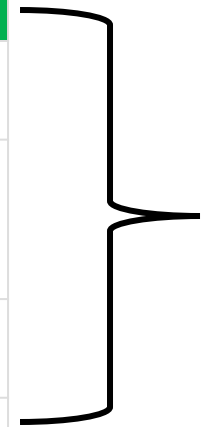




# Announcements

- **Announcements:** Recordings available on Canvas (under Kaltura tab)
- Class roadmap:

Tuesday Sept. 16	Architectures II: Subquadratic Architectures
Thursday Sept. 18	Language Models I
Tuesday Sept. 23	Language Models II
Thursday Sept. 25	Prompting
Tuesday Sept. 30	Specialization & Adaptation



Mostly Language Models

# Outline

- **Review Attention**

- Notions of attention, self-attention, basic attention layer, QKV setup and intuition, positional encodings

- **Transformers**

- Architecture, encoder and decoder setups

- **Subquadratic Models**

- Basic ideas. Examples: S4, Mamba.

# Self-Attention: Retrieval Intuition

- How should we design the interactions?

- Analogy: **search**

“Which restaurants near me are open at 9:00 pm?”

Query

Key

Value

Objects:

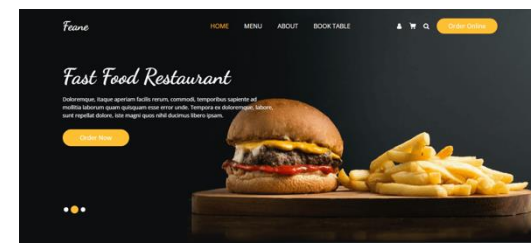
Query

Key

Value

Score 0.3

Score 0.7



# Self-Attention: Query, Key, Value Vectors

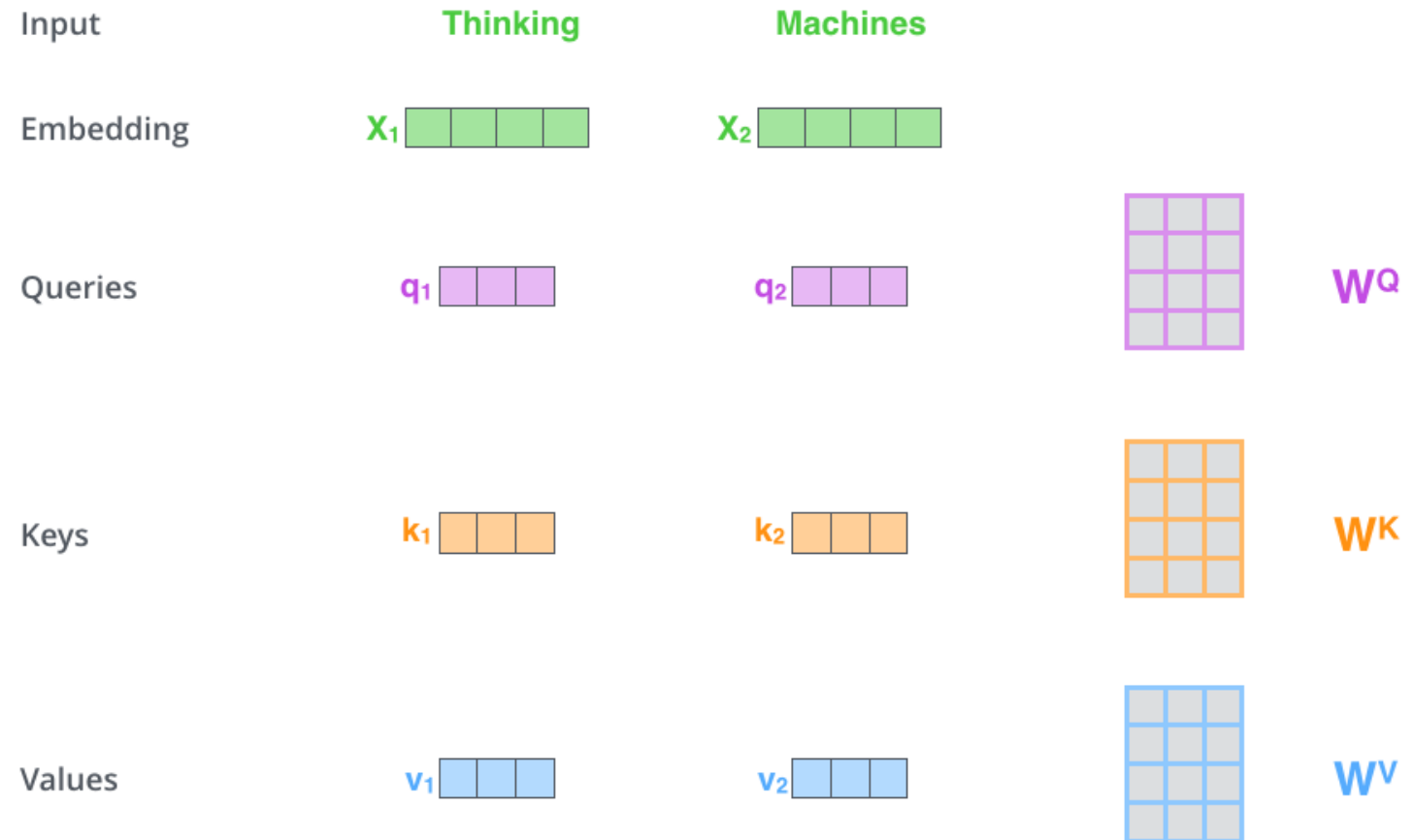
- *Transform* incoming word vectors,
- Enable *interactions* between words
- Get our **query**, **key**, **value** vectors via weight matrices: linear transformations!

**Objects:**

**Query**

**Key**

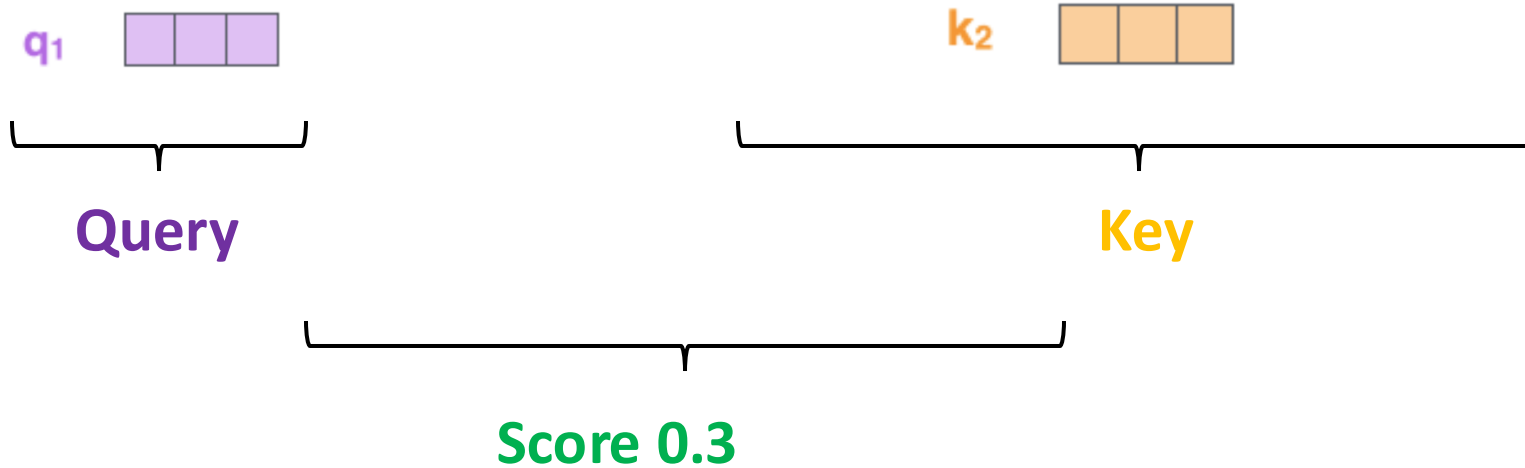
**Value**



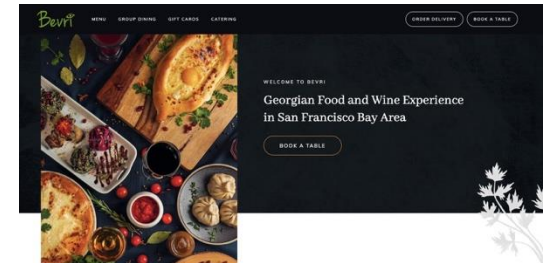
# Self-Attention: Score Functions

Have **query**, **key**, **value** vectors

- Next, get our **score**



- Lots of things we could do --- **simpler** is usually better!
- Dot product  $q_1 \cdot k_2 = 96$
- Then we'll do **softmax**



# Self-Attention: Scoring and Scaling

- *Transform* incoming word vectors,
- Enable *interactions* between words
- Get our **query**, **key**, **value** vectors via weight matrices: linear transformations!
- Compute scores

**Objects:**

**Query**

**Key**

**Value**

Input

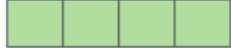
Embedding

Queries

Keys

Values

Thinking

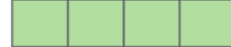
$x_1$  

$q_1$  

$k_1$  

$v_1$  

Machines

$x_2$  

$q_2$  

$k_2$  

$v_2$  

# Self-Attention: Putting it Together

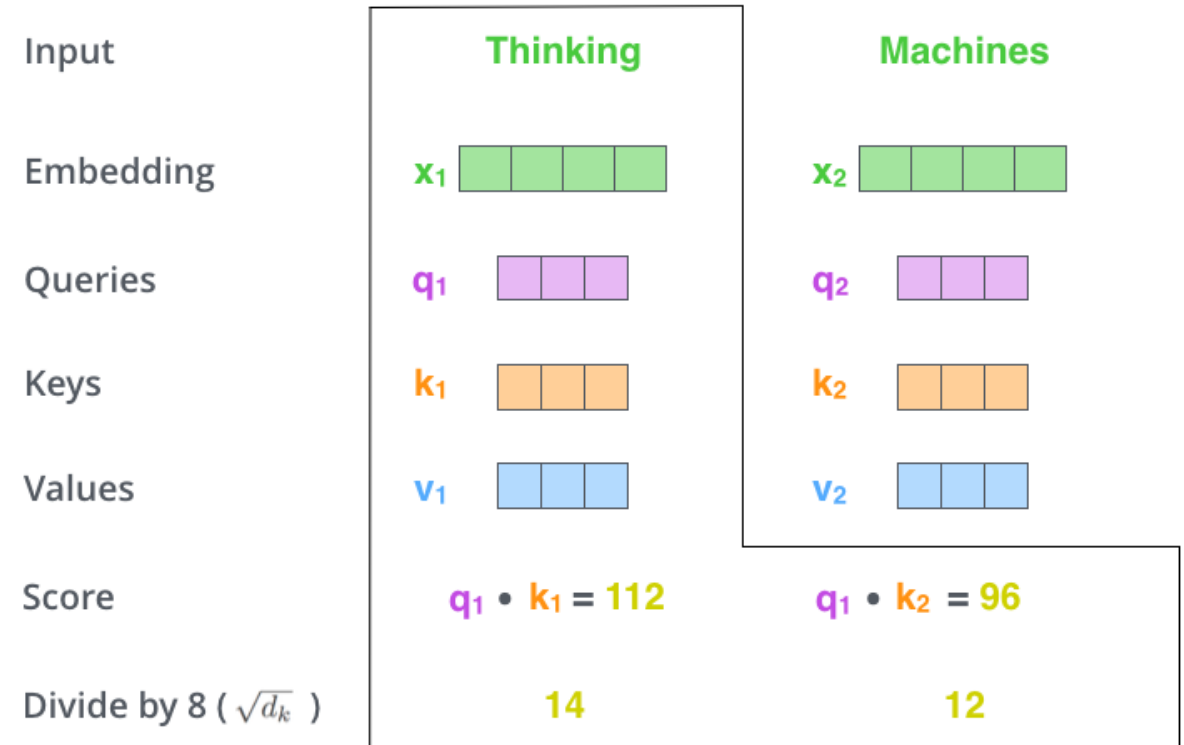
- Have **query**, **key**, **value** vectors via weight matrices: linear transformations!
- Have softmax score outputs (**focus**)
- Add up the values!

**Objects:**

**Query**

**Key**

**Value**





# Self-Attention: Matrix Formulas

- Have **query**, **key**, **value** vectors via weight matrices: linear transformations!
- Have softmax score outputs (**focus**)
- Add up the values!

Objects:

Query

Key

Value

$$Q = XW_Q, K = XW_K, V = XW_V$$

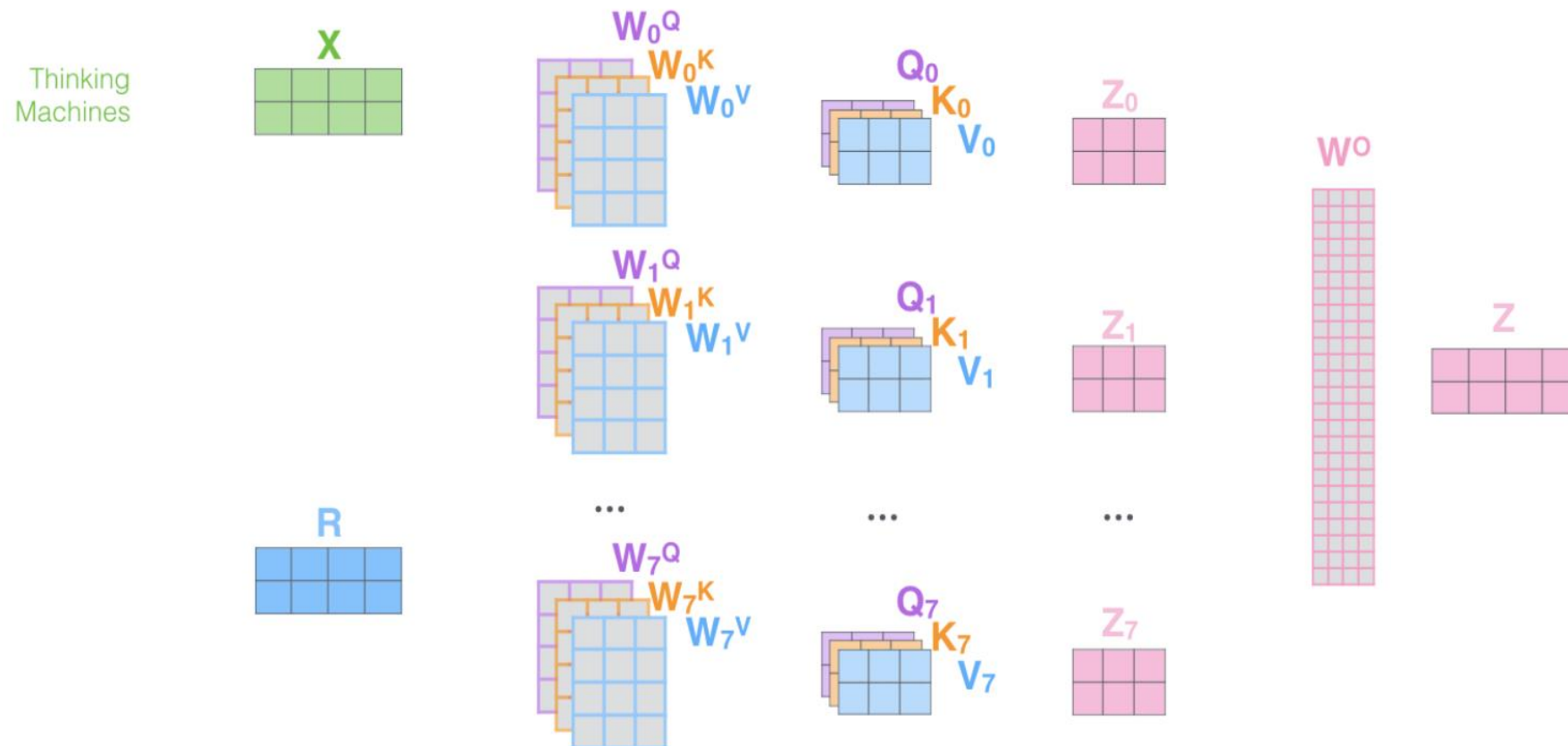
$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

$$\text{Attention}(Q, K, V) = \text{softmax} \left( X \frac{W_Q W_K^T}{\sqrt{d_k}} X^T \right) V$$

# Self-Attention: Multi-head

This is great but will we capture everything in one?

- Do we use just 1 kernel in CNNs? **No!**
- Do it many times in parallel: **multi-headed attention**. Concatenate outputs



# Self-Attention: Positional Encodings

Almost have a full layer designed.

- One annoying issue: so far, order of words (**position**) **doesn't matter!**
- Solution: add positional encodings

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



Location index



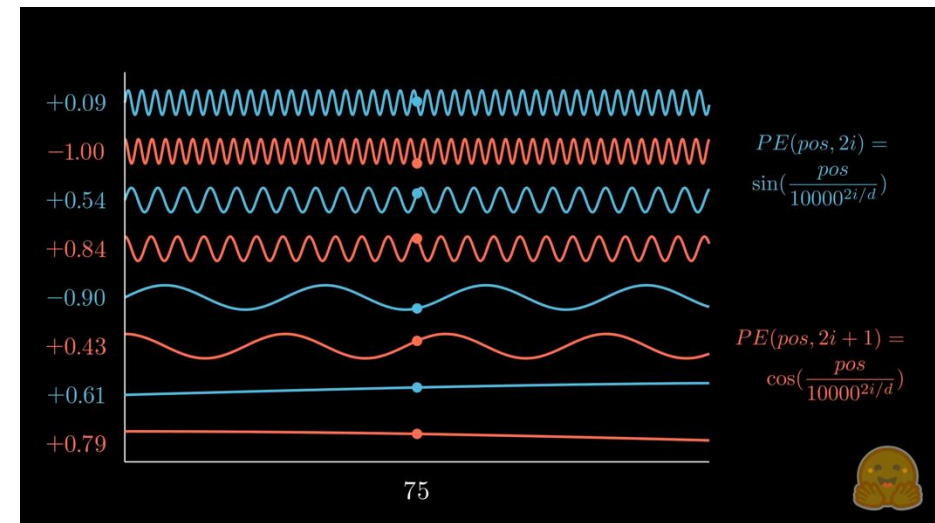
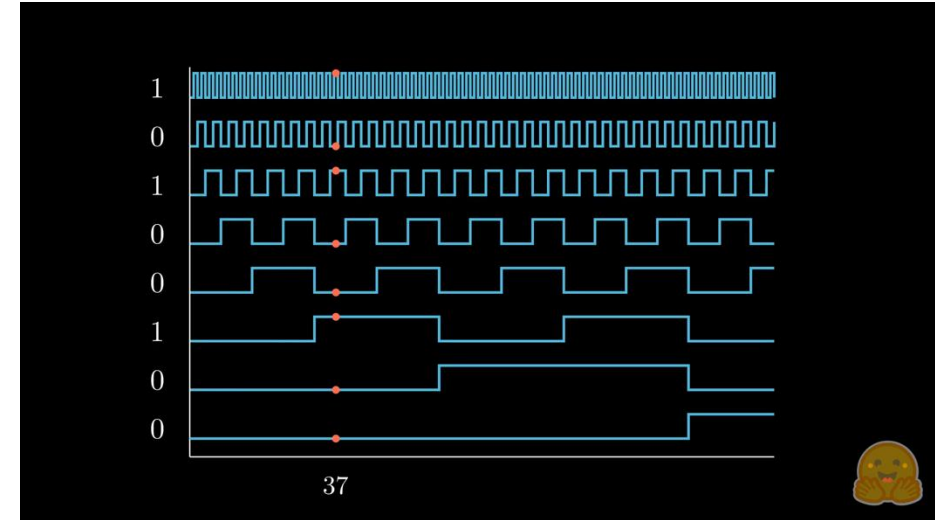
# Self-Attention: Positional Encodings

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Why these **mysterious formulas**? Want properties:

- Consistent encoding
- Smooth
- Linearity across positions
  - Alternating sin and cos: can multiply by rotation matrix to obtain shifts



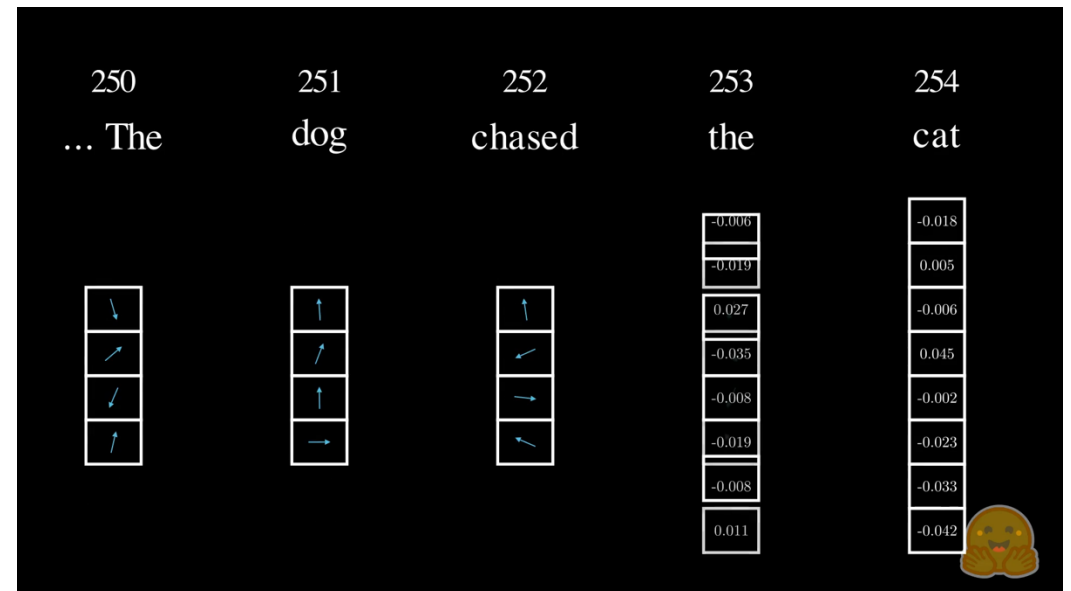
# Self-Attention: Modern Positional Encodings

These *sinusoidal* embeddings were defined in the original Transformers paper,

- Added once (as we saw) prior to the first layer

Many new variants of positional encodings that behave slightly differently

- Example: *multiplicative* instead of *additive*
- Popular: **Rotary Positional Encoding (RoPE)**
- Note: perform in every attention layer



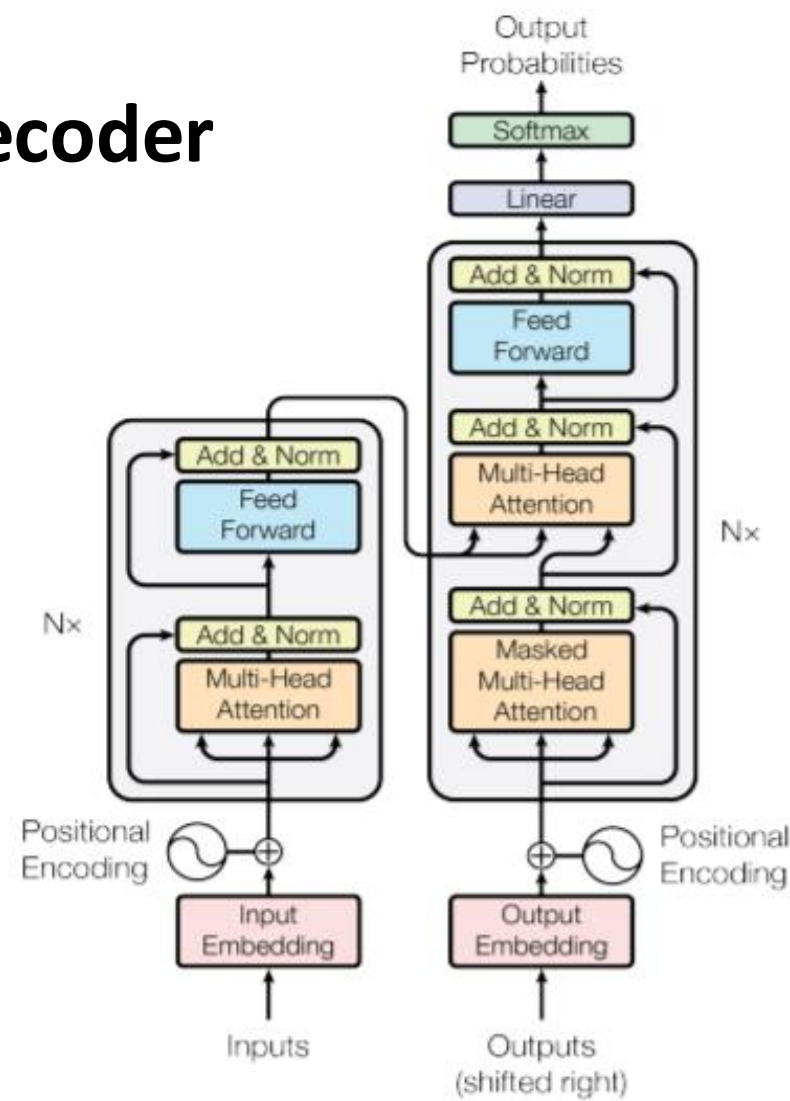




# Break & Questions

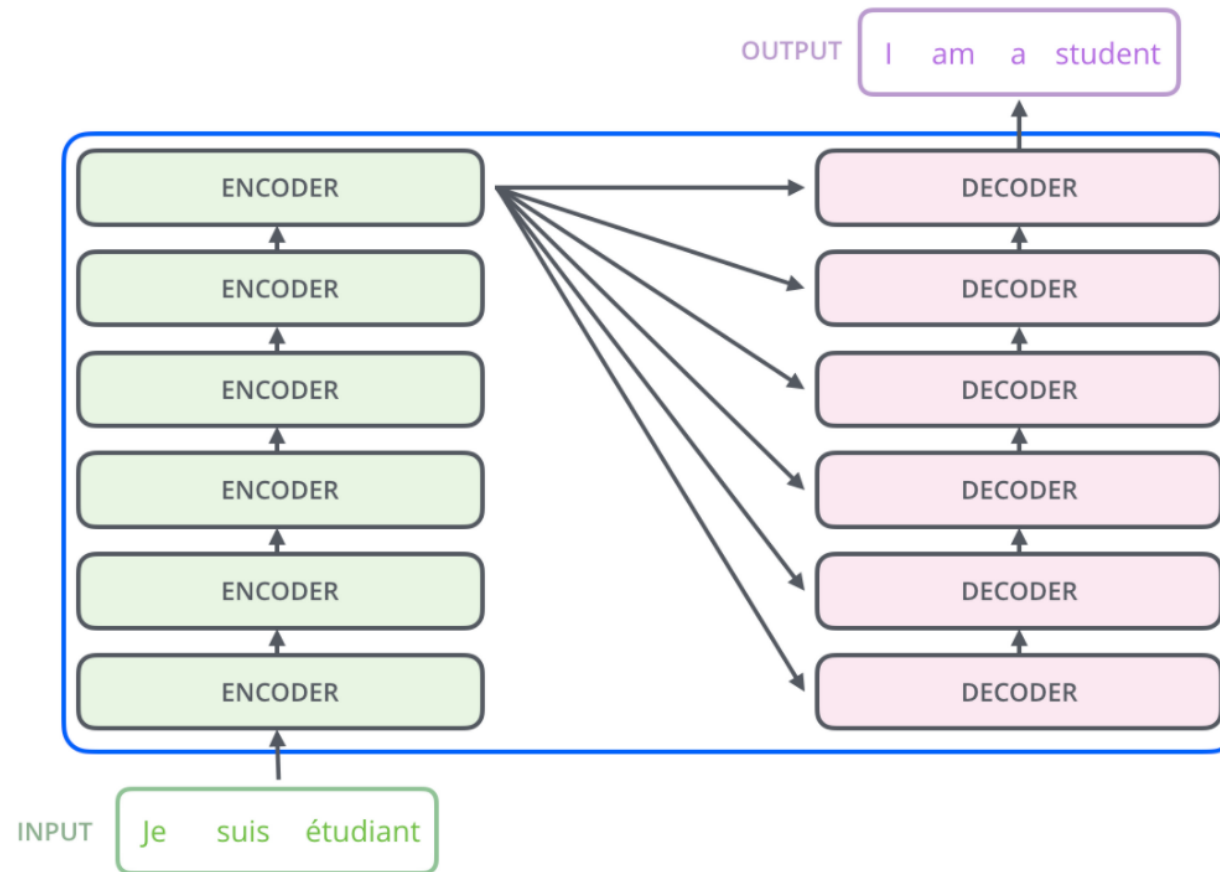
# Transformers: Model Architecture

- Initial goal for an architecture: **encoder-decoder**
  - Get **rid of recurrence**
  - Replace with **self-attention**
- Architecture
  - The famous picture you've seen
  - Centered on self-attention blocks



# Transformers: Architecture

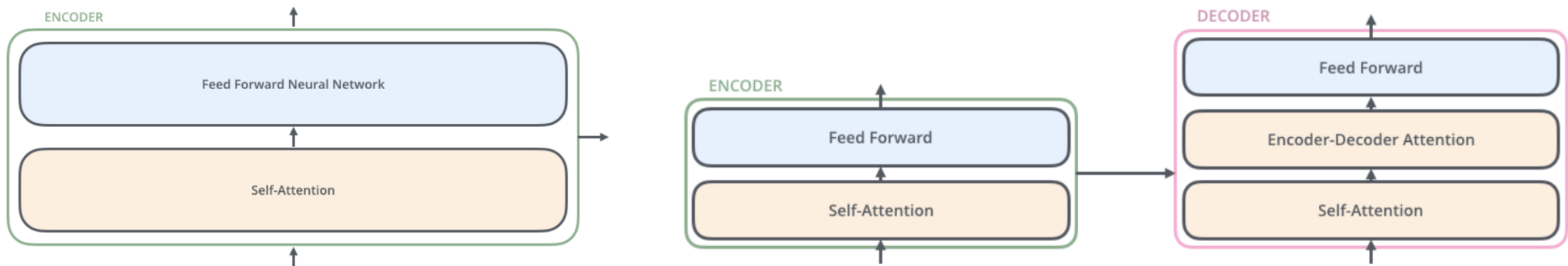
- **Sequence-sequence** model with **stacked** encoders/decoders:
  - For example, for French-English translation:





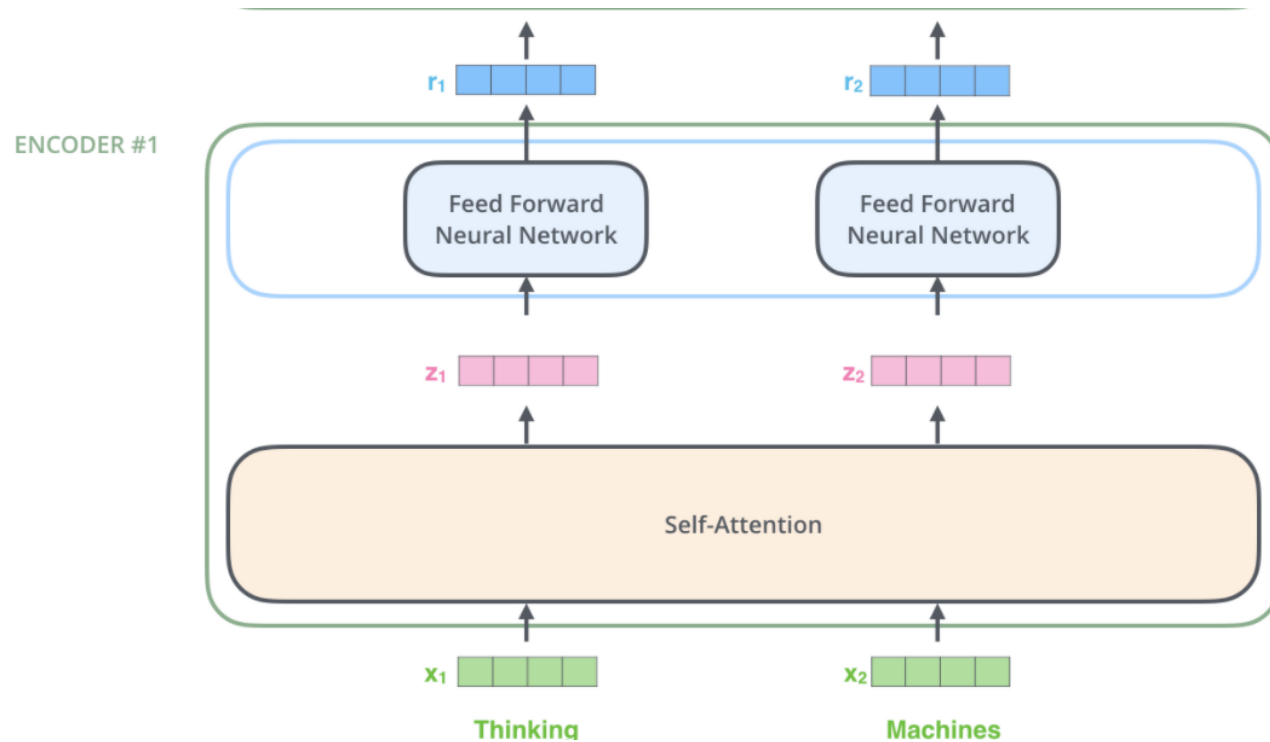
# Transformers: Architecture

- Sequence-sequence model with **stacked** encoders/decoders:
  - What's inside each encoder/decoder unit?
- Focus encoder first: **pretty simple!** 2 components:
  - Self-attention block
  - Fully-connected layers (i.e., an MLP)



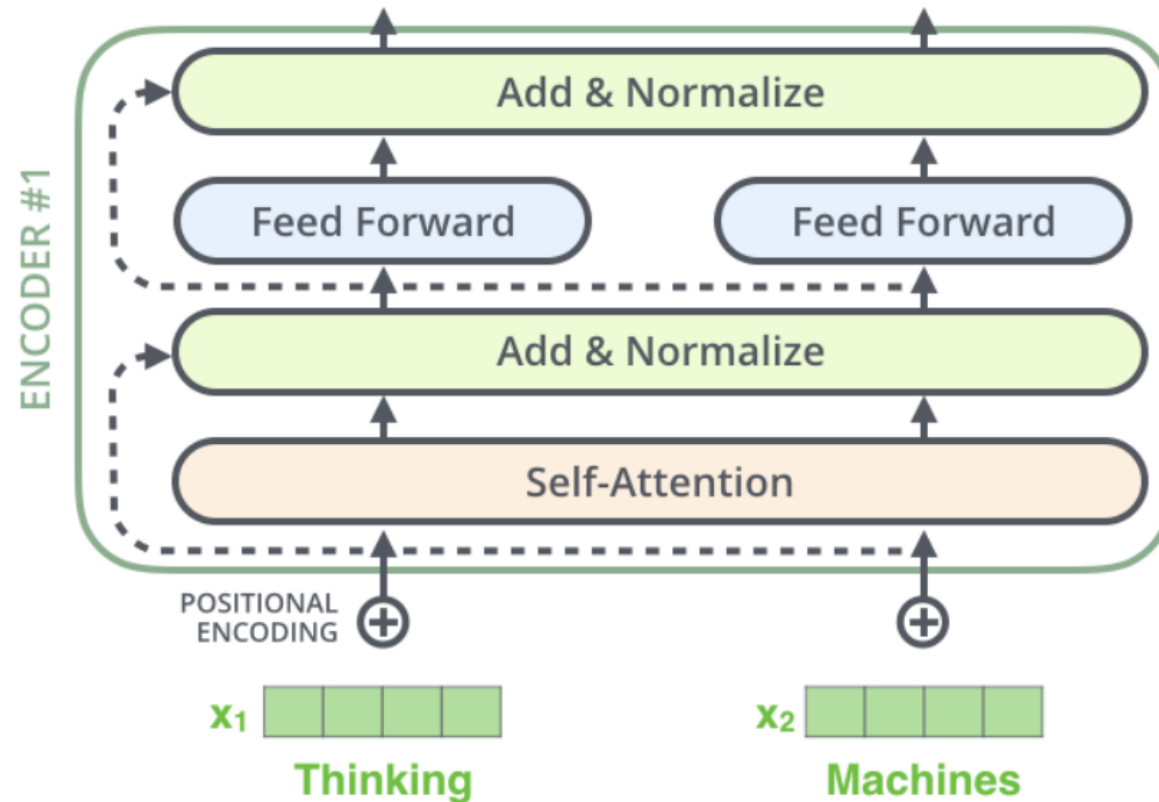
# Transformers: Inside an Encoder

- Let's take a look at the encoder. Two components:
  - 1. **Self-attention** layer (covered this)
  - 2. “Independent” **feedforward nets** for each head



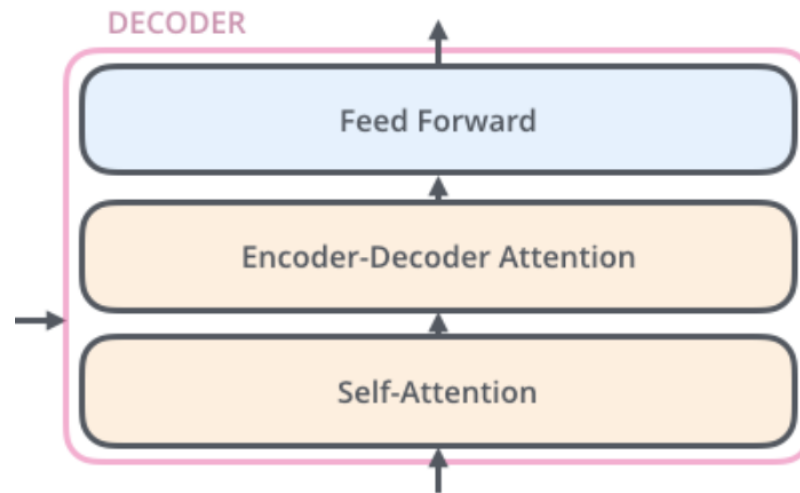
# Transformers: More Tricks

- Recall a big innovation for ResNets: residual connections
  - And also layer normalizations
  - Apply to our encoder layers



# Transformers: Inside a Decoder

- Let's take a look at the decoder. Three components:
  - 1. **Self-attention** layer (covered this)
  - 2. Encoder-decoder attention (same, but K, V come from encoder)
  - 3. “Independent” **feedforward nets** for each head



# Transformers: Last Layers

- Next let's look at the end. Similar to a CNN,

- 1. Linear layer
- 2. Softmax

Get probabilities of words

Which word in our vocabulary is associated with this index?

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

am

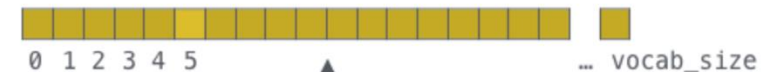
5

log\_probs



Softmax

logits



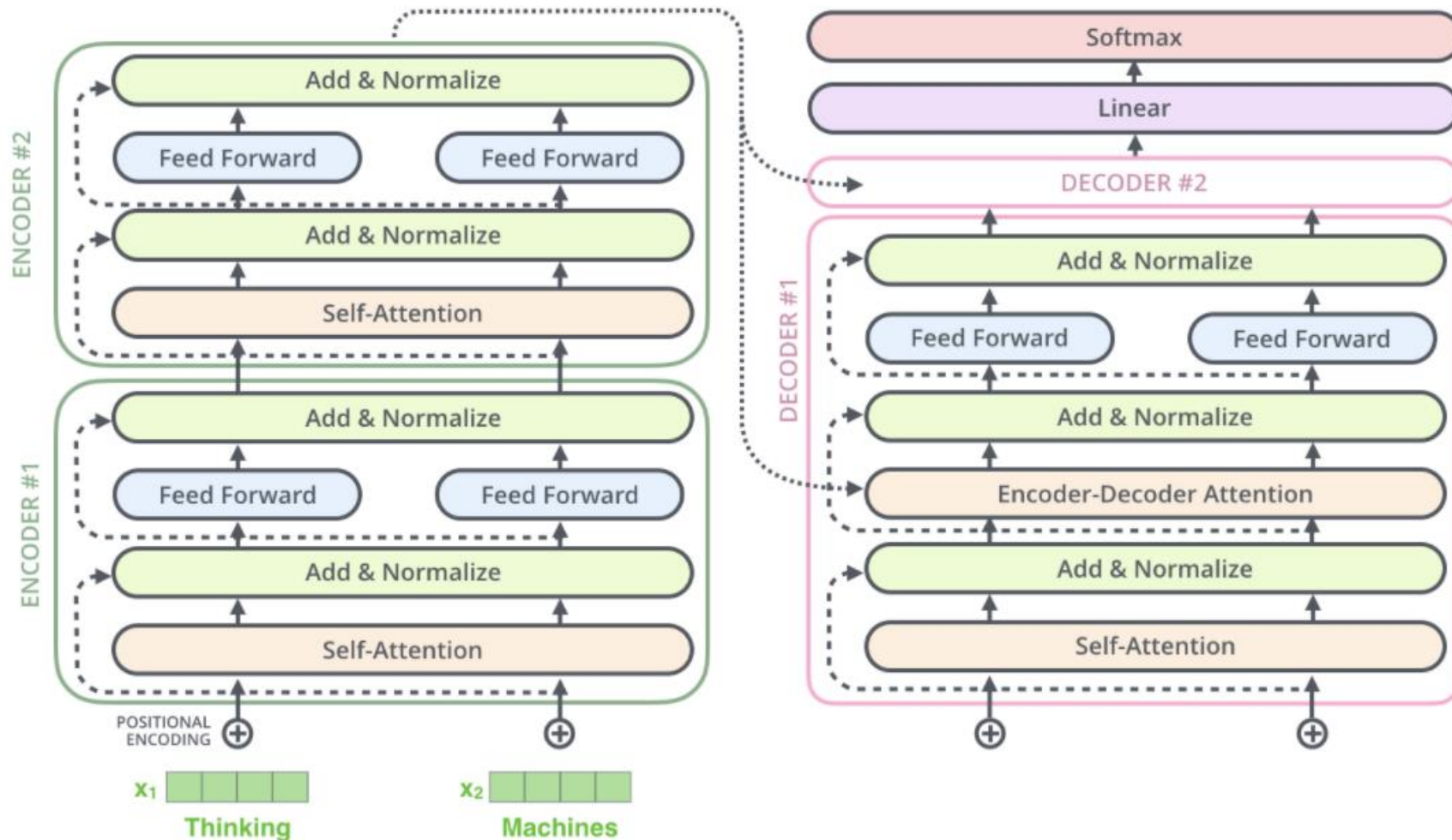
Linear

Decoder stack output



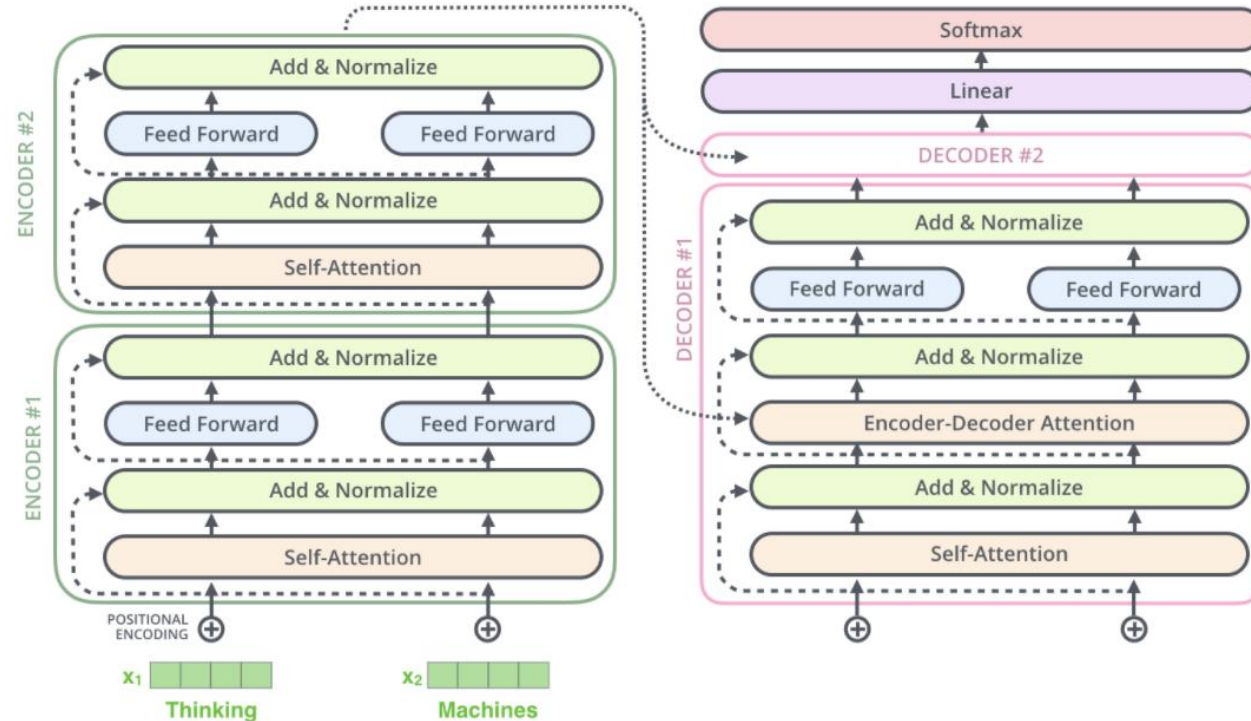
# Transformers: Putting it All Together

- What does the full architecture look like?



# Transformers: The Rest

- Next time: we'll talk about
  - How to **use** it (i.e., outputs)
  - How to **train** it
  - How to **rip** it apart and build other models with it.



# Transformers: The Rest

- Next time: we'll talk about
  - How to use it (i.e., outputs)
  - **How to train it**
  - How to rip it apart and build other models with it.

## 5.1 Training Data and Batching

We trained on the standard WMT 2014 English-German dataset consisting of about 4.5 million sentence pairs. Sentences were encoded using byte-pair encoding [3], which has a shared source-target vocabulary of about 37000 tokens. For English-French, we used the significantly larger WMT 2014 English-French dataset consisting of 36M sentences and split tokens into a 32000 word-piece vocabulary [38]. Sentence pairs were batched together by approximate sequence length. Each training batch contained a set of sentence pairs containing approximately 25000 source tokens and 25000 target tokens.



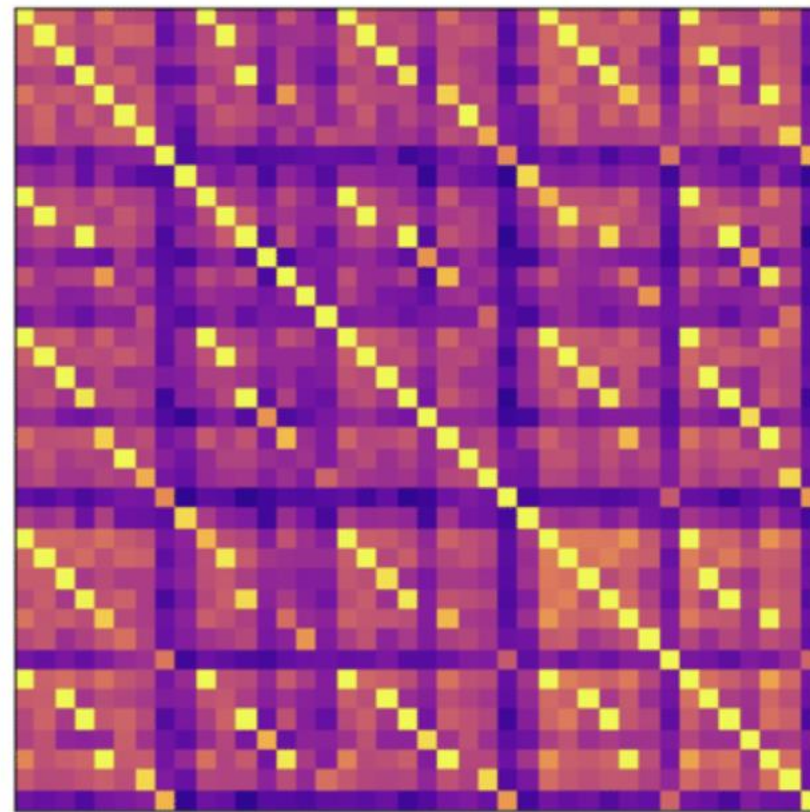


**Break & Questions**

# Attention Alternatives?

- One annoying thing: if the sequence length is  $L$ , we're doing a  $O(L^2)$  operation.
- This can be quite limiting for long sequences...

I.e., 4000 tokens is fine, but  $10^6$  tokens is not.



# Attention Alternatives?

Recently, lots of different approaches that attempt to get rid of this quadratic dependency

- Sometimes called **sub-quadratic** models.
- We'll briefly study a few.
- Step 1: let's get inspired by something RNN-like (well, fully linear for now). Borrow from continuous models:

$$x'(t) = \mathbf{A}x(t) + \mathbf{B}u(t)$$

$$y(t) = \mathbf{C}x(t) + \mathbf{D}u(t)$$

# State-Space Model

Step 1: let's get inspired by something RNN-like (well, fully linear for now). Borrow from continuous models:

$$\begin{array}{c} \text{State} \quad \quad \quad \text{Input} \\ \downarrow \quad \quad \quad \downarrow \\ x'(t) = \mathbf{A}x(t) + \mathbf{B}u(t) \\ \text{Output} \rightarrow y(t) = \mathbf{C}x(t) + \mathbf{D}u(t) \end{array}$$

- Can ignore the “D” (think of this as a skip connection).
- Inputs, outputs are 1-D, state is higher dimensional.

# State-Space Model: Discrete Form

Step 2: let's make this a discrete function

$$\begin{array}{ccc} & \text{State} & \text{Input} \\ & \downarrow & \downarrow \\ x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k \\ \text{Output} \rightarrow & y_k = \overline{\mathbf{C}}x_k & \end{array}$$

- Ignored D
- Can create approximations of A,B,C through discretizing.
- Looks a lot like an RNN! (or, a linear version of one)

# State-Space Model: Convolutional Form

Step 3: let's unroll the recursion

$$\begin{array}{lll} x_0 = \overline{B}u_0 & x_1 = \overline{A}\overline{B}u_0 + \overline{B}u_1 & x_2 = \overline{A}^2\overline{B}u_0 + \overline{A}\overline{B}u_1 + \overline{B}u_2 \\ y_0 = \overline{C}\overline{B}u_0 & y_1 = \overline{C}\overline{A}\overline{B}u_0 + \overline{C}\overline{B}u_1 & y_2 = \overline{C}\overline{A}^2\overline{B}u_0 + \overline{C}\overline{A}\overline{B}u_1 + \overline{C}\overline{B}u_2 \end{array}$$

$$y_k = \overline{C}\overline{A}^k\overline{B}u_0 + \overline{C}\overline{A}^{k-1}\overline{B}u_1 + \cdots + \overline{C}\overline{A}\overline{B}u_{k-1} + \overline{C}\overline{B}u_k$$

- In general,  $y = \overline{K} * u.$

- This is a **convolution**!

# State-Space Model: Convolutional Form

Step 3: let's unroll the recursion

- Convolution

$$y_k = \overline{CA}^k \overline{B} u_0 + \overline{CA}^{k-1} \overline{B} u_1 + \cdots + \overline{CA} \overline{B} u_{k-1} + \overline{CB} u_k$$

$$y = \overline{K} * u.$$

- But a weird one. It's a very **long** convolution.
  - Kernel as long as the input sequence (say, L).
  - Naively, is this better than attention?
  - Let's do **something else** instead.

# Interlude: Time & Frequency Domains

Back to Signals and Systems class,

- Convolution in the time-domain is element-wise multiplication in the frequency domain
- So low-complexity.
- But, need to convert to frequency domain
- Solution: **FFT**.  $O(L \log L)$  (and also for iFFT, to invert back).
- So, can compute fast and use during training!

$$y_k = \overline{CA}^k \overline{B} u_0 + \overline{CA}^{k-1} \overline{B} u_1 + \cdots + \overline{CAB} u_{k-1} + \overline{CB} u_k$$
$$y = \overline{K} * u.$$



# Back to SSM Picture

Back to the formula

$$x_k = \overline{A}x_{k-1} + \overline{B}u_k$$

$$y_k = \overline{C}x_k$$

- Just directly making all of these trainable parameters doesn't work so well.
  - Similar issues as in RNNs: stuff blowing up
  - Instead, various models propose approaches

S4 (Structured State Space Models) Gu et al' 22

- Build A with a special fixed transition matrix that is good at memorization
- Couple with a particular parametrization to get the discretization.

# Using SSMs as Layers

Back to the formula

$$x_k = \overline{\mathbf{A}}x_{k-1} + \overline{\mathbf{B}}u_k$$
$$y_k = \overline{\mathbf{C}}x_k$$

S4 (Structured State Space Models) Gu et al' 22

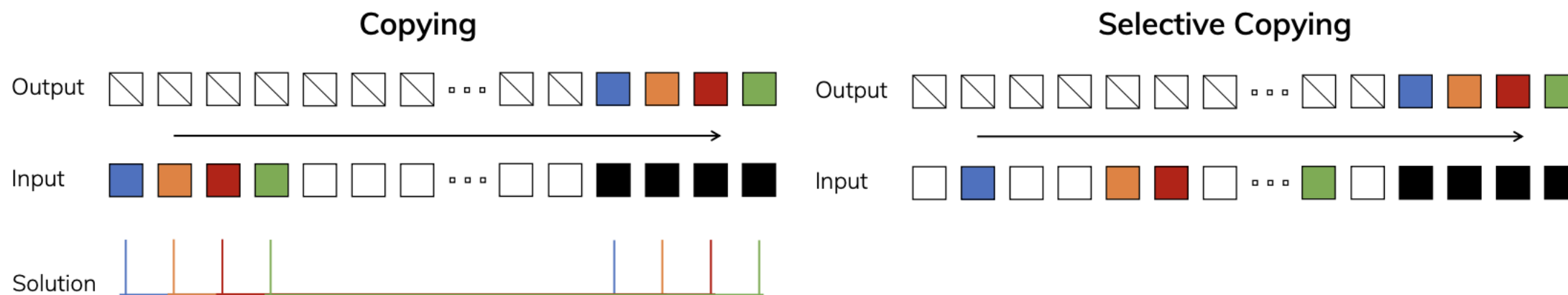
- Special A state transition matrix
- Special parametrization/choice of trainable parameters
- How to actually use these? Need to define a layer,
  - Stack H of them together (similar to conv layers, multihead attn)
  - Mix with linear layer, place activation function at the end

# S4 Results: The Good and the Bad

Models like S4 can address **very long sequences**

- “S4 solves the **Path-X task**, an extremely challenging task that involves reasoning about LRDs over sequences of length ... 16384. All previous models have failed...”

- But, can struggle with “selective” tasks.



# S4 Results: The Good and the Bad

Solution: need some type of context-aware approach

## • Mamba Model

- Gu and Dao '23, "Mamba: Linear-Time Sequence Modeling with Selective State Spaces"

**Algorithm 1** SSM (S4)

**Input:**  $x : (B, L, D)$

**Output:**  $y : (B, L, D)$

1:  $\mathbf{A} : (\mathbf{D}, \mathbf{N}) \leftarrow \text{Parameter}$ 

- Represents structured  $N \times N$  matrix

2:  **$B$**  :  $(D, N) \leftarrow \text{Parameter}$ 3:  $C : (D, N) \leftarrow \text{Parameter}$ 4:  $\Delta : (D) \leftarrow \tau_{\Delta}(\text{Parameter})$ 5:  $\overline{A}, \overline{B} : (\mathcal{D}, \mathcal{N}) \leftarrow \text{discretize}(\Delta, A, B)$ 6:  $y \leftarrow \text{SSM}(\overline{A}, \overline{B}, C)(x)$ 

- ▷ Time-invariant: recurrence or convolution

7: **return**  $y$ 

### Algorithm 2 SSM + Selection (S6)

**Input:**  $x : (B, L, D)$

**Output:**  $y : (B, L, D)$

1:  $\mathbf{A} : (\mathbf{D}, \mathbf{N}) \leftarrow \text{Parameter}$ 

- Represents structured  $N \times N$  matrix

2: **B** :  $(B, L, N) \leftarrow s_B(x)$ 3:  $\mathbf{C} : (\mathbf{B}, \mathbf{L}, \mathbf{N}) \leftarrow s_C(x)$ 4:  $\Delta : (B, L, D) \leftarrow \tau_{\Delta}(\text{Parameter} + s_{\Delta}(x))$ 5:  $\bar{A}, \bar{B} : (\mathbf{B}, \mathbf{L}, \mathbf{D}, \mathbf{N}) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$ 6:  $y \leftarrow \text{SSM}(\overline{A}, \overline{B}, C)(x)$ 

- **Time-varying:** recurrence (*scan*) only

7: **return**  $y$

# Lots of Related Approaches & Variations

- **Linear attention.** “Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention”. Katharopoulos et al, ‘20
- **RWKV.** “RWKV: Reinventing RNNs for the Transformer Era”, Peng et al ‘23

We’ll see more as we go!