# CS 839: Foundation Models
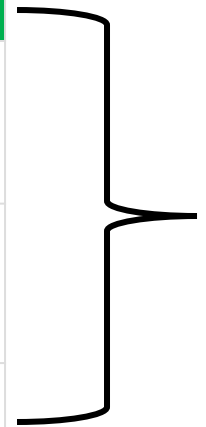## **Transformers and Attention**

Fred Sala

University of Wisconsin-Madison

**Sept. 11, 2025**

# Announcements

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

| Thursday Sept. 11 | **Architectures I: Transformers & Attention** |
|---|---|
| Tuesday Sept. 16 | Architectures II: Subquadratic Architectures |
| Thursday Sept. 18 | Language Models I |
| Tuesday Sept. 23 | Language Models II |
| Thursday Sept. 25 | Prompting |

Mostly Language Models

# Outline

- **Basic Attention**
  - Notions of attention, self-attention, basic attention layer, QKV setup and intuition
- **Additional Elements**
  - Multi-head attention, positional encodings
- **Transformers**
  - Architecture, encoder and decoder setups

# Outline

- **Basic Attention**
  - Notions of attention, self-attention, basic attention layer, QKV setup and intuition
- **Additional Elements**
  - Multi-head attention, positional encodings
- **Transformers**
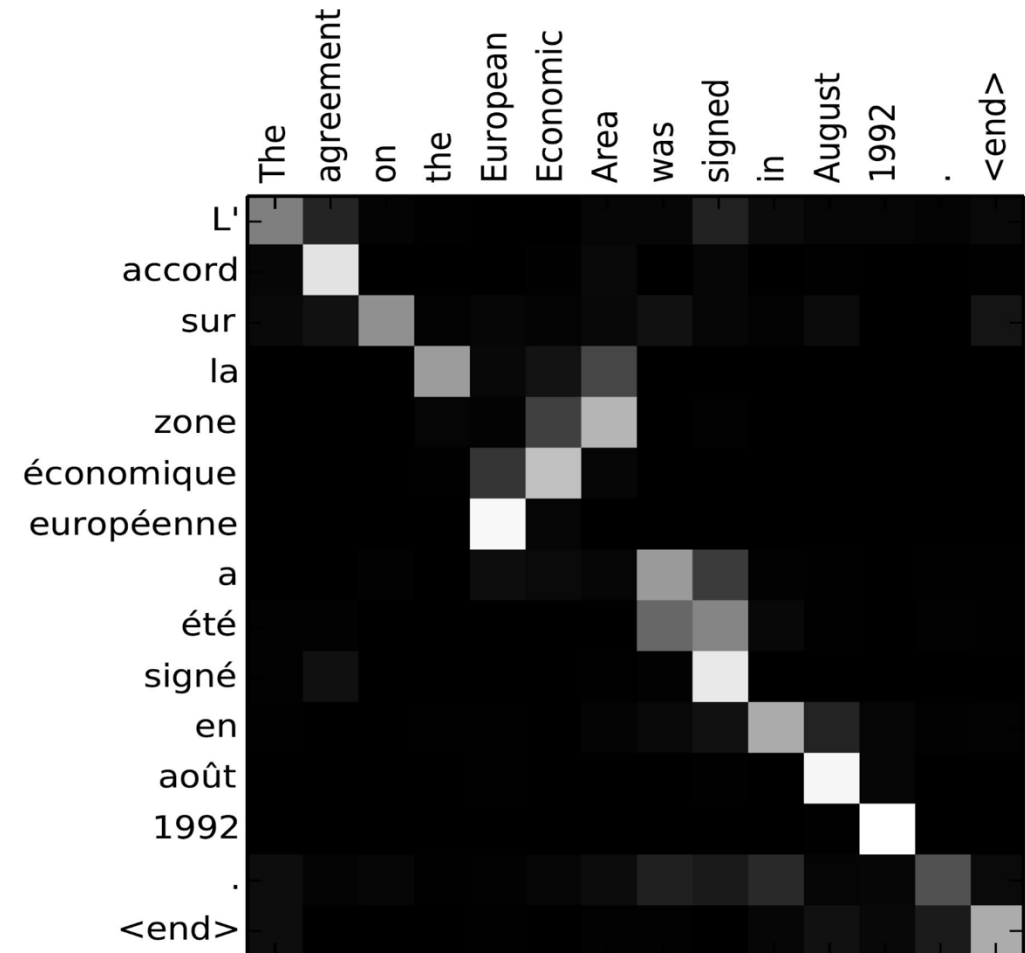  - Architecture, encoder and decoder setups

# History of Attention

Basic motivation: in NLP *fixed* context vector **not** enough

- Why?
  - Words depend on each other
  - Dependencies are complex

- Need: mechanism to help model **focus** on the right "part"

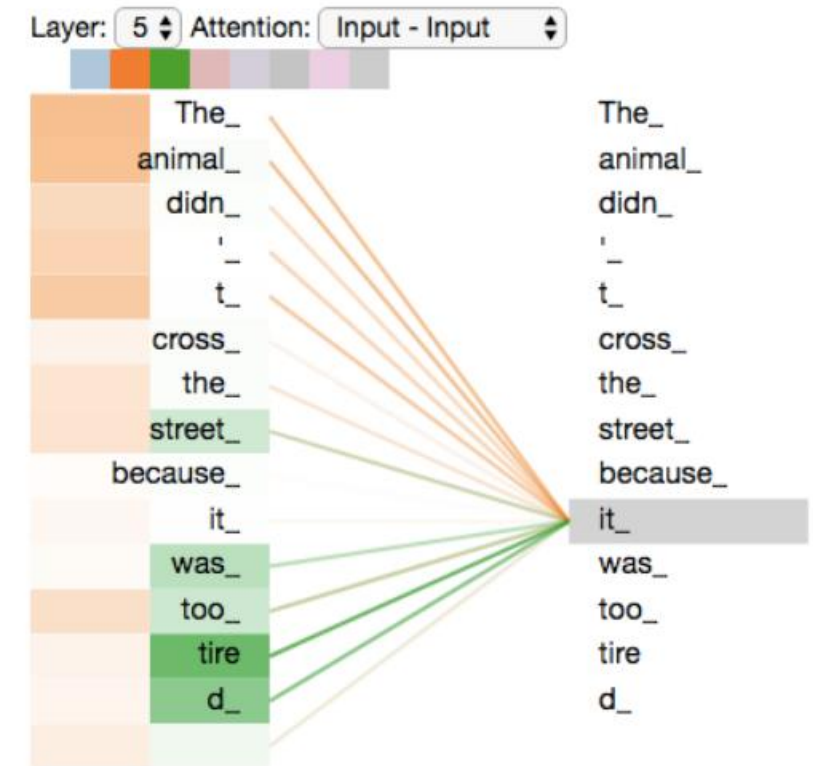Lots of approaches from 2014 on
  - Bahdanau et al, 2014

Bahdanau et al, 2014

# Self-Attention: Motivation

Popularized from 2017 on…

- From bottom-up. Let's design a basic layer.
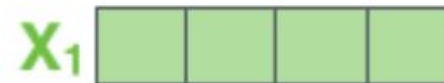  - Intuition: dependencies **within** same sentence



Cheng et al, 2016

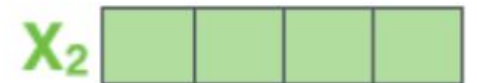Jay Alammar

# **Self-Attention:** Goals and Inputs

From bottom-up. Let's design a basic layer.

- Two criteria
  - *Transform* incoming word vectors,
  - Enable *interactions* between words

- Input: vectors for words

**Thinking**

$X_1$

**Machines**

$X_2$

**Note:** All visualizations are due to Jay Alammar

Excellent resource: https://jalammar.github.io/illustrated-transformer/

# **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**
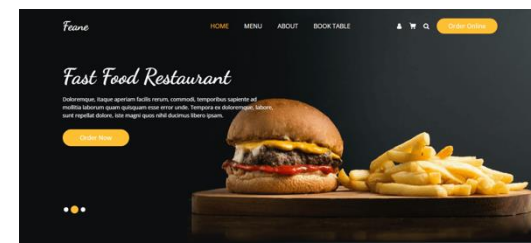
Objects:

Query
Key
Value

Score 0.3

Score 0.7

# Self-Attention: Query, Key, Value Vectors
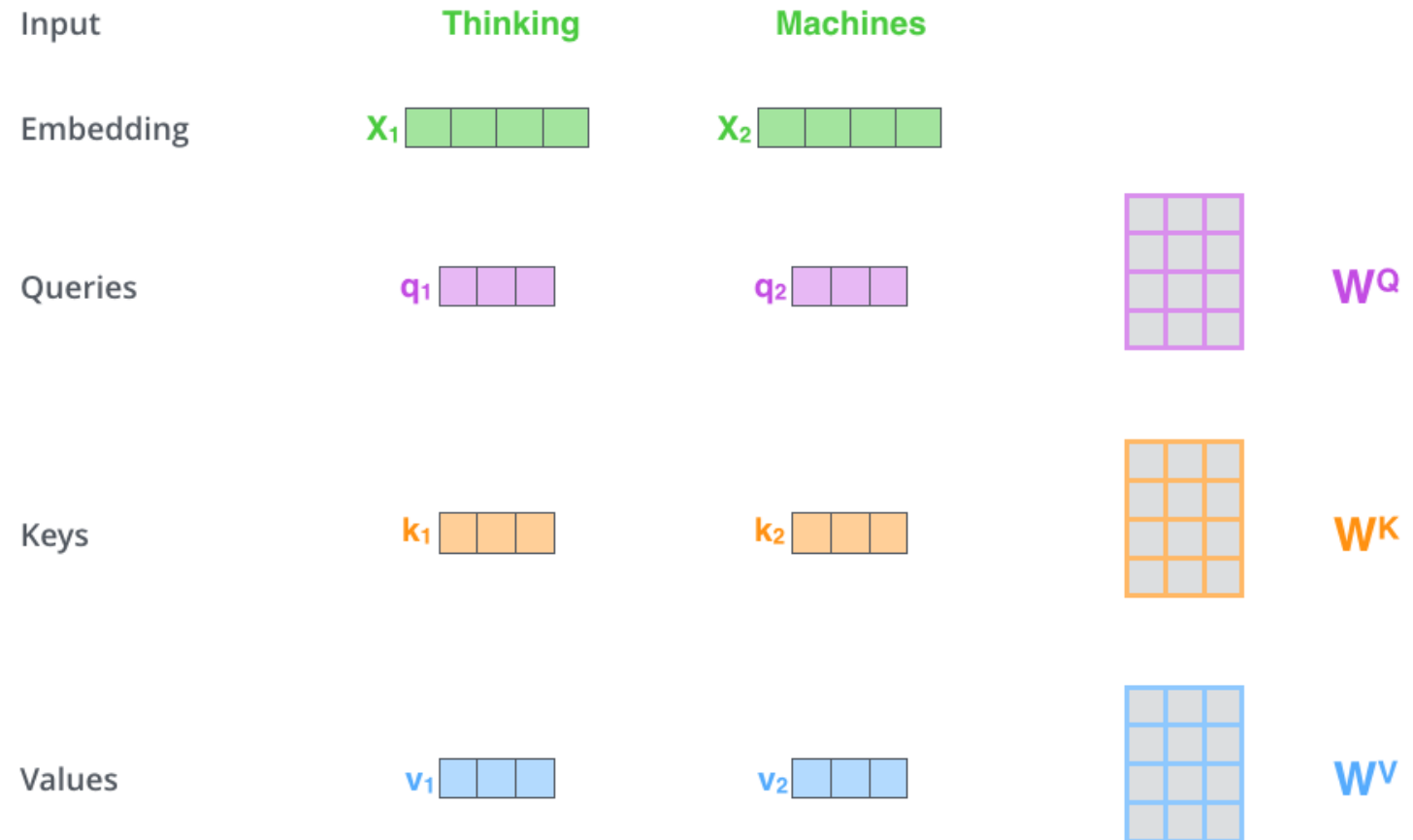
- *Transform* incoming word vectors,
- Enable *interactions* between words

- Get our **query, key, value** vectors via weight matrices: linear transformations!
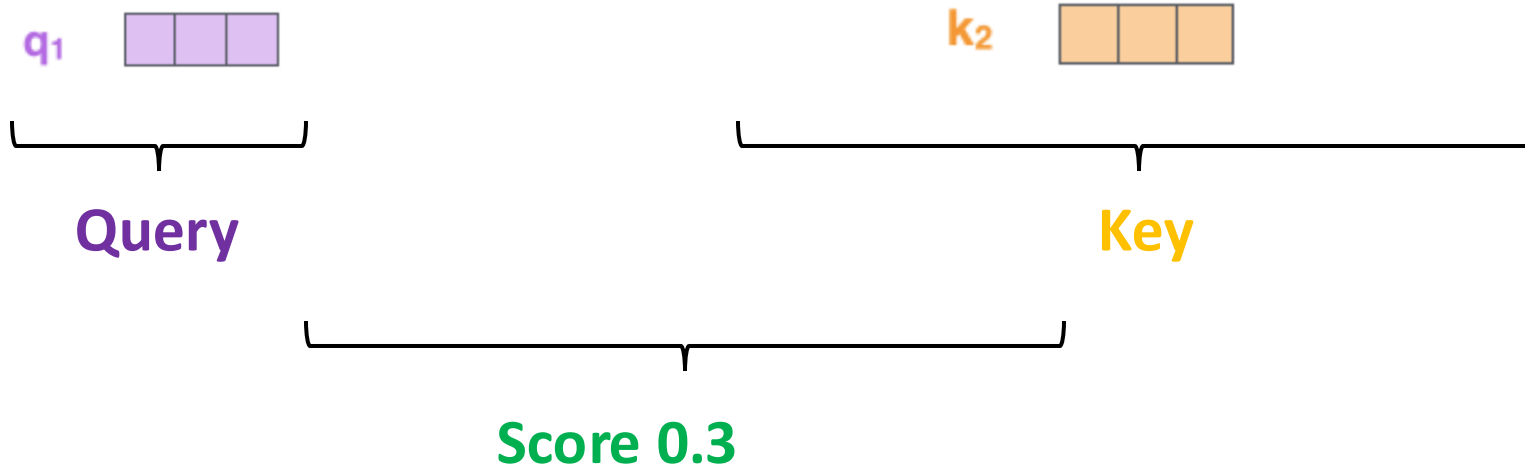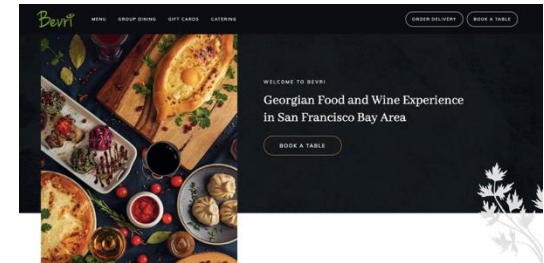
**Objects:**

**Query**
**Key**
**Value**

| Input | Thinking | Machines |
|---|---|---|
| Embedding | $X_1$ | $X_2$ |
| Queries | $q_1$ | $q_2$ | $W^Q$ |
| Keys | $k_1$ | $k_2$ | $W^K$ |
| Values | $v_1$ | $v_2$ | $W^V$ |

# **Self-Attention:** Score Functions

Have **query, key, value** vectors

- Next, get our **score**

$q_1$ [ ][ ][ ]    $k_2$ [ ][ ][ ]

**Query**    **Key**

**Score 0.3**

- 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!
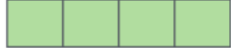- Compute scores

**Objects:**

**Query**
**Key**
**Value**

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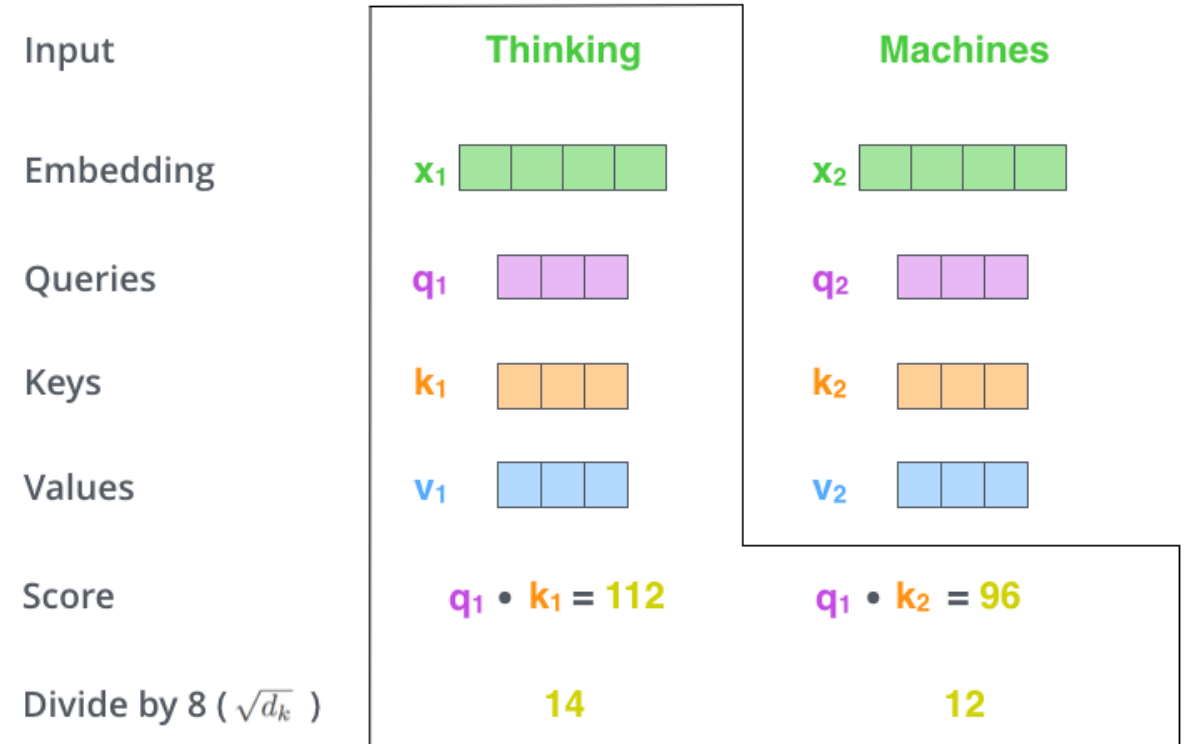
- Have **query, key, value** vectors via weight matrices: linear transformations!
- Have softmax score outputs (**focus**)
- Add up the values!

**Objects:**

**Query**
**Key**
**Value**

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |

# 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**

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

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

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

# Break & Questions

# Outline

- **Basic Attention**
  - Notions of attention, self-attention, basic attention layer, QKV setup and intuition
- **Additional Elements**
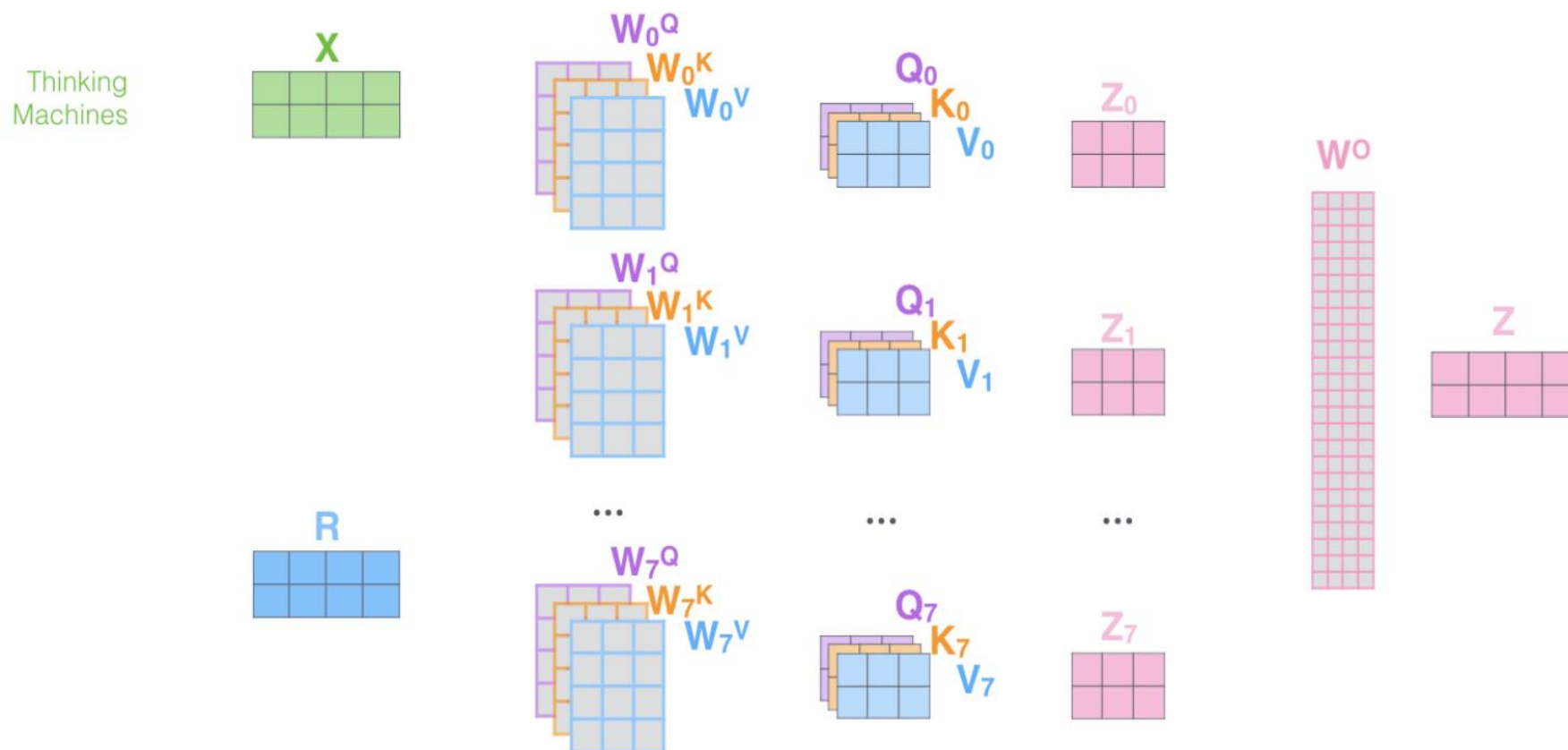  - Multi-head attention, positional encodings
- **Transformers**
  - Architecture, encoder and decoder setups

# 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

| POSITIONAL ENCODING | 0 | 0 | 1 | 1 | | 0.84 | 0.0001 | 0.54 | 1 | | 0.91 | 0.0002 | -0.42 | 1 |

+ + +

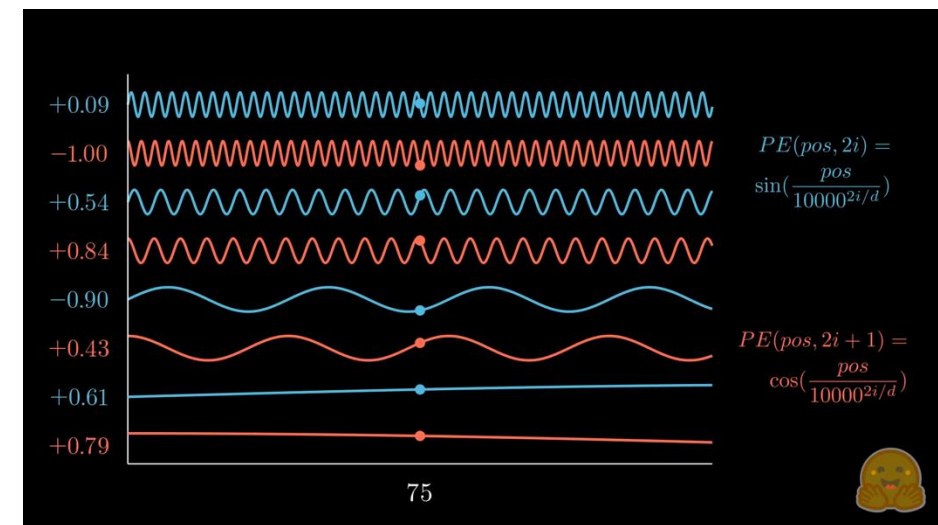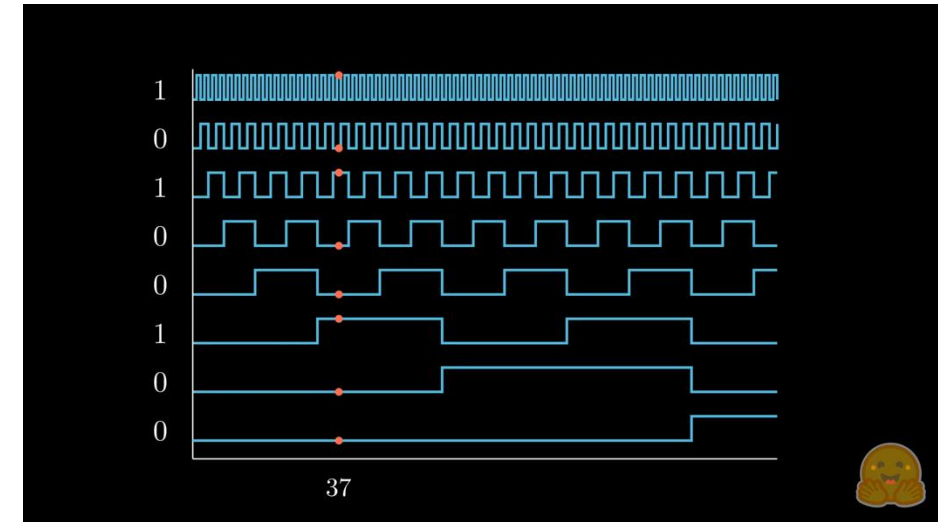| EMBEDDINGS | $x_1$ | | $x_2$ | | $x_3$ |

| INPUT | Je | suis | étudiant |

# 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





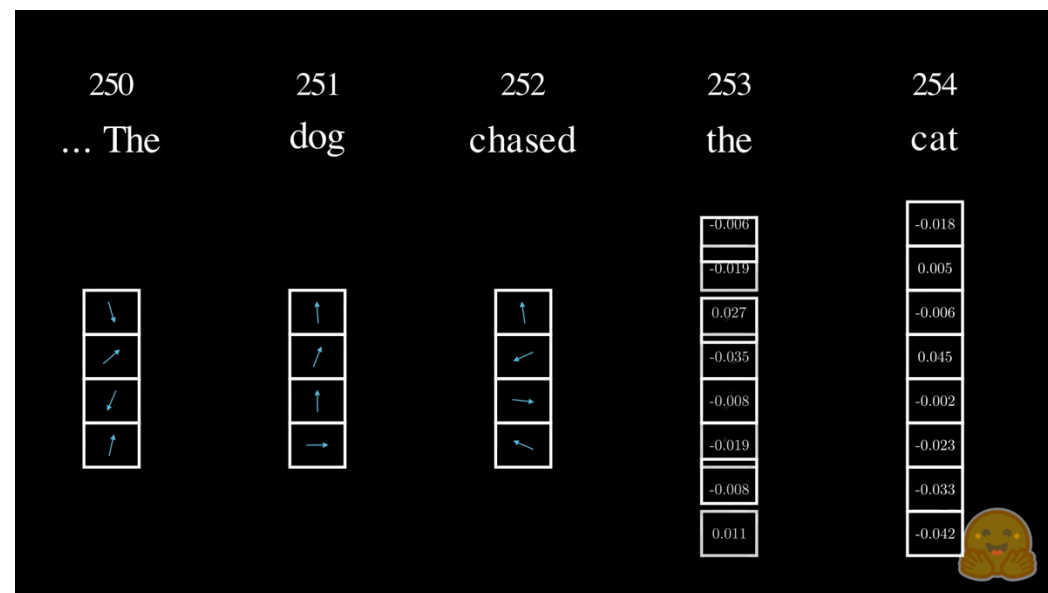https://huggingface.co/blog/designing-positional-encoding

# 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



https://huggingface.co/blog/designing-positional-encoding

# Break & Questions

# Outline

- **Basic Attention**
  - Notions of attention, self-attention, basic attention layer, QKV setup and intuition
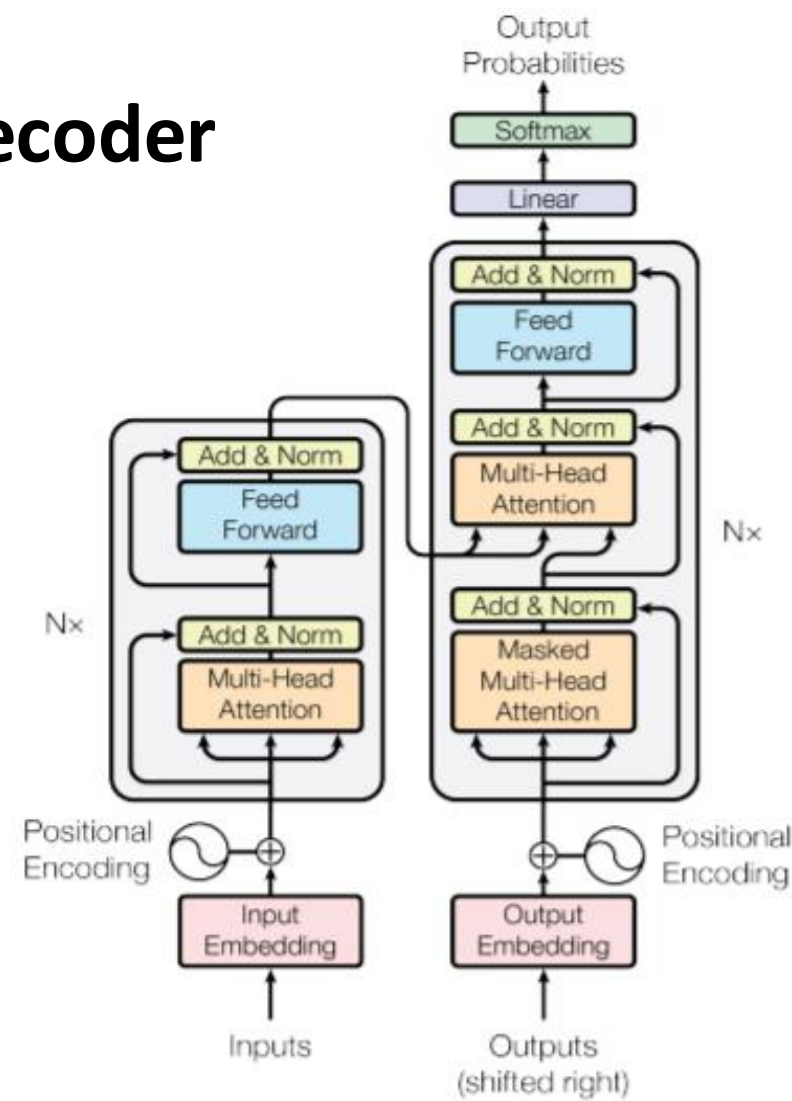- **Additional Elements**
  - Multi-head attention, positional encodings
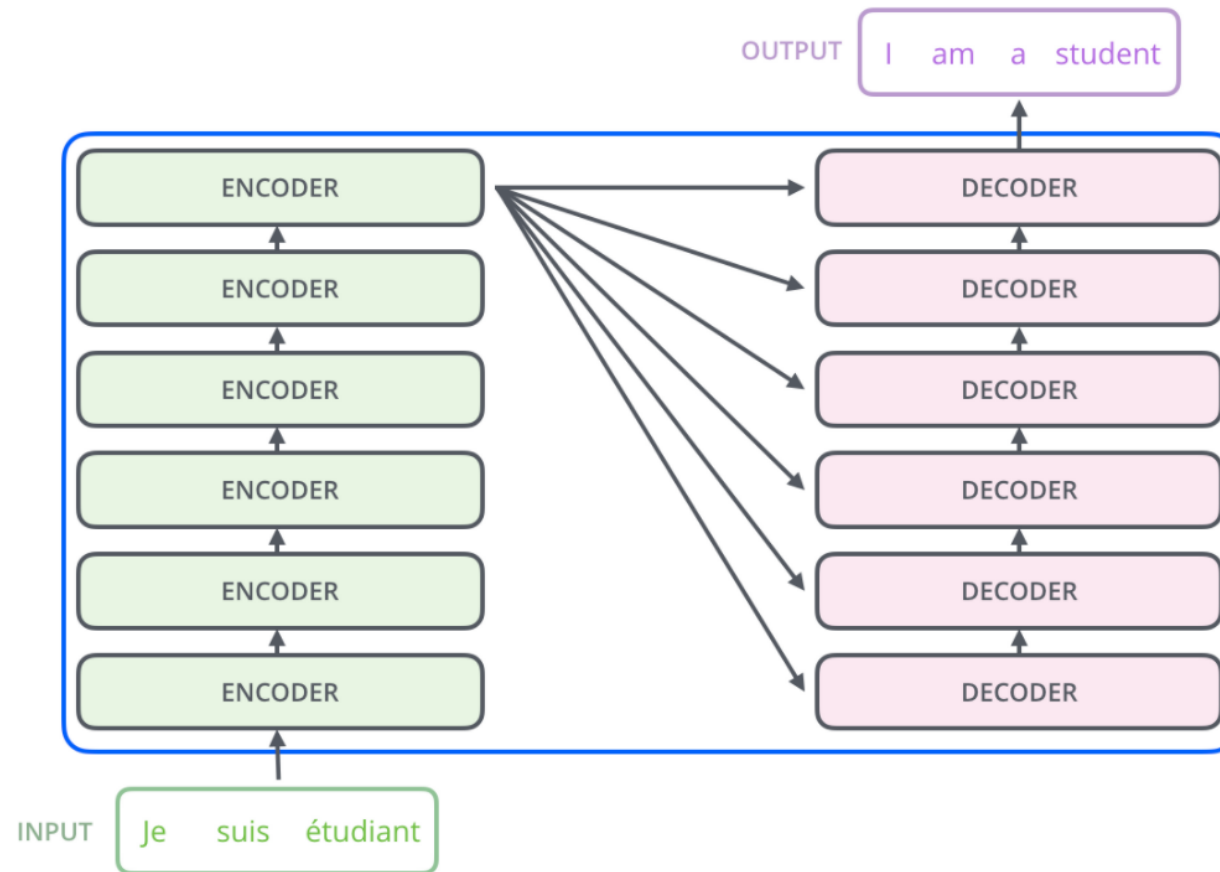- **Transformers**
  - Architecture, encoder and decoder setups

# **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
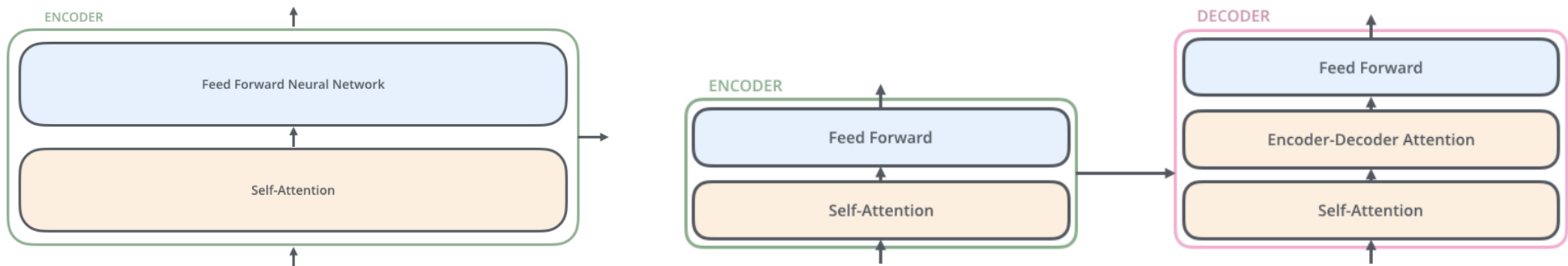  - Centered on self-attention blocks

Vaswani et al. '17

# **Transformers**: Architecture

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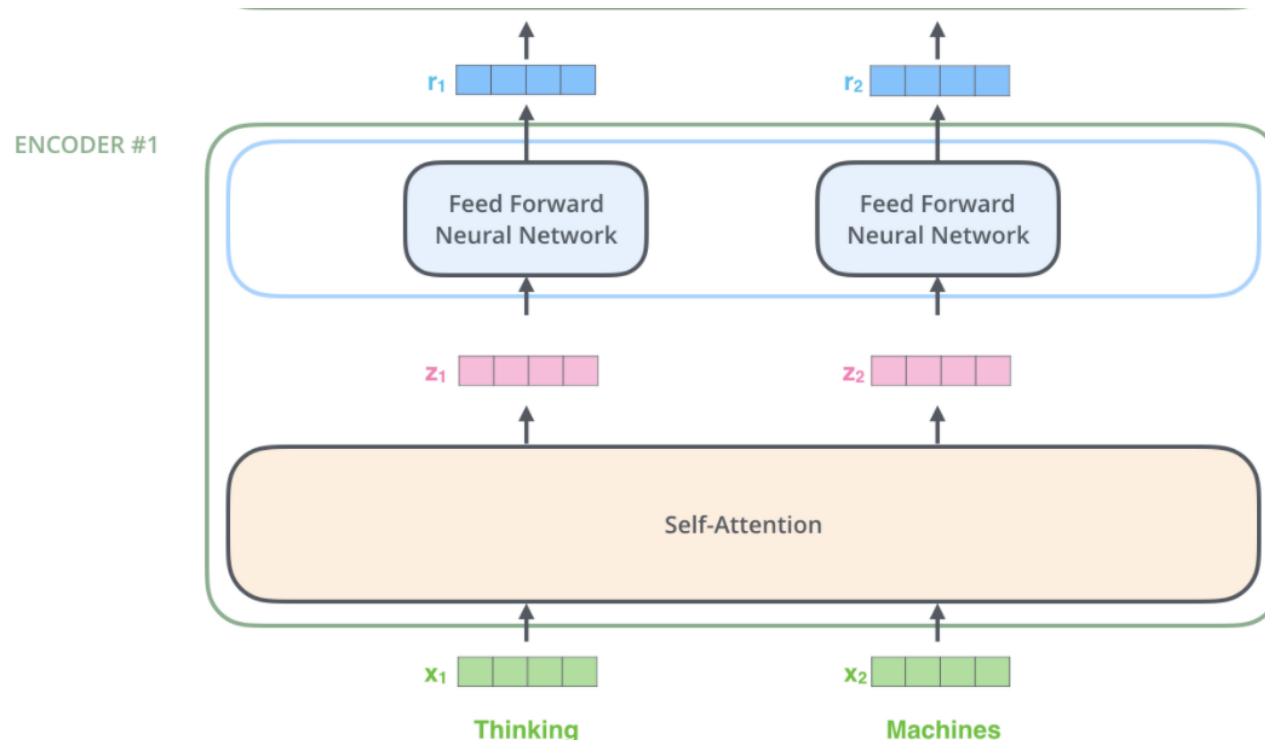  - For example, for French-English translation:

# **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
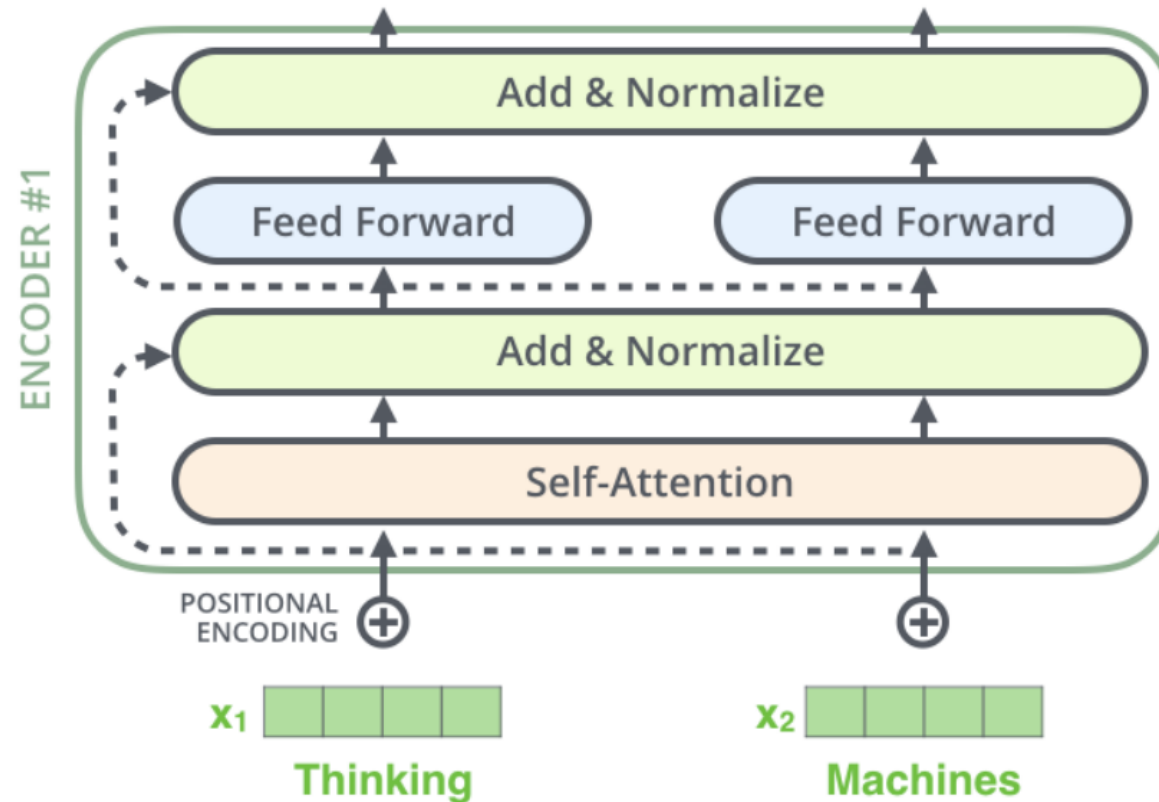    - Fully-connected layers (i.e., an MLP)

# **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
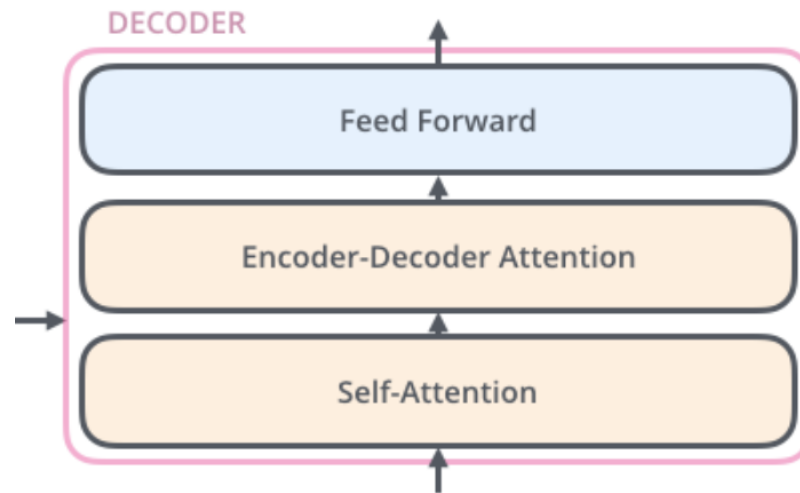
# **Transformers**: More Tricks

- 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?

am

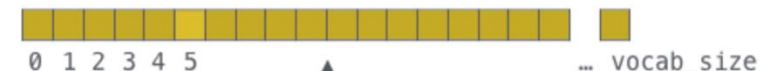Get the index of the cell with the highest value (argmax)

5

log_probs

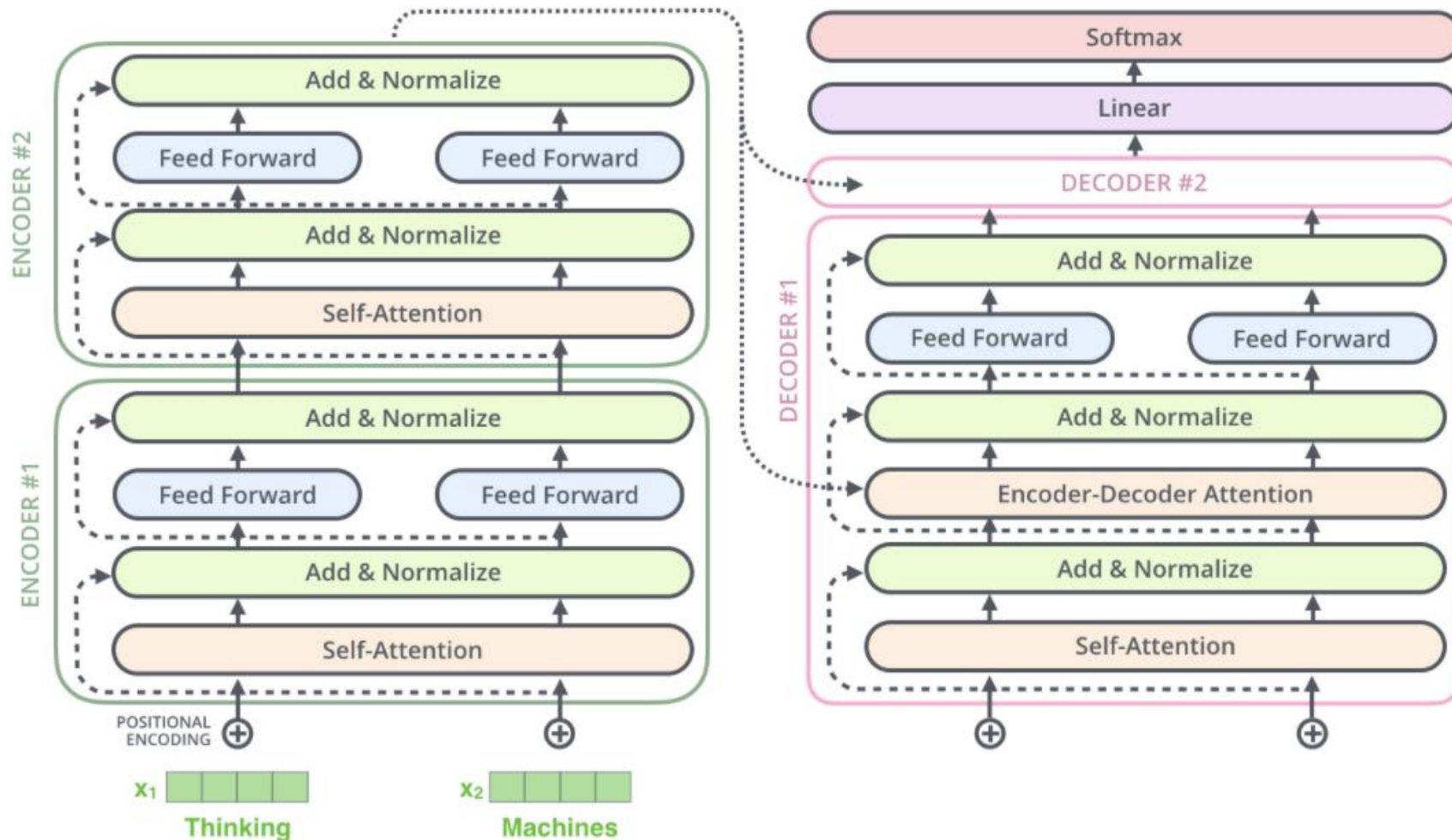0 1 2 3 4 5 ... vocab_size

Softmax

logits

0 1 2 3 4 5 ... vocab_size
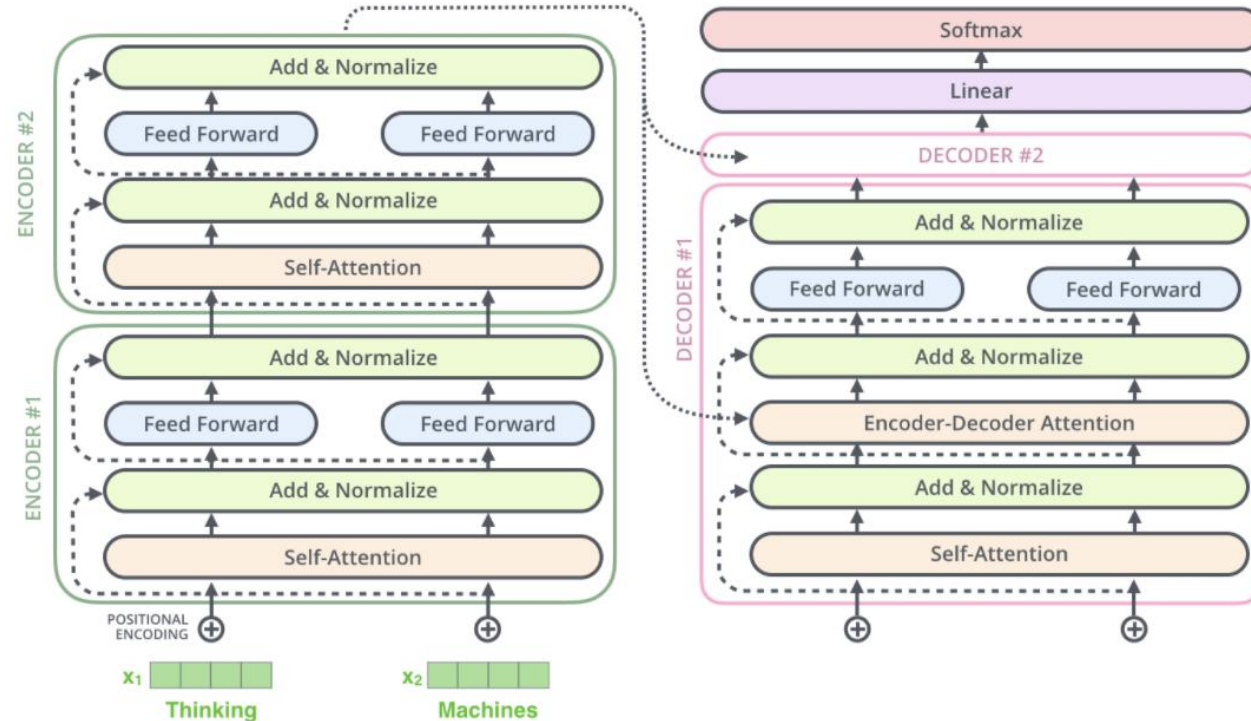
Linear

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.

# Thank You!