$0 \left(\frac{2 \cdot 0/d}{}\right)$  $\left(\frac{2 \cdot 1/d}{}\right)$  $\left(2 \cdot \frac{d-1}{2}\right)/d$

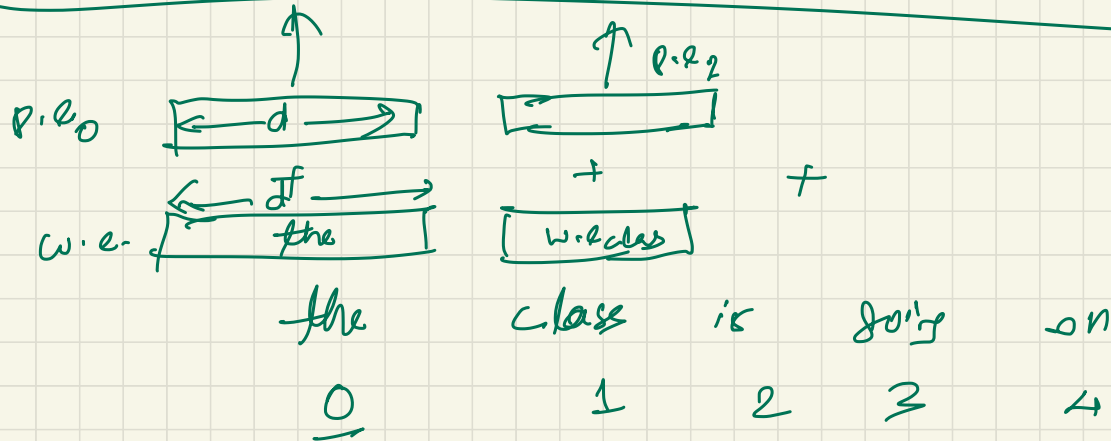Sin   Cos   2 dim   Sin      $i$        Sin   Cos   d

O

pos$^n$

$K_i$

sinusoidal

= not learnable

$\alpha = 512$
$d = 512$

$\alpha - 1$

$\alpha \times d$   values

Pos$^n$ embeddy

$$P(K, \ell) = Sin\left(\frac{K}{10000^{2 \cdot 1/d}}\right)$$

Encoder

$p.e_0$

$d$

$p.e_2$

$+$ $+$

$w.e$ $f$

the w.e class

the class is going on

0 1 2 3 4

# Multi-headed Self Attention Block



1) This is our input sentence*

2) We embed each word*

Thinking Machines

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one
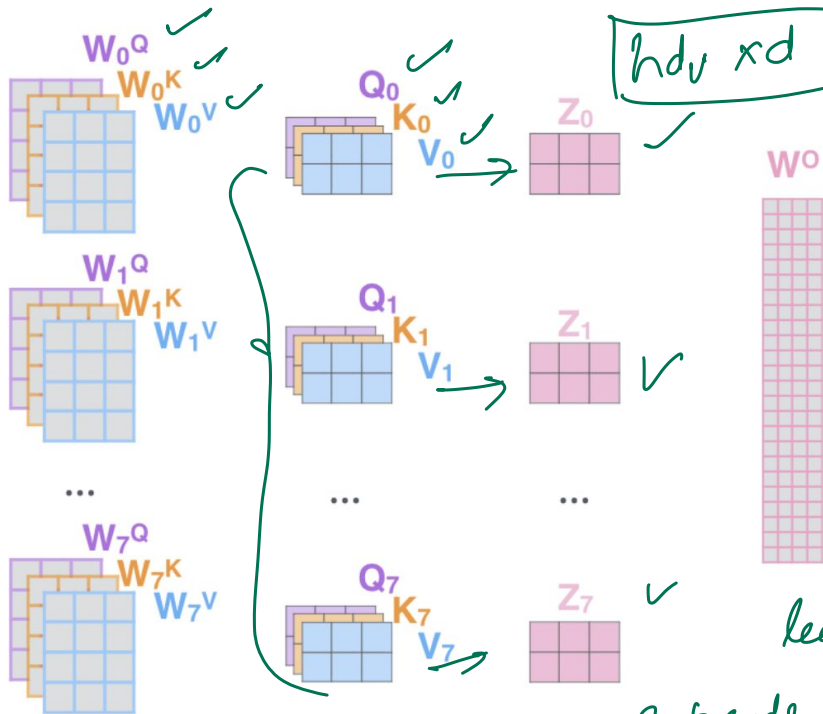
R

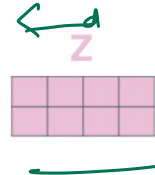3) Split into 8 heads. We multiply $X$ or $R$ with weight matrices

$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

...

$W_7^Q$
$W_7^K$
$W_7^V$

learnable

4) Calculate attention using the resulting $Q$/$K$/$V$ matrices

$Q_0$ $K_0$ $V_0$

$Q_1$ $K_1$ $V_1$

...

$Q_7$ $K_7$ $V_7$

8 heads

5) Concatenate the resulting $Z$ matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

$Z_0$

$Z_1$

...

$Z_7$

$W^O$

$Z$

learnable

$4d^2$

$hd_v \times d$

$\underline{1\ m}$ self-attn

$d = 512$

$h = 8$

$d_q = 64$

Figure: Jay Alammar

# Encoder Block



Figure: Jay Alammar

# Transformer with encoders and decoders
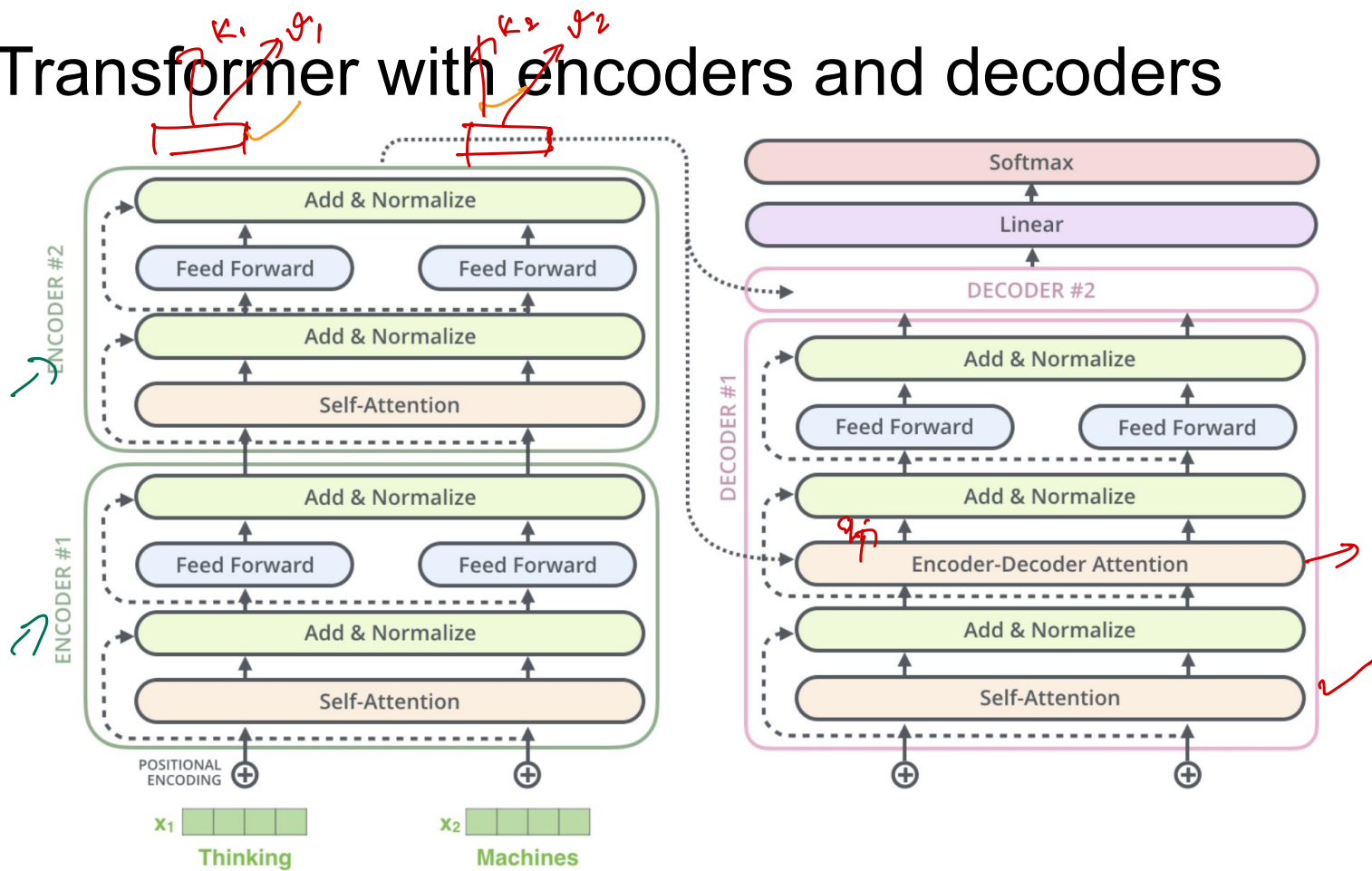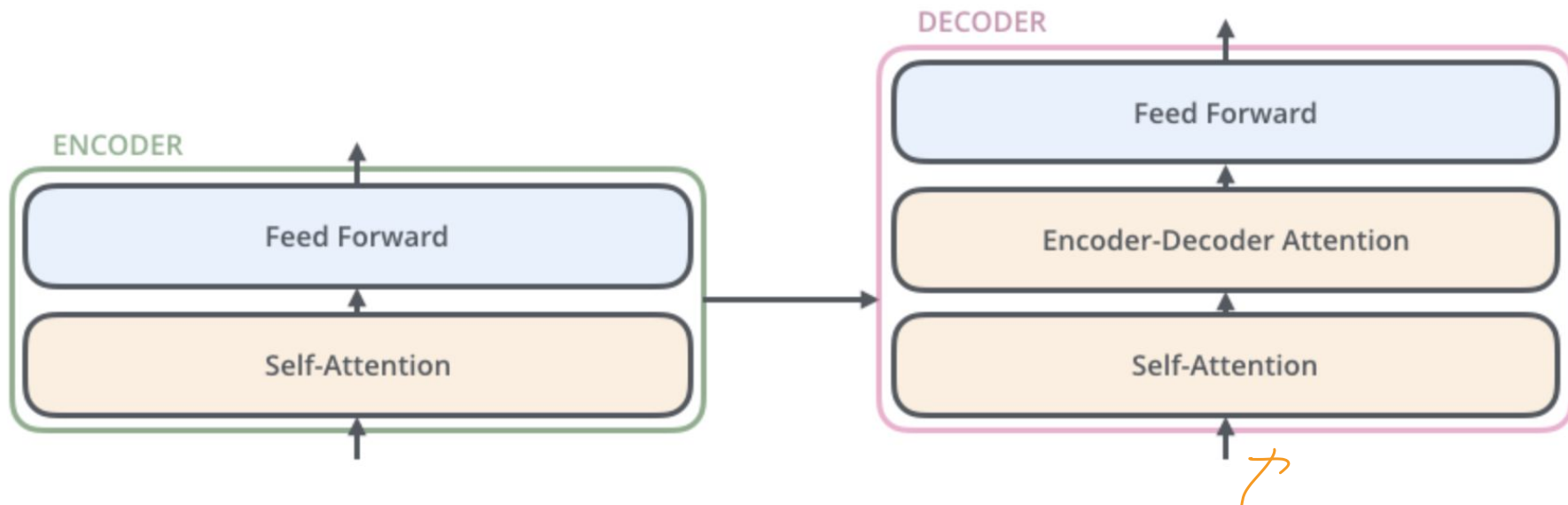


Figure: [Jay Alammar](https://jalammar.github.io/)

# How is the decoder different?



Figure: Jay Alammar

# Masked Self-attention

In the **decoder**, the self-attention layer is only allowed to attend to **earlier positions** in the output sequence. Otherwise, we'd be cheating!

This is done by **masking future positions (setting the dot product score for future positions to -inf)** before the Softmax step in the self-attention calculation.

$$e^{0} = 1 \qquad \boxed{e^{-\infty} = 0}$$

**+ M**

-inf

in general, upper triangle set to -inf

$$Z = \text{softmax}\left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$
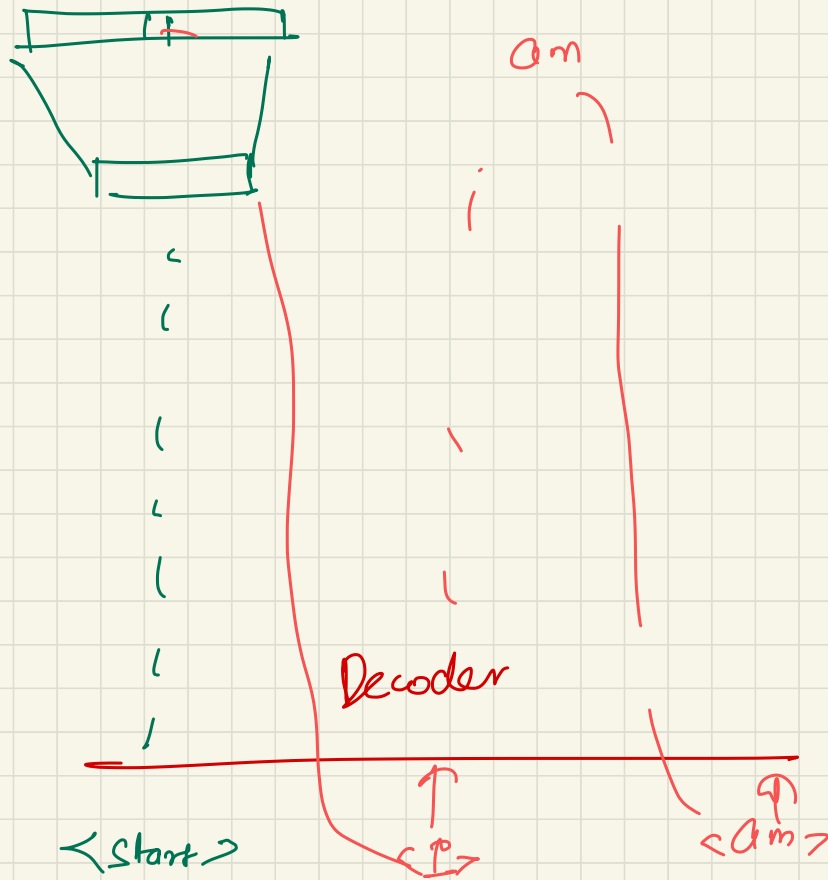
# Masked Self-attention for decoder
## (to avoid seeing the future tokens)



|  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | q1·k1 | −∞ | −∞ | −∞ | −∞ |
|  | q2·k1 | q2·k2 | −∞ | −∞ | −∞ |
| N | q3·k1 | q3·k2 | q3·k3 | −∞ | −∞ |
|  | q4·k1 | q4·k2 | q4·k3 | q4·k4 | −∞ |
|  | q5·k1 | q5·k2 | q5·k3 | q5·k4 | q5·k5 |

N

Masked

Self- attn
matrix

https://web.stanford.edu/~jurafsky/slp3/

am

No masking
at inference
time

Decoder

<start>

<p>

decode ✓

last
decoder

<Start> I am going to class

# Encoder-Decoder Attention



Decoding time step: 1 2 3 4 5 6        OUTPUT

$K_{encdec}$   $V_{encdec}$

Linear + Softmax

ENCODER

ENCODER

DECODER

DECODER

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant
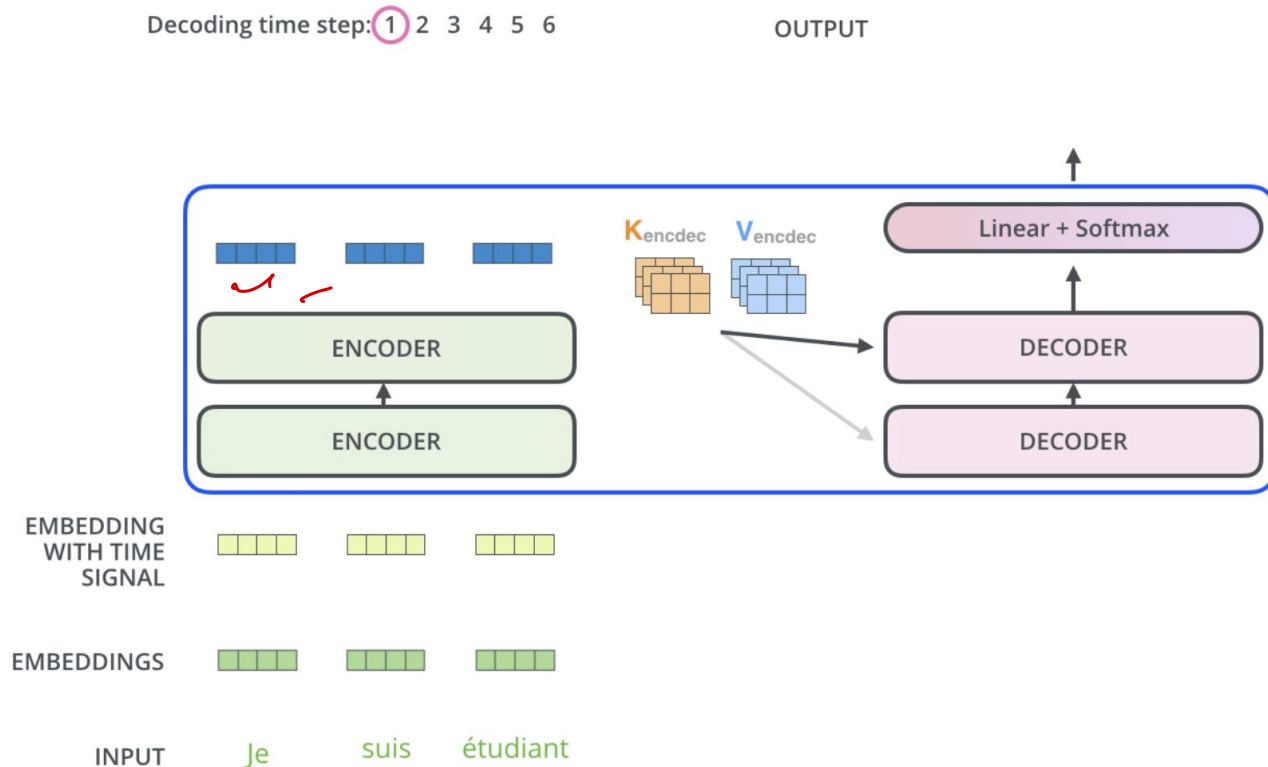
The output of the top encoder is transformed into a set of attention vectors K and V.

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:

Figure: Jay Alammar

# Decoding



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.

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

Figure: Jay Alammar

# Converting decoder stack output to words
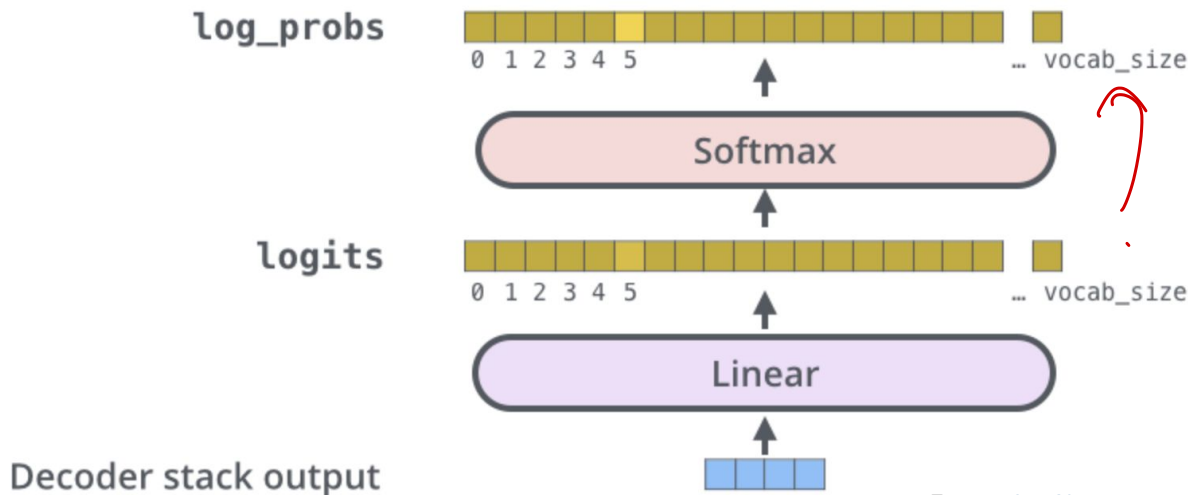
Which word in our vocabulary is associated with this index?

am

Get the index of the cell with the highest value (argmax)

5

That's the job of the final **Linear layer** which is followed by a **Softmax Layer**.

Encoders

Encoder-Decoder

Decoders?

log_probs
0 1 2 3 4 5 ... vocab_size

Softmax

logits
0 1 2 3 4 5 ... vocab_size

Linear

Decoder stack output

Figure: Jay Alammar

# Transformers as Language models (Decoder)

**Language Model Head**

takes $h^L_N$ and outputs a distribution over vocabulary V

| y1 | y2 | ... | y|V| | Word probabilities | 1 x |V| |

Softmax over vocabulary V

| u1 | u2 | ... | u|V| | Logits | 1 x |V| |

Unembedding layer = $E^T$

Unembedding layer   d x |V|

| $h^L_1$ | $h^L_2$ | | $h^L_N$ | 1 x d |

Layer L Transformer Block

...

| w1 | w2 | | $w_N$ |

- Masked self-attention
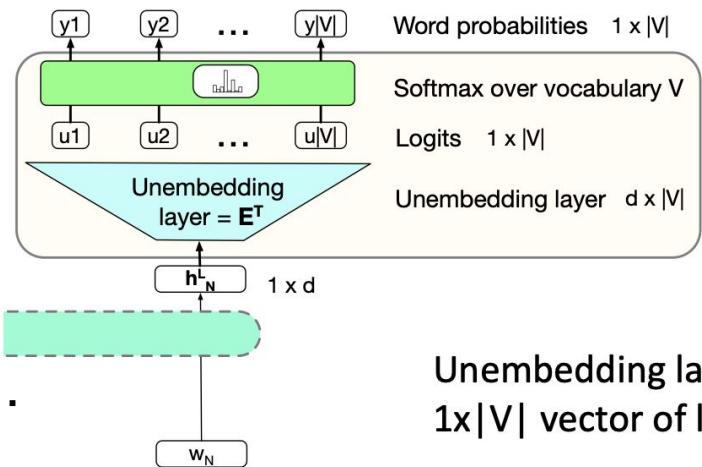- No encoder-decoder attention

$|V \times d|$

E — Embedding matrix

1-hot

(learnable with the model)

# Transformers as Language models

**Unembedding layer**: linear layer projects from $h^L_N$ (shape $[1 \times d]$) to logit vector



Why "unembedding"? **Tied** to $E^T$

**Weight tying**, we use the same weights for two different matrices

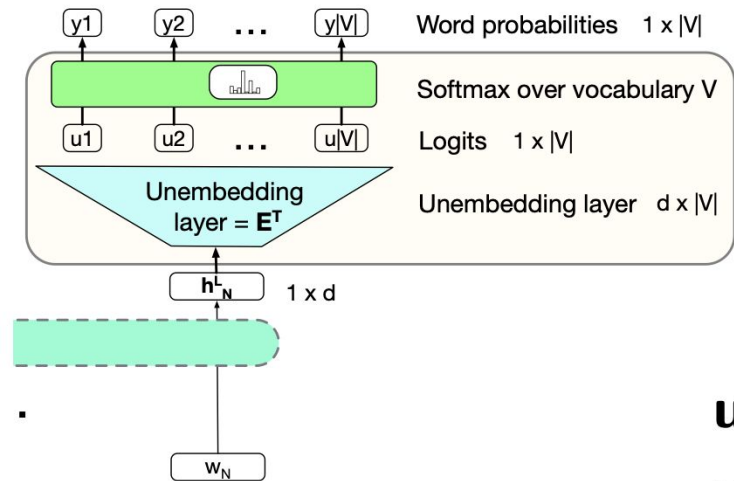Unembedding layer maps from an embedding to a 1x|V| vector of logits

https://web.stanford.edu/~jurafsky/slp3/

# Language modeling head

**Logits**, the score vector u

One score for each of the $|V|$ possible words in the vocabulary $V$. Shape $1 \times |V|$.

**Softmax** turns the logits into probabilities over vocabulary. Shape $1 \times |V|$.



| | | |
|---|---|---|
| y1 | y2 ... y|V| | Word probabilities 1 x |V| |

Softmax over vocabulary V

u1 u2 ... u|V|  Logits 1 x |V|

Unembedding layer = $\mathbf{E^T}$   Unembedding layer d x |V|

$\mathbf{h^L_N}$   1 x d

$w_N$

$$\mathbf{u} = \mathbf{h^L_N} \, \mathbf{E^T}$$

$$\mathbf{y} = \mathrm{softmax}(\mathbf{u})$$

https://web.stanford.edu/~jurafsky/slp3/

# Transformers LM

Token probabilities

$y1$  $y2$  $\cdots$  $y|V|$

$w_{i+1}$

Sample token to generate at position i+1

Language Modeling Head

softmax

logits    $u1$  $u2$  $\cdots$  $u|V|$

U

$h^L_i$

feedforward
layer norm
attention
layer norm

Layer L

$h^{L-1}_i = x^L_i$

$\cdots$  $h^2_i = x^3_i$

feedforward
layer norm
attention
layer norm

Layer 2

$h^1_i = x^2_i$

feedforward
layer norm
attention
layer norm

Layer 1

$x^1_i$

Input Encoding

$i$

E

pos$^n$ encoding (Encoder/decoder)

https://web.stanford.edu/~jurafsky/slp3/

# Number of parameters in encoder-decoder?

6  Layers — Encoders

$\qquad$ Decoders

$d = 512$

$r = 40K$

# heads = 8 (h)  $\qquad$ $d_K = d_q = 64$  $\qquad$ $FFN_{up} = 2048$

Embedds: $\underline{512 \times 40K}$  (= Unembeddy) $\approx$ $\boxed{20m}$ (weight tying)

Encoder: $\underline{Self-att^n} :\Rightarrow h^{\infty} \underbrace{\left( 3 \times d \times \dfrac{d}{h} \right)}_{W_q,\ W_K,\ W_v} + \underbrace{h \times \dfrac{d}{h} \times d}_{W_o} = \underline{4d^2}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad +$

$\qquad\qquad \underline{FFN} :\Rightarrow \qquad\qquad d \times 4d + 4d \times d \qquad = \underline{8d^2}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underline{12d^2}$

$\underline{\dfrac{12 \times 512 \times 512 \times 6}{3m \times 6 = 18m}}$  $\qquad$ $\boxed{12d^2 \times 2}) \rightarrow encoders$

# Number of parameters in encoder-decoder?

Decoder

self-attn: $4d^2$

feed-forward: $8d^2$

encoder-decoder: $4d^2$

$16d^2$

assumes a default conf

$h \times d_v = d$! — ①

Up-proj = $4d$ — ②

$16d^2 \times 6$

$4m \times 6 \quad = 24m$

$18m + 24m + (20m) emb \quad \approx 62m$

# Number of parameters in the decoder only Transformer?

Same as encoder

( 18m )

+

embeddg /unembeddg

GPT-3   ( 175 B )

$12 d^2 \times \ell$

$d = 12288$     $\ell = 96$

Try this one

# Decoding: Greedy and Beam Search are deterministic!

- Greedy decoding as well as Beam Search decoding will give a "deterministic" output
- Other common decoding algorithms involve "sampling", and bring in some degree of "randomness"

$x \sim p(x) \rightarrow$ choose $x$ by sampling from the distribution $p(x)$

## Random Sampling

$$i \leftarrow 1$$
$$w_i \sim p(w)$$
$$\textbf{while } w_i \mathrel{!}= \text{EOS}$$
$$\quad i \leftarrow i + 1$$
$$\quad w_i \sim p(w_i \mid w_{<i})$$

# Quality vs Diversity trade-off

Various sampling methods enable trading off two important factors in generation: *quality* and *diversity*.

> ### *Quality vs Diversity trade-off*
>
> - Methods that emphasize the most probable words tend to produce more coherent and accurate generations but also tend to be repetitive and boring
>
> - Methods that give bit more weight to the middle probability words tend to be more creative and diverse, but likely to be incoherent and less factual

# Random Sampling with Temperature

## Intuition from thermodynamics

A system at a *high temperature* is flexible and can explore various states, while a system at a *low temperature* is likely to explore a subset of lower energy (better) states

## How is this implemented

Divide the logits by a temperature parameter $\tau \in (0, 1]$ before passing it through softmax

Random sampling: $y = softmax(u)$

Random sampling with temperature: $y = softmax(u/\tau)$

$\tau = 0.2$

$u = [3, 1, -1]$

$e^3 \quad e \quad e^{-1}$

$u/\tau = [15, 5, -5]$

$e^{15} e^5 \quad e^{-5}$

$0.99$

# Why does this work?



*Without scaling, random sampling often generates incoherent gibberish*

$$p(y_i \mid x) = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

Higher T: softens probabilities.
Lower T: sharpens probabilities.

https://utah-cs6340-nlp.notion.site/Natural-Language-Processing-bd1a2ca290fc44f69556908ad8d25c70