



# CS 760: Machine Learning Language Models

Misha Khodak

University of Wisconsin-Madison  
27 October 2025

# Announcements

- Homework 3 due November 5<sup>th</sup>
- Reminder: no class November 26<sup>th</sup>

# Outline

- **Language Models & NLP**
  - RNNs, word embeddings, attention
- **Transformer Model**
  - Properties, architecture breakdown
- **Transformer-based Models**
  - BERT, GPTs, Foundation Models

# Outline

- **Language Models & NLP**
  - RNNs, word embeddings, attention
- **Transformer Model**
  - Properties, architecture breakdown
- **Transformer-based Models**
  - BERT, GPTs, Foundation Models

# Language Models: Word Embeddings

- One way to encode words: one-hot vectors
  - Want something smarter...

**Distributional semantics:** account for relationships

- Representations should be close/similar to other words that appear in a similar context

Dense vectors:

$$\text{dog} = [0.13 \quad 0.87 \quad -0.23 \quad 0.46 \quad 0.87 \quad -0.31]^T$$

$$\text{cat} = [0.07 \quad 1.03 \quad -0.43 \quad -0.21 \quad 1.11 \quad -0.34]^T$$

AKA **word embeddings**



# Training Word Embeddings

Many approaches (very popular 2010-present)

- Word2vec: a famous approach
- Write out a likelihood

$$L(\theta) = \prod_{t=1}^T \prod_{-a \leq j \leq a} P(w_{t+j} | w_t, \theta)$$

Our word vectors (weights)

All positions

Windows of length  $2a$



# Training Word Embeddings

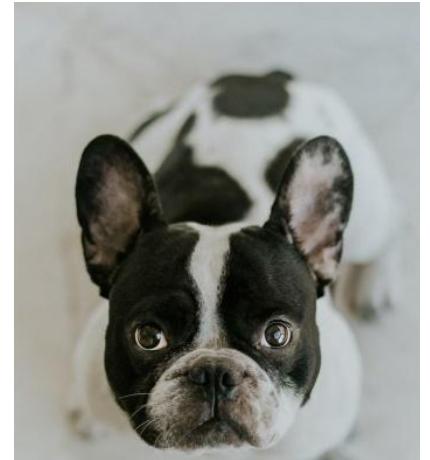
Word2vec likelihood

$$L(\theta) = \prod_{t=1}^T \prod_{-a \leq j \leq a} P(w_{t+j} | w_t, \theta)$$

- Expression for the probability:

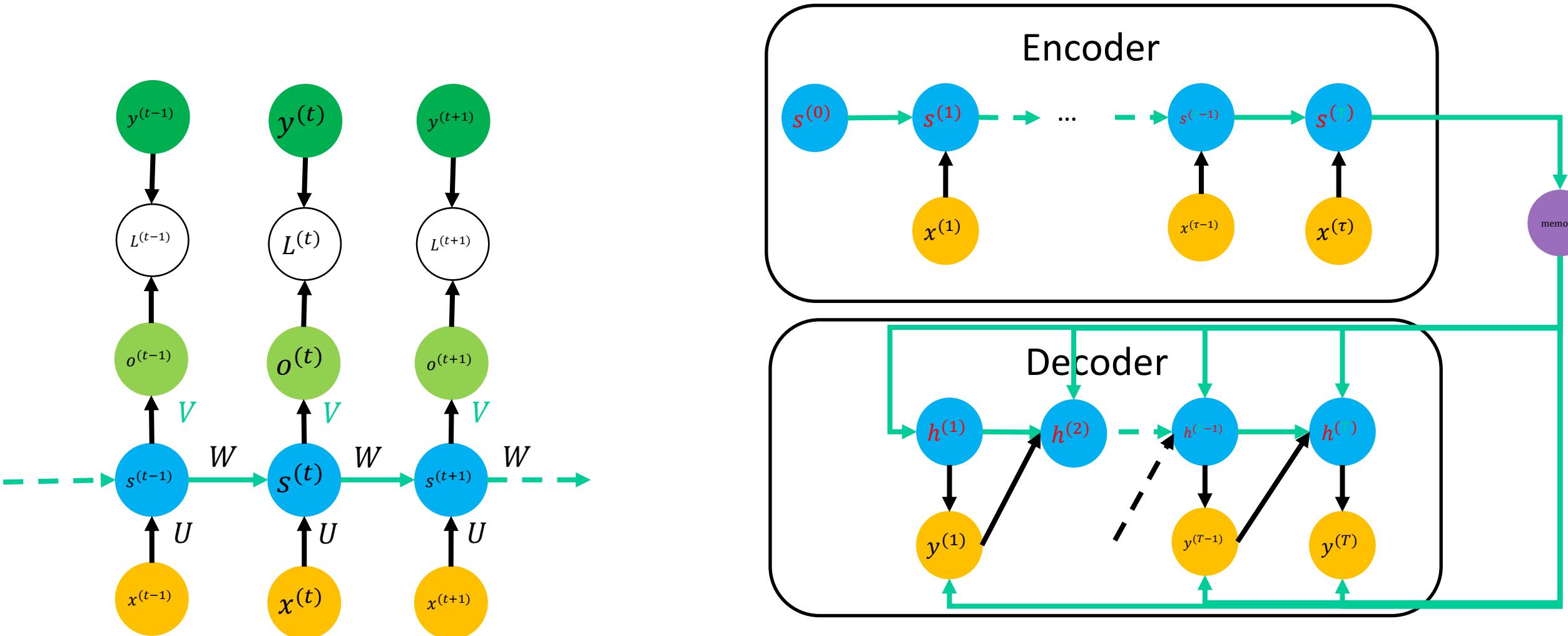
$$P(w' | w, \theta) = \frac{\exp((\theta_{w',o})^\top \theta_{w,c})}{\sum_{v \in V} \exp((\theta_{v,o})^\top \theta_{w,c})}$$

- $\theta_{w,o}$ : occurrence vector for word  $w$
- $\theta_{w,c}$ : context vector for word  $w$



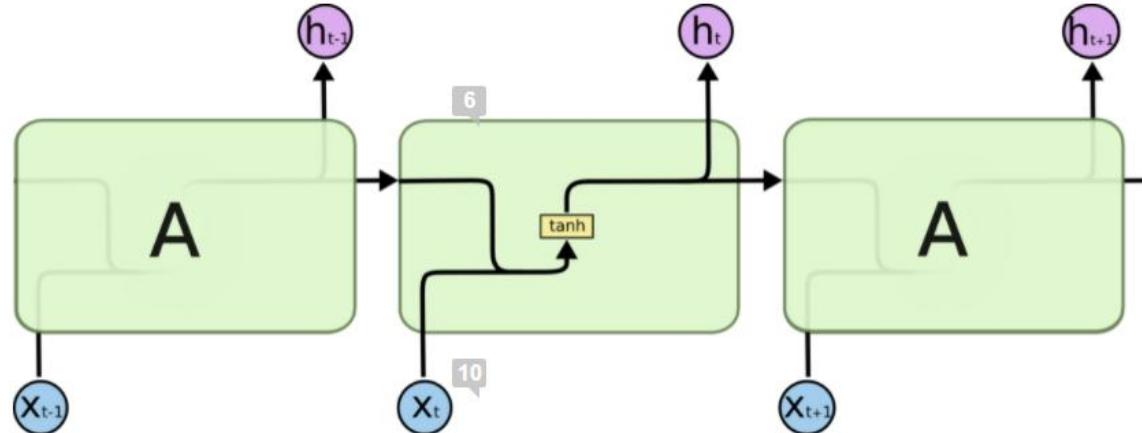
# Language Models: RNN Review

- Classical RNN model / Encoder-Decoder variant:

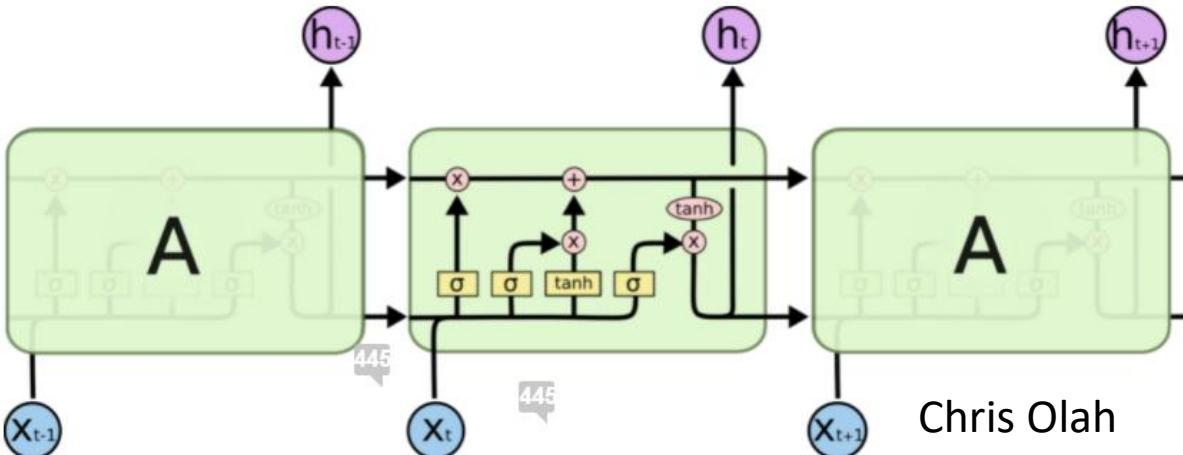


# Language Models: LSTM Review

- RNN: can write structure as:

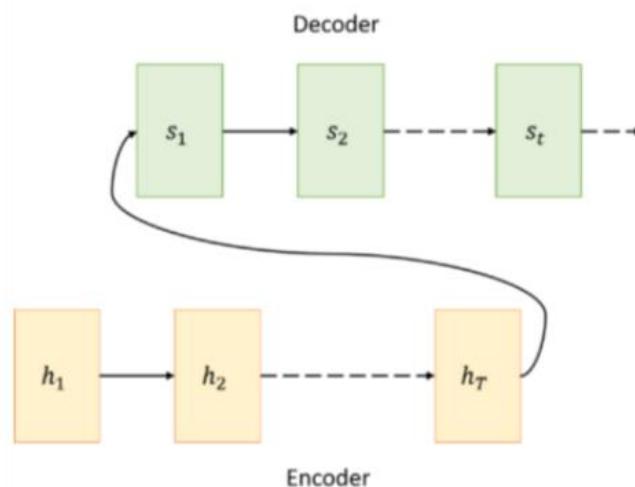


- Long Short-Term Memory: deals with problem. Cell:

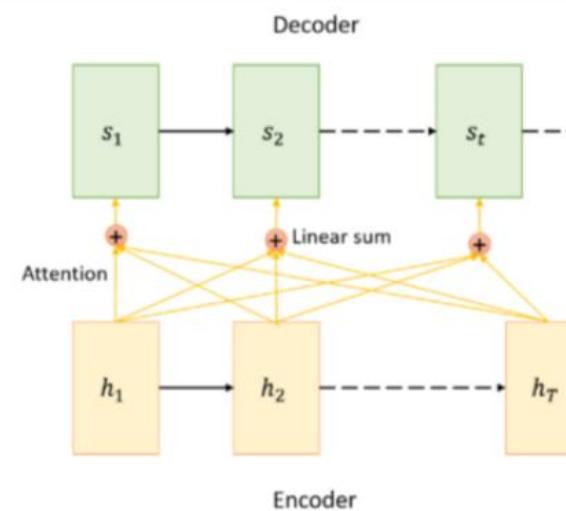


# Language Models: Attention

- One challenge: dealing with the hidden state
  - Everything gets compressed there
  - Might lose information
- Solution: **attention** mechanism
  - Similar to residual connections in ResNets!



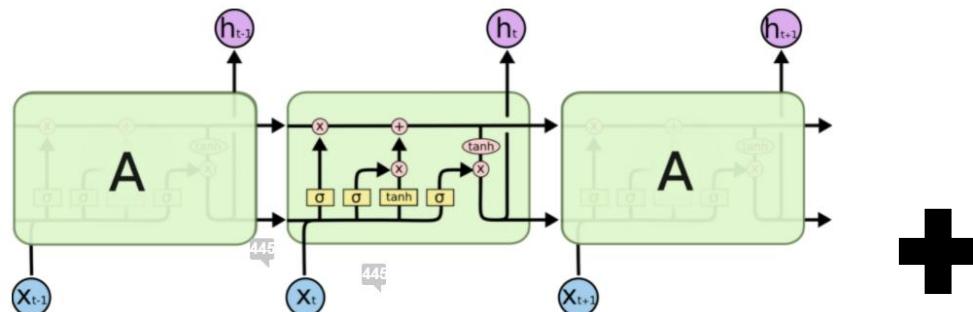
(a) Vanilla Encoder Decoder Architecture



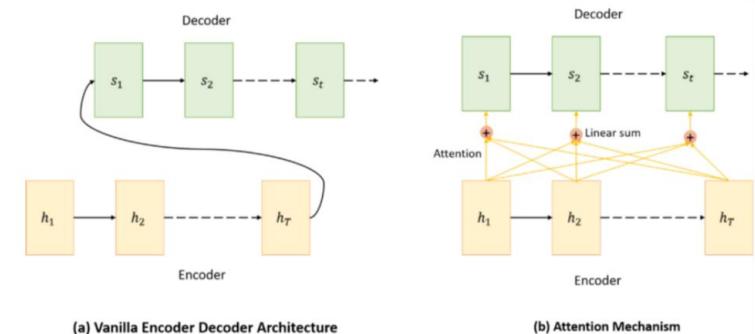
(b) Attention Mechanism

# Language Models: Putting it All Together

- Before 2017: best language models
  - Use encoder/decoder architectures based on **RNNs**
  - Use **word embeddings** for word representations
  - Use attention mechanisms



$$\text{dog} = [0.13 \quad 0.87 \quad -0.23 \quad 0.46 \quad 0.87 \quad -0.31]^T$$

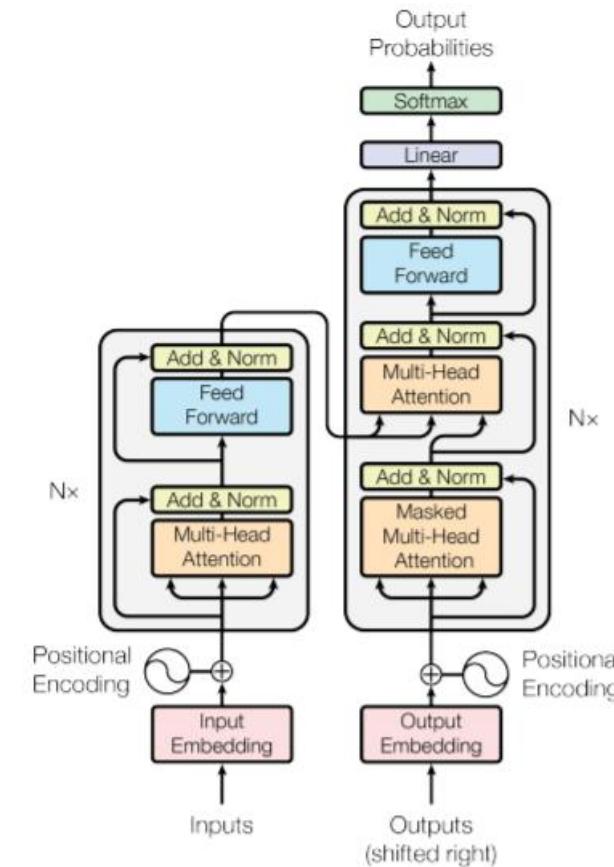


# Outline

- **Language Models & NLP**
  - k-gram models, RNN review, word embeddings, attention
- **Transformer Model**
  - Properties, architecture breakdown
- **Transformer-based Models**
  - BERT, GPTs, Foundation Models

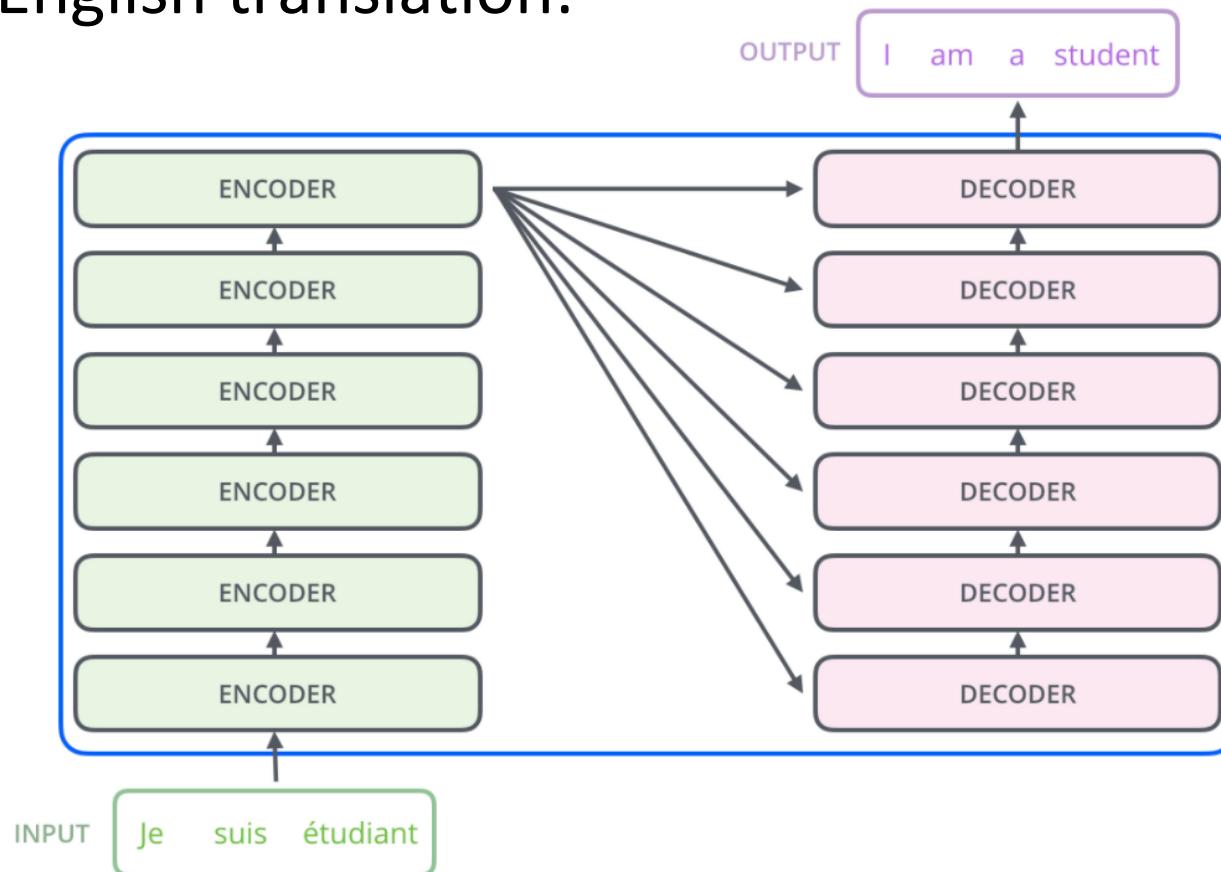
# Transformers: Idea

- Initial goal for an architecture: encoder-decoder
  - Get **rid of recurrence**
  - Replace with **self-attention**



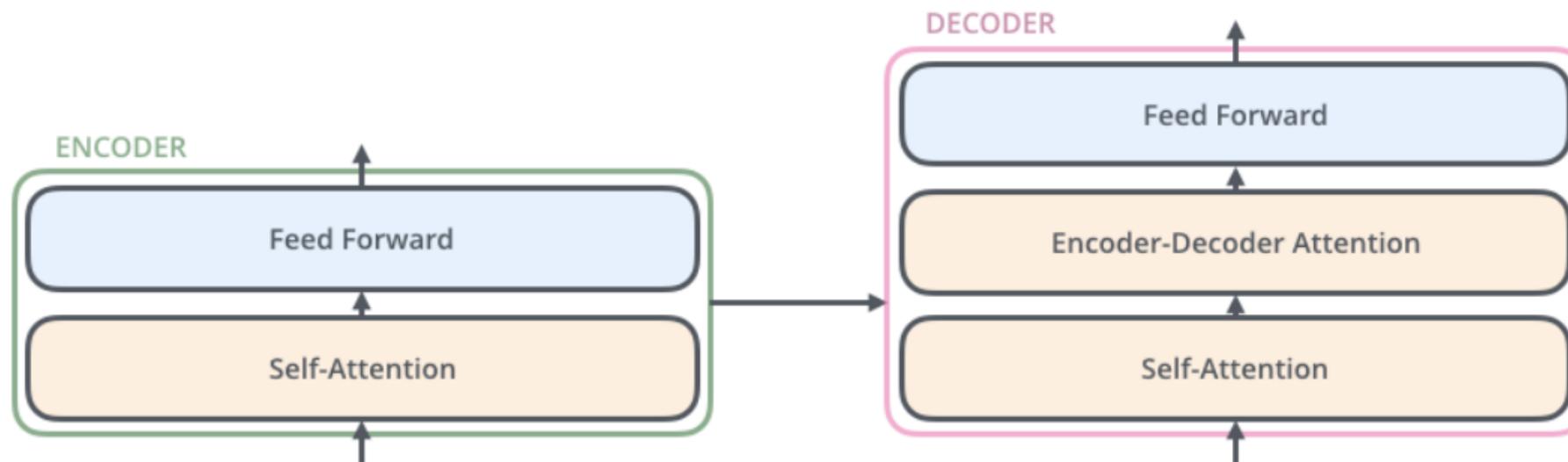
# Transformers: Architecture

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



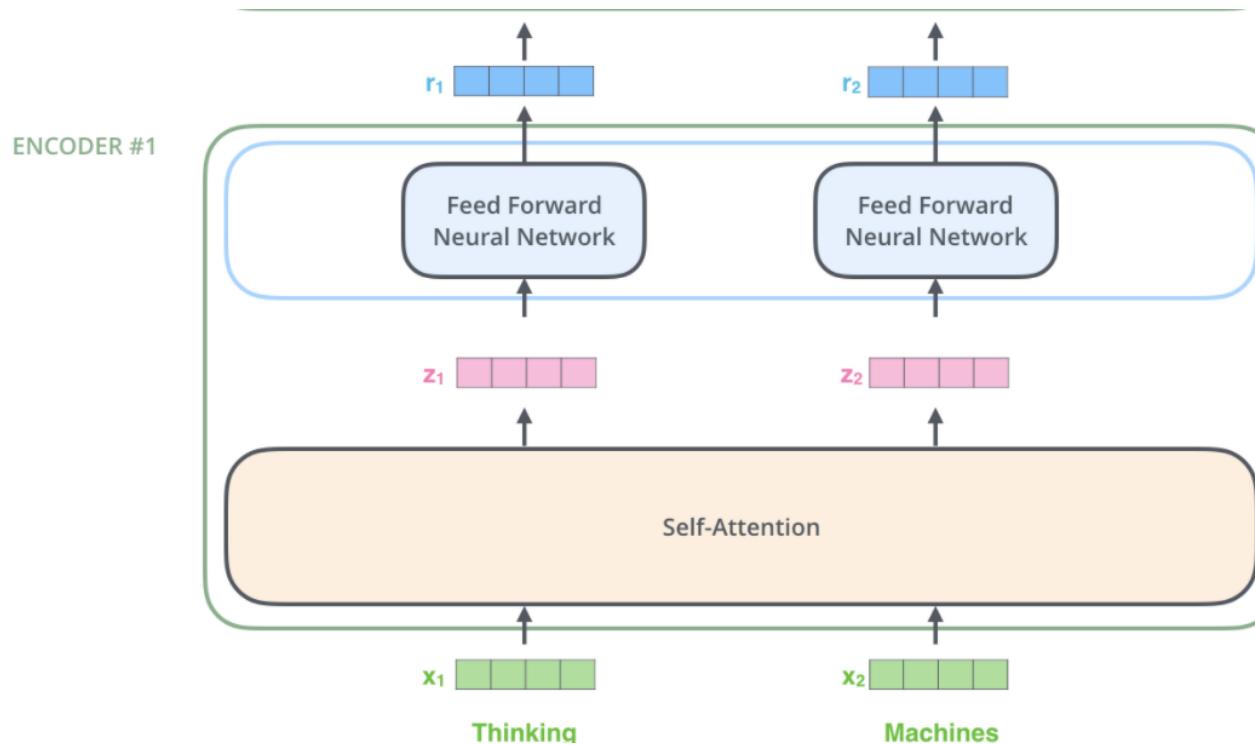
# Transformers: Architecture

- Sequence-sequence model with **stacked** encoders/decoders:
  - What's inside each encoder/decoder unit?



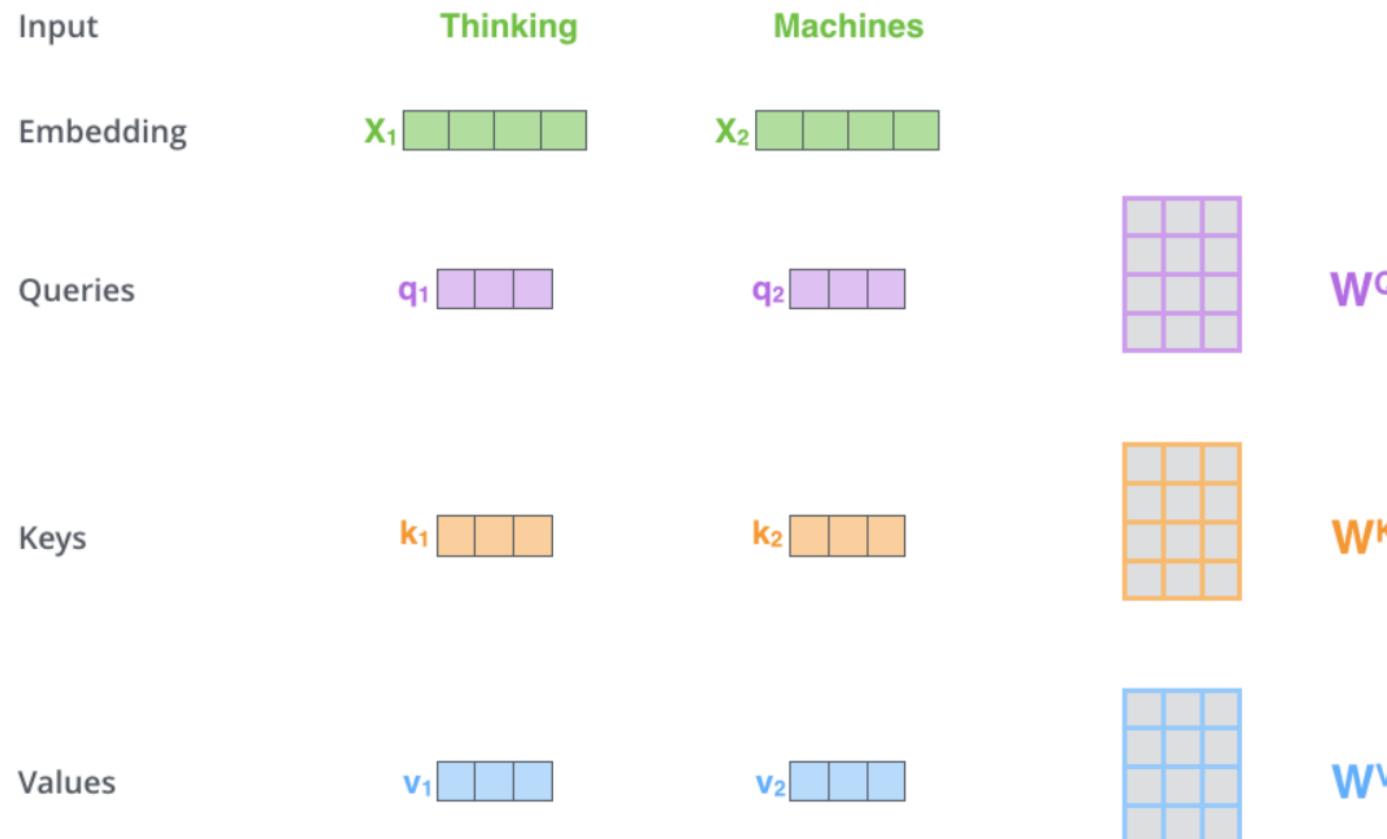
# Transformers: Inside an Encoder

- Let's take a look at the encoder. Two components:
  - 1. **Self-attention** layer
  - 2. **Feedforward nets**



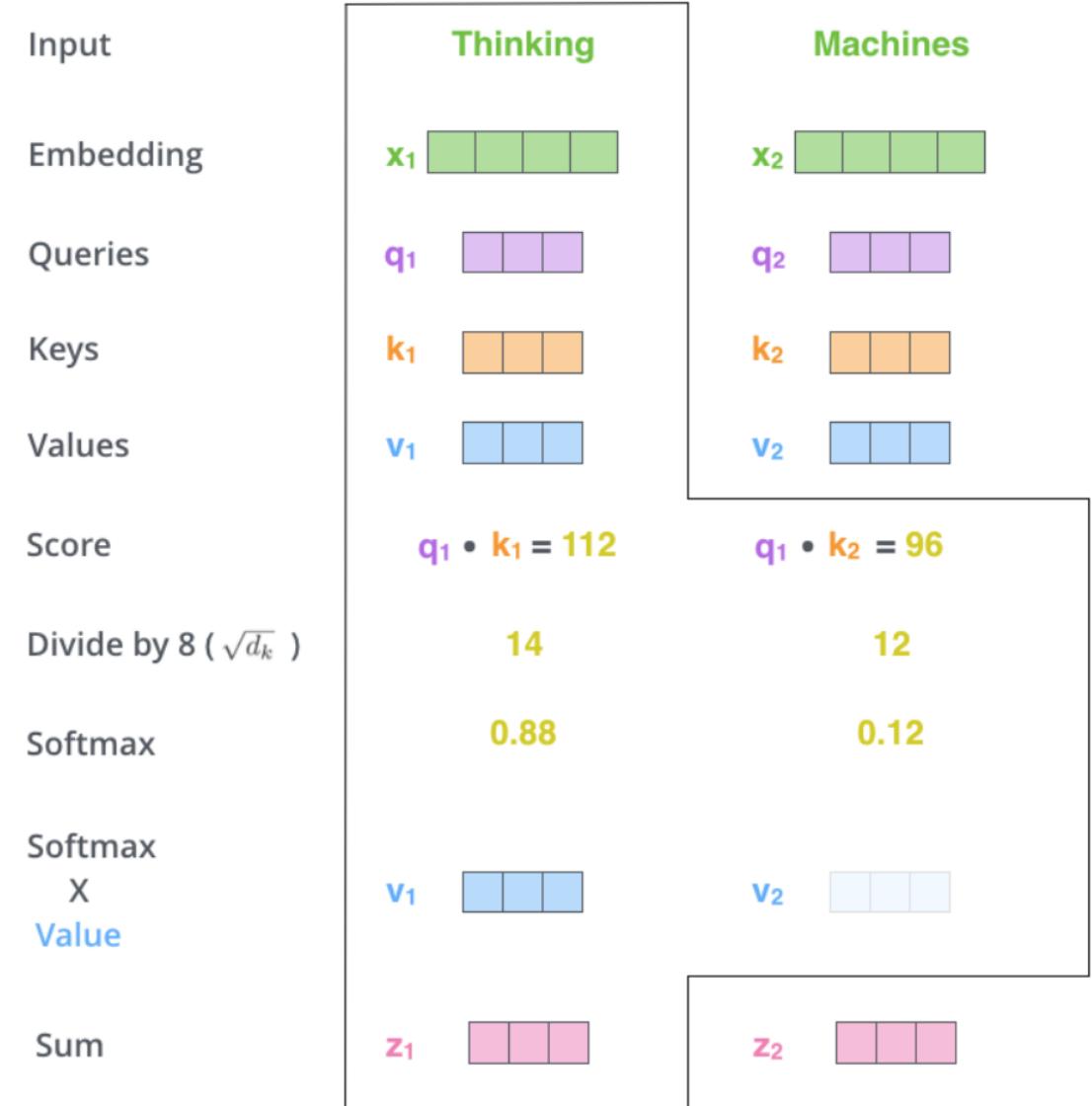
# Transformers: Self-Attention

- Self-attention is the key layer in a transformer stack
  - Get 3 vectors for each embedding: **Query**, **Key**, **Value**



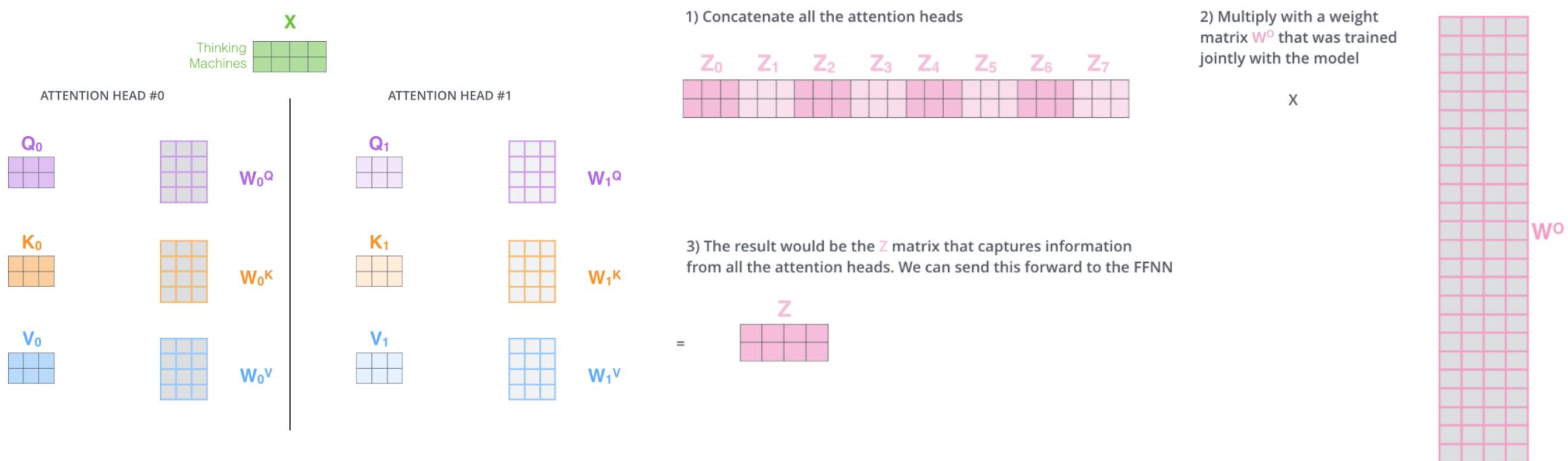
# Transformers: Self-Attention

- Self-attention is the key layer in a transformer stack
  - Illustration. Recall the three vectors for each embedding: **Query**, **Key**, **Value**
  - The sum values are the outputs of the self-attention layer
  - Send these to feedforward NNs
  - Highly parallelizable!



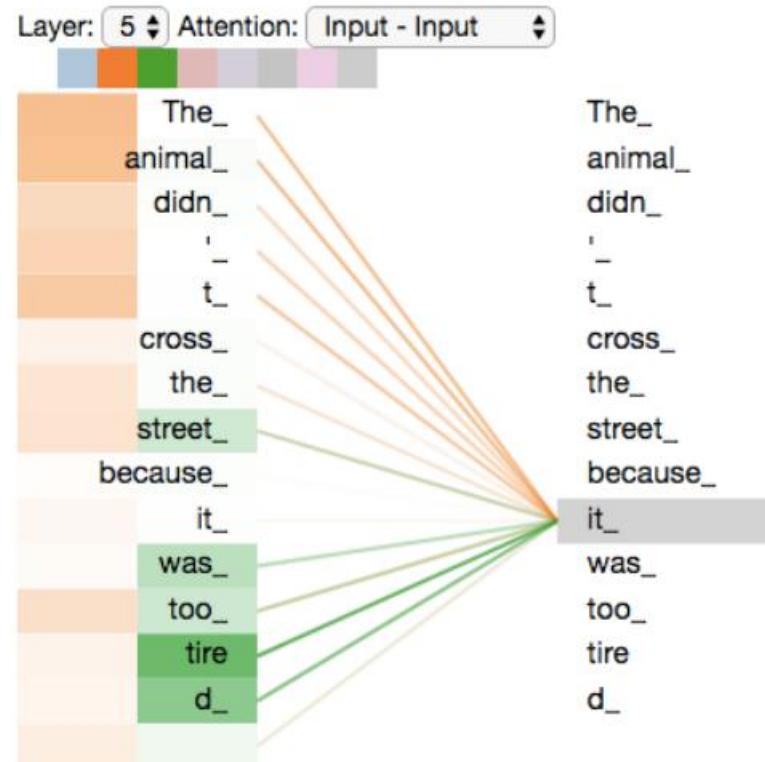
# Transformers: Multi-Headed Attention

- We can do this multiple times in parallel
  - Called multiple heads
  - Need to combine the resulting output sums



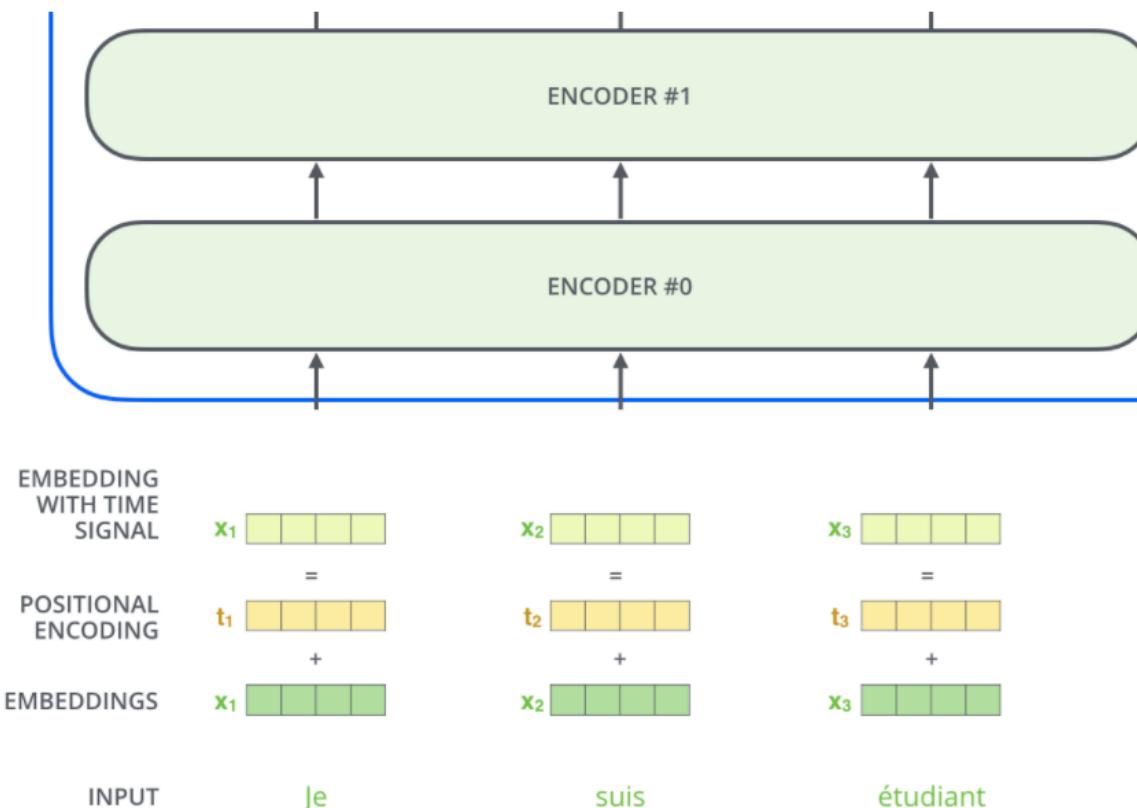
# Transformers: Attention Visualization

- Attention tells us where to focus the information
  - Illustration for a sentence:



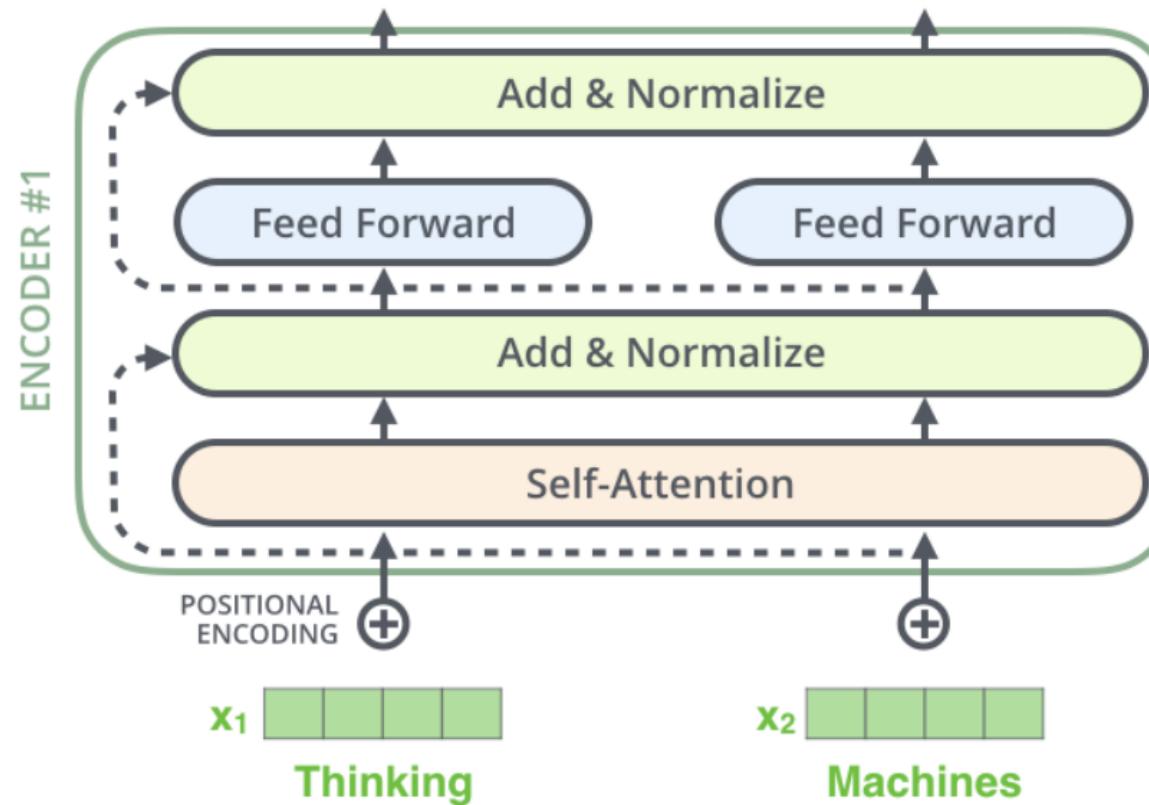
# Transformers: Positional Encodings

- One thing we haven't discussed: the order of the symbols/elements in the sequence
  - Add a vector containing a special positional formula's embedding



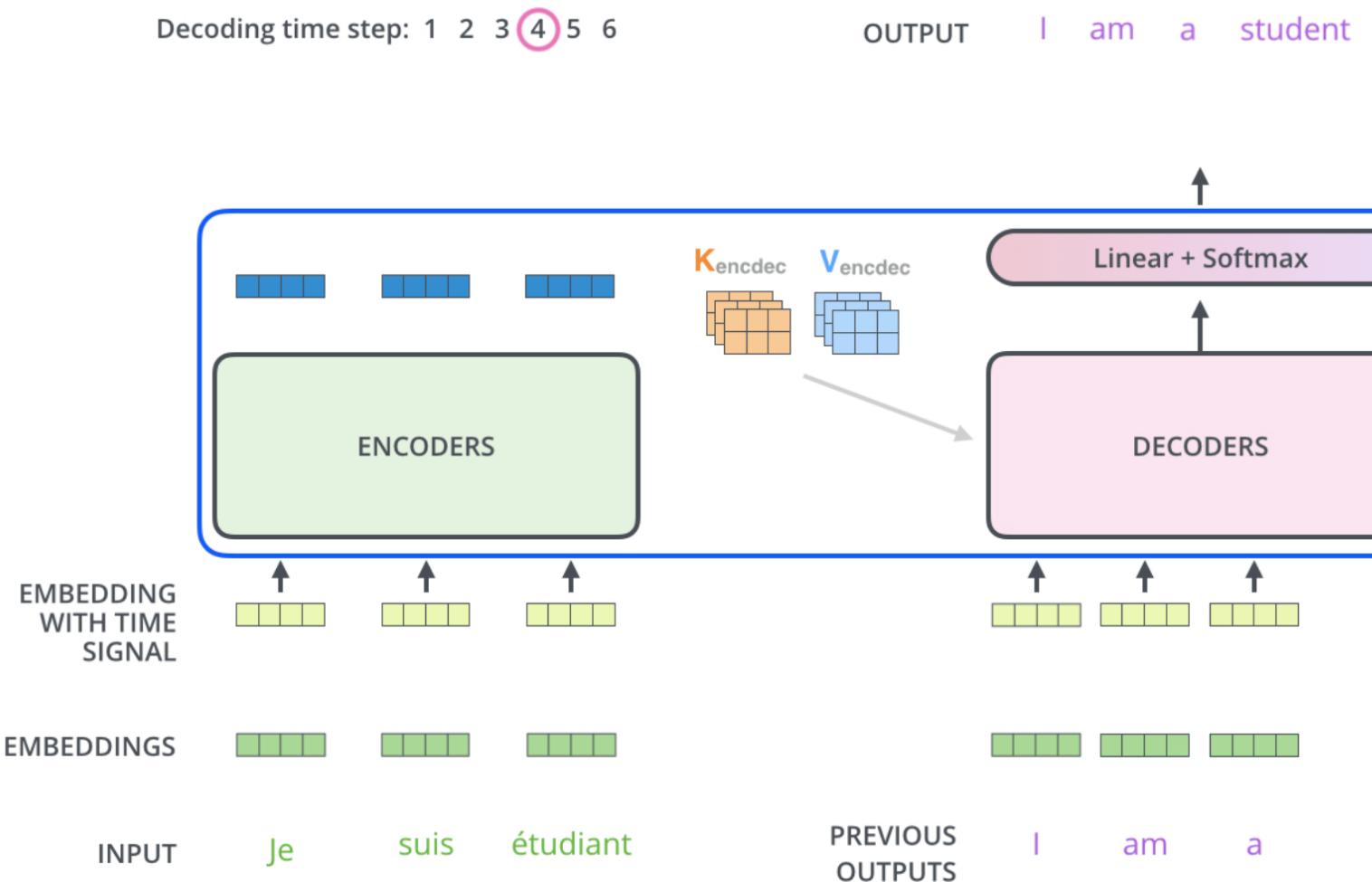
# Transformers: More Tricks

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



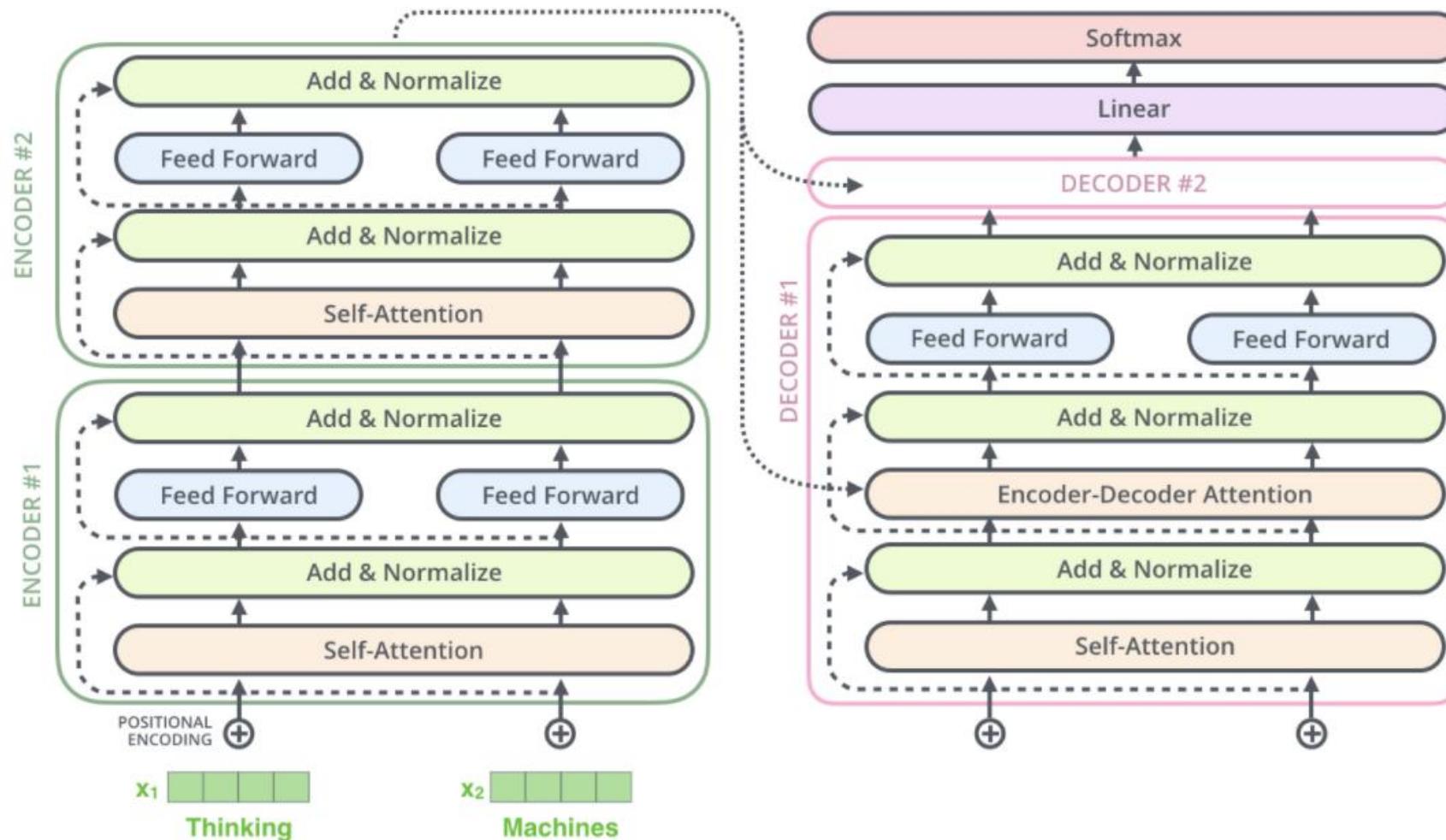
# Transformers: Decoder

- Similar to encoders (see blog post for more details).
- E.g. Generating a translation



# Transformers: Putting it All Together

- What does the full architecture look like?



# Outline

- **Language Models & NLP**
  - k-gram models, RNN review, word embeddings, attention
- **Transformer Model**
  - Properties, architecture breakdown
- **Transformer-based Models**
  - BERT, GPTs, Foundation Models

# Transformer-Based Models: BERT

- Semi-supervised learning + Transformers
  - Semi-supervised learning to learn embeddings in encoder

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.

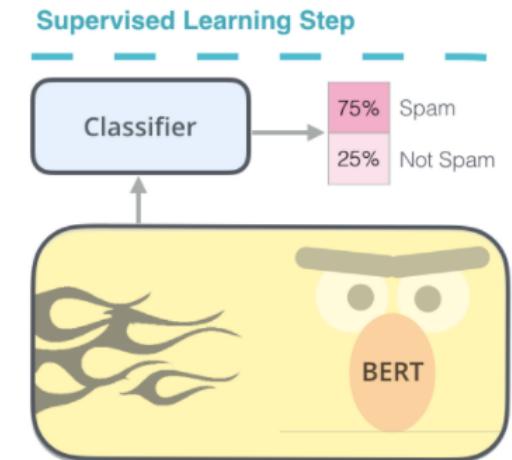


Dataset:



Objective:  
Predict the masked word  
(language modeling)

2 - **Supervised** training on a specific task with a labeled dataset.



Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

# BERT: Concepts

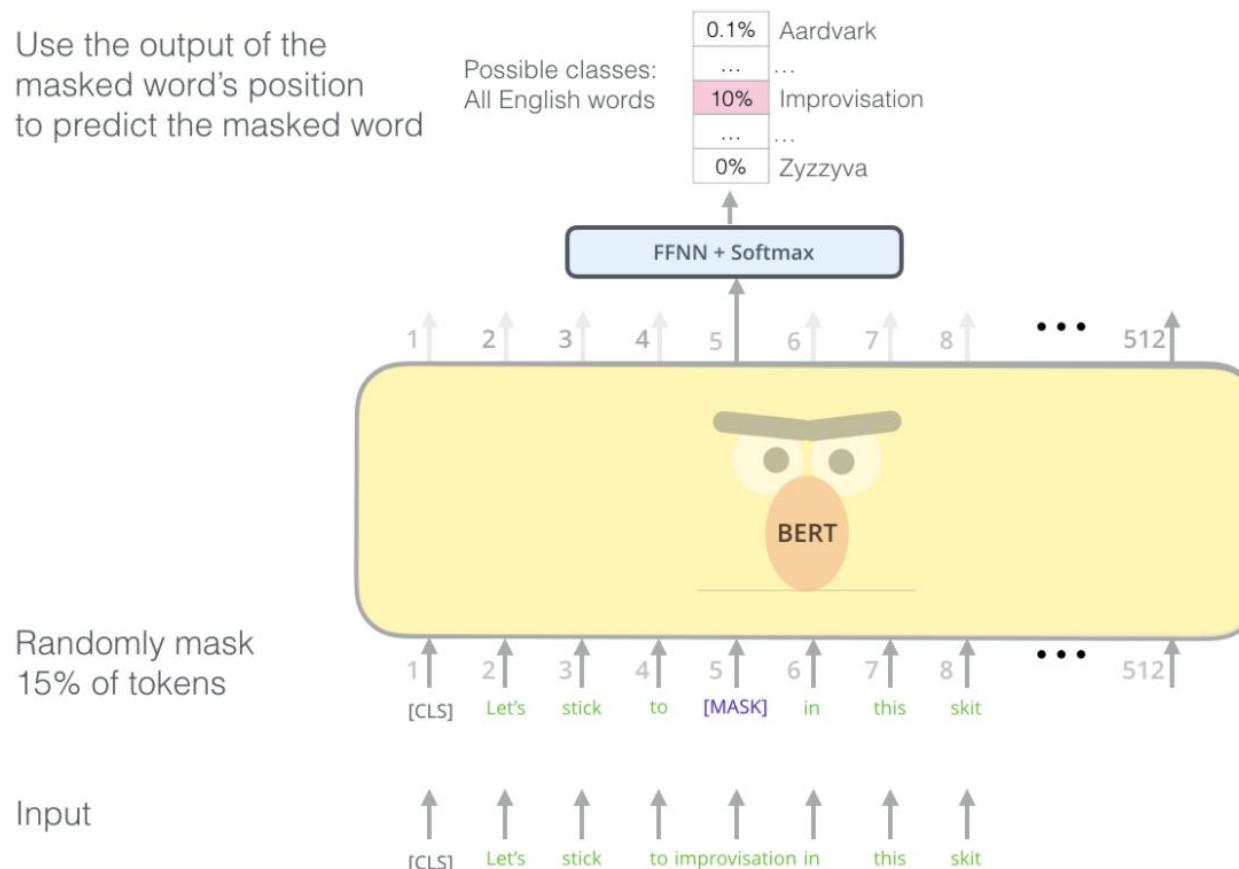
- What makes BERT work? A bunch of ideas:
  - 1. Use the **Transformer** architecture
  - 2. **Pre-training** on corpora using pretext tasks
    - Then fine-tune for a particular task
  - 3. **Scale**: BERT-Large has 340 million parameters

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

Results: Devlin et al, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

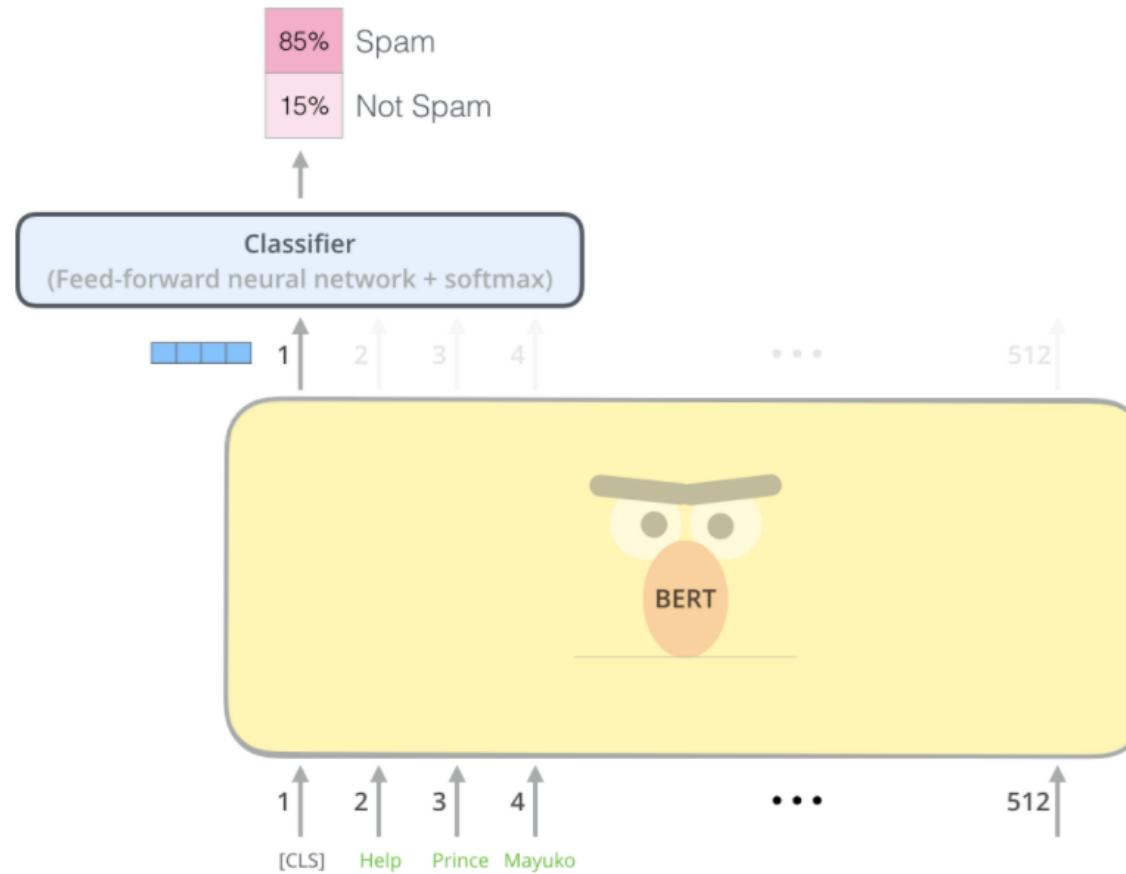
# BERT: Training

- BERT is trained on a simple tasks on a huge amount of data:
  - Recall our pretext tasks in self-supervised learning
  - **Masked word prediction:**



# BERT: Classification

- Then, fine-tune on a particular task
  - Example: **binary classification**, spam VS not spam

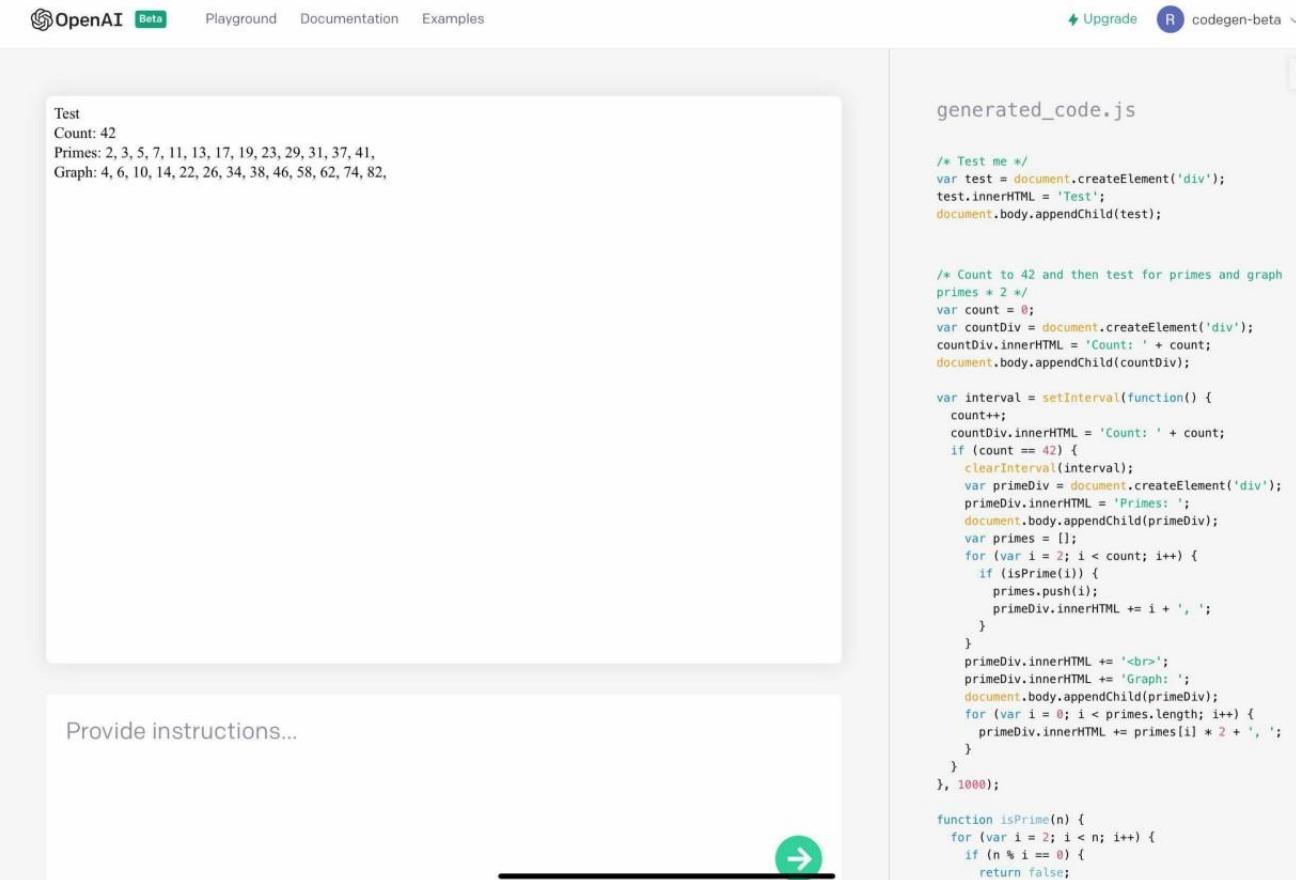


# GPT Series of Models

- **GPT: Generative Pre-trained Transformer**
  - Also built on top of transformer model architecture
  - Essentially the decoder part only
- Goal: generate text (possibly from a **prompt**)
- Scale: huge!
  - GPT-3: 175 billion parameters

# Codex

- Codex: a variant of GPT-3 based on source code
  - Outputs code. Ex: show primes



The image shows a screenshot of the OpenAI Codex interface. On the left, a text input field contains the instruction: "Provide instructions...". On the right, the generated code is displayed in a code editor window titled "generated\_code.js". The code is a JavaScript script that first creates a test div with the text "Test". It then counts to 42 and tests for primes and graphs. For each prime number, it creates a prime div, adds the prime to an array, and appends the prime to the prime div's inner HTML. Finally, it appends the prime div to the body. The code includes a check for primality using a helper function "isPrime".

```
/* Test me */
var test = document.createElement('div');
test.innerHTML = 'Test';
document.body.appendChild(test);

/* Count to 42 and then test for primes and graph
primes * 2 */
var count = 0;
var countDiv = document.createElement('div');
countDiv.innerHTML = 'Count: ' + count;
document.body.appendChild(countDiv);

var interval = setInterval(function() {
  count++;
  countDiv.innerHTML = 'Count: ' + count;
  if (count == 42) {
    clearInterval(interval);
    var primeDiv = document.createElement('div');
    primeDiv.innerHTML = 'Primes: ';
    document.body.appendChild(primeDiv);
    var primes = [];
    for (var i = 2; i < count; i++) {
      if (isPrime(i)) {
        primes.push(i);
        primeDiv.innerHTML += i + ', ';
      }
    }
    primeDiv.innerHTML += '<br>';
    primeDiv.innerHTML += 'Graph: ';
    document.body.appendChild(primeDiv);
    for (var i = 0; i < primes.length; i++) {
      primeDiv.innerHTML += primes[i] * 2 + ', ';
    }
  }
}, 1000);

function isPrime(n) {
  for (var i = 2; i < n; i++) {
    if (n % i == 0) {
      return false;
    }
  }
}
```

# DALL-E

- Create images from text
  - Prompt: “an armchair in the shape of an avocado. . .”



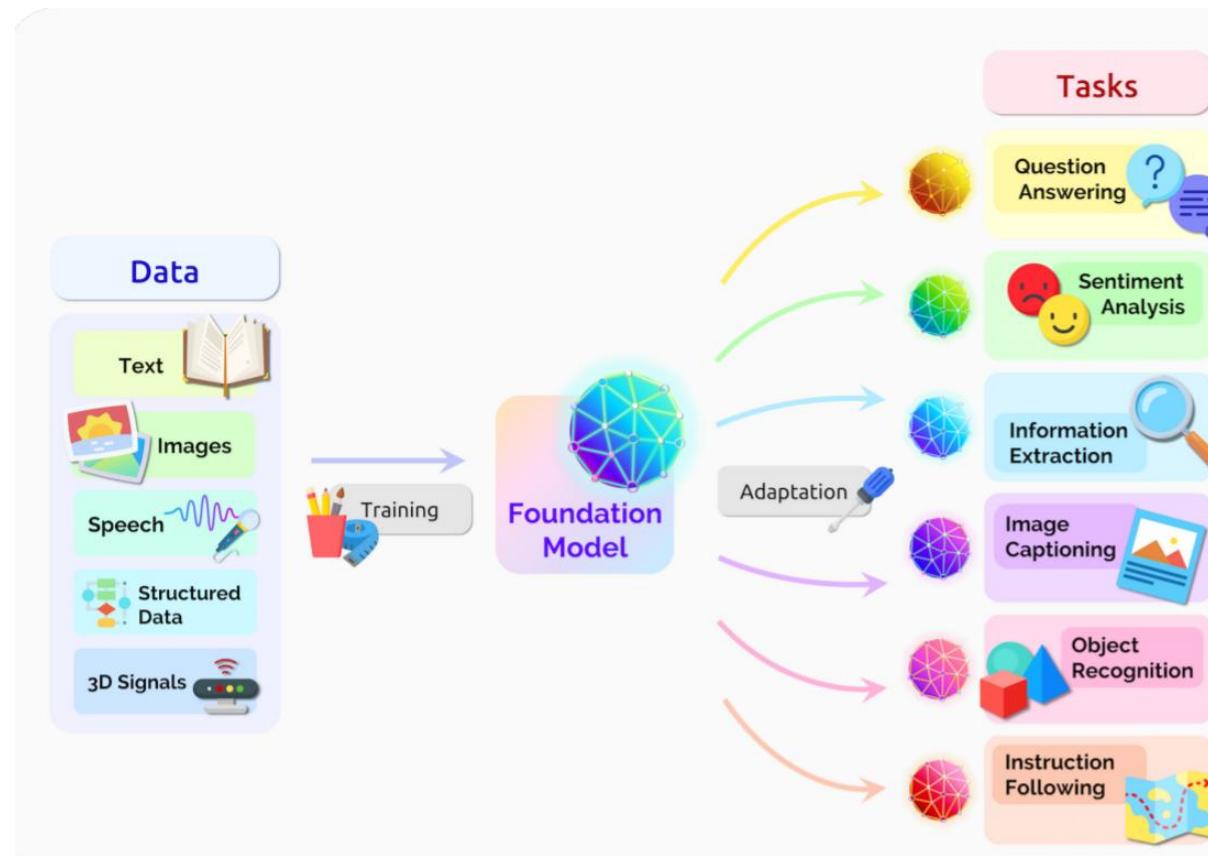
<https://openai.com/blog/dall-e/>

- Note: several online demos. Try it yourself!

<https://huggingface.co/stabilityai/stable-diffusion-3.5-large>

# Foundation Models

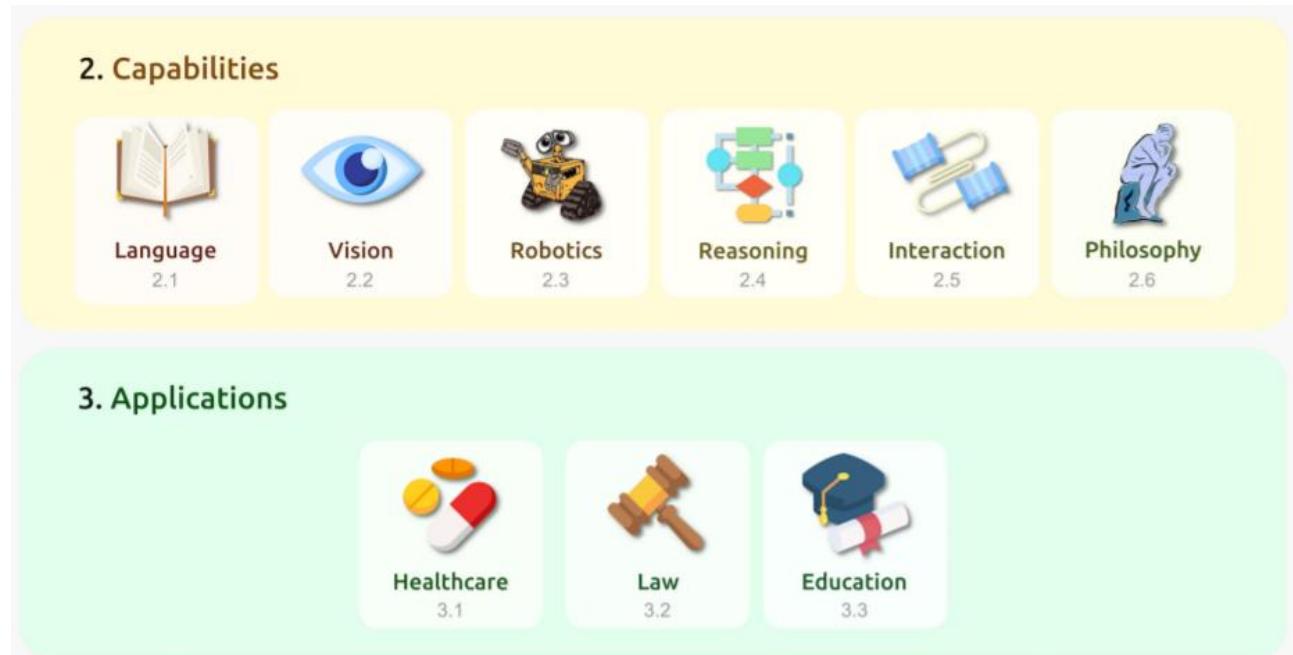
- Many more large scale models
  - Not just focused on text



Bommasani et al, "On the Opportunities and Risks of Foundation Models"

# Conclusion

- “Foundation” models based on transformers and beyond
  - Huge, expensive to train, challenging in various ways... but
  - Remarkably powerful for a vast number of tasks.
  - AGI??



Bommasani et al, “On the Opportunities and Risks of Foundation Models”



# Thanks Everyone!

Some of the slides in these lectures have been adapted/borrowed from materials developed by Mark Craven, David Page, Jude Shavlik, Tom Mitchell, Nina Balcan, Elad Hazan, Tom Dietterich, Pedro Domingos, Jerry Zhu, Yingyu Liang, Volodymyr Kuleshov, Jay Alammar, Fred Sala, Kirthi Kandasamy, Tengyang Xie