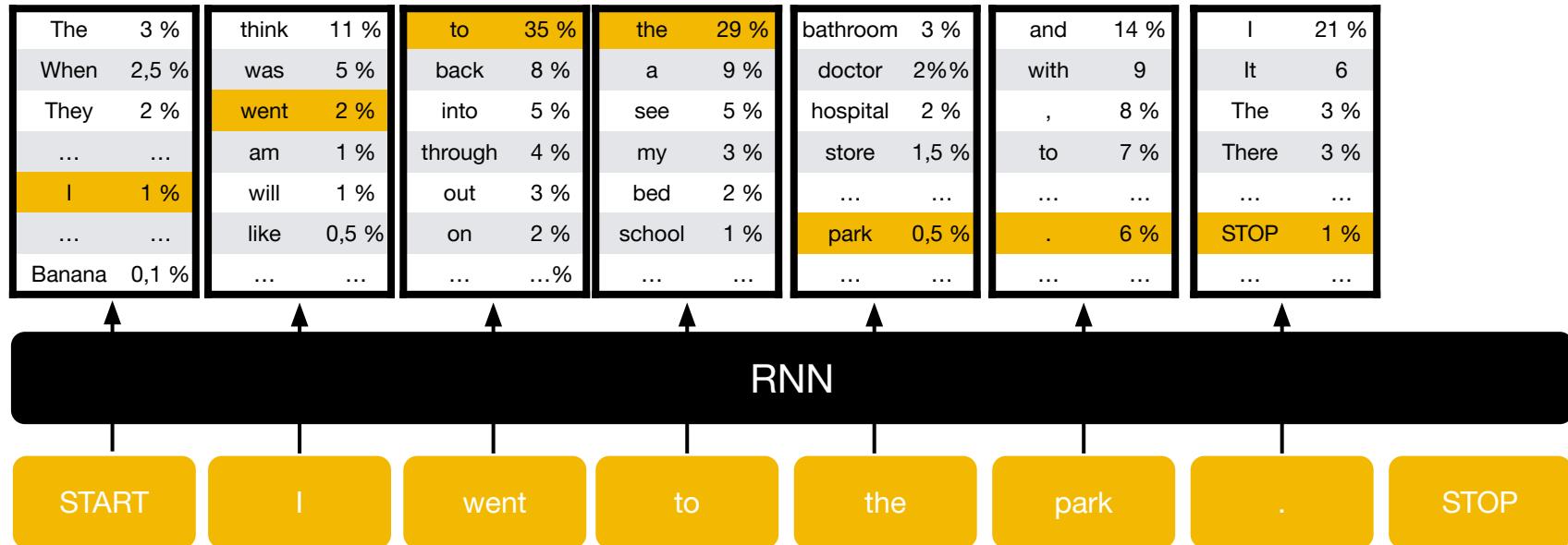
$y_1 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

(Actual observed word)

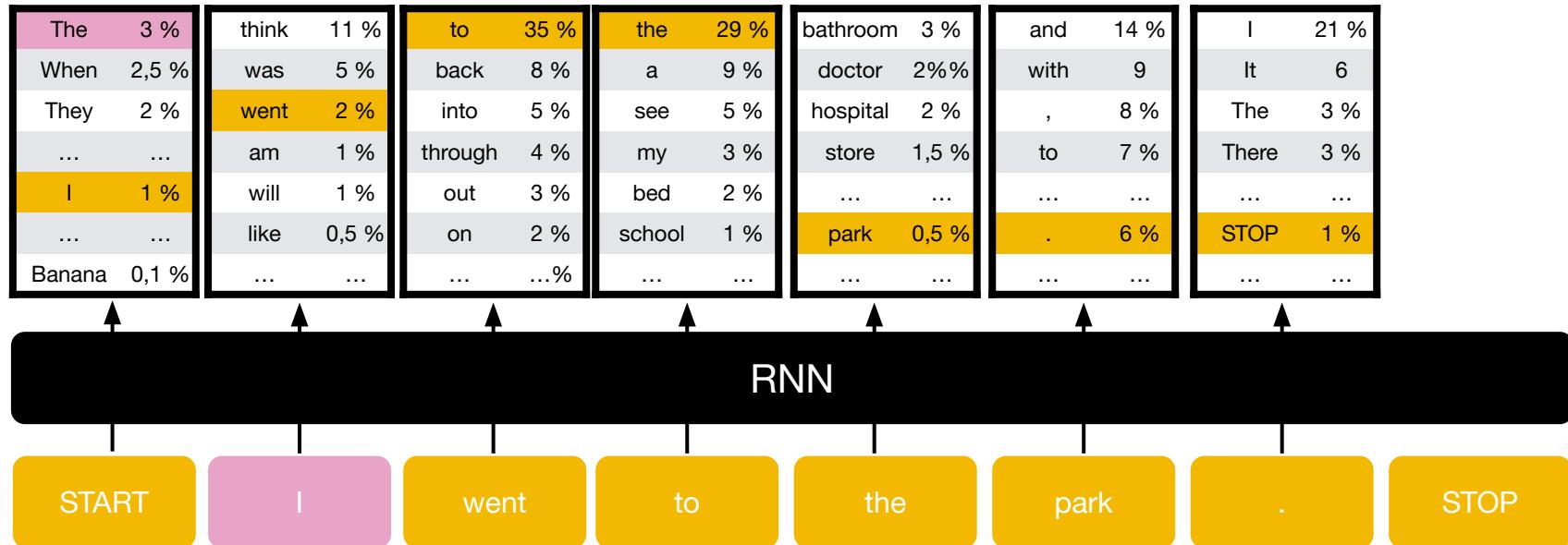


3. Calculate loss compared to true  
 $y$ 's (Cross Entropy Loss)

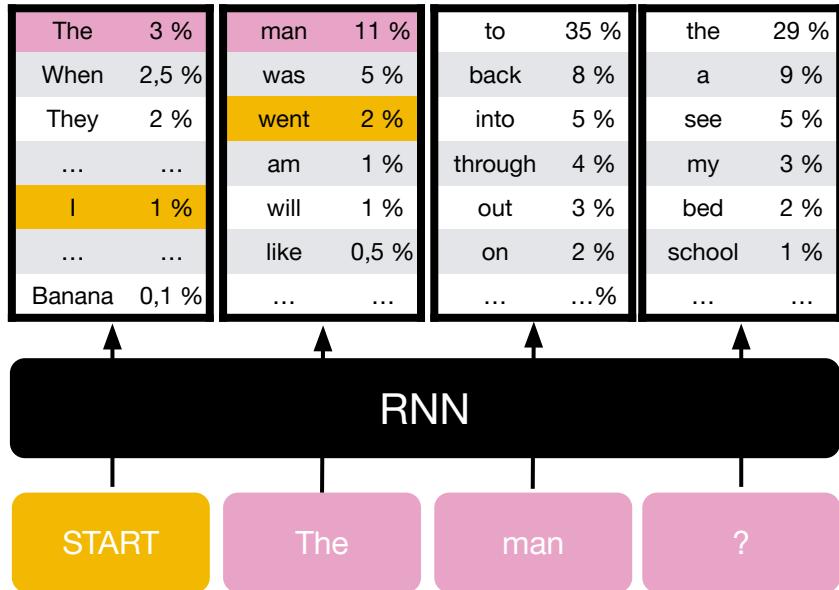
# Teacher Forcing



# Train-test mismatch: exposure bias



# Train-test mismatch: exposure bias

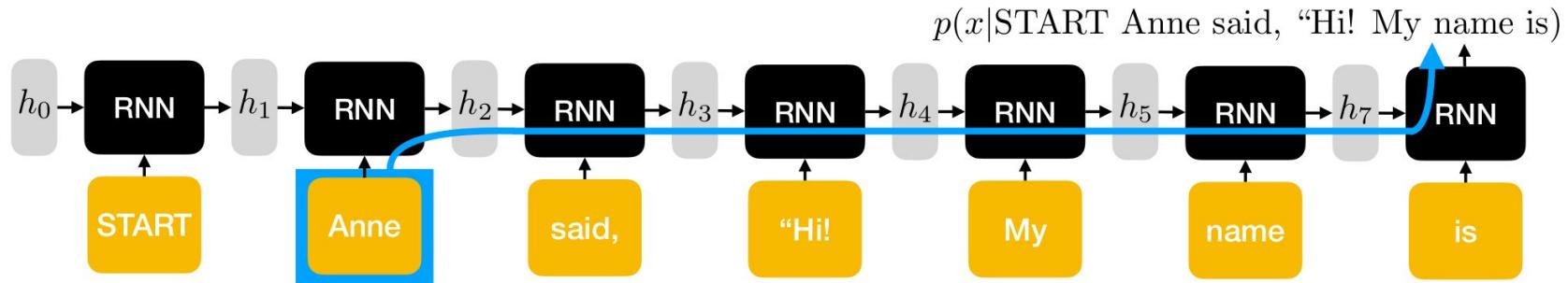


- Keep in mind: with RNN it's trivial to train the model with its own generation instead of teacher forcing
- reduces exposure bias

# RNN limit 1: vanishing gradients

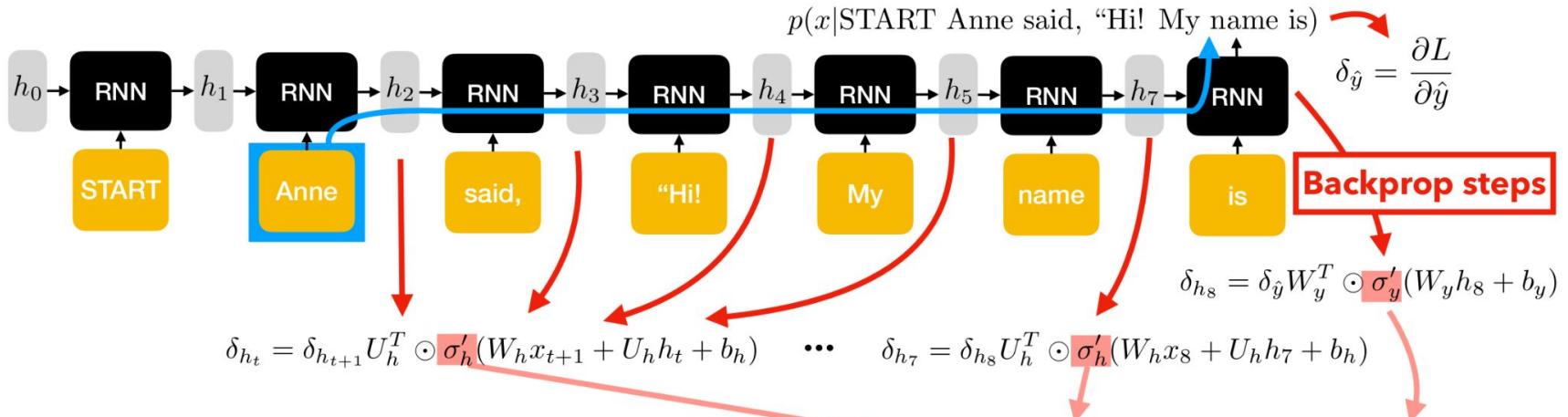
- What word is likely to come next for this sequence?

Anne said, "Hi! My name is



- Need relevant information to flow across many time steps
- When we backpropagate, we want to allow the relevant information to flow

# RNN limit 1: vanishing gradients



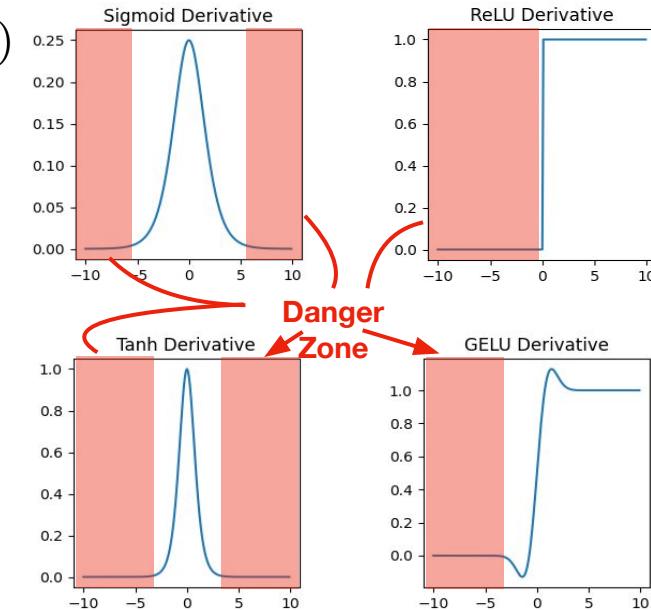
However, when we backprop, it involves multiplying a chain of computations from time  $t_1$  to time  $t_7$

If any of the terms are close to zero, the whole gradient goes to zero (vanishes!)

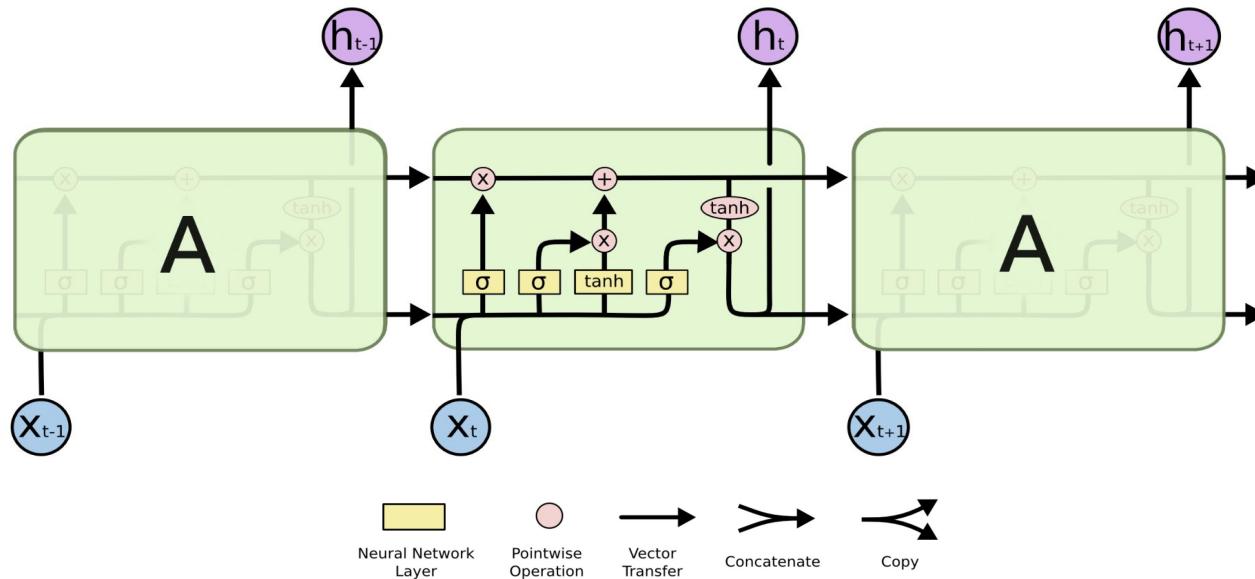
# RNN limit 1: vanishing gradients

$$\delta_{h_t} = \delta_{h_{t+1}} U_h^T \odot \sigma'_h(W_h x_{t+1} + U_h h_t + b_h)$$

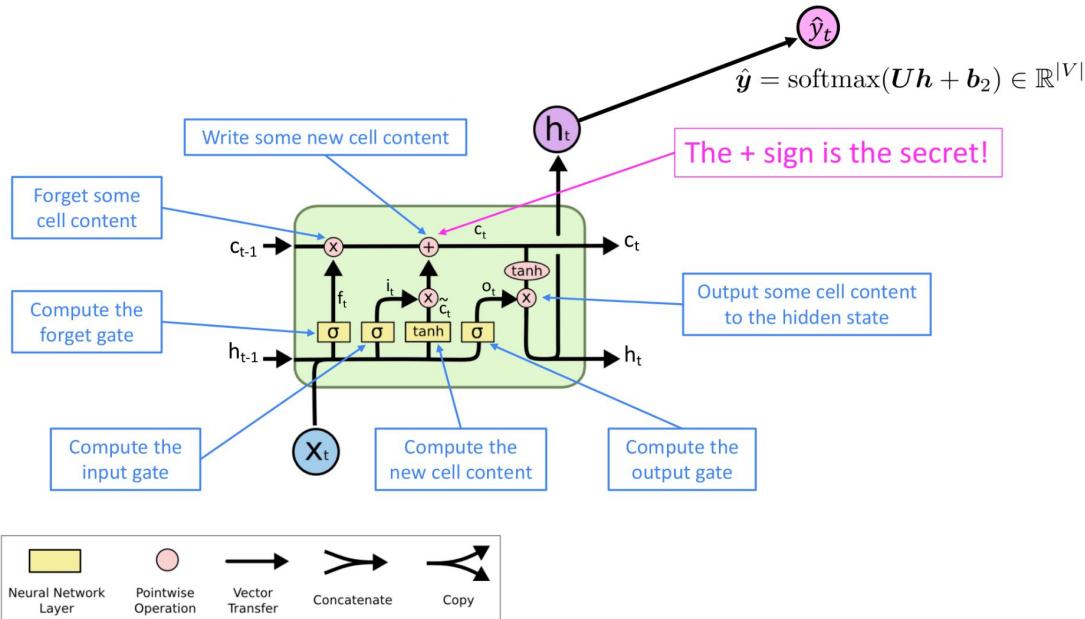
- If any of the terms are close to zero, the whole gradient goes to zero (vanishes!)
- This happens often for many activation functions... the gradient is close to zero when outputs get very large or small
- The more time steps back, the more chances for a vanishing gradient



# v1.1: Long-Short Term Memory (LSTM)

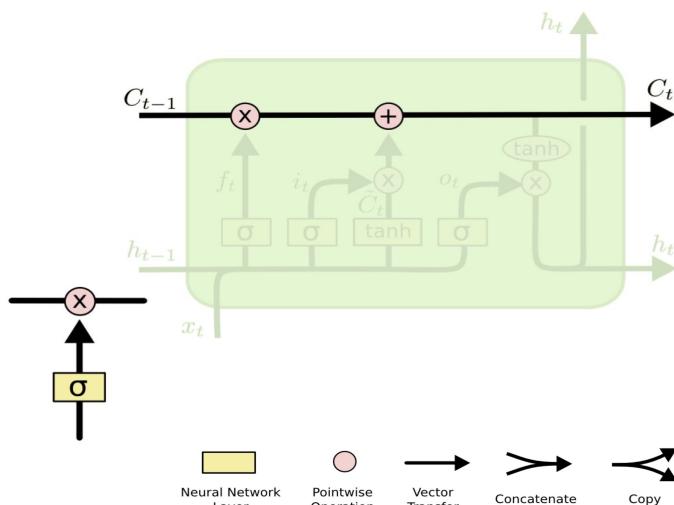


# v1.1: Long-Short Term Memory (LSTM)



# v1.1: Long-Short Term Memory (LSTM)

Cell state (**long-term memory**): allows information to flow with only small, **linear interactions** (good for gradients!)

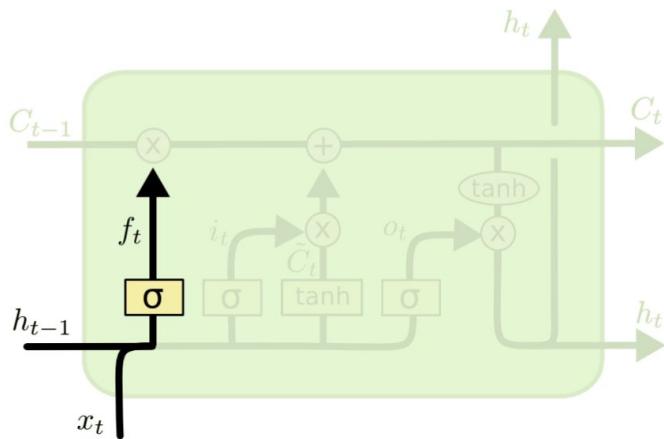


$$\begin{aligned}
 f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \\
 i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \\
 o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \\
 \tilde{c}_t &= \sigma_c(W_c x_t + U_c h_{t-1} + b_c) \\
 c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\
 h_t &= o_t \odot \sigma_h(c_t)
 \end{aligned}$$

- $x_t \in \mathbb{R}^d$ : input vector to the LSTM unit
- $f_t \in (0, 1)^h$ : forget gate's activation vector
- $i_t \in (0, 1)^h$ : input/update gate's activation vector
- $o_t \in (0, 1)^h$ : output gate's activation vector
- $h_t \in (-1, 1)^h$ : hidden state vector also known as output vector of the LSTM unit
- $\tilde{c}_t \in (-1, 1)^h$ : cell input activation vector
- $c_t \in \mathbb{R}^h$ : cell state vector

# v1.1: Long-Short Term Memory (LSTM)

Input Gate Layer: Decide what information to “**forget**”



$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \sigma_h(c_t)$$

$x_t \in \mathbb{R}^d$ : input vector to the LSTM unit

$f_t \in (0, 1)^h$ : forget gate's activation vector

$i_t \in (0, 1)^h$ : input/update gate's activation vector

$o_t \in (0, 1)^h$ : output gate's activation vector

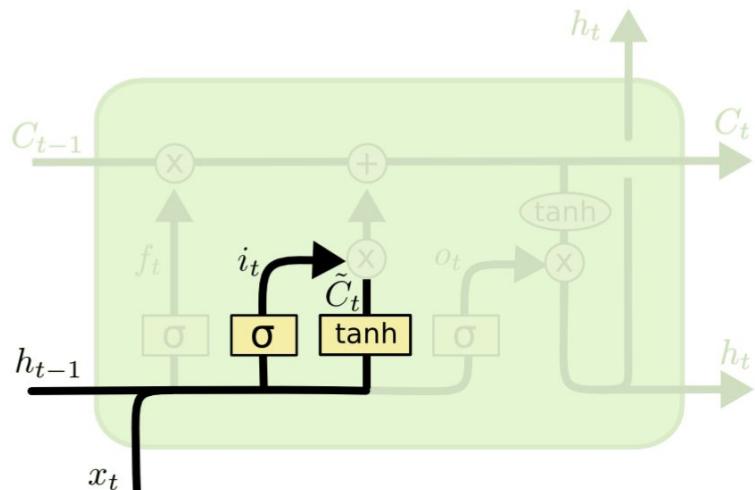
$h_t \in (-1, 1)^h$ : hidden state vector also known as output vector of the LSTM unit

$\tilde{c}_t \in (-1, 1)^h$ : cell input activation vector

$c_t \in \mathbb{R}^h$ : cell state vector

# v1.1: Long-Short Term Memory (LSTM)

Candidate state values: Extract candidate information to put into the cell vector: "**remember**"



$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \sigma_h(c_t)$$

$x_t \in \mathbb{R}^d$ : input vector to the LSTM unit

$f_t \in (0, 1)^h$ : forget gate's activation vector

$i_t \in (0, 1)^h$ : input/update gate's activation vector

$o_t \in (0, 1)^h$ : output gate's activation vector

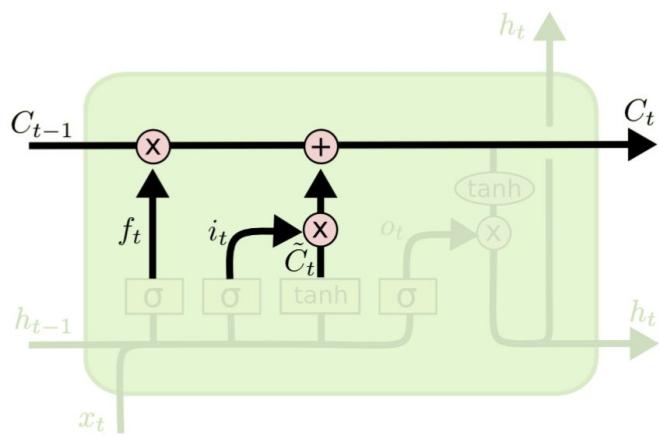
$h_t \in (-1, 1)^h$ : hidden state vector also known as output vector of the LSTM unit

$\tilde{c}_t \in (-1, 1)^h$ : cell input activation vector

$c_t \in \mathbb{R}^h$ : cell state vector

# v1.1: Long-Short Term Memory (LSTM)

Update cell: “Forget” the information we decided to forget and update with new candidate information



If  $f\_t$  is:

- High: we remember more previous info
- Low: we “forget” more info

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \sigma_h(c_t)$$

If  $i\_t$  is:

- High: we add more new info
- Low: we add less new info

$x_t \in \mathbb{R}^d$ : input vector to the LSTM unit

$f_t \in (0, 1)^h$ : forget gate’s activation vector

$i_t \in (0, 1)^h$ : input/update gate’s activation vector

$o_t \in (0, 1)^h$ : output gate’s activation vector

$h_t \in (-1, 1)^h$ : hidden state vector also known as output vector of the LSTM unit

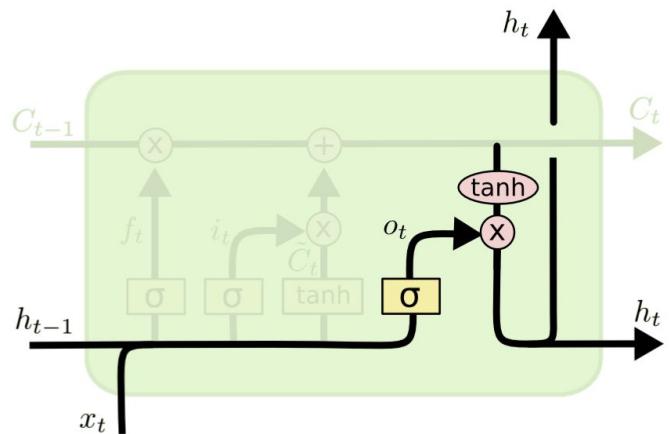
$\tilde{c}_t \in (-1, 1)^h$ : cell input activation vector

$c_t \in \mathbb{R}^h$ : cell state vector

# v1.1: Long-Short Term Memory (LSTM)

Output/**Short-term Memory** (as in RNN):

Pass information onto the next state/for use in output  
(e.g., probabilities)



$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \sigma_h(c_t)$$

$x_t \in \mathbb{R}^d$ : input vector to the LSTM unit

$f_t \in (0, 1)^h$ : forget gate's activation vector

$i_t \in (0, 1)^h$ : input/update gate's activation vector

$o_t \in (0, 1)^h$ : output gate's activation vector

$h_t \in (-1, 1)^h$ : hidden state vector also known as output vector of the LSTM unit

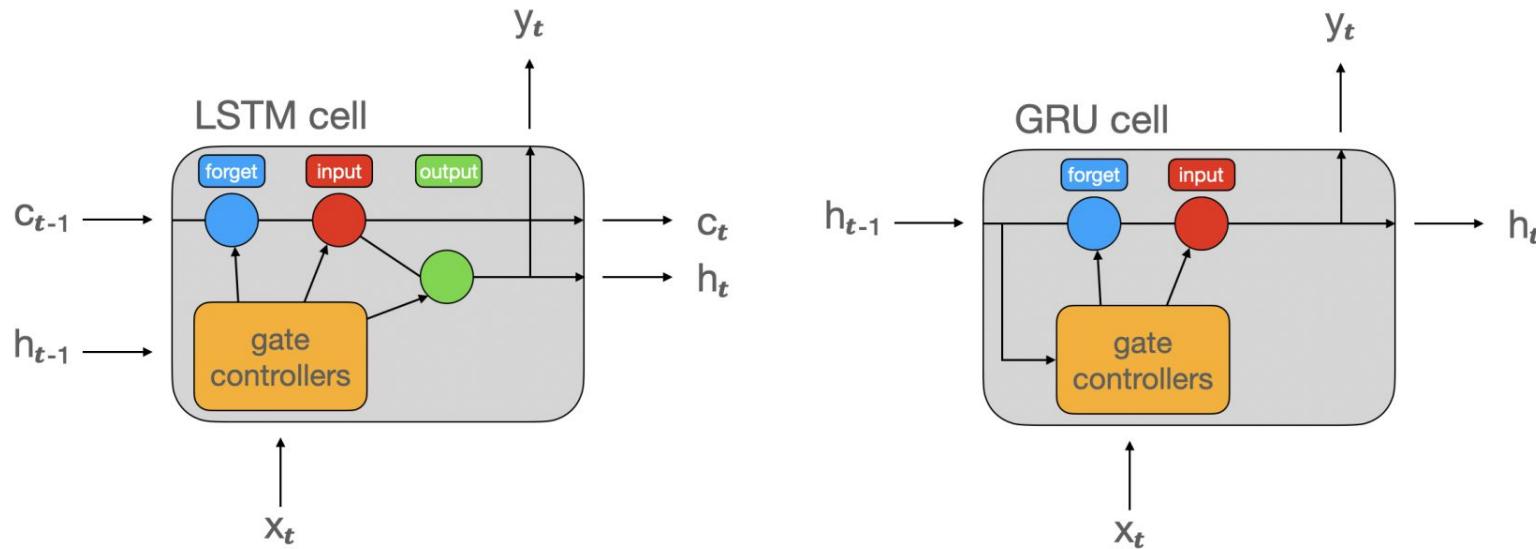
$\tilde{c}_t \in (-1, 1)^h$ : cell input activation vector

$c_t \in \mathbb{R}^h$ : cell state vector

# LSTM for Machine Translation

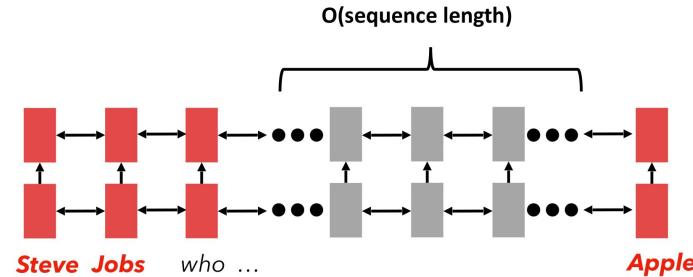
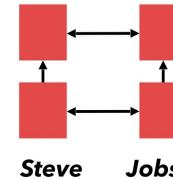
Type	Sentence
<b>Our model</b>	Ulrich UNK , membre du conseil d' administration du constructeur automobile Audi , affirme qu' il s' agit d' une pratique courante depuis des années pour que les téléphones portables puissent être collectés avant les réunions du conseil d' administration afin qu' ils ne soient pas utilisés comme appareils d' écoute à distance .
<b>Truth</b>	Ulrich Hackenberg , membre du conseil d' administration du constructeur automobile Audi , déclare que la collecte des téléphones portables avant les réunions du conseil , afin qu' ils ne puissent pas être utilisés comme appareils d' écoute à distance , est une pratique courante depuis des années .
<b>Our model</b>	“ Les téléphones cellulaires , qui sont vraiment une question , non seulement parce qu' ils pourraient potentiellement causer des interférences avec les appareils de navigation , mais nous savons , selon la FCC , qu' ils pourraient interférer avec les tours de téléphone cellulaire lorsqu' ils sont dans l' air ” , dit UNK .
<b>Truth</b>	“ Les téléphones portables sont véritablement un problème , non seulement parce qu' ils pourraient éventuellement créer des interférences avec les instruments de navigation , mais parce que nous savons , d' après la FCC , qu' ils pourraient perturber les antennes-relais de téléphonie mobile s' ils sont utilisés à bord ” , a déclaré Rosenker .

# Gated Recurrent Units (GRU)

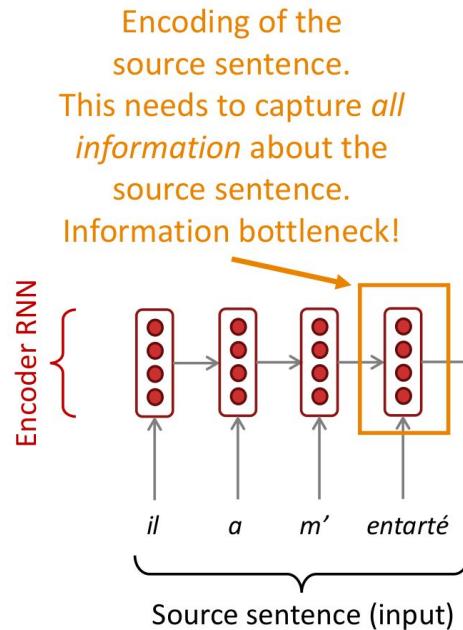


# RNN limit 2: Linear Interaction Distance

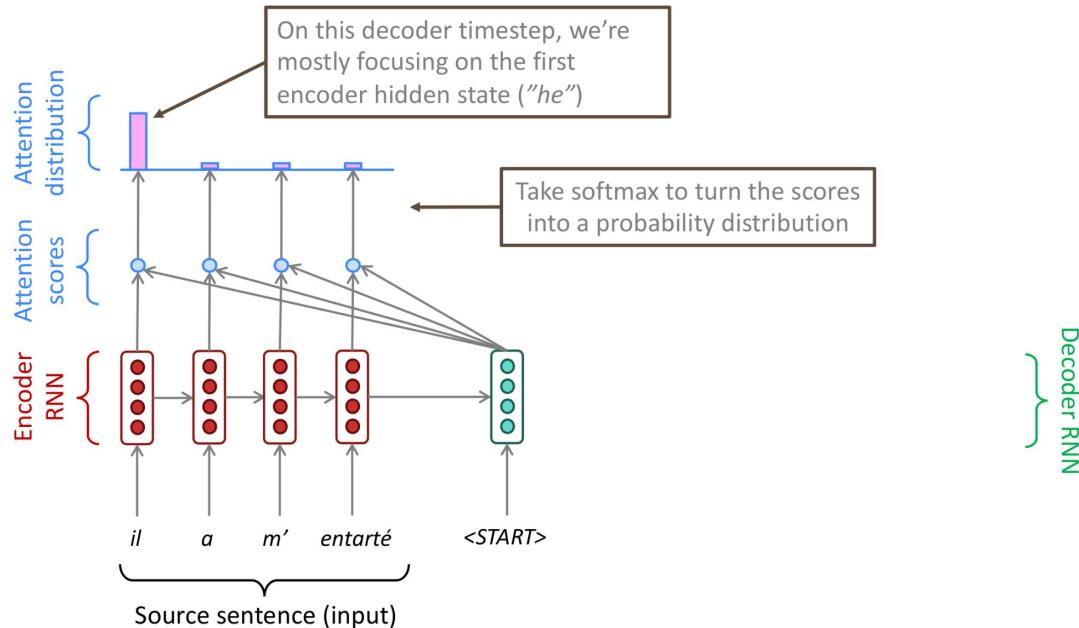
- RNNs are unrolled left-to-right.
- Linear locality is a useful heuristic: nearby words often affect each other's meaning!
- Better with LSTM than RNN, but still:
  - Failing to capture long-term dependencies
  - Vanishing gradient



# Why Attention?

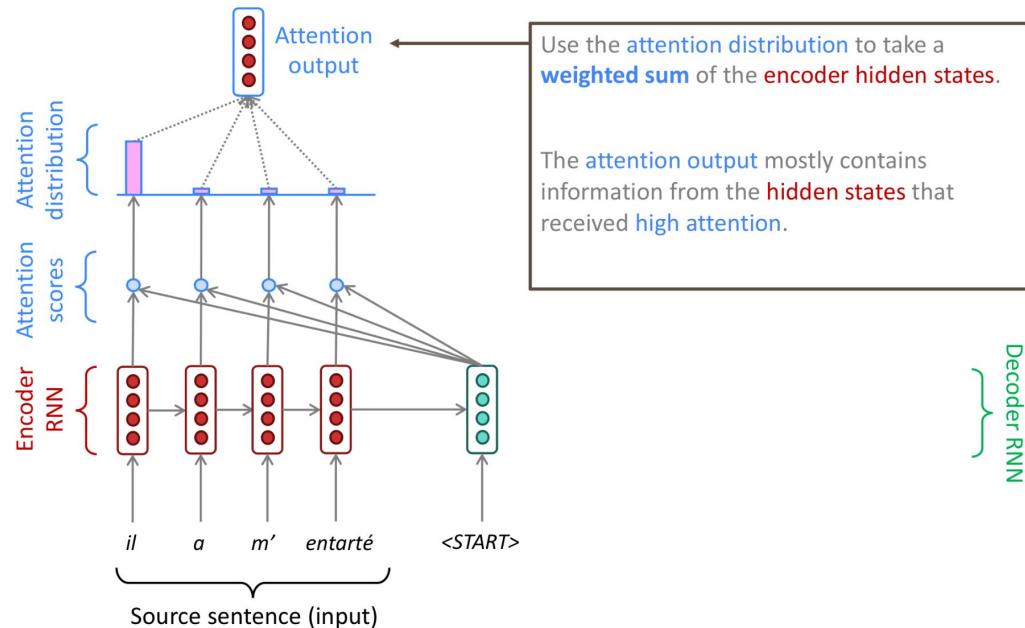


# LSTM with Attention



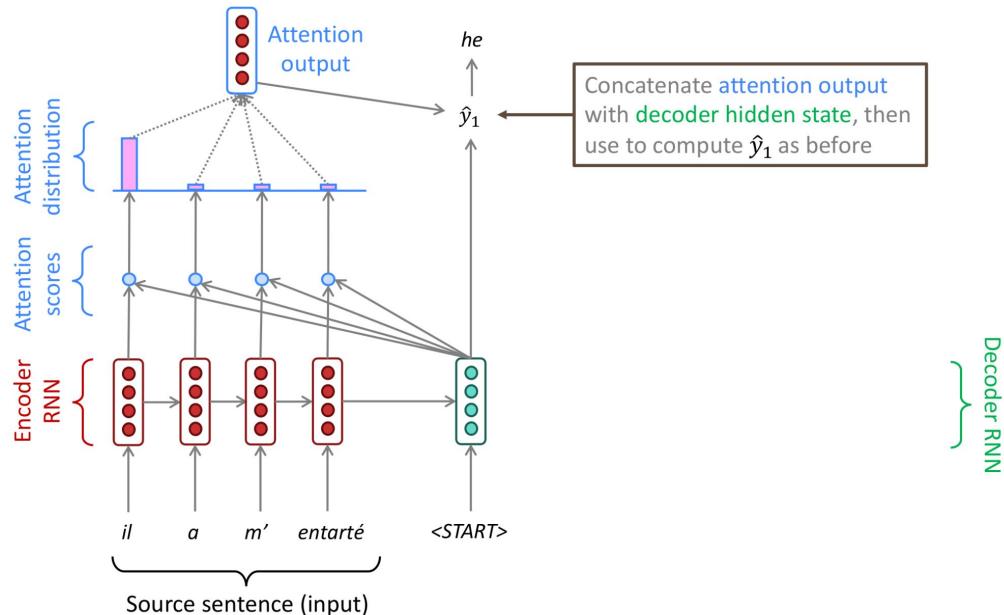
On each step of the decoder, use **direct connection** to the encoder to focus on a particular part of the source sequence

# LSTM with Attention



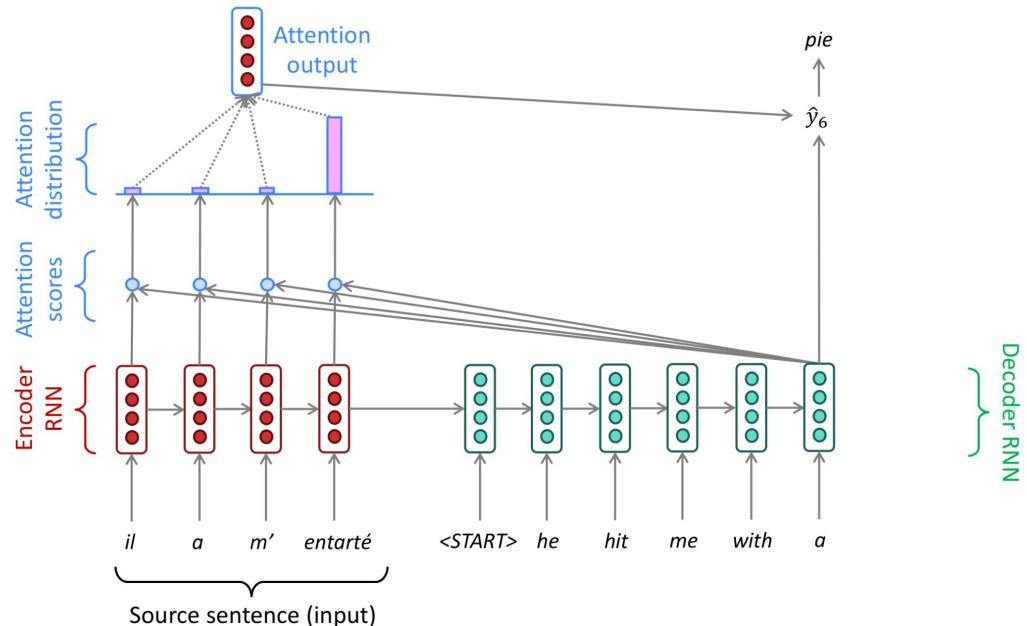
On each step of the decoder, use **direct connection** to the encoder to focus on a particular part of the source sequence

# LSTM with Attention



On each step of the decoder, use **direct connection** to the encoder to focus on a particular part of the source sequence

# LSTM with Attention



On each step of the decoder, use **direct connection** to the encoder to focus on a particular part of the source sequence

# Why Attention?

- Attention significantly improves Translation performance:  
It's very useful to allow decoder to focus on certain parts of the source
- Attention provides a more “human-like” model of the MT process:  
You can look back at the source sentence while translating, rather than needing to remember it all
- Attention solves the bottleneck problem:  
allows decoder to look directly at source
- Attention helps with the vanishing gradient problem:  
Provides shortcut to faraway states

# GRU+Attention for Machine Translation

Source: An admitting privilege is the right of a doctor to admit a patient to a hospital or a medical centre to carry out a diagnosis or a procedure, based on his status as a health care worker at a hospital.

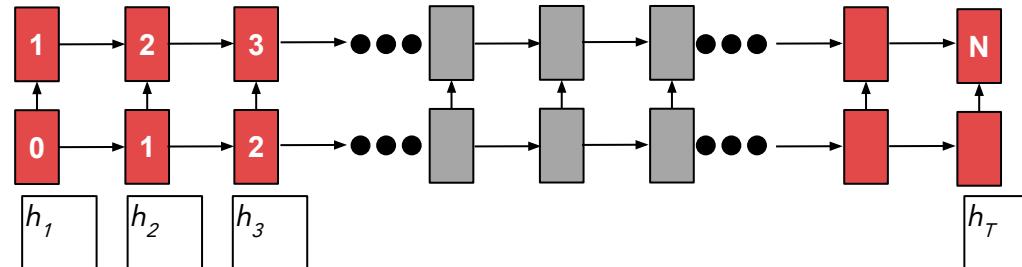
Reference: Le privilège d'admission est le droit d'un médecin, en vertu de son statut de membre soignant d'un hôpital, d'admettre un patient dans un hôpital ou un centre médical afin d'y délivrer un diagnostic ou un traitement.

RNNenc-50: Un privilège d'admission est le droit d'un médecin de reconnaître un patient à l'hôpital ou un centre médical d'un diagnostic ou de prendre un diagnostic en fonction de son état de santé.

RNNsearch-50: Un privilège d'admission est le droit d'un médecin d'admettre un patient à un hôpital ou un centre médical pour effectuer un diagnostic ou une procédure, selon son statut de travailleur des soins de santé à l'hôpital.

# RNN limit 3: Parallelization

- Forward and backward passes have  $O(\text{sequence length})$  unparallelizable operations
- GPUs can perform many independent computations (like addition) at once!
- But future RNN hidden states can't be computed in full before past RNN hidden states have been computed.
- Training and inference are slow; inhibits on very large datasets!

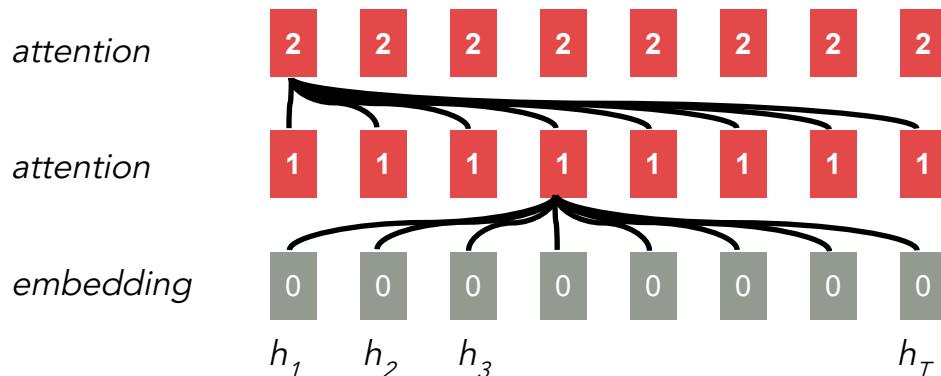


Numbers indicate min # of steps before a state can be computed

# Break for questions and "appel"

# "Attention is all you need": Transformers

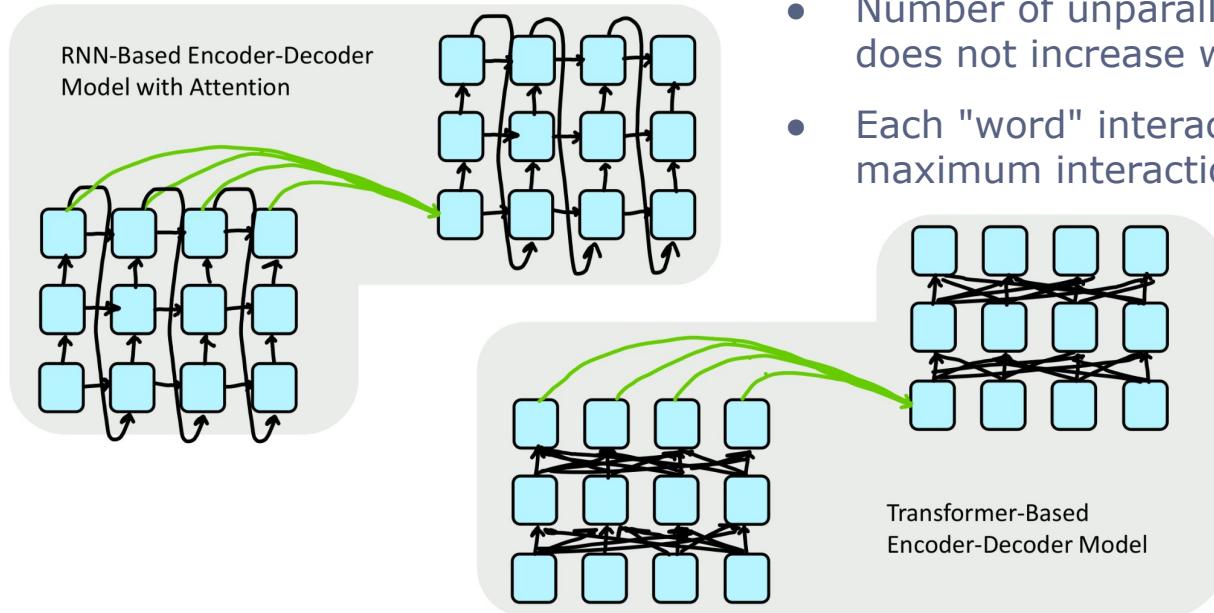
- Keeps only the attention mechanism, removes RNN
- Attention treats each token's representation as a query to access and incorporate information from a set of values.
- Number of unparallelizable operations does NOT increase with sequence length.



Maximum interaction distance:  $O(1)$ , since all tokens interact at every layer!

All tokens attend to all tokens in previous layer; most arrows here are omitted

# Recurrence vs. Attention



- Number of unparallelizable operations does not increase with sequence length.
- Each "word" interacts with each other, so maximum interaction distance is  $O(1)$ .

# Deep Learning is made out of GPUs

## Zuckerberg's Meta Is Spending Billions to Buy 350,000 Nvidia H100 GPUs

In total, Meta will have the compute power equivalent to 600,000 Nvidia H100 GPUs to help it develop next-generation AI, says CEO Mark Zuckerberg.



By Michael Kan January 18, 2024 [f](#) [X](#) [g](#) ...



(David Paul Morris/Bloomberg via Getty Images)

Paul Lerner – February 2025

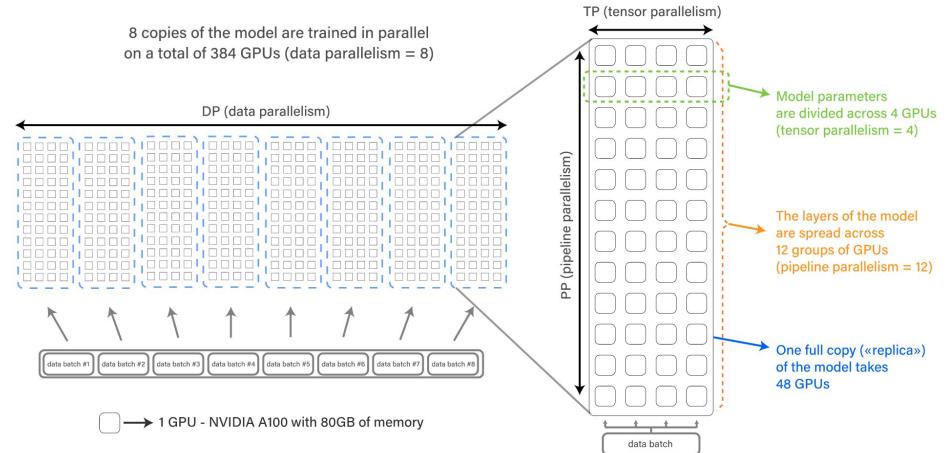


# Deep Learning is made out of GPUs

Nvidia



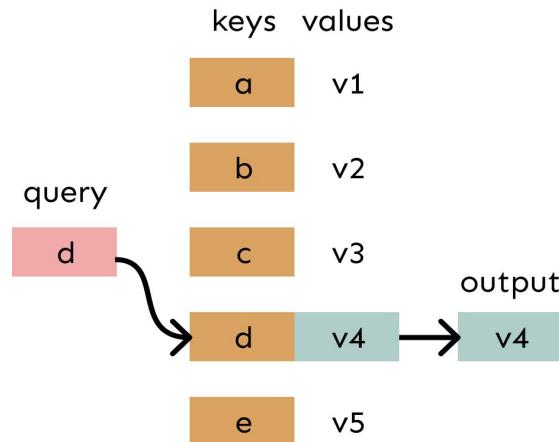
+38,400% in 12 years



# Attention as a soft, averaging lookup table

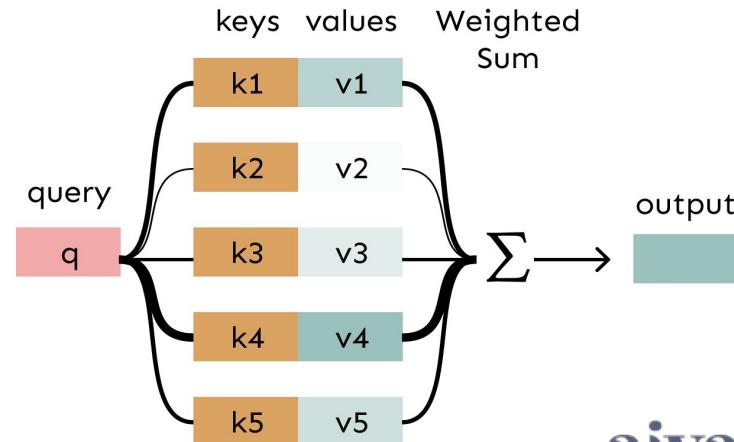
We can think of attention as performing fuzzy lookup in a key-value store.

In a lookup table, we have a table of **keys** that map to **values**. The **query** matches one of the keys, returning its value.



Paul Lerner – February 2025

In attention, the **query** matches all **keys** *softly*, to a weight between 0 and 1. The keys' **values** are multiplied by the weights and summed.



# Self-Attention: Basic Concepts

Each vector receives three representations (“roles”)

$$[W_Q] \times [ ] = [ ]$$

**Query:** vector **from** which the attention is looking

*“Hey there, do you have this information?”*

$$[W_K] \times [ ] = [ ]$$

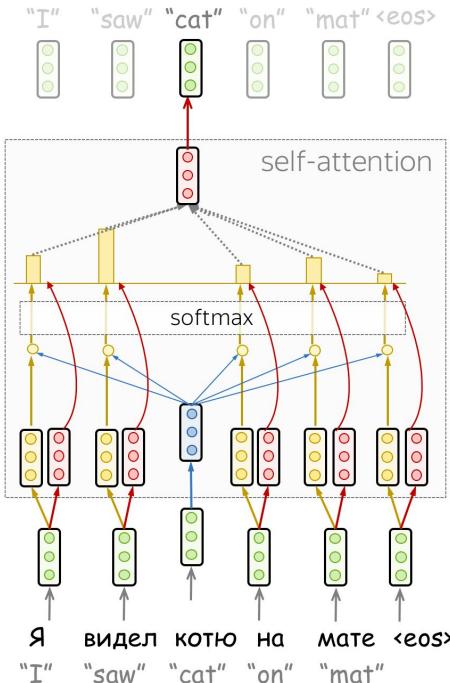
**Key:** vector **at** which the query looks to compute weights

*“Hi, I have this information – give me a large weight!”*

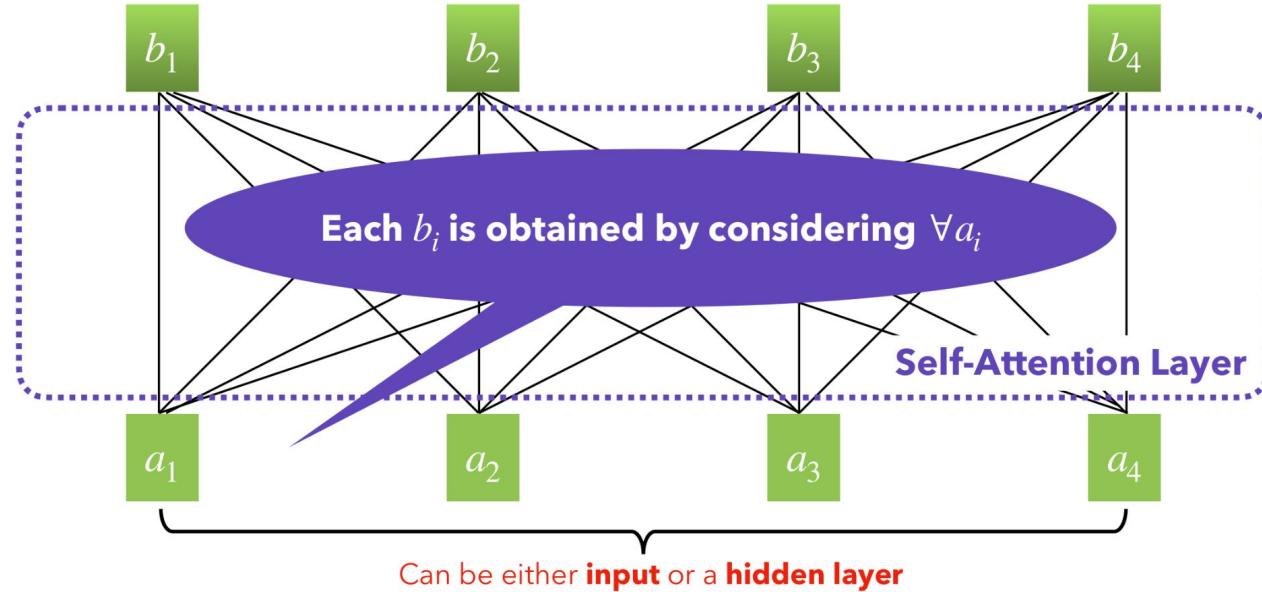
$$[W_V] \times [ ] = [ ]$$

**Value:** their weighted sum is attention output

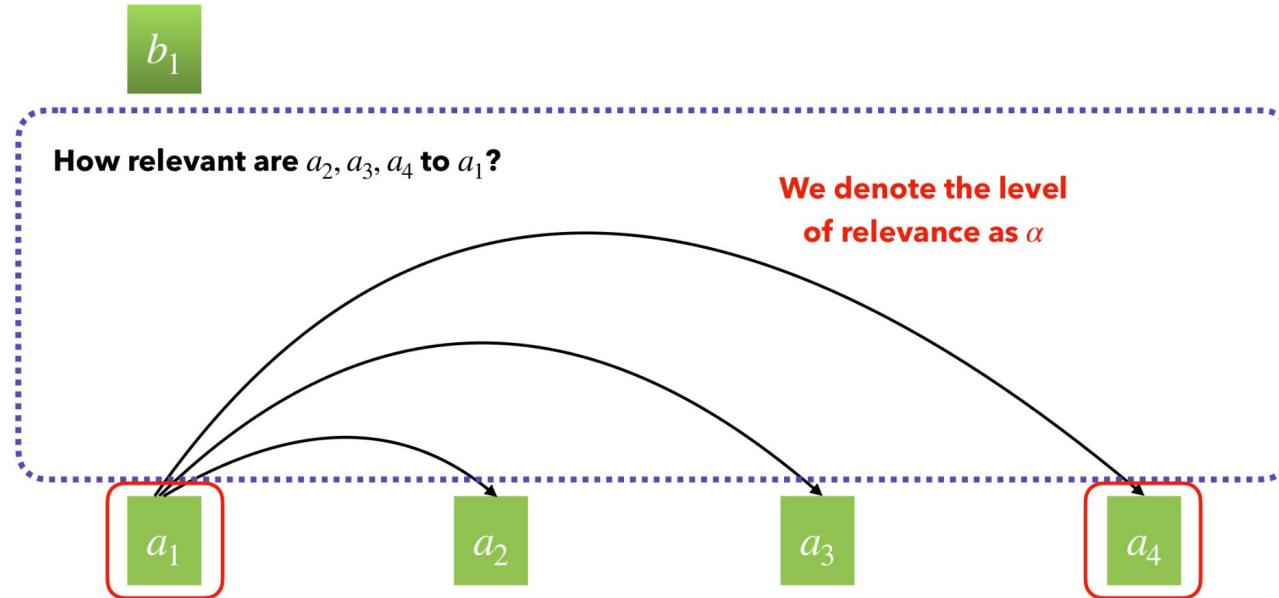
*“Here’s the information I have!”*



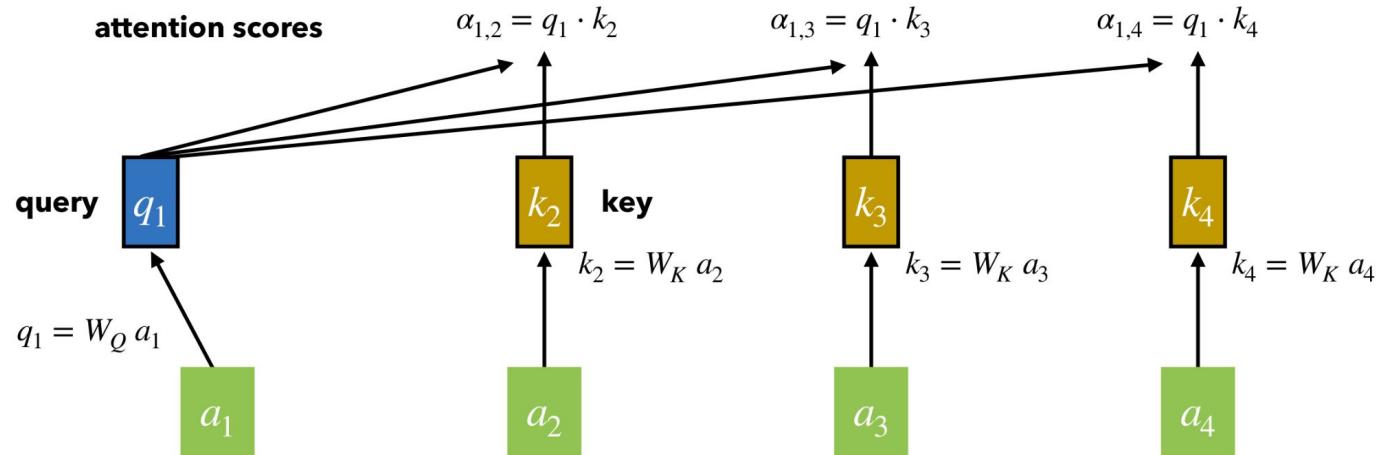
# Self-Attention: Walk-through



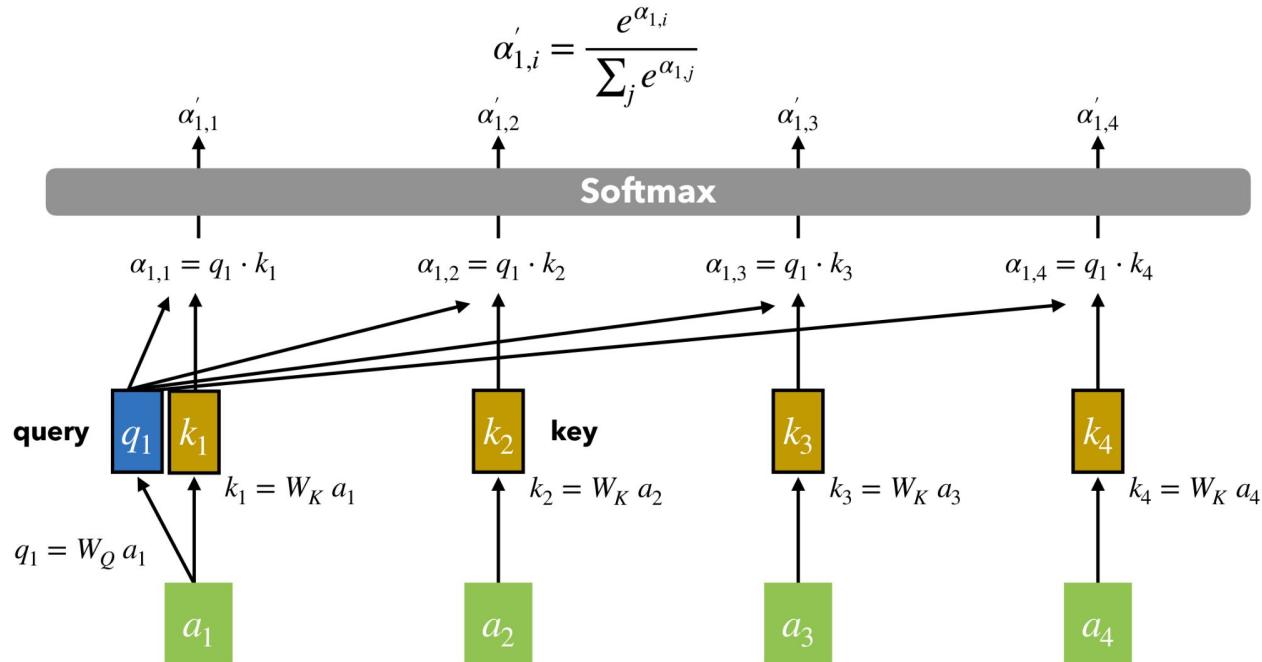
# Self-Attention: Walk-through



# Self-Attention: Walk-through



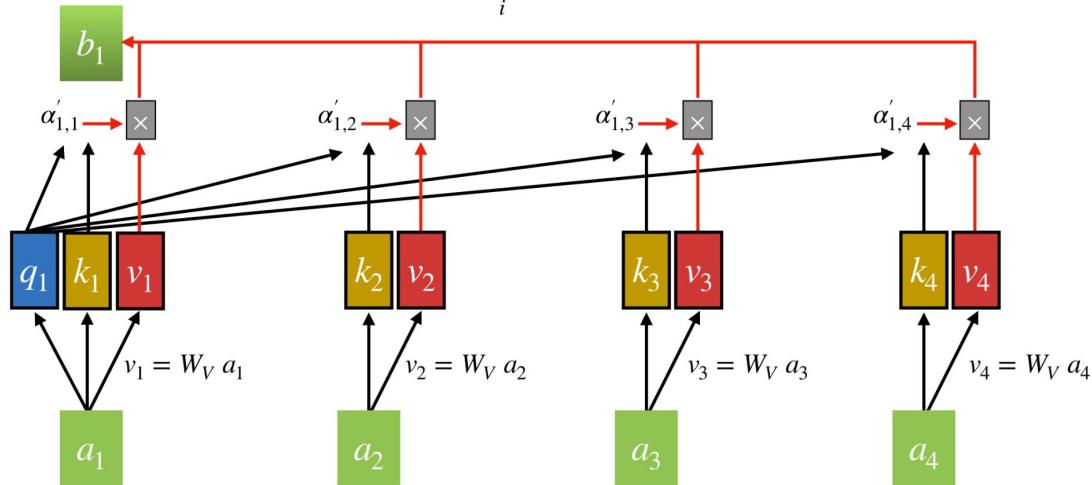
# Self-Attention: Walk-through



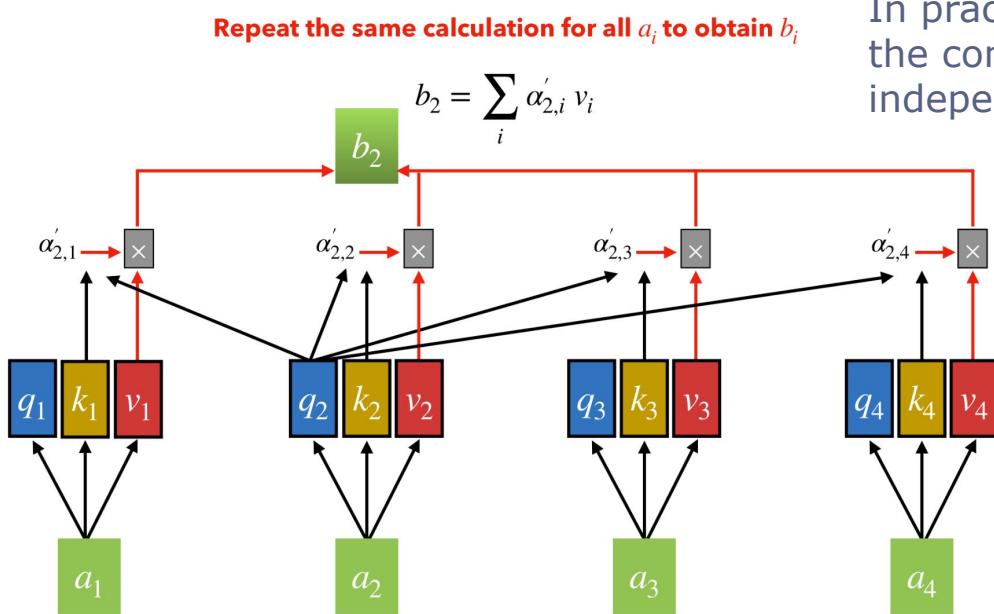
# Self-Attention: Walk-through

Use attention scores to extract information

$$b_1 = \sum_i \alpha'_{1,i} v_i$$

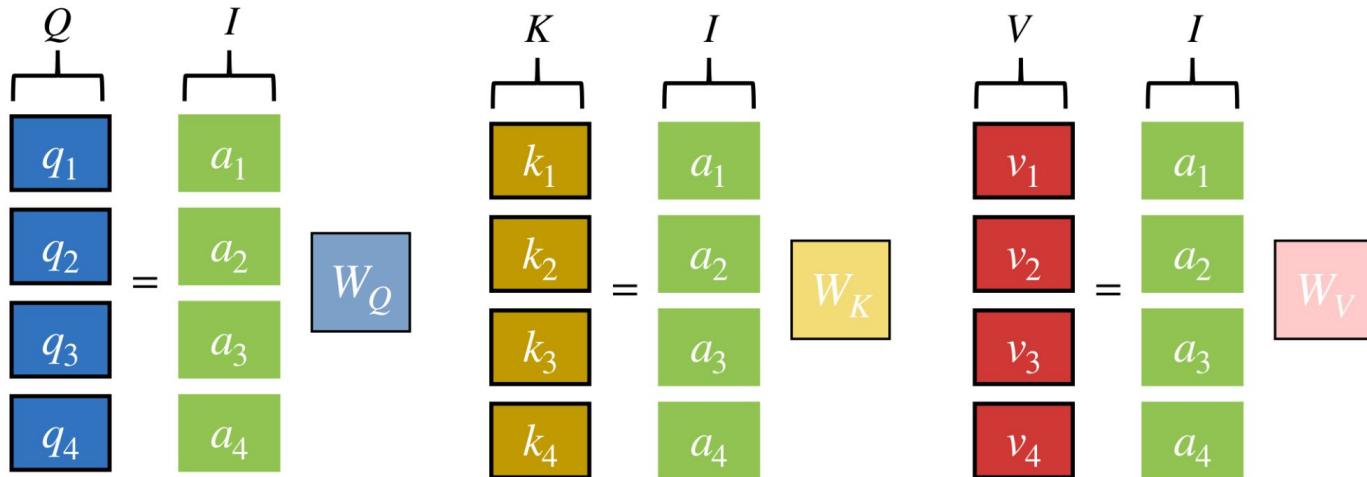


# Self-Attention: Walk-through

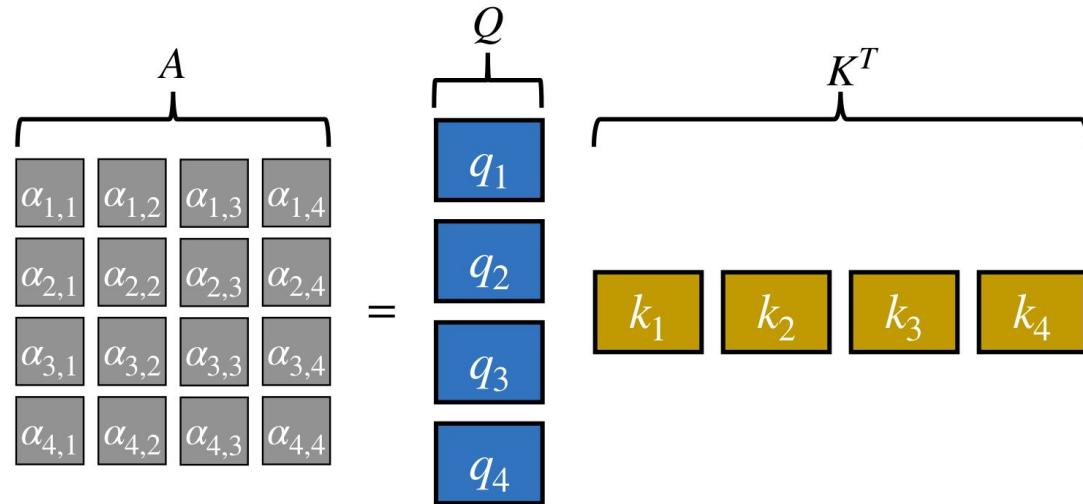


In practice this is done in **parallel**:  
the computation of  $b_1$  is  
independent of  $b_2$

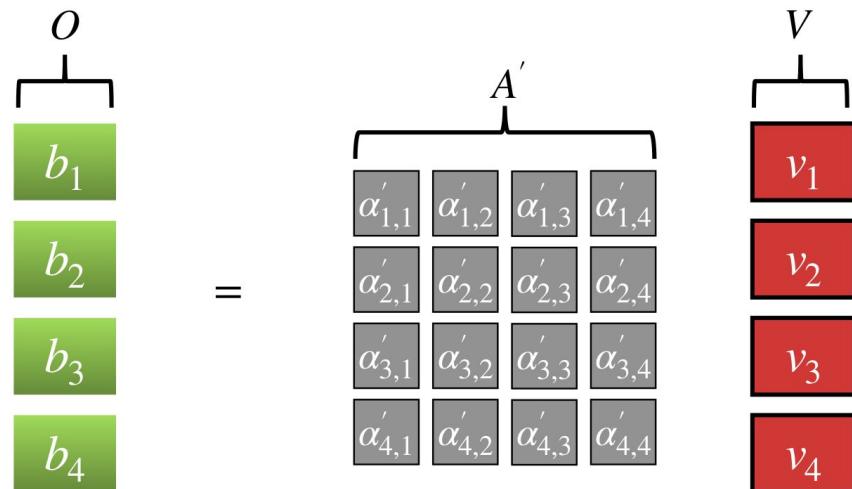
# Self-Attention: in parallel



# Self-Attention: in parallel



# Self-Attention: in parallel

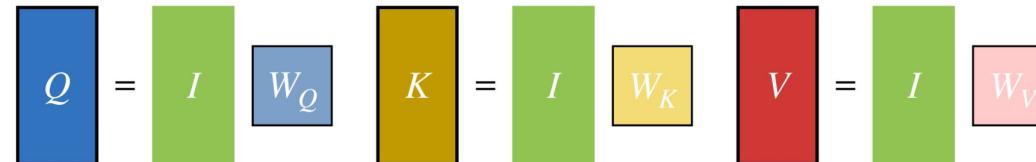


# Self-Attention: in parallel

$$Q = I W_Q$$

$$K = I W_K$$

$$V = I W_V$$

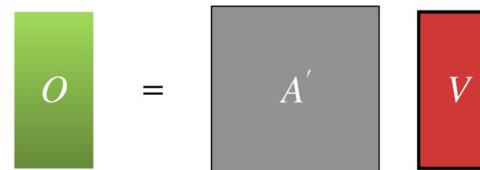


$$A = Q K^T$$

$$A = I W_Q (I W_K)^T = I W_Q W_K^T I^T$$

$$A' = \text{softmax}(A)$$

$$O = A' V$$



# Self-Attention: formally

$$Q = I W_Q$$

$$K = I W_K$$

$$V = I W_V$$

$$\left. \begin{array}{l} I = \{a_1, \dots, a_n\} \in \mathbb{R}^{n \times d}, \text{ where } a_i \in \mathbb{R}^d \\ W_Q, W_K, W_V \in \mathbb{R}^{d \times d} \\ Q, K, V \in \mathbb{R}^{n \times d} \end{array} \right\}$$

$$A = Q K^T$$

$$A = I W_Q (I W_K)^T = I W_Q W_K^T I^T$$

$$A' = \text{softmax}(A)$$

$$\left. \begin{array}{l} A', A \in \mathbb{R}^{n \times n} \end{array} \right\}$$

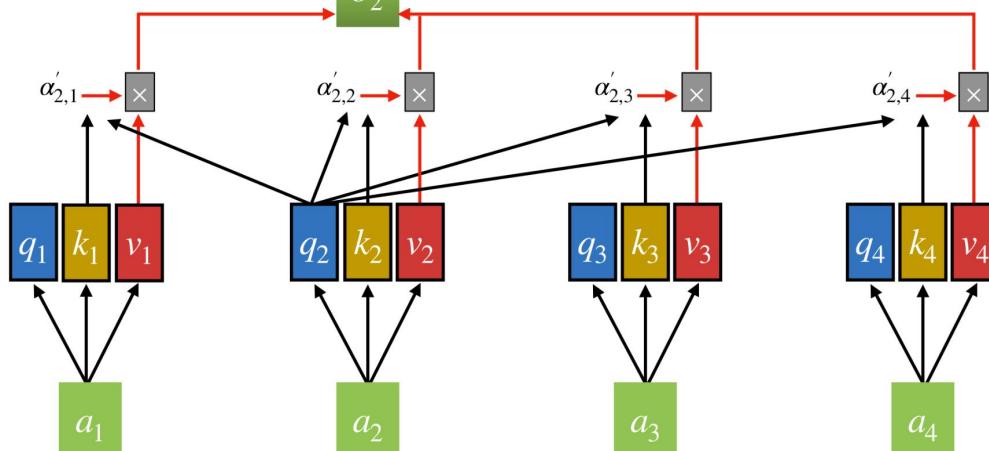
$$O = A' V$$

$$\left. \begin{array}{l} O \in \mathbb{R}^{n \times d} \end{array} \right\}$$

# Permutation-invariant: Transformer = Bag of Word?

**b\_2** sums over all **a\_i**, does not matter if **a\_1** comes before/after **a\_2**

$$b_2 = \sum_i \alpha'_{2,i} v_i$$



$x_1$ : yes , we have no bananas

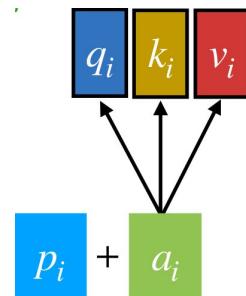
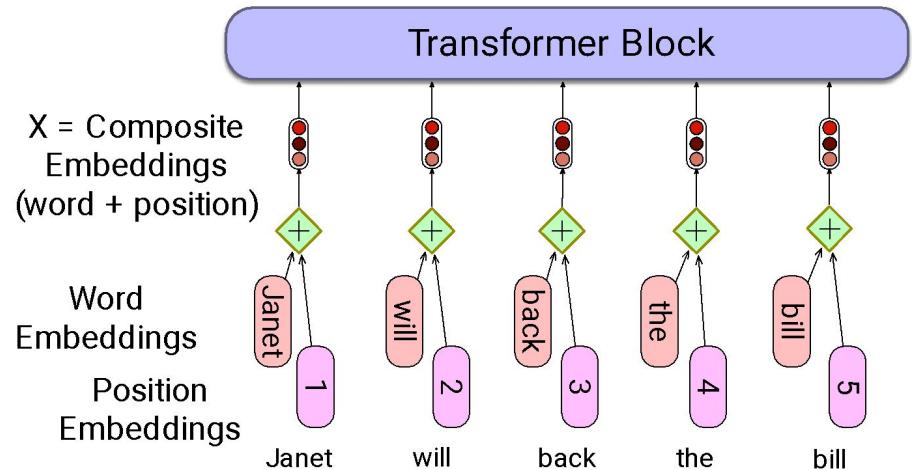
$x_2$ : say yes for bananas

$x_3$ : no bananas , we say

	1	2	3
,	1	0	1
bananas	1	1	1
for	0	1	0
have	1	0	0
no	1	0	1
say	0	1	1
we	1	0	1
yes	1	1	0

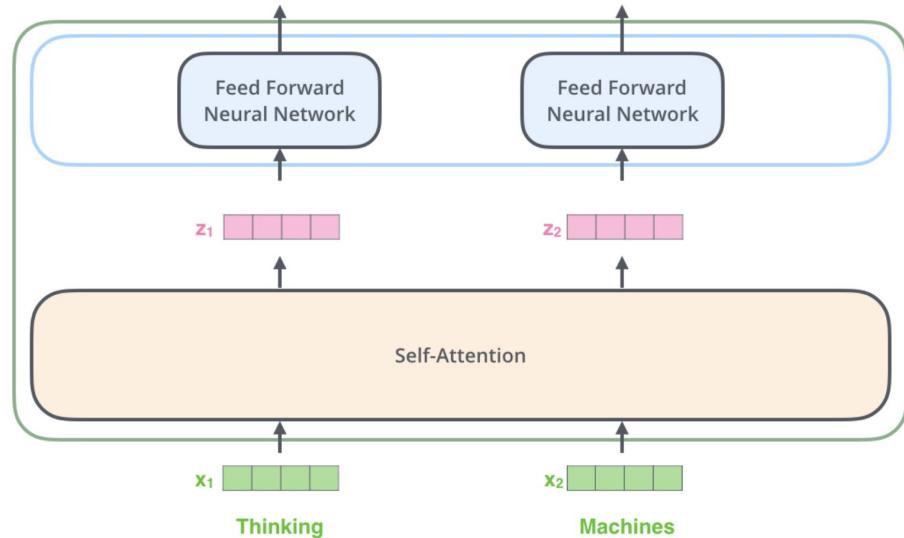
# Position Encoding

- Most basic method: At the 1st layer, add an embedding of the position to the word embedding (BERT, GPT-3)
- Typically initialized randomly and learned like any other parameter of the model
- Despite adding position, several papers argue that Transformer models are **permutation invariant** / do not model the order of words
- More recent methods we won't cover modify self-attention (RoPE, ALiBi)



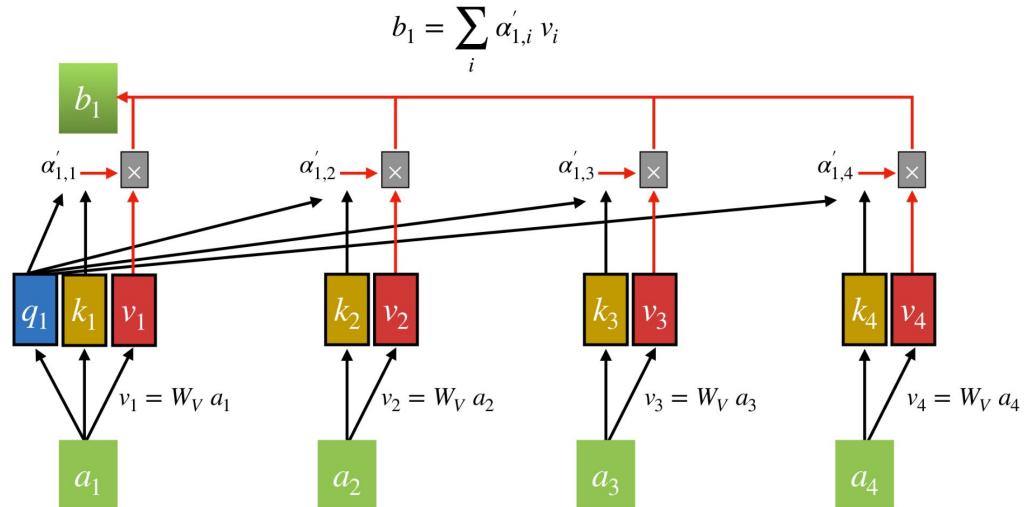
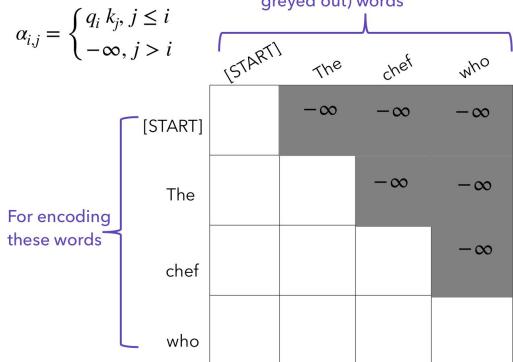
# Attention is almost all you need

- Since there are no element-wise non-linearities, self-attention is simply performing a re-averaging of the value vectors
- Apply a feedforward layer to the output of attention, providing non-linear activation (and additional expressive power)



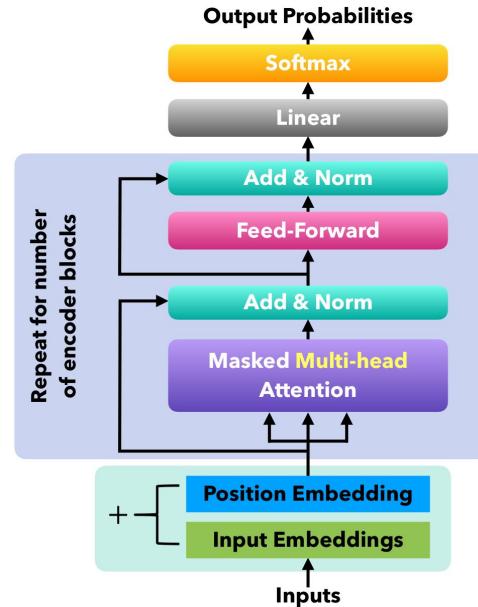
# Self-attention? What about causality?

- $b_1$  depends on  $a_2$  and  $a_3$ ... but the goal is the **generate  $a_2$**  and  $a_3$
- Mask** attention scores of future words



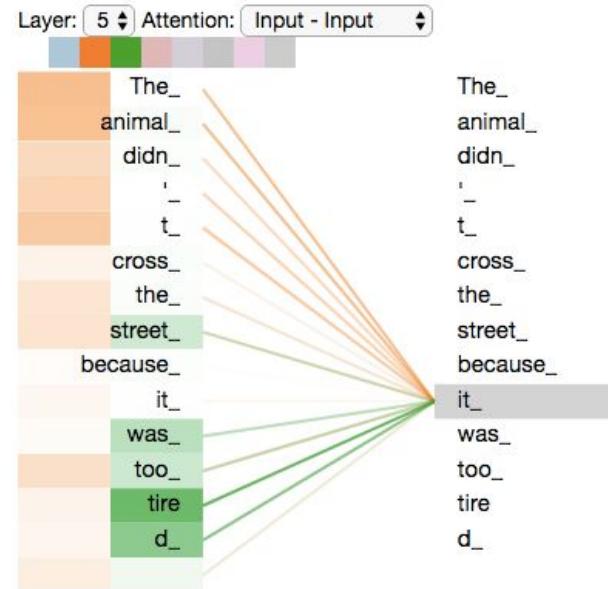
# Putting the pieces together

- Positional Encoding: otherwise permutation-invariant
- (*Multi-head*) Self-attention: essential part to model relations between words
  - masked for causality
- *Residual Connection*: for stable training/deeper networks
- Feedforward for nonlinearity/expressiveness
- Linear/softmax: output layer back to vocabulary dimension

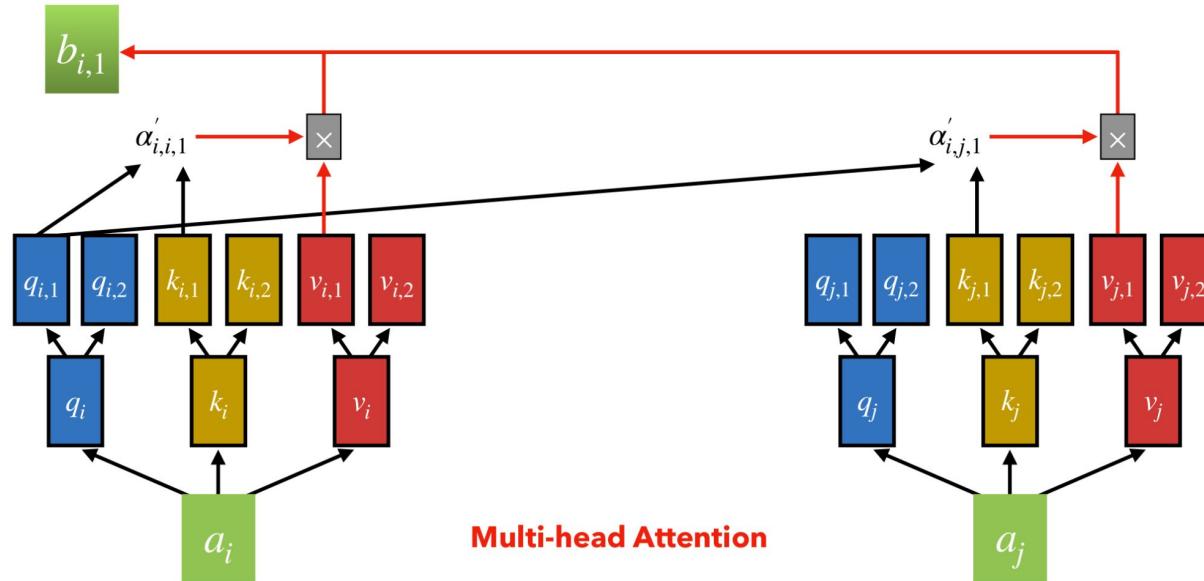


# Why Multi-Head Attention?

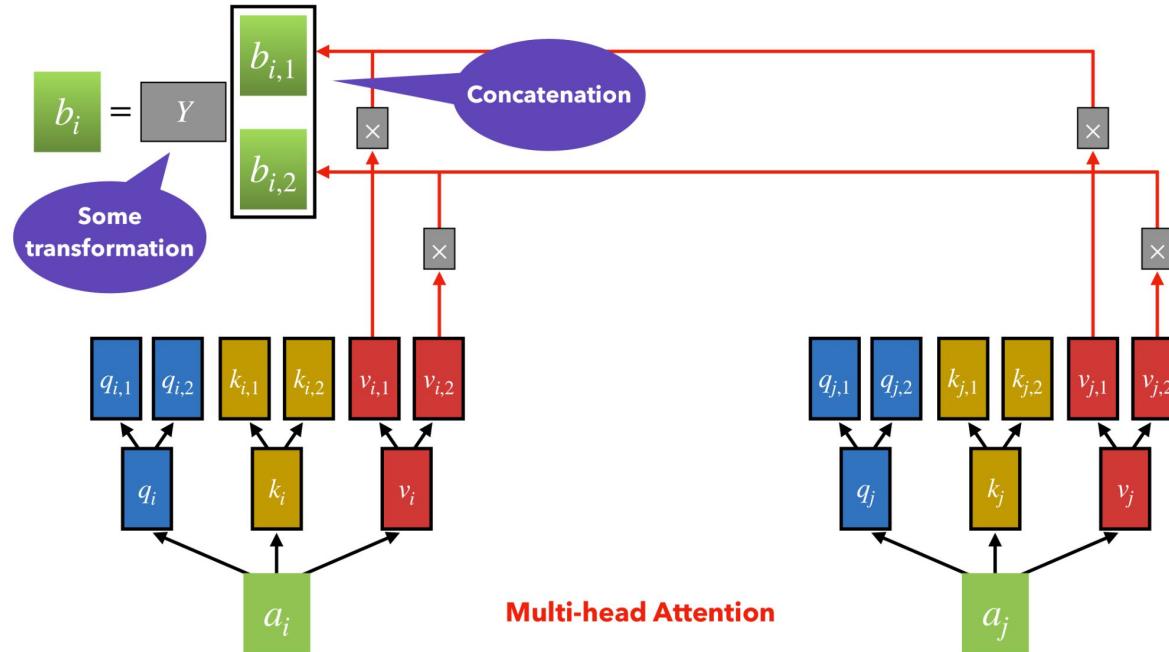
- What if we want to look in multiple places in the sentence at once?
- For word  $i$ , self-attention “looks” where  $x_i Q_i K_j$  is high, but maybe we want to focus on different  $j$  for different reasons?
- orange head for the coreference "The animal"
- green head for the past participle "tired"



# Multi-Head Attention: Walk-through



# Multi-Head Attention: Walk-through



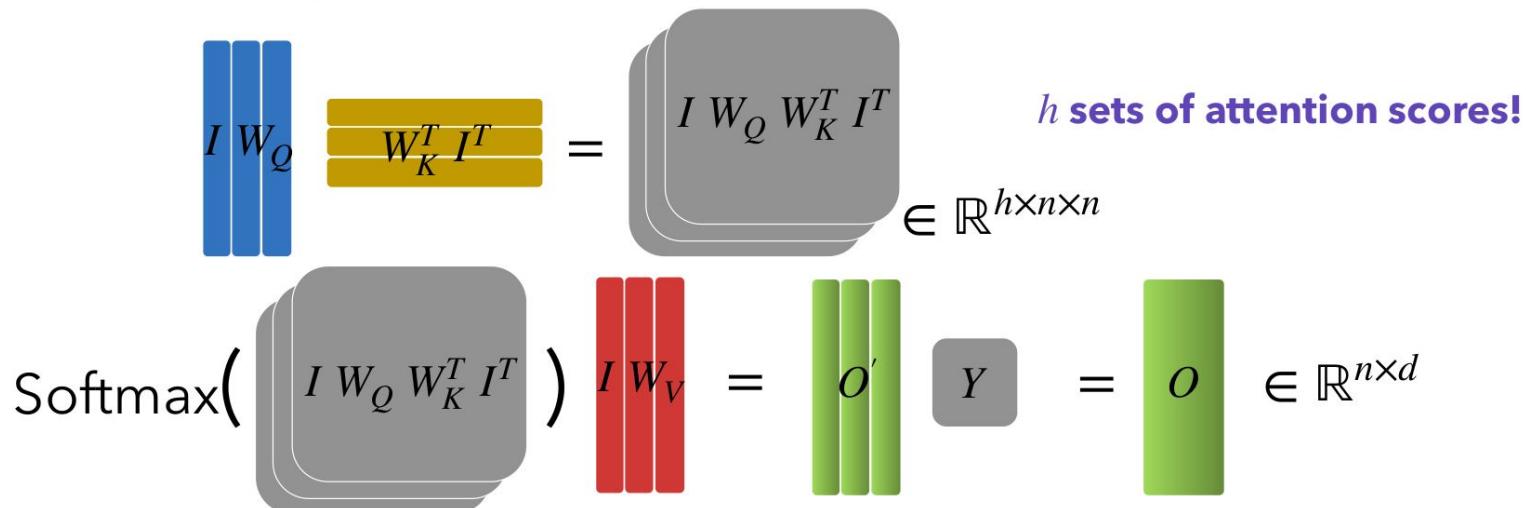
# Multi-Head Attention: formally

$$\begin{aligned}
 Q^l &= I W_Q^l & I &= \{a_1, \dots, a_n\} \in \mathbb{R}^{n \times d}, \text{ where } a_i \in \mathbb{R}^d \\
 K^l &= I W_K^l & W_Q^l, W_K^l, W_V^l &\in \mathbb{R}^{d \times \frac{d}{h}} & \text{Multiple attention "heads" can be} \\
 V^l &= I W_V^l & Q^l, K^l, V^l &\in \mathbb{R}^{n \times \frac{d}{h}} & \text{defined via multiple } \mathbf{W}^* \text{ matrices} \\
 A^l &= Q^l K^l{}^T & A^{l'}, A^l &\in \mathbb{R}^{n \times n} \\
 A^{l'} &= \text{softmax}(A^l) & O^l &\in \mathbb{R}^{n \times \frac{d}{h}} & \text{Each attention head performs attention independently} \\
 O^l &= A^{l'} V^l & Y &\in \mathbb{R}^{d \times d} \\
 O &= [O^1; \dots; O^h] Y & [O^1; \dots; O^h] &\in \mathbb{R}^{n \times d} & \text{Their results are concatenated} \\
 && O &\in \mathbb{R}^{n \times d}
 \end{aligned}$$

# Multi-Head Attention: in parallel (as always)

compute  $I W_Q \in \mathbb{R}^{n \times d}$ , and then reshape to  $\mathbb{R}^{n \times h \times \frac{d}{h}}$

Then we transpose to  $\mathbb{R}^{h \times n \times \frac{d}{h}}$ ; **now the head axis is like a batch axis**



# Residual Connections for stable training

- Residual connections are a trick to help models train better.
  - Instead of  $X^{(i)} = \text{Layer}(X^{(i-1)})$  (where  $i$  represents the layer)

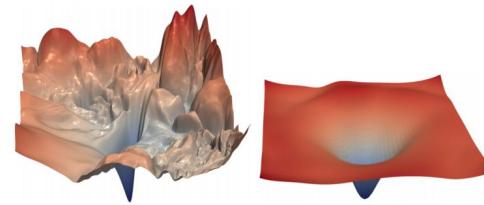


- We let  $X^{(i)} = X^{(i-1)} + \text{Layer}(X^{(i-1)})$  (so we only have to learn “the residual” from the previous layer)



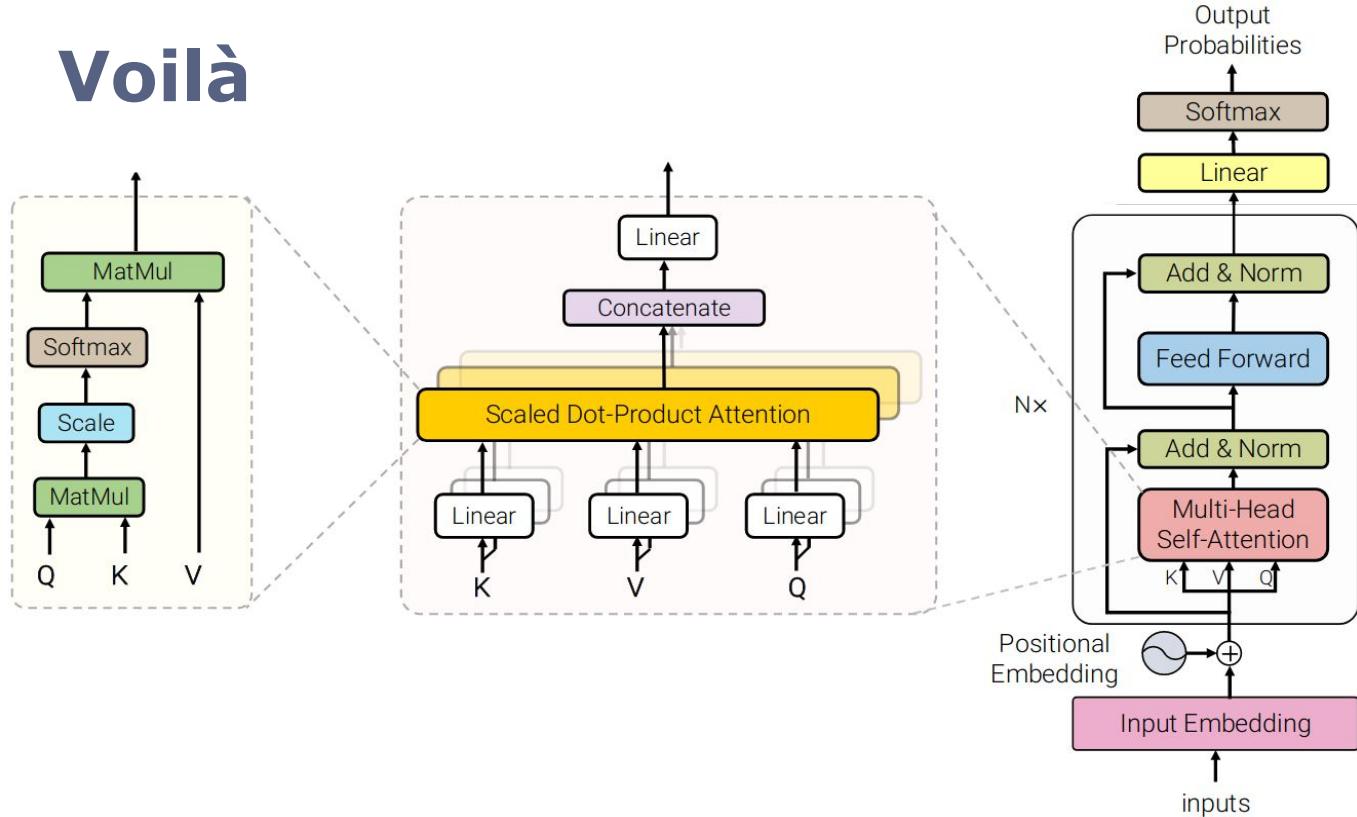
- Gradient is great through the residual connection; it's 1!
- Bias towards the identity function!

Remember the Cell state in LSTM



[Loss landscape visualization,  
Li et al., 2018, on a ResNet]

# Voilà



# Transformer for Machine Translation

Source: An admitting privilege is the right of a doctor to admit a patient to a hospital or a medical centre to carry out a diagnosis or a procedure, based on his status as a health care worker at a hospital.

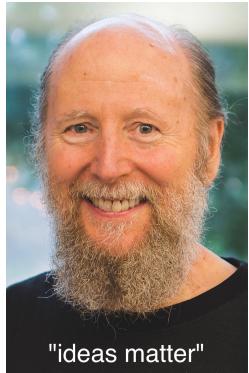
Reference: Le privilège d'admission est le droit d'un médecin, en vertu de son statut de membre soignant d'un hôpital, d'admettre un patient dans un hôpital ou un centre médical afin d'y délivrer un diagnostic ou un traitement.

RNNsearch-50: Un privilège d'admission est le droit d'un médecin d'admettre un patient à un hôpital ou un centre médical pour effectuer un diagnostic ou une procédure, selon son statut de travailleur des soins de santé à l'hôpital.

Transformer (fairseq wmt14.en-fr): Un privilège d'admission est le droit d'un médecin d'admettre un patient dans un hôpital ou un centre médical pour y effectuer un diagnostic ou une intervention, en fonction de son statut de travailleur de la santé dans un hôpital.

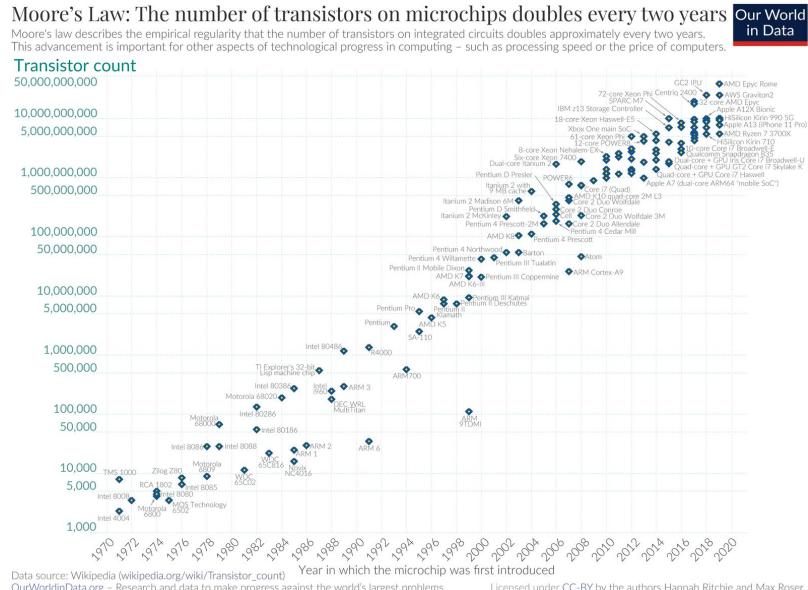
# Scaling Transformers (the bitter lesson)

**"The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin"**

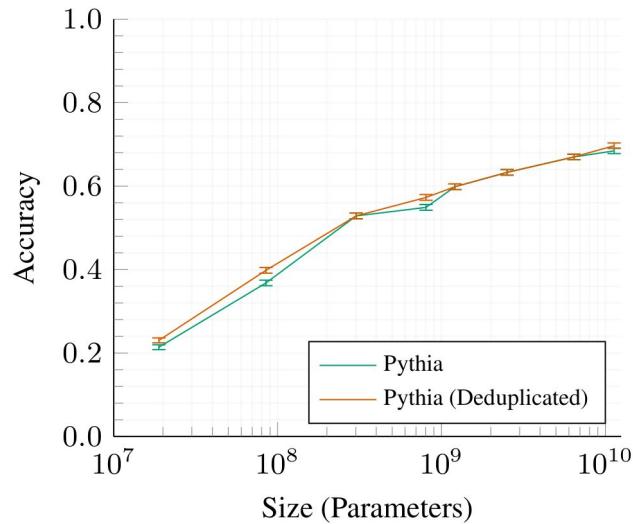


**Rich Sutton**

Paul Lerner – February 2025

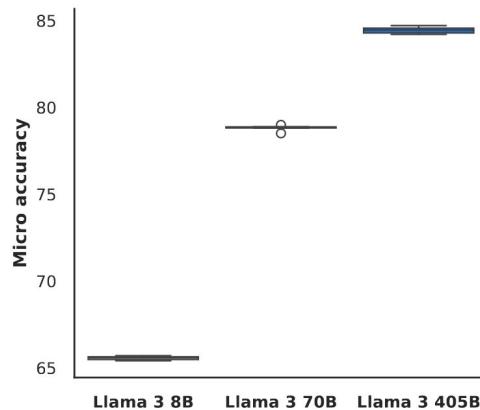


# Scaling Transformers (the bitter lesson)

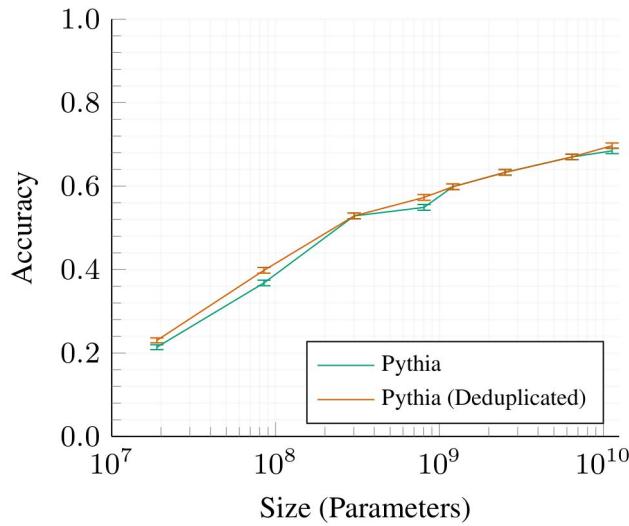


(a) LAMBADA (OpenAI)

There seems to be no limit, from millions to billions to trillions of parameters

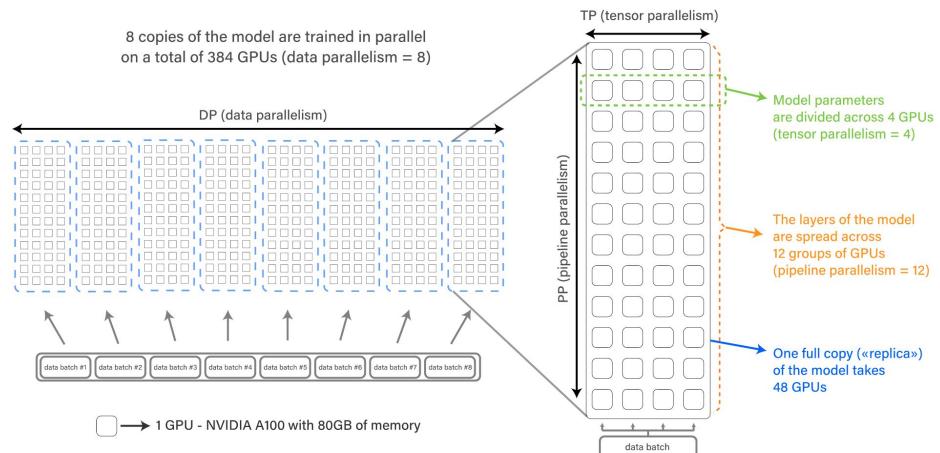


# Scaling Transformers: buy more GPUs

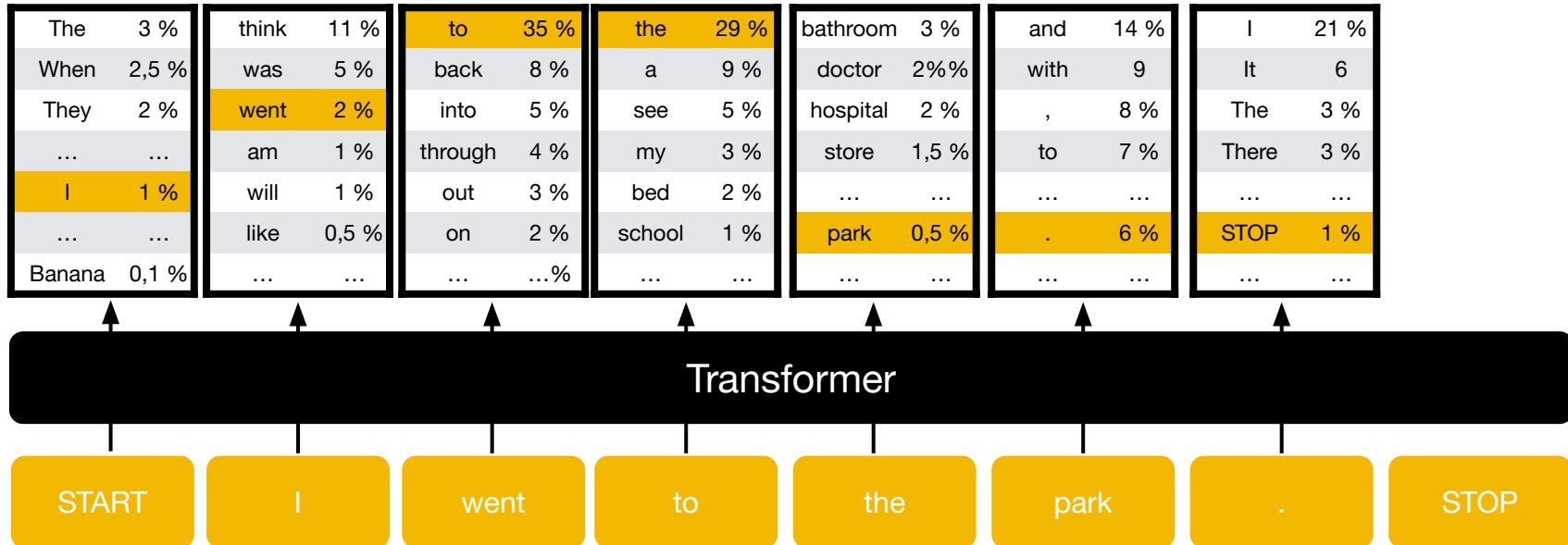


(a) LAMBADA (OpenAI)

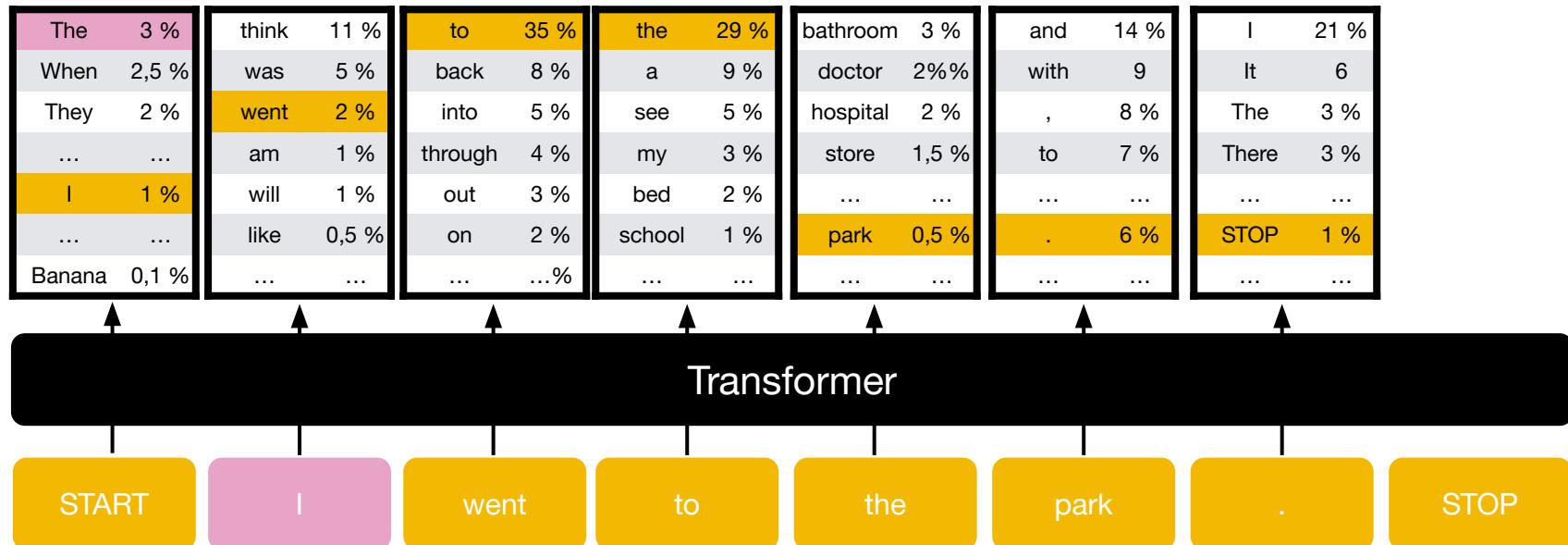
From  $8 \times 32\text{GB}$  V100 GPUs (RoBERTa)  
to  
 $16,384 \times 94\text{GB}$  H100 GPUs (Llama 3)



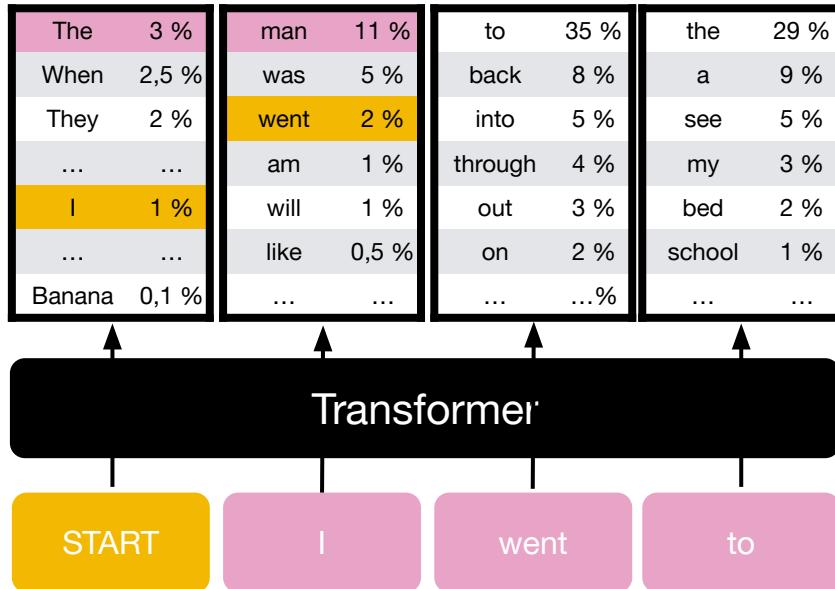
# Limit: Teacher Forcing



# Train-test mismatch: exposure bias



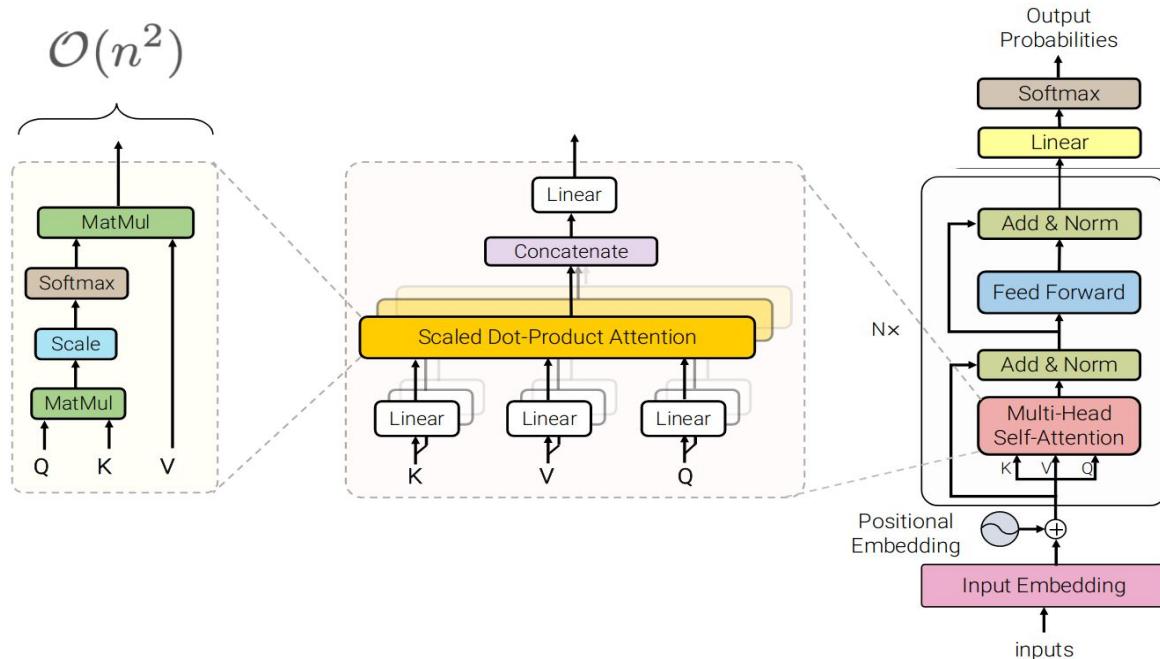
# Not easily solved, unlike for RNN



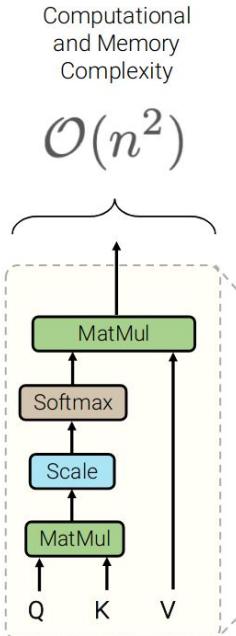
- Because computation is done in parallel, we cannot access the generation of the model
- Teacher forcing is done systematically, models are subject to exposure bias

Mihaylova, T., & Martins, A. F. T. (2019).  
Scheduled Sampling for Transformers.

# Limit: Quadratic complexity



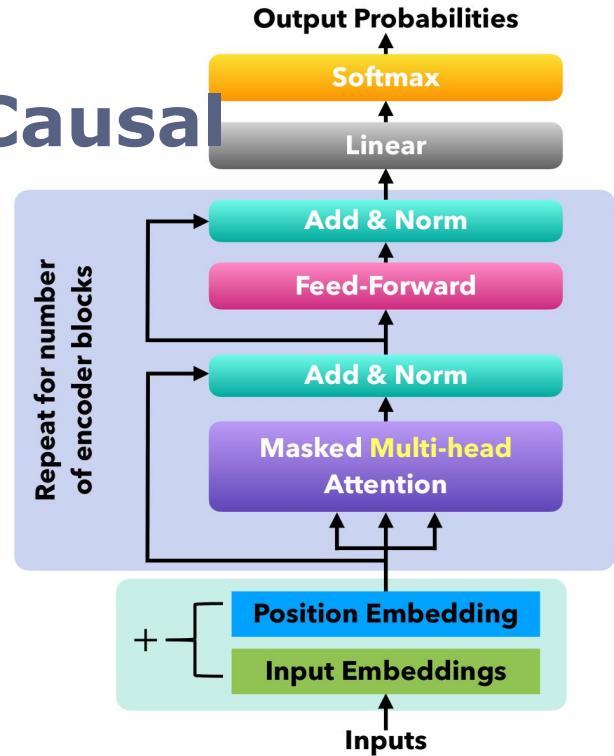
# Limit: Quadratic complexity



- Computing all pairs of interactions means our computation grows **quadratically** with the sequence length!
- for recurrent models, it only grew **linearly**
- Large body of work on this question (Tay et al., 2020)  
"Efficient Transformers: A Survey"
- But vanilla Transformer still used in state-of-the-art LLMs

# Transformer A: (Autoregressive) Decoder/Causal

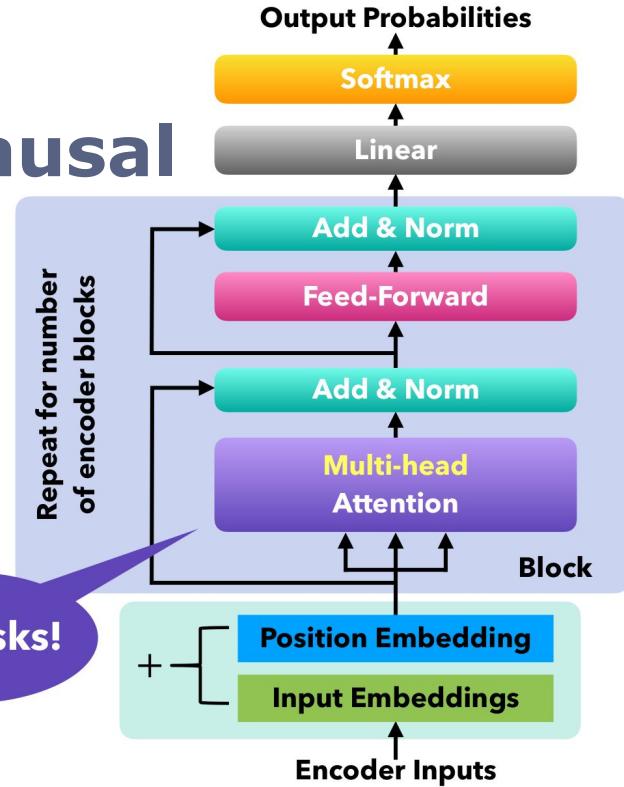
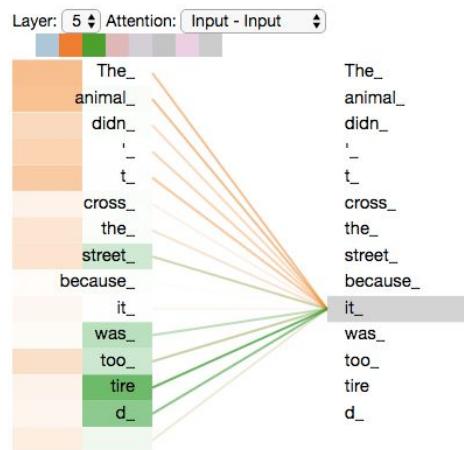
- Described previously: main architecture for LLMs (**GPT-3**, Llama-\*, and many many many more)
- **Causal/unidirectional mask:**  
can only see past words
- First purpose: **Language Modeling** / autoregressive generation
- But now every task of NLP is cast as Language Modeling, even classification



**Transformer B:  
(Bidirectional) Encoder/non-causal**

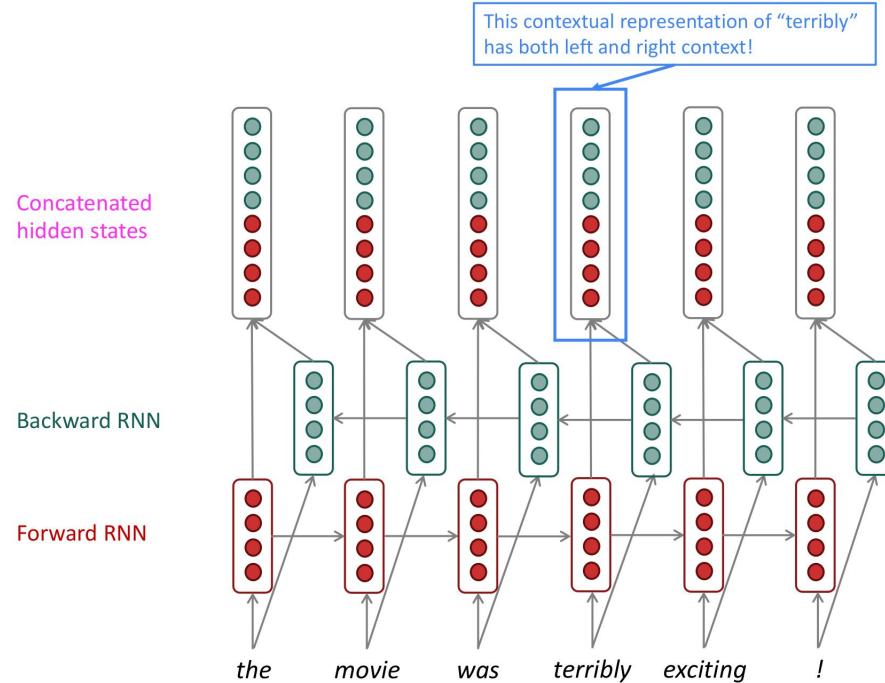
- Removes the mask from self-attention: now every word can see future and past
  - Use for classification (but now words have a better context, unlike bag of words)
  - Famous examples:  
**BERT, mBERT,**  
**RoBERTa,**  
**DeBERTa,**  
**CamemBERT, ...**

The diagram illustrates the attention mechanism at layer 5. It shows two sequences: an input sequence ('The animal didn't cross the street') and an output sequence ('The animal didn't - cross the street'). The input sequence is represented by colored squares: orange for 'The', green for 'animal', grey for 'didn't', light blue for 'cross', and light orange for 'the'. The output sequence is also represented by colored squares. Orange lines connect the input words to their corresponding output words. The word 'street' in the input has a green line connecting it to the same word in the output, indicating it has the highest attention weight to itself.

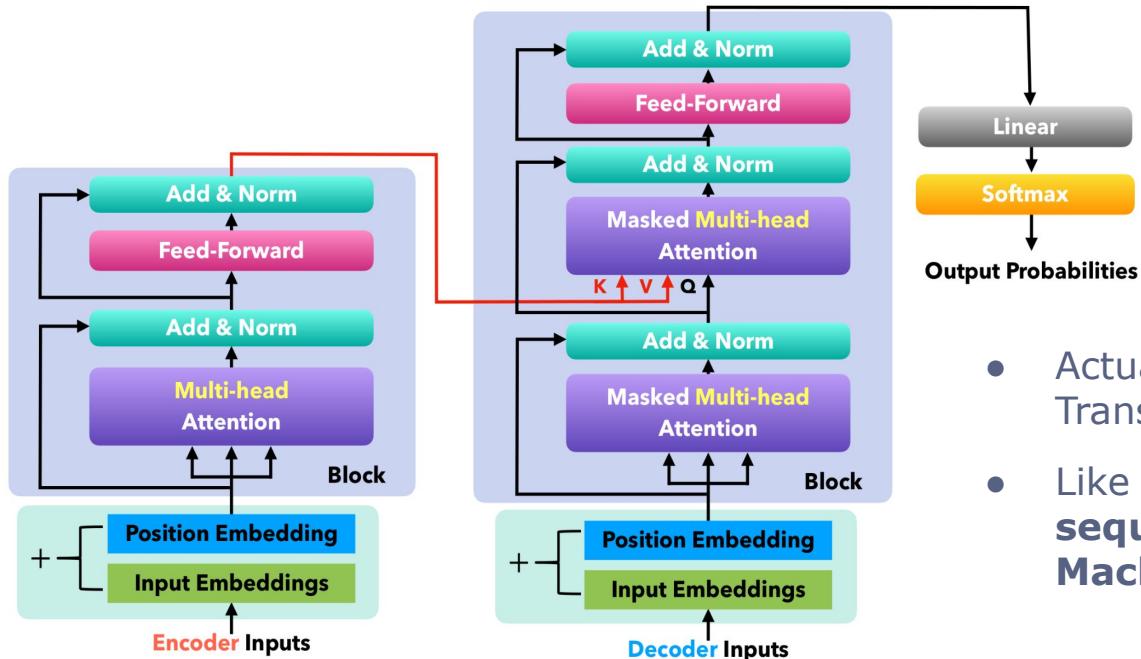


# Note on Bidirectional RNNs

- RNNs could also be bidirectional but you then needed two!
- → long sequential (unparallelizable) operations, although still  $O(n)$  theoretically



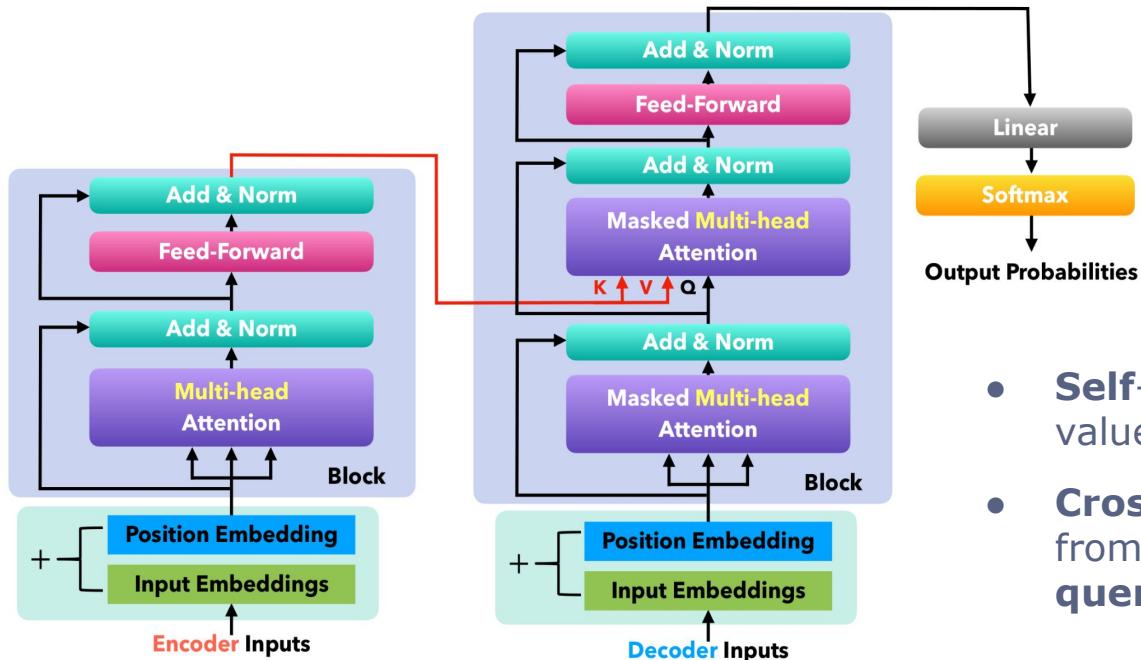
# Transformer C: Encoder-Decoder



- Famous examples:  
T5, BART,  
BARTez, ...

- Actually the first variant proposed for Translation by Vaswani et al. (2017)
- Like an RNN Encoder-Decoder, use for **sequence-to-sequence** tasks like **Machine Translation**

# Cross-Attention in Encoder-Decoder

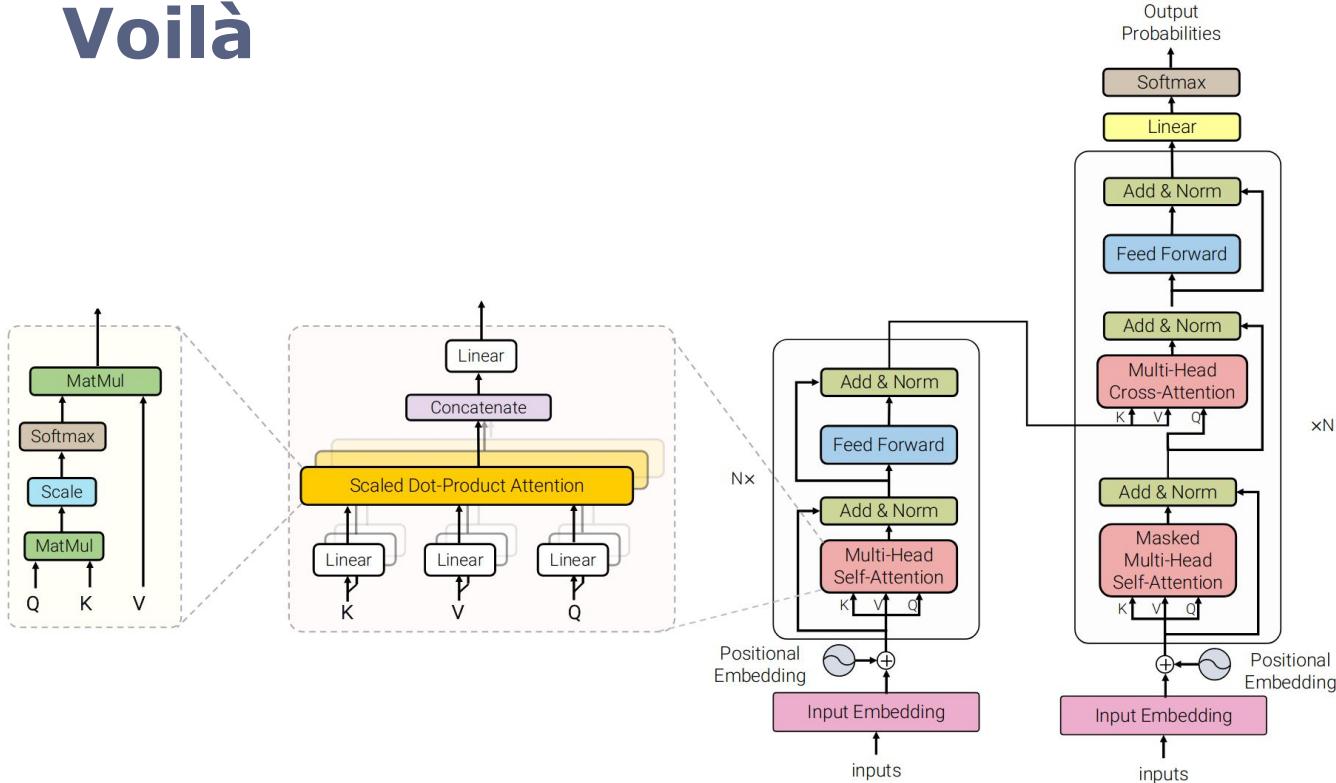


- **Self-attention:** queries, keys, and values come **from the same source**
- **Cross-Attention:** keys and values are from *Encoder* (like a memory); **queries** are from **Decoder**

# Summarizing

- Use Transformers for:
  - Sequence-to-sequence (e.g. Translation)
  - Natural Language Generation/Language Modeling
  - Actually anything: encode text (with context) then classify
- Attention was invented for RNNs and Translation, to focus on parts of the source sentence: relieves information bottleneck
- Transformer leverages attention in a parallelized way:  
scales as long as you buy enough GPUs

# Voilà

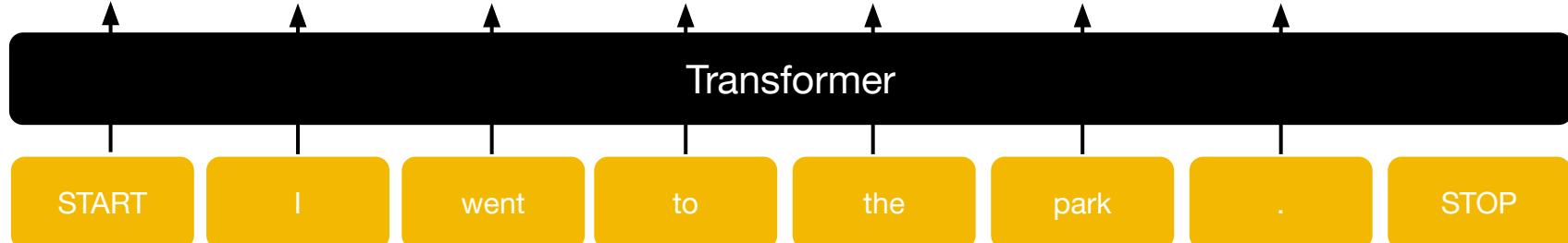


# Homework!

# Next class: (*Pre-*)Training Transformers (LLMs)

$p(x|\text{START})p(x|\text{START I})p(x|\dots\text{went}) \ p(x|\dots\text{to}) \ p(x|\dots\text{the}) \ p(x|\dots\text{park}) \ p(x|\text{START I went to the park.})$

The	3 %	think	11 %	to	35 %	the	29 %	bathroom	3 %	and	14 %	I	21 %
When	2,5 %	was	5 %	back	8 %	a	9 %	doctor	2 %	with	9	It	6
They	2 %	went	2 %	into	5 %	see	5 %	hospital	2 %	,	8 %	The	3 %
...	...	am	1 %	through	4 %	my	3 %	store	1,5 %	to	7 %	There	3 %
I	1 %	will	1 %	out	3 %	bed	2 %	...	...	...	...	...	...
...	...	like	0,5 %	on	2 %	school	1 %	park	0,5 %	.	6 %	STOP	1 %
Banana	0,1 %	...	...	...	...	...	...	...	...	...	...	...	...



# Acknowledgements

This class directly builds upon:

- **Jurafsky, D., & Martin, J. H.** (2024). Speech and Language Processing : An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models (3rdéd.).
- **Eisenstein, J.** (2019). Natural Language Processing. 587.
- **Yejin Choi.** (Winter 2024). CSE 447/517: Natural Language Processing (University of Washington - Paul G. Allen School of Computer Science & Engineering)
- **Noah Smith.** (Winter 2023). CSE 447/517: Natural Language Processing (University of Washington - Paul G. Allen School of Computer Science & Engineering)
- **Benoit Sagot.** (2023-2024). *Apprendre les langues aux machines* (Collège de France)
- **Chris Manning.** (Spring 2024). Stanford CS224N: Natural Language Processing with Deep Learning
- Classes where I was/am Teacher Assistant:
  - **Christopher Kermorvant.** Machine Learning for Natural Language Processing (ENSAE)
  - **François Landes** and **Kim Gerdes.** Introduction to Machine Learning and NLP (Paris-Saclay)

Also inspired by:

- My PhD thesis: *Répondre aux questions visuelles à propos d'entités nommées* (2023)
- **Noah Smith** (2023): Introduction to Sequence Models (LxMLS)
- **Kyunghyun Cho:** Transformers and Large Pretrained Models (LxMLS 2023), Neural Machine Translation (ALPS 2021)
- My former PhD advisors **Olivier Ferret** and **Camille Guinaudeau** and postdoc advisor **François Yvon**
- My former colleagues at LISN



# aivancity

PARIS-CACHAN

**advancing education  
in artificial intelligence**