



# CS 760: Machine Learning **Large Language Models**

Fred Sala

University of Wisconsin-Madison

**Dec. 9, 2021**

# Announcements

- **Logistics:**

- HW 8 / Project Due Dec. 14
- Final: Dec. 20<sup>th</sup>
- **Surveys** --- thank you!

- **Class roadmap:**

<b>Thurs., Dec 9</b>	<b>Large Language Models</b>
Tues., Dec 14	Fairness & Ethics
Monday, Dec 20	<b>Final Exam</b>



# Training Human Intelligence Updates

- 6 months of training
  - ~100 billion neurons
  - Compare to today's model parameter counts



# Outline

- **Language Models & NLP**

- k-gram models, RNN review, word embeddings, attention

- **Transformer Model**

- Properties, architecture breakdown

- **Transformer-based Models**

- BERT, GPTs, Foundation Models

# Outline

- **Language Models & NLP**

- k-gram models, RNN review, word embeddings, attention

- **Transformer Model**

- Properties, architecture breakdown

- **Transformer-based Models**

- BERT, GPTs, Foundation Models

# Language Models: History

- Pre-date neural networks
- Basic idea: use probabilistic models to **assign a probability to a sentence**

$$P(W) = P(w_1, w_2, \dots, w_n) \text{ or } P(w_{\text{next}} | w_1, w_2 \dots)$$

- Goes back to Shannon
  - Information theory: letters

Zero-order approximation	XFOML RXKHRJFFJUJ ALPWXFVJXYJ FFJEYVJCSGHYD QPAAMKBZAACIBZLKJQD
First-order approximation	OCRO HLO RGWR NMIELWIS EU LL NBNESEBYA TH EEI ALHENHTTPA OObTTVA NAH BRL
Second-order approximation	ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY ACHIN D ILONASIVE TUOOWE AT TEASONARE FUSO TIZIN ANDY TOBE SEACE CTISBE
Third-order approximation	IN NO IST LAT WHEY CRATICT FROURE BIRS GROCID PONDENOME OF DEMONSTURES OF THE REPTAGIN IS REGOACTIONA OF CRE
First-order word approximation	REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME CAN DIFFERENT NATURAL HERE HE THE A IN CAME THE TO OF TO EXPERT GRAY COME TO FURNISHES THE LINE MESSAGE HAD BE THESE

# Language Models: History

- Classic approach: make Markov-type assumptions:

$$P(w_i | w_{i-1} w_{i-2} \dots w_1) = P(w_i | w_{i-1} w_{i-2} \dots w_{i-k})$$

- Ex: k=1: Bigram model

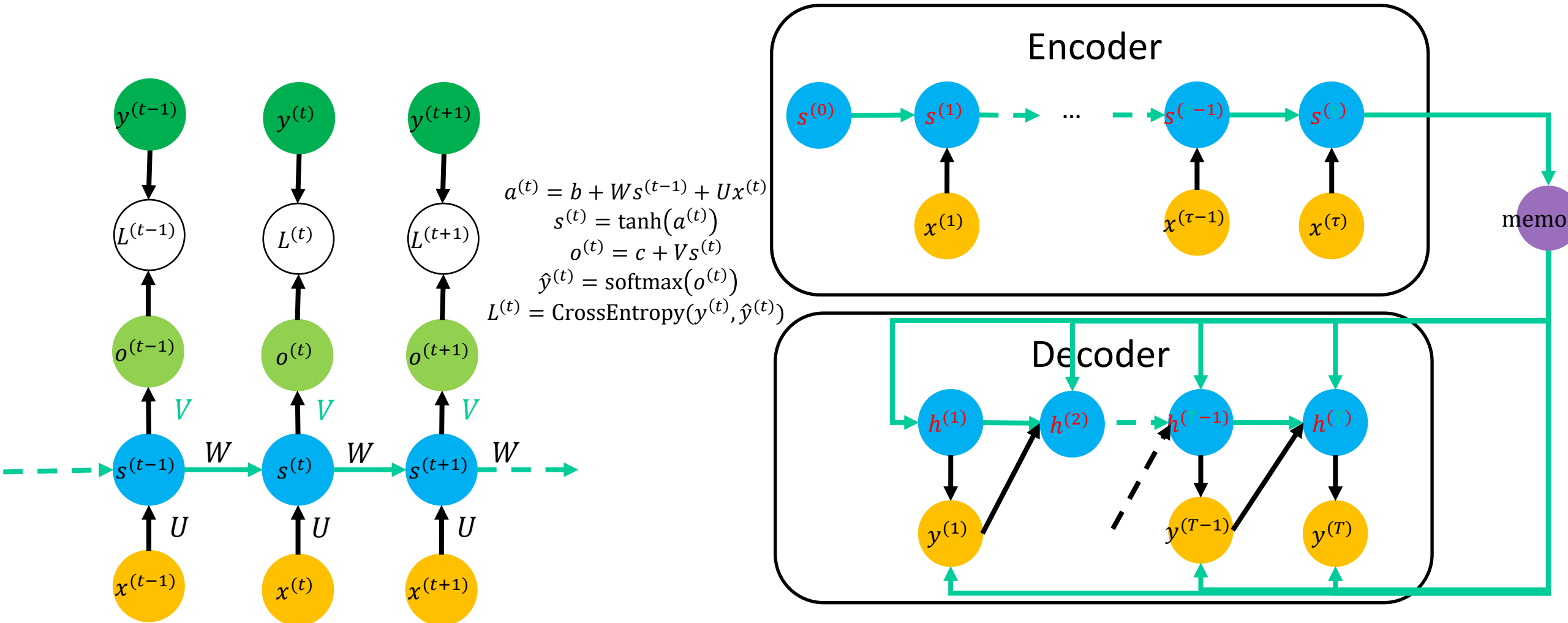
$$P(w_1, w_2, \dots, w_n) = P(w_1) P(w_2 | w_1) P(w_3 | w_2) \dots P(w_n | w_{n-1})$$

- Example:

texaco, rose, one, in, this, issue, is, pursuing, growth, in, a, boiler,  
house, said, mr., gurria, mexico, 's, motion, control, proposal, without,  
permission, from, five, hundred, fifty, five, yen outside, new, car, parking,  
lot, of, the, agreement, reached this, would, be, a, record, november

# Language Models: RNN Review

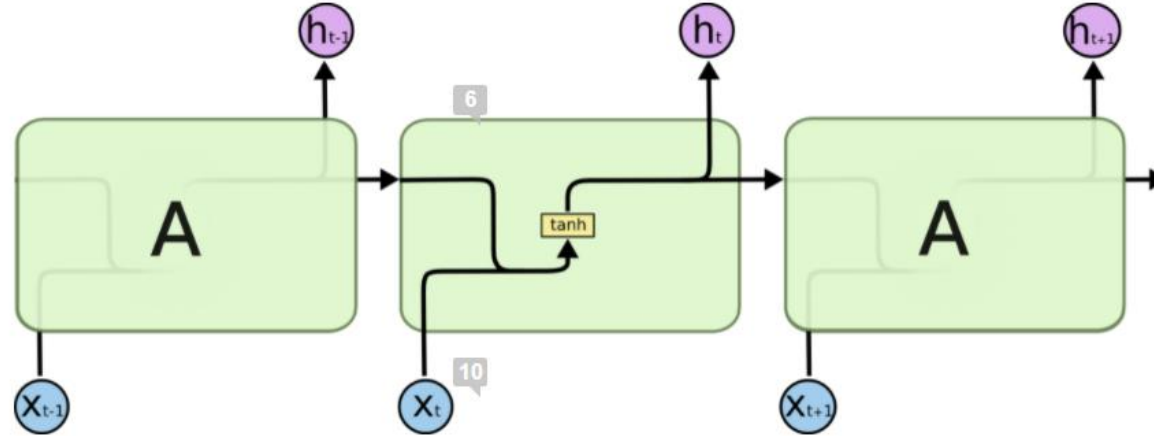
- Classical RNN model / Encoder-Decoder variant:



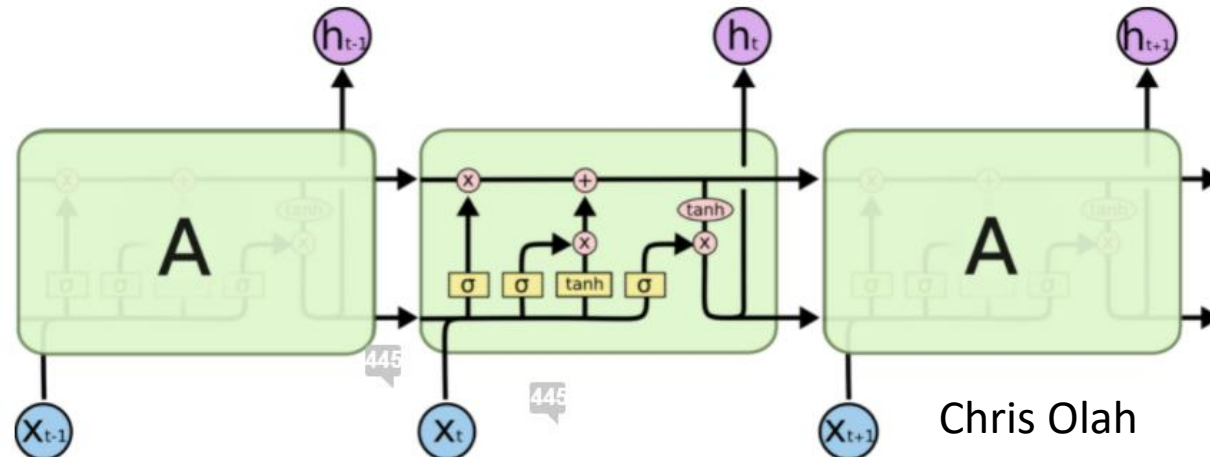


# Language Models: LSTM Review

- RNN: can write structure as:



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



# Language Models: 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 (super popular 2010-present)

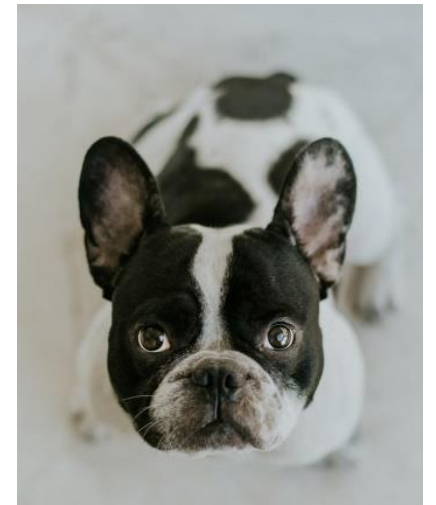
- Word2vec: a famous approach
- What's our likelihood?

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

Windows of length  $2a$

Our word vectors (weights)

All positions



# Training Word Embeddings

Word2vec likelihood

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

- Maximize this; what's the probability?
- Two vectors per word.  $v_w$ ,  $u_w$  for center/context (o is context word, c is center)

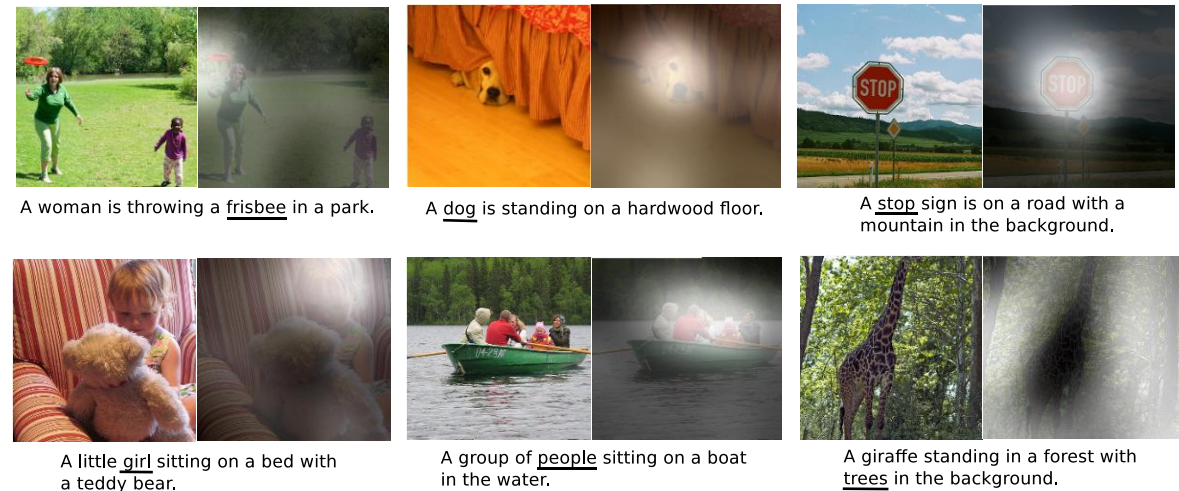
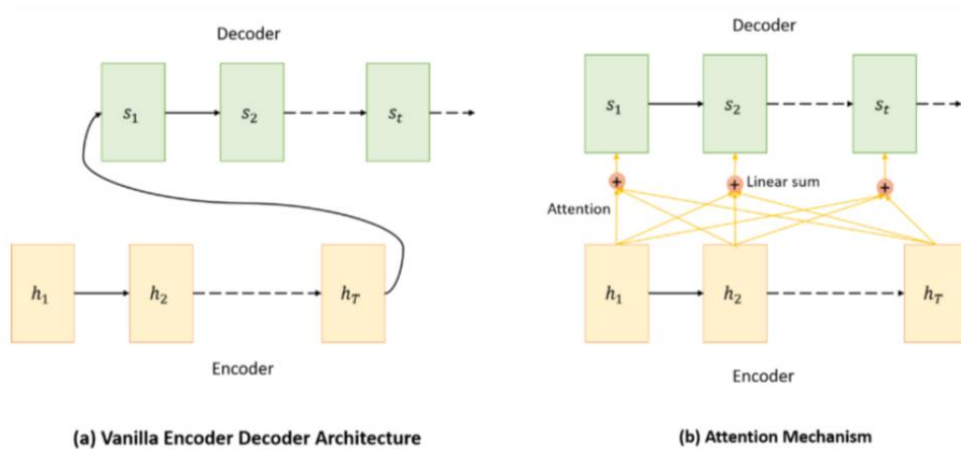
Softmax

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$



# Language Models: Attention

- One challenge: dealing with the hidden state
  - Everything gets compressed there
  - Might lose information
- Solution: **attention** mechanism
  - Tells us where to focus. Note: usable in other models too!

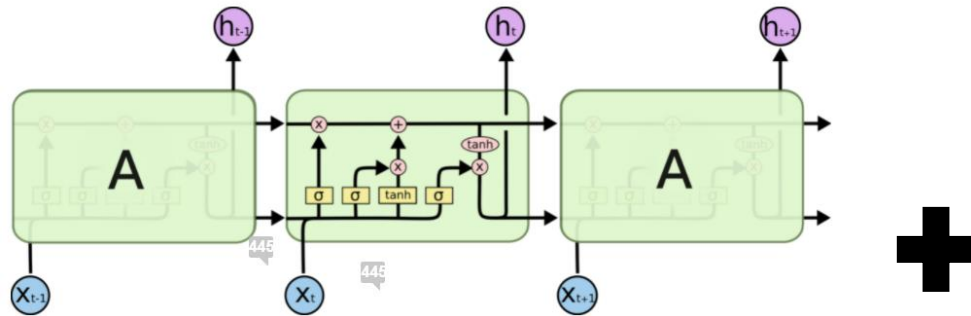


Xu et al, Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

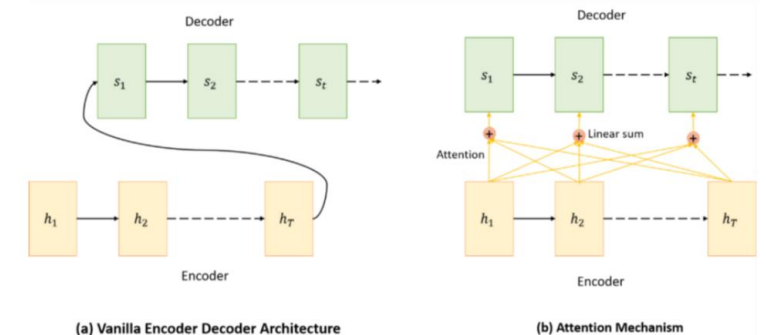


# Language Models: Putting it All Together

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



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





**Break & Quiz**

# Outline

- Language Models & NLP

- k-gram models, RNN review, word embeddings, attention

- **Transformer Model**

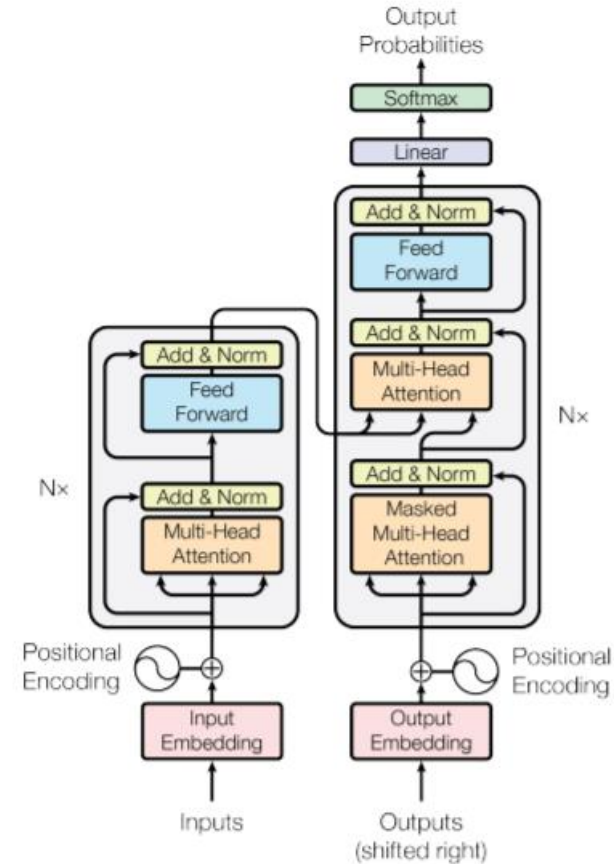
- Properties, architecture breakdown

- Transformer-based Models

- BERT, GPTs, Foundation Models

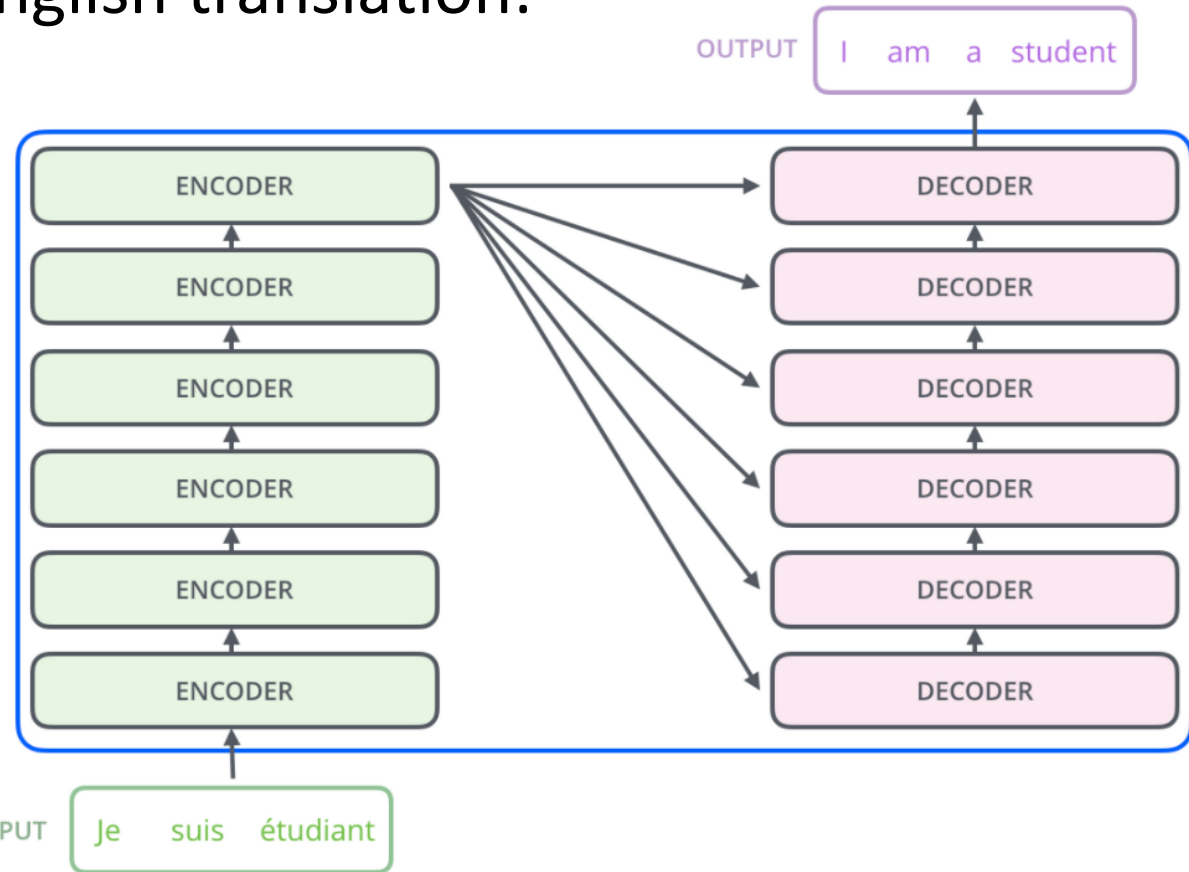
# Transformers: Idea

- Initial goal for an architecture: encoder-decoder
  - Get rid of recurrence
  - Replace with **self-attention**
- Architecture
  - We'll cover this step-by-step
- Results:
  - Best results on translation tasks.



# Transformers: Architecture

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



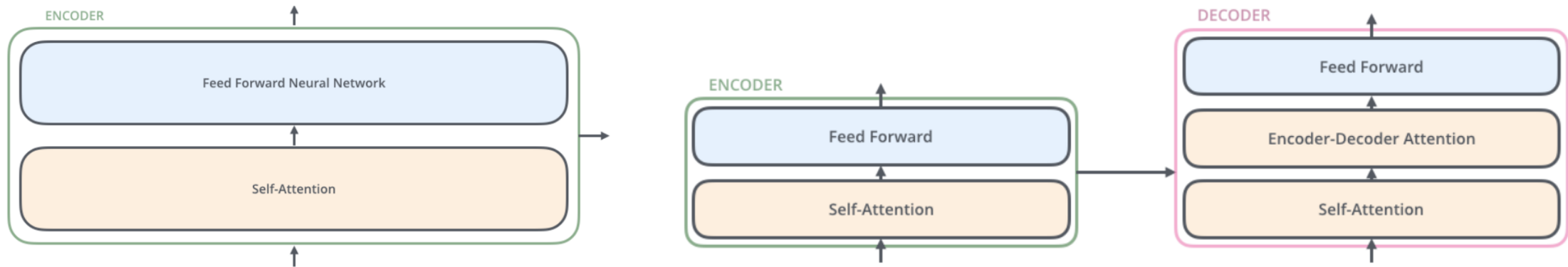
**Note:** All visualizations are due to Jay Alammar

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



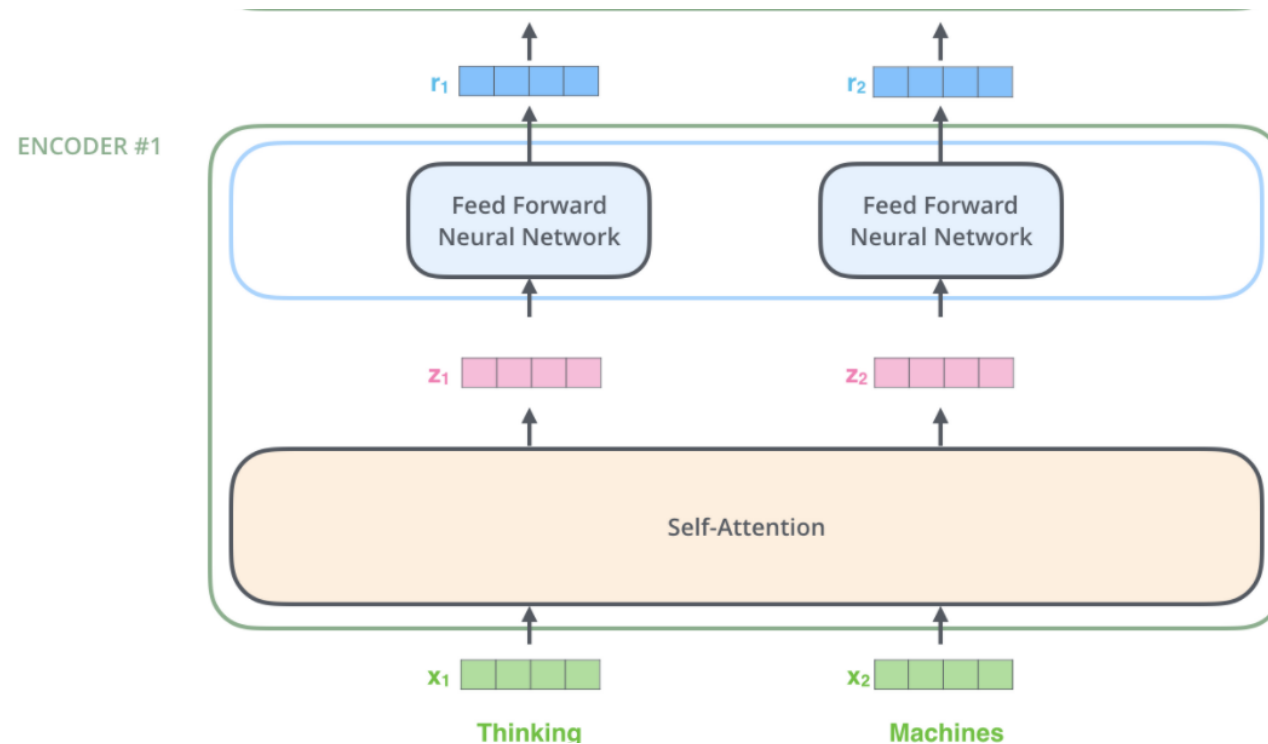
# Transformers: Architecture

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



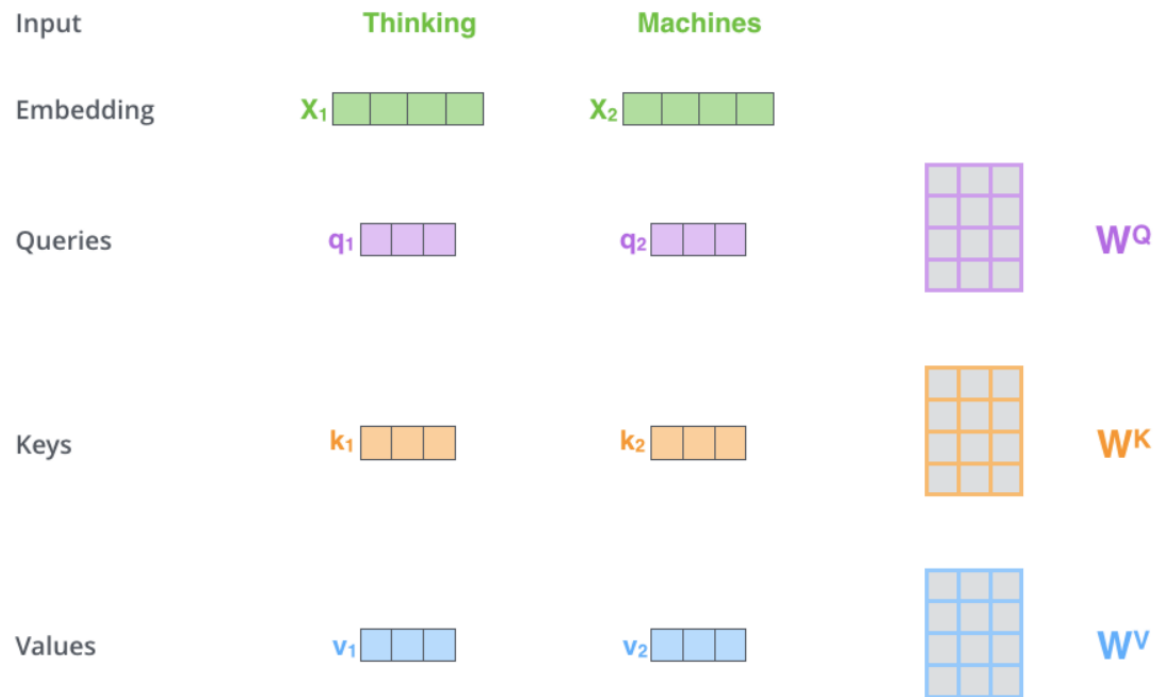
# Transformers: Insider an Encoder

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



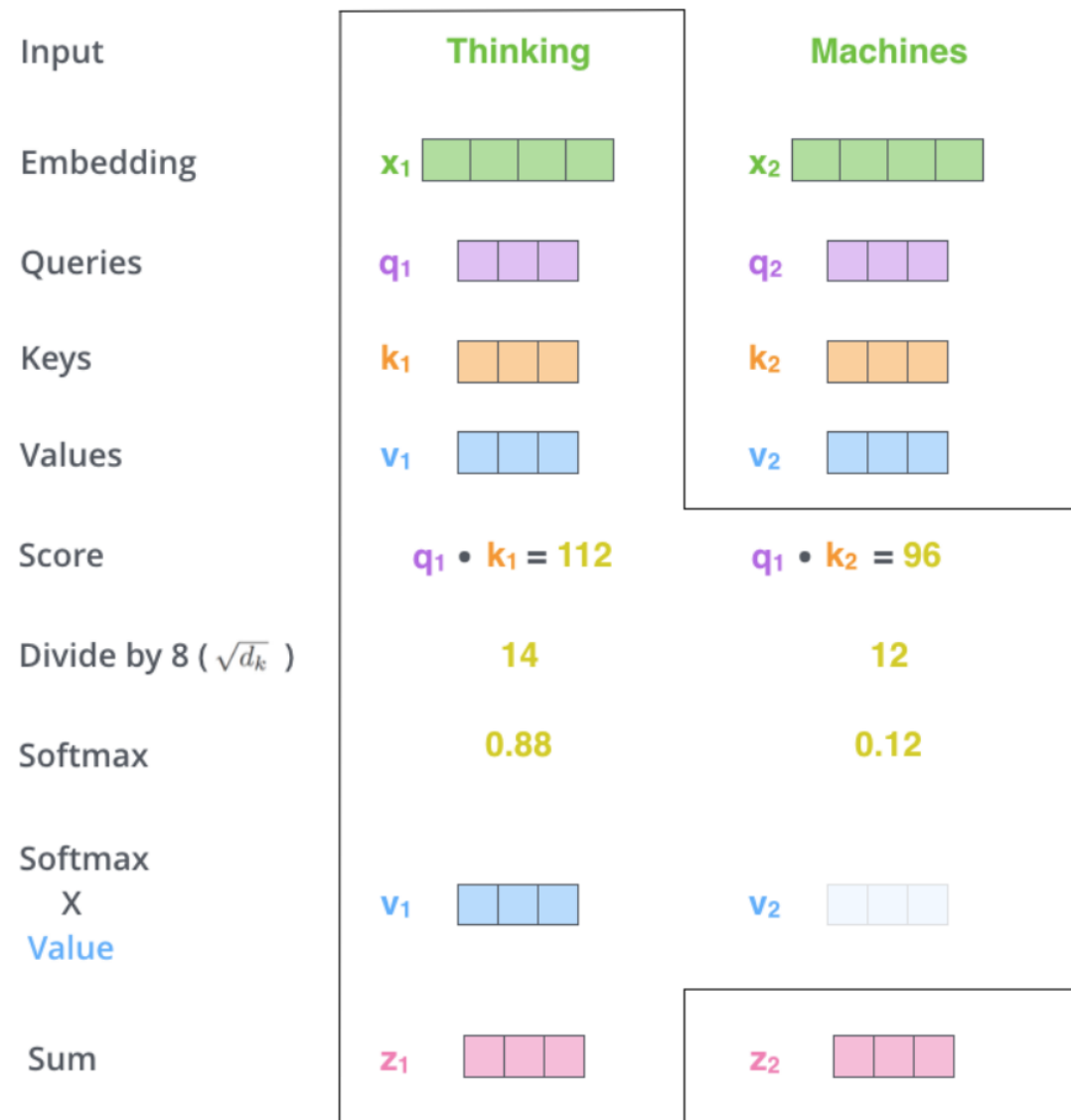
# Transformers: Self-Attention

- Self-attention is the key layer in a transformer stack
  - Get 3 vectors for each embedding: **Query**, **Key**, **Value**
  - We'll combine query at  $i$  with each key  $j$  ( $j=1,\dots$ ), and run softmax
  - Compute values vectors by adding weights (from softmax)



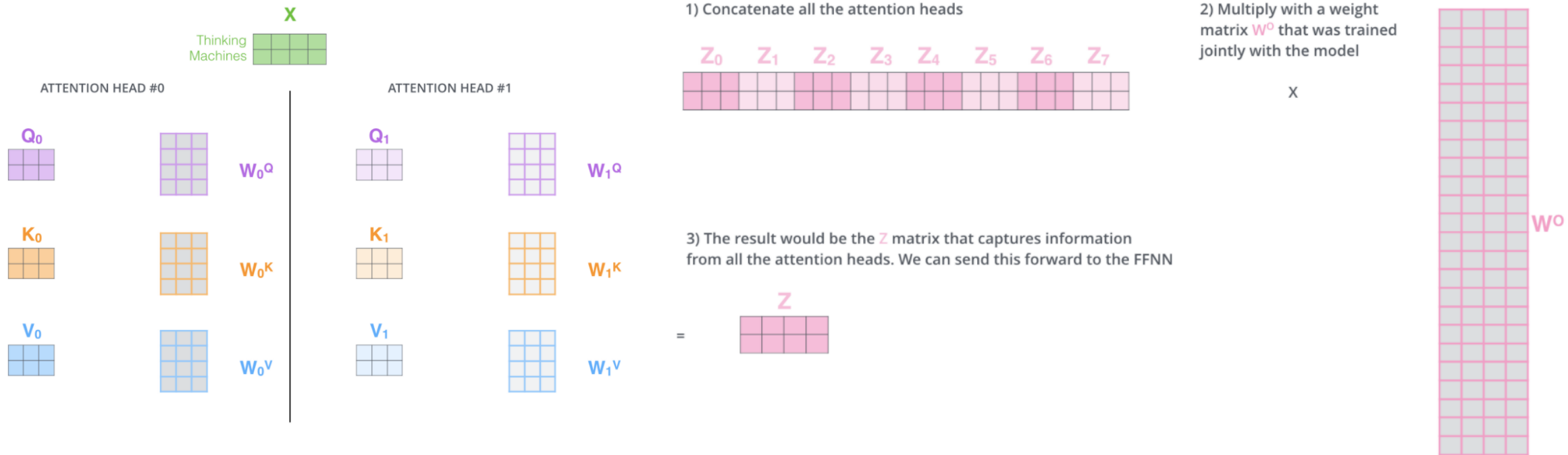
# Transformers: Self-Attention

- 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 for this layer
  - Note: matrix versions in practice
- Send these to feedforward NNs



# Transformers: Multi-Headed Attention

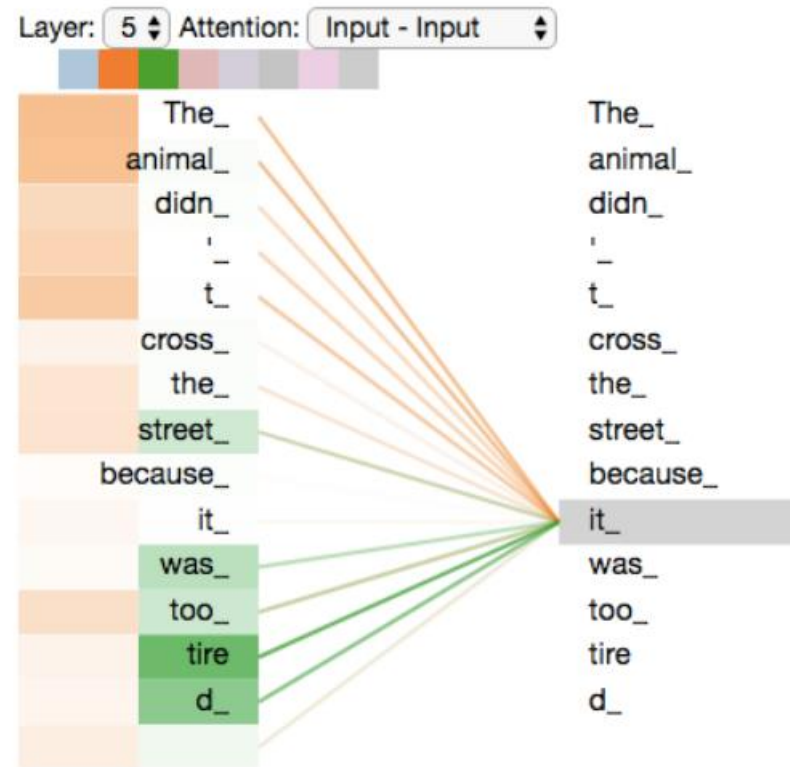
- We can do this multiple times in parallel
  - Called multiple heads
  - Need to combine the resulting output sums





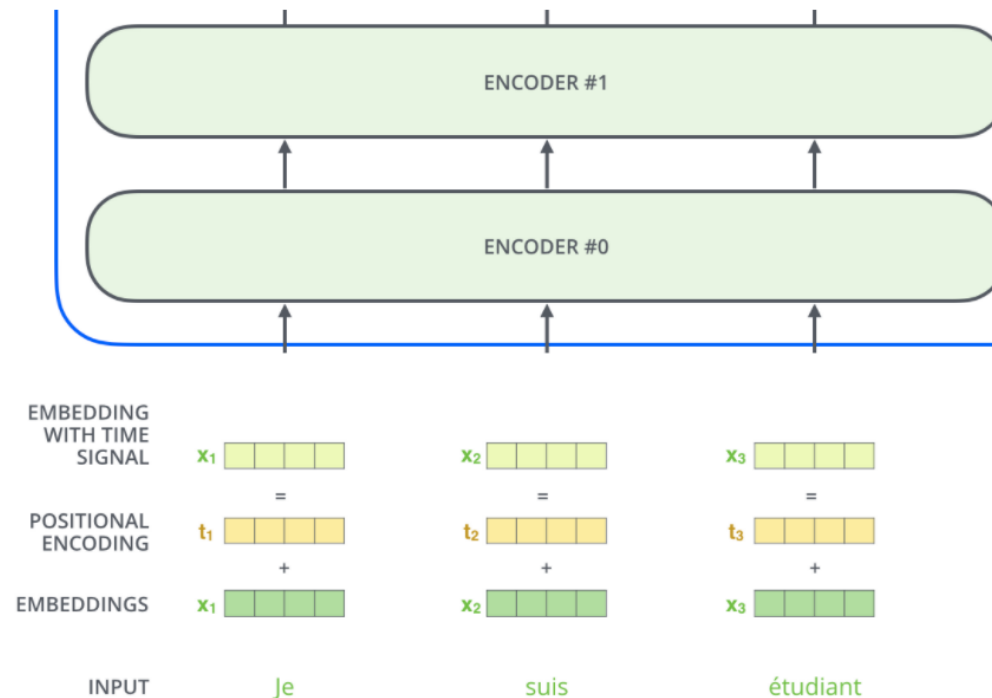
# Transformers: Attention Visualization

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



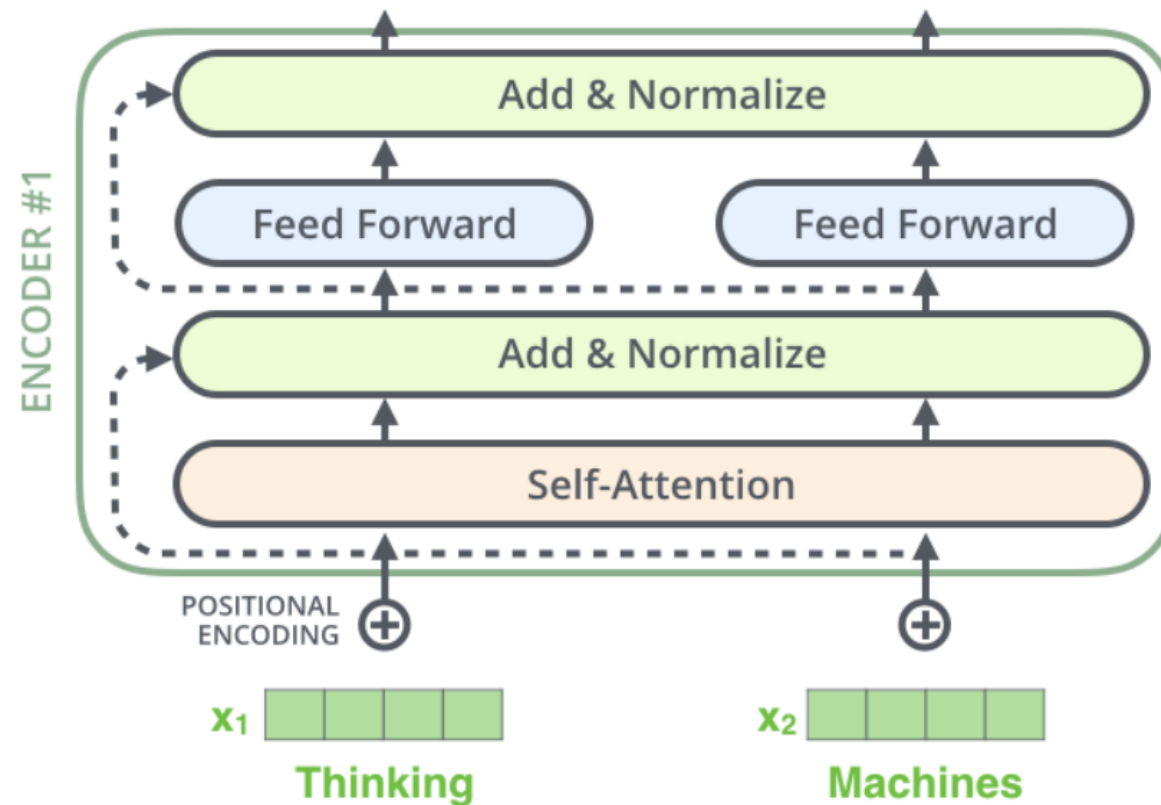
# Transformers: Positional Encodings

- One thing we haven't discussed: the order of the symbols/elements in the sequence
  - Clearly important. **Q:** How to include?
  - **A:** add a vector containing a special positional formula's embedding



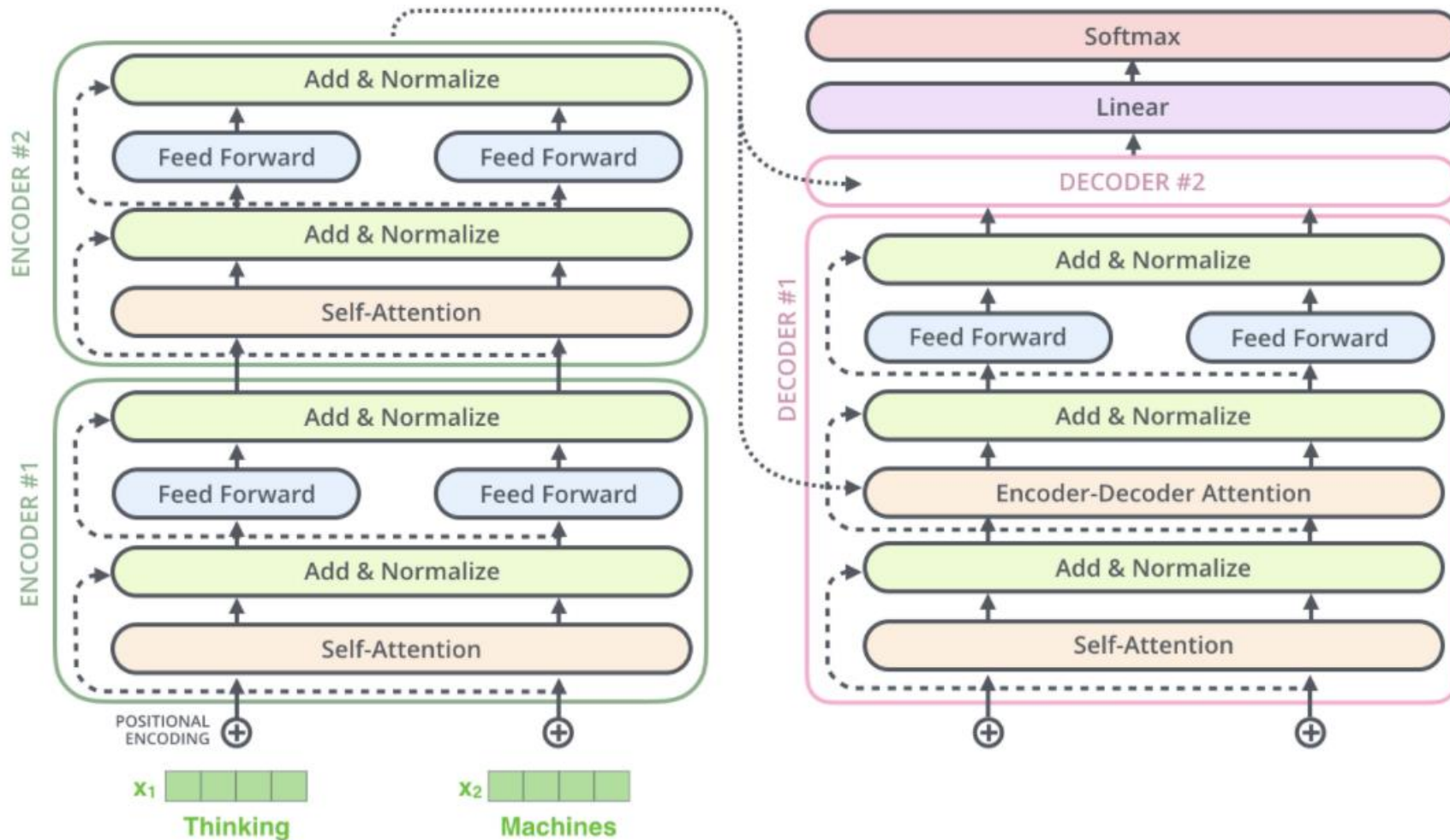
# Transformers: More Tricks

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



# Transformers: Putting it All Together

- What does the full architecture look like?





# Break & Quiz



# Outline

- Language Models & NLP

- k-gram models, RNN review, word embeddings, attention

- Transformer Model

- Properties, architecture breakdown

- **Transformer-based Models**

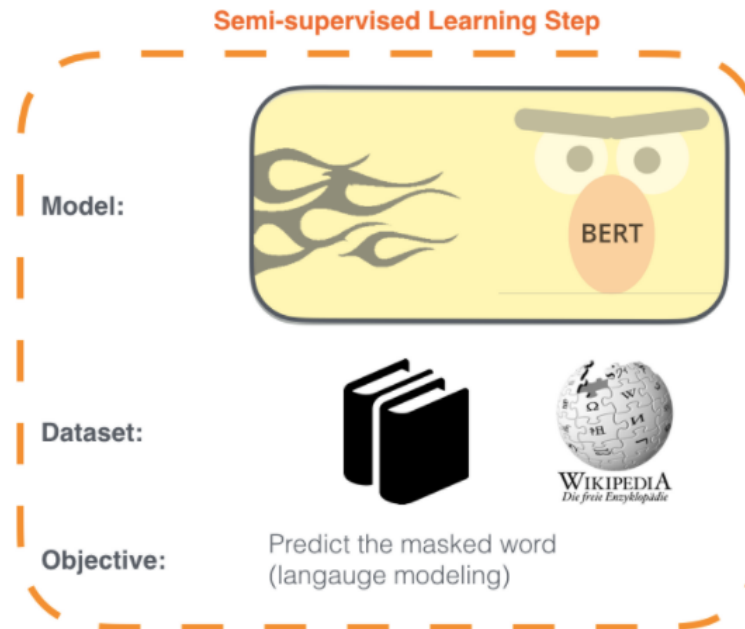
- BERT, GPTs, Foundation Models

# Transformer-Based Models: BERT

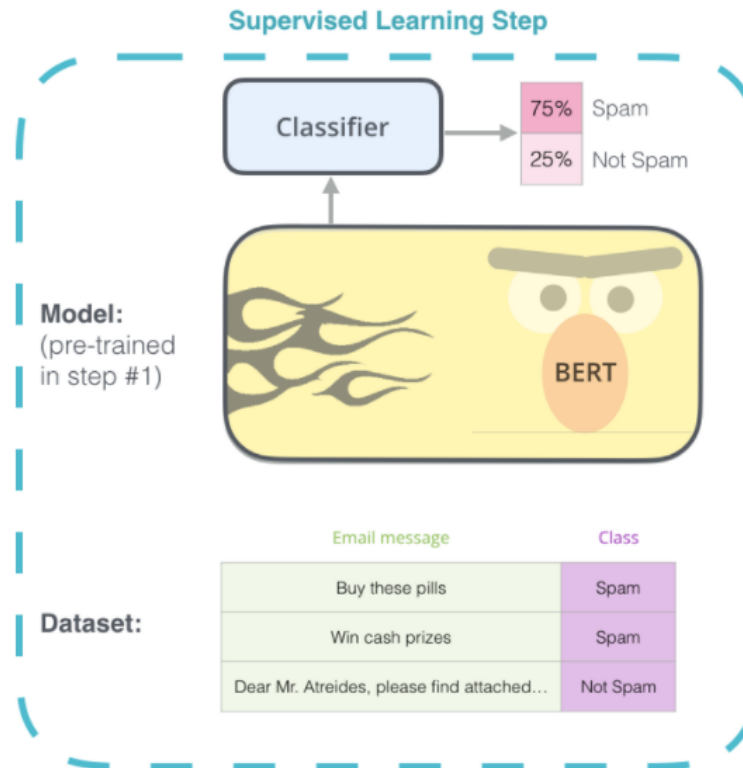
- The transformer model works very well...
  - Why not apply it to **everything** in NLP?
  - BERT does this. Two-step process to use for a particular task

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.



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



# BERT: Concepts

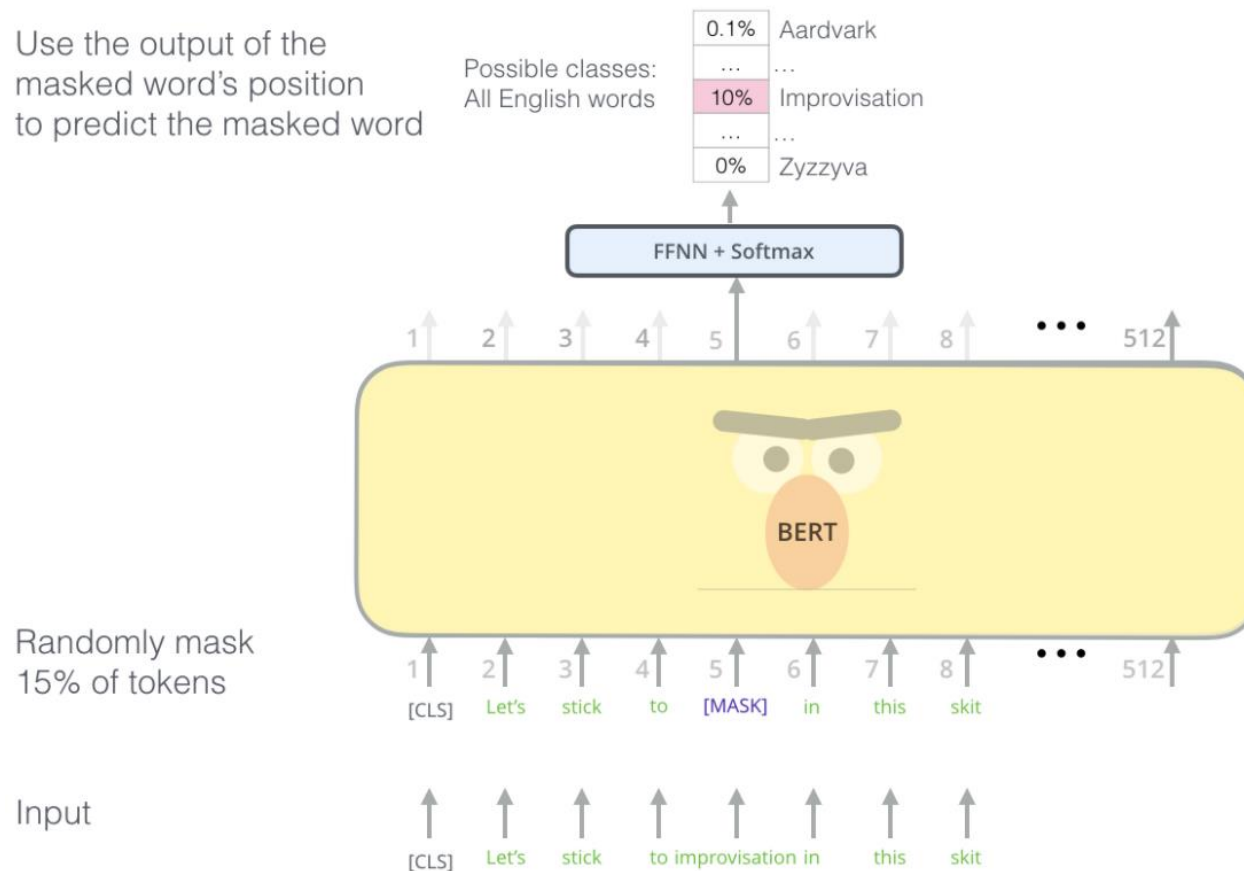
- What makes BERT work? A bunch of ideas:
  - 1. Use the **Transformer** architecture
    - Encoder stacks in particular (not doing sequence-sequence)
  - 2. **Pre-training** on corpora
    - Then fine-tune for a particular task
  - 3. **Scale**: BERT-Large has 340 million parameters

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

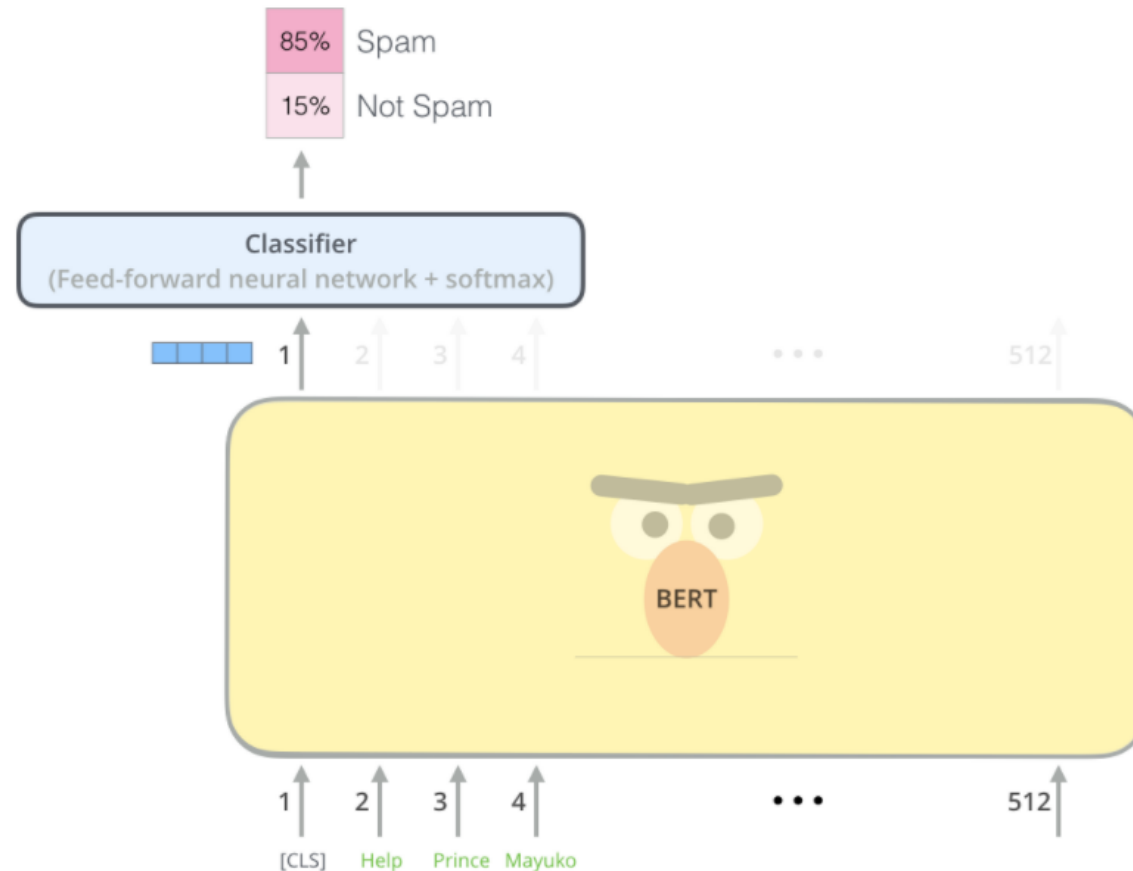
# BERT: Training

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



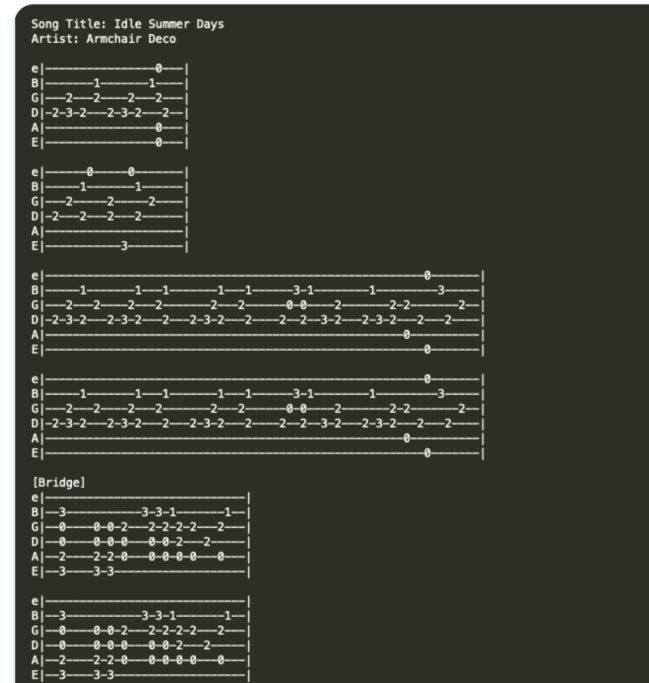
# BERT: Classification

- Then, fine-tune on a particular task
  - Example: **binary classification**, spam VS not spam
  - Also use for word embeddings. Note: contextual, unlike word2Vec



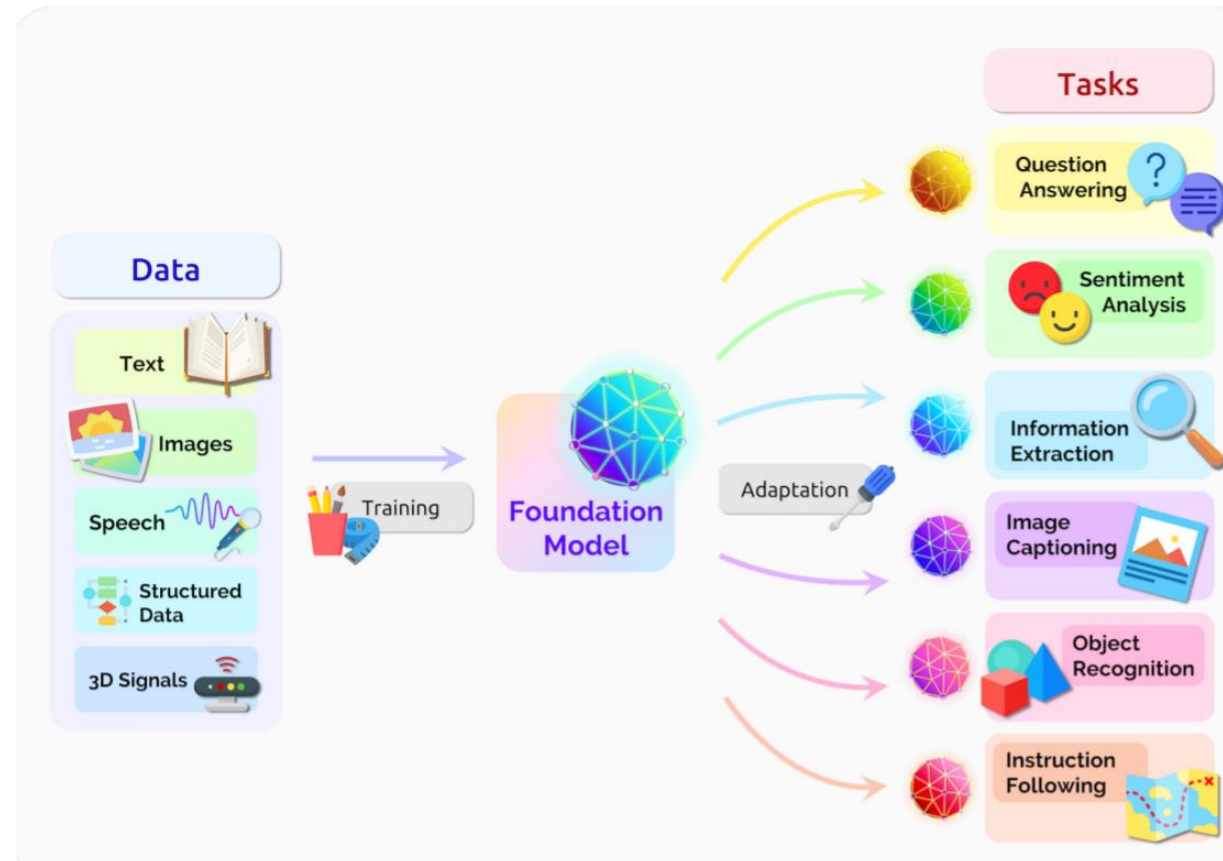
# GPT Series of Models

- **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**)
  - Not fine-tuning for a particular task.
- Scale: huge!
  - GPT-3: 175 billion parameters



# Foundation Models

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

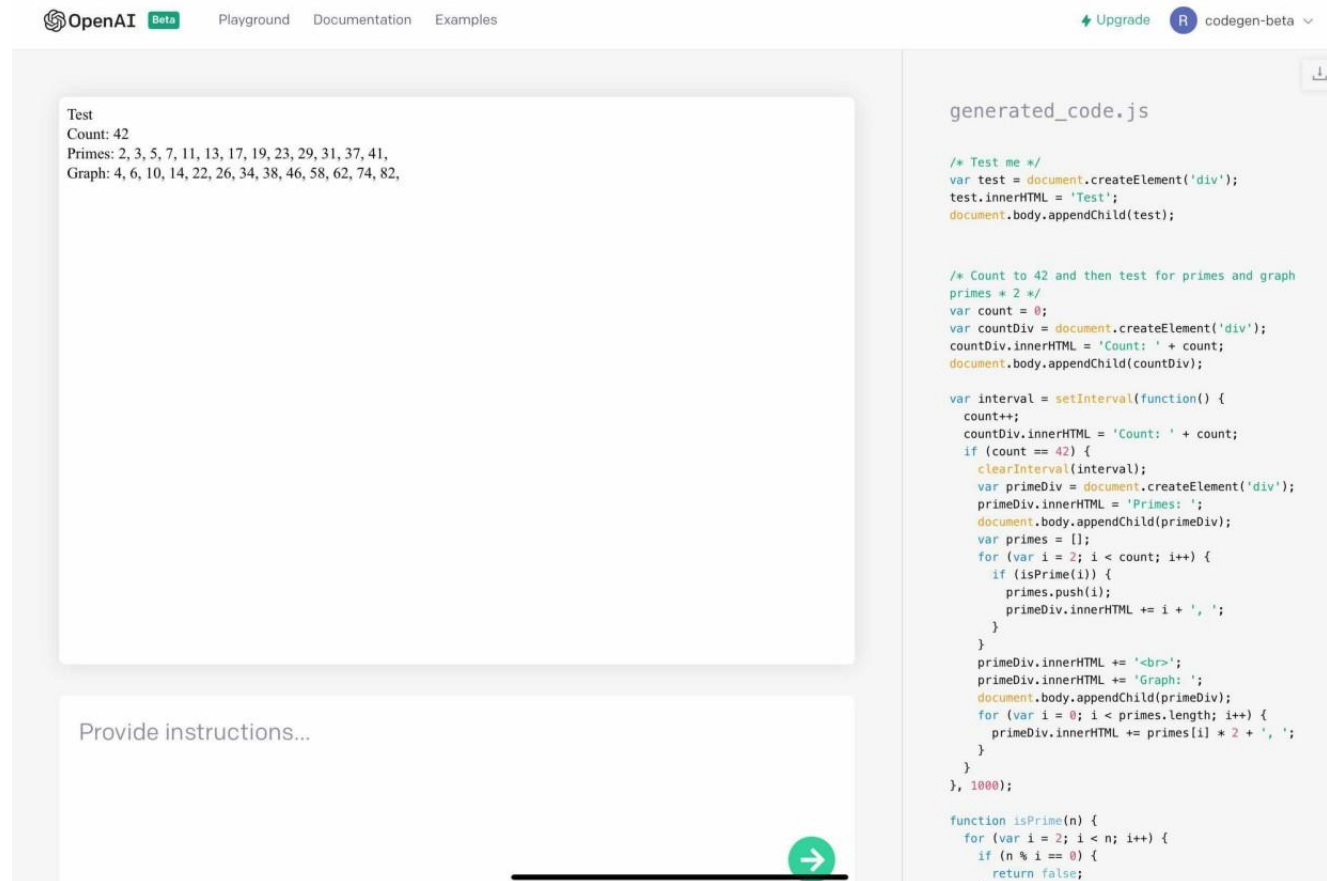


Bommasani et al, "On the Opportunities and Risks of Foundation Models"



# Codex

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



The screenshot shows the OpenAI Playground interface. At the top, there is a navigation bar with "OpenAI Beta", "Playground", "Documentation", and "Examples". On the right, there are "Upgrade" and "codegen-beta" options. The main area is split into two panels. The left panel contains a text input field with the prompt: "Test", "Count: 42", "Primes: 2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, 41,", and "Graph: 4, 6, 10, 14, 22, 26, 34, 38, 46, 58, 62, 74, 82,". Below the input field is a "Provide instructions..." field and a green arrow button. The right panel shows the generated JavaScript code in a file named "generated\_code.js". The code includes comments for each section and implements a prime number generator and a graph visualization.

```
generated_code.js

/* 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. . . .”

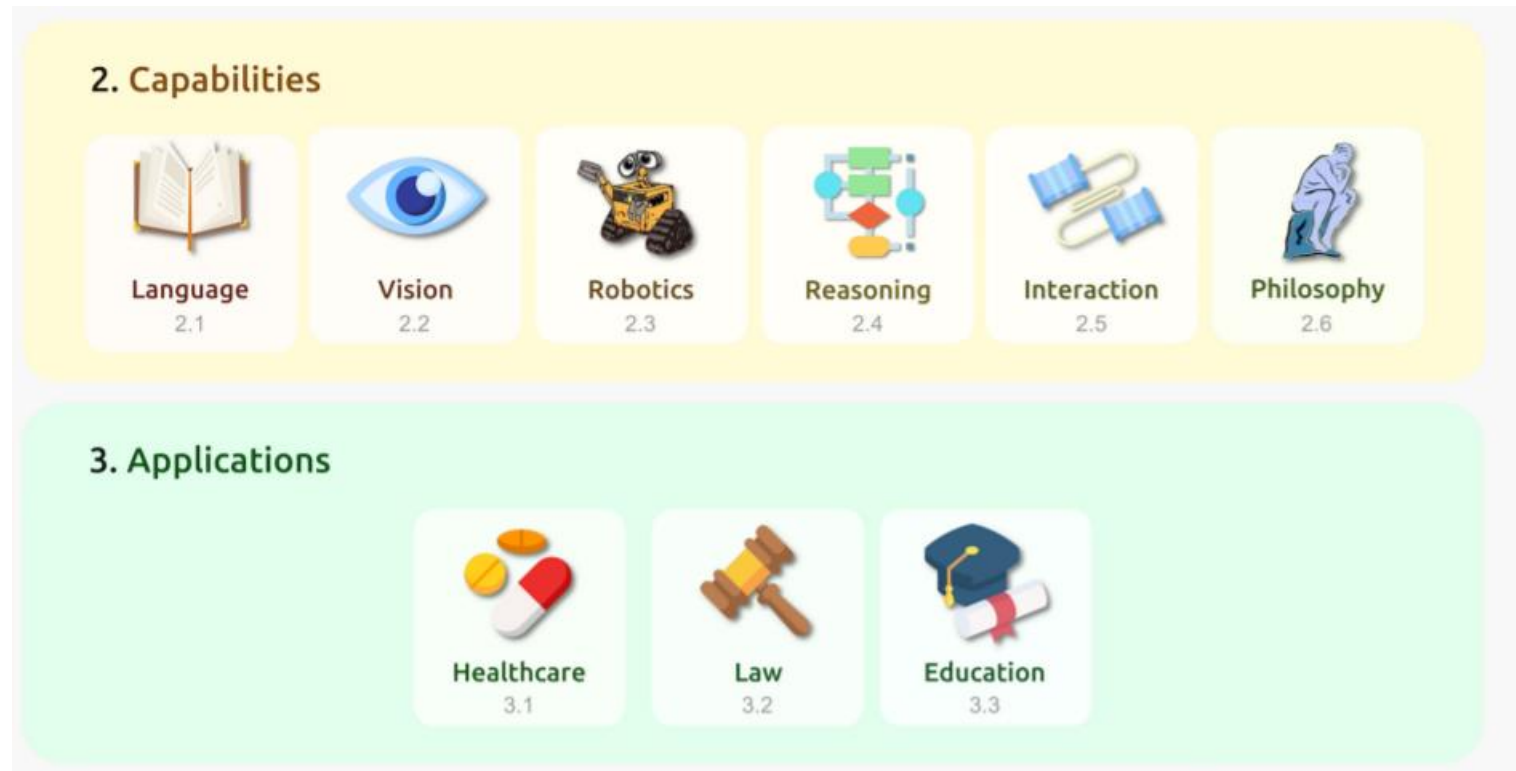


<https://openai.com/blog/dall-e/>

- Note: several online demos. Try it yourself!

# 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.
  - The future of AI??





# 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 Alammarr