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

Kirthi Kandasamy

University of Wisconsin-Madison

**May 1, 2023**

# Announcements

- **Course evaluations**
  - **33** have already filled out. Thank you!
  - Please fill out if you haven't already

- **Homework 7** was due this morning

- **Finals**
  - 5/12/2023, Friday 7:45am - 9:45am
  - MICROBIAL SCIENCES BLDG 1520

# 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:

$$\mathrm{dog} = \begin{bmatrix} 0.13 & 0.87 & -0.23 & 0.46 & 0.87 & -0.31 \end{bmatrix}^T$$

$$\mathrm{cat} = \begin{bmatrix} 0.07 & 1.03 & -0.43 & -0.21 & 1.11 & -0.34 \end{bmatrix}^T$$

AKA **word embeddings**

# Training **Word Embeddings**

Many approaches (very popular 2010-present)
- Word2vec: a famous approach
- Write out a likelihood

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

Windows of length 2a

Our word vectors (weights)

All positions

# Training **Word Embeddings**

Word2vec likelihood

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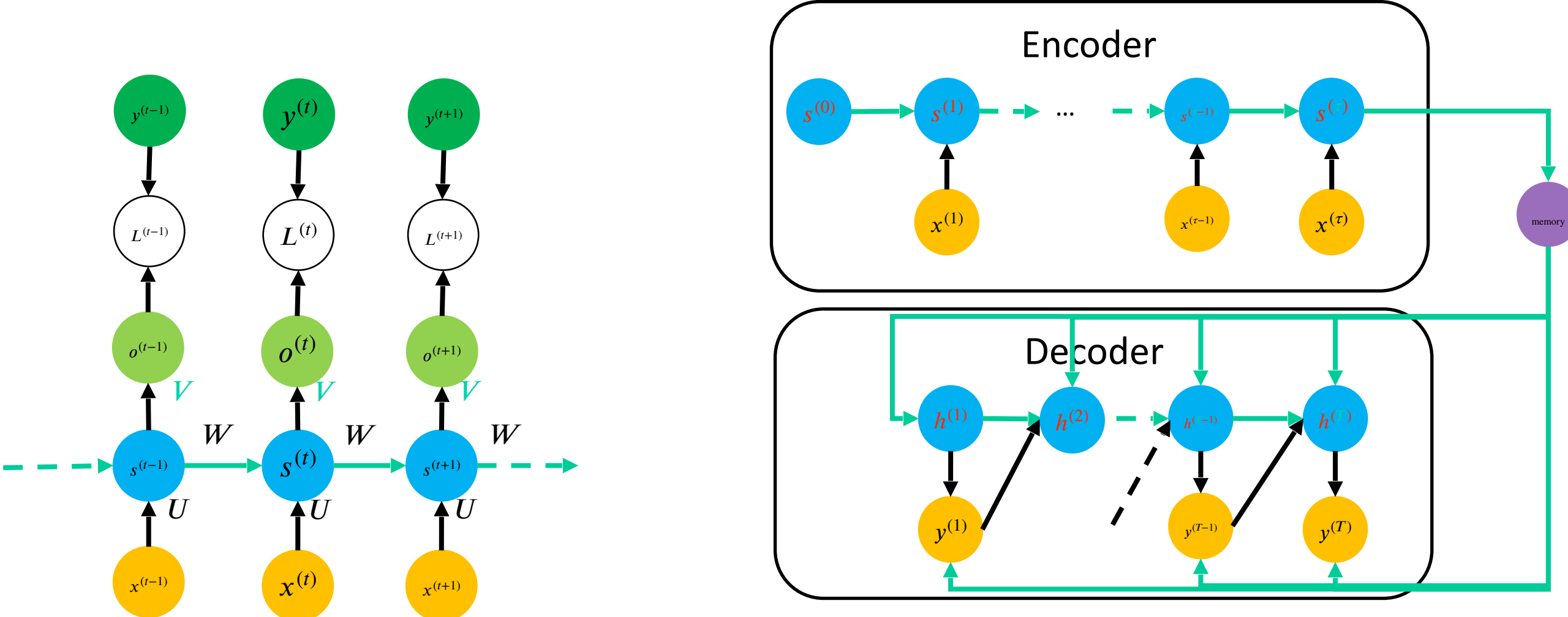
- Expression for the probability:

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

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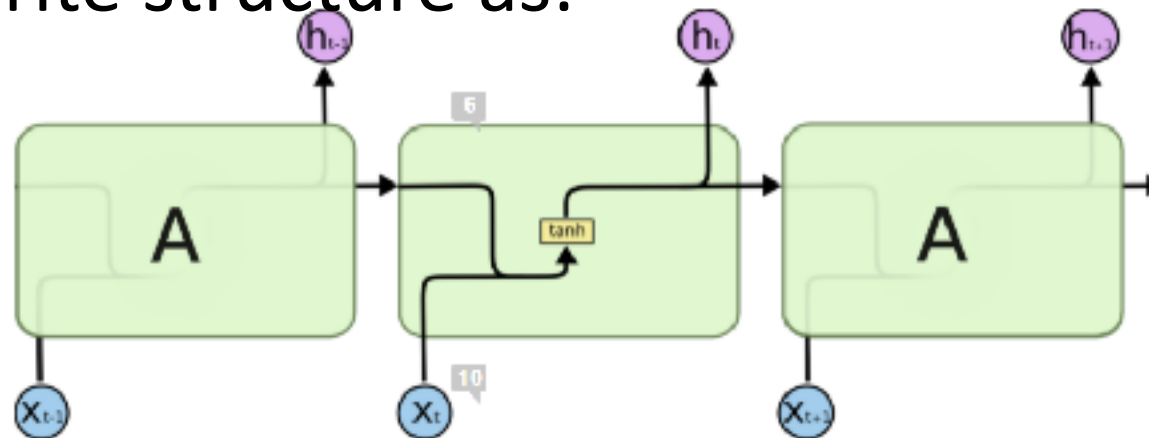- $\theta_{w,c}$: context vector for word $w$

# **Language Models**: RNN Review
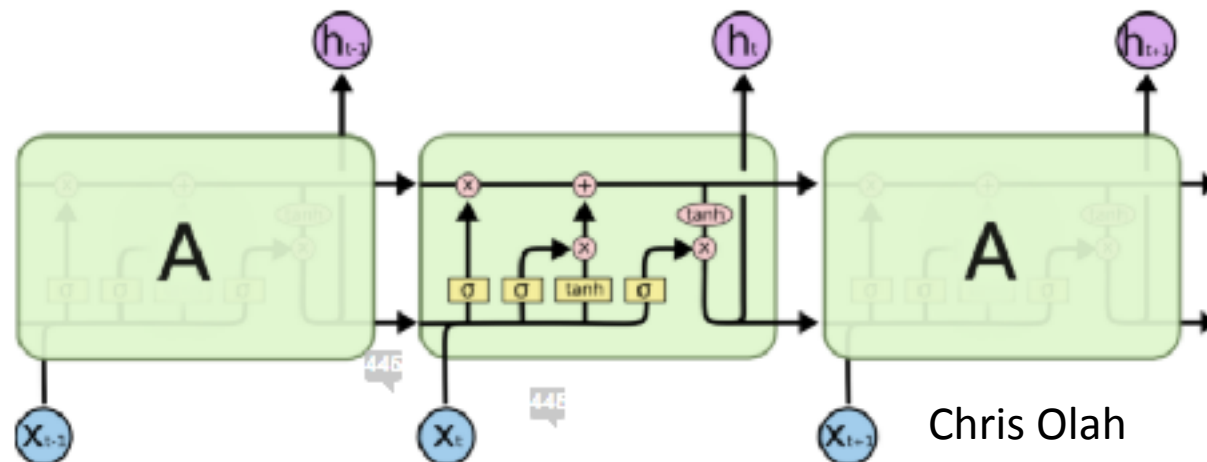
- Classical RNN model / Encoder-Decoder variant:

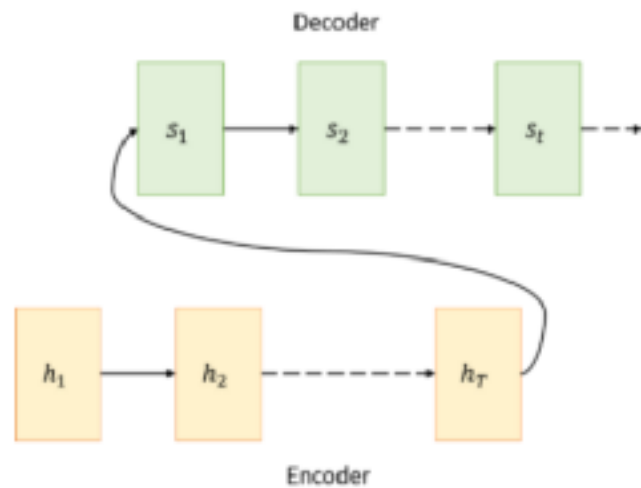# Language Models: LSTM Review

- RNN: can write structure as:



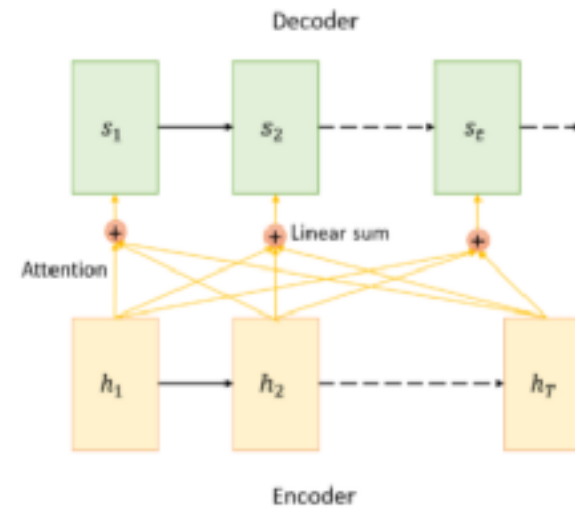- Long Short-Term Memory: deals with problem. Cell:



Chris Olah

# **Language Models**: Attention

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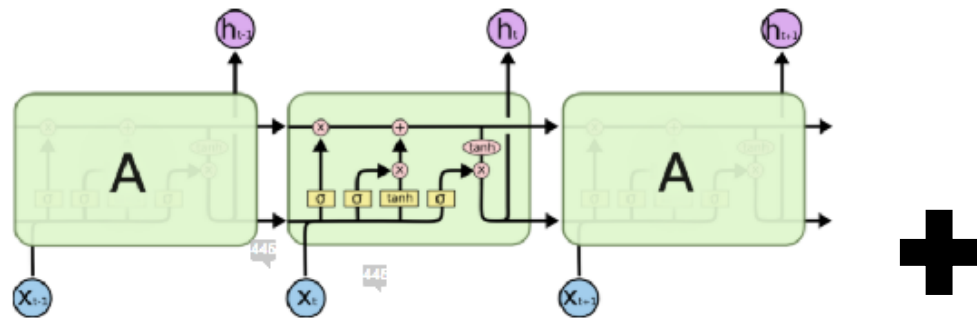  - Similar to residual connections in ResNets!


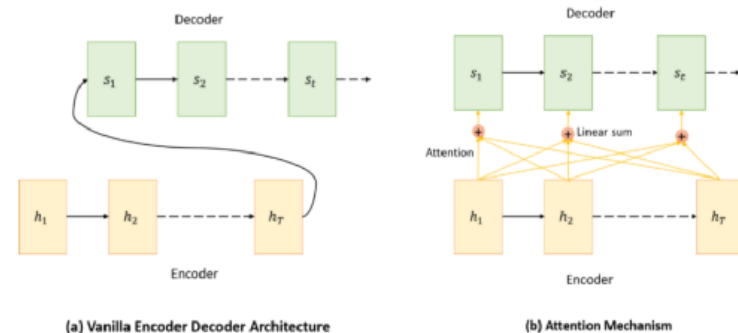
(a) Vanilla Encoder Decoder Architecture

(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



$$\text{dog} = \begin{bmatrix} 0.13 & 0.87 & -0.23 & 0.46 & 0.87 & -0.31 \end{bmatrix}^{T}$$



(a) Vanilla Encoder Decoder Architecture

(b) Attention Mechanism

# Outline

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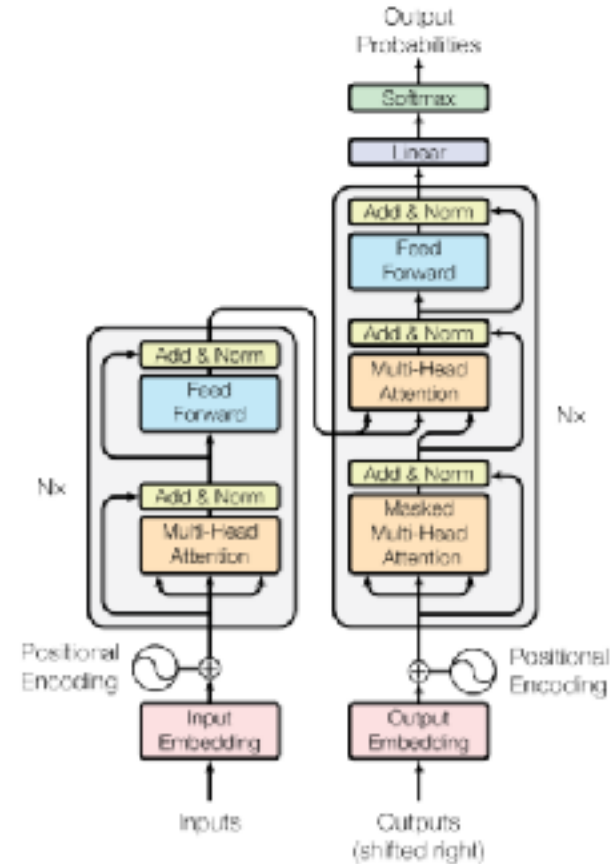- **Transformer Model**
  - Properties, architecture breakdown
- **Transformer-based Models**
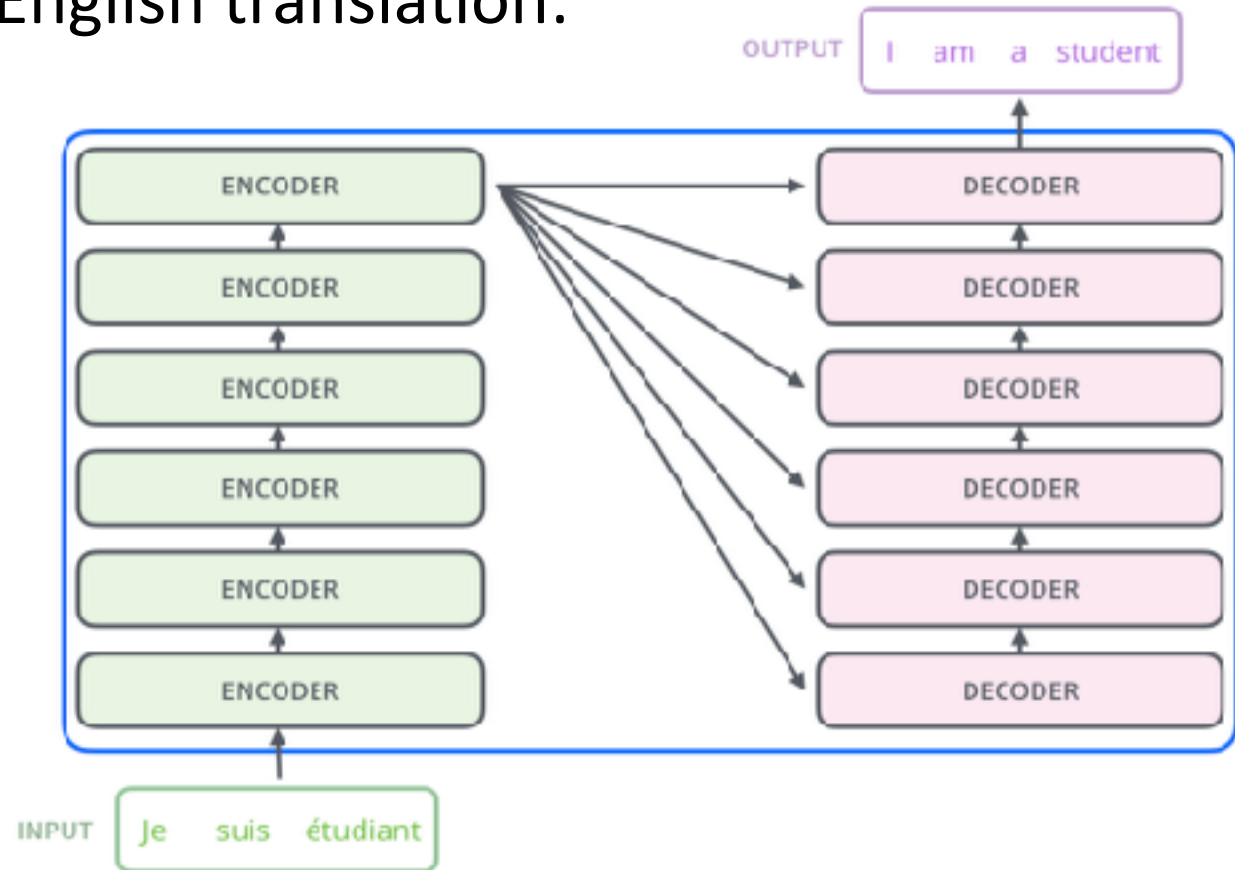  - BERT, GPTs, Foundation Models

# **Transformers**: Idea

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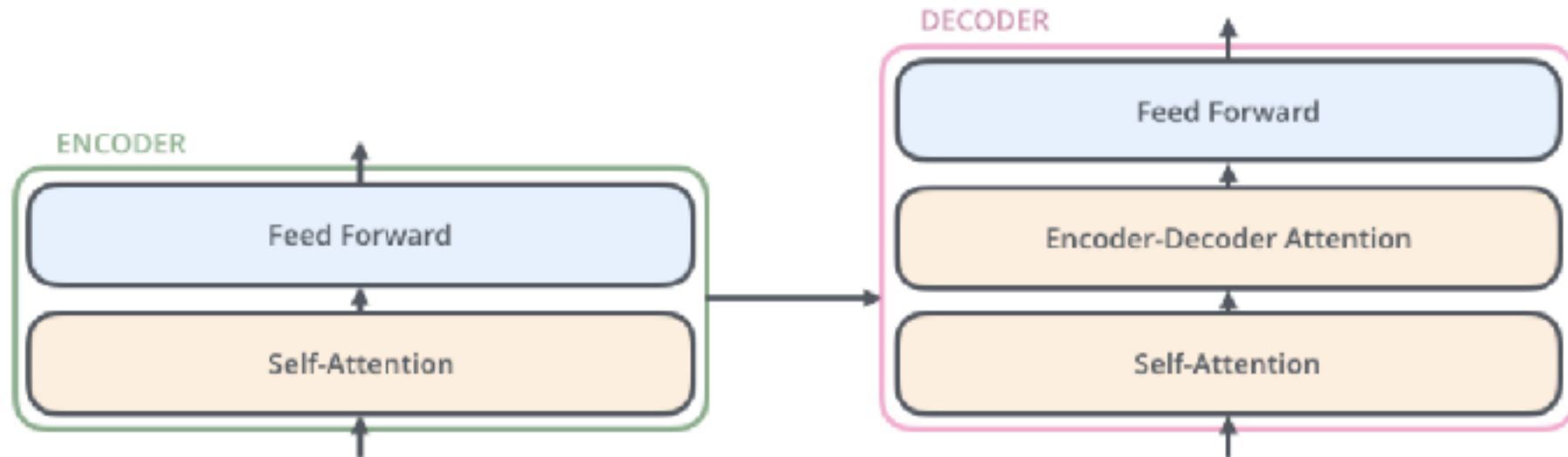  - Replace with **self-attention**



Vaswani et al. '17

# **Transformers**: Architecture

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



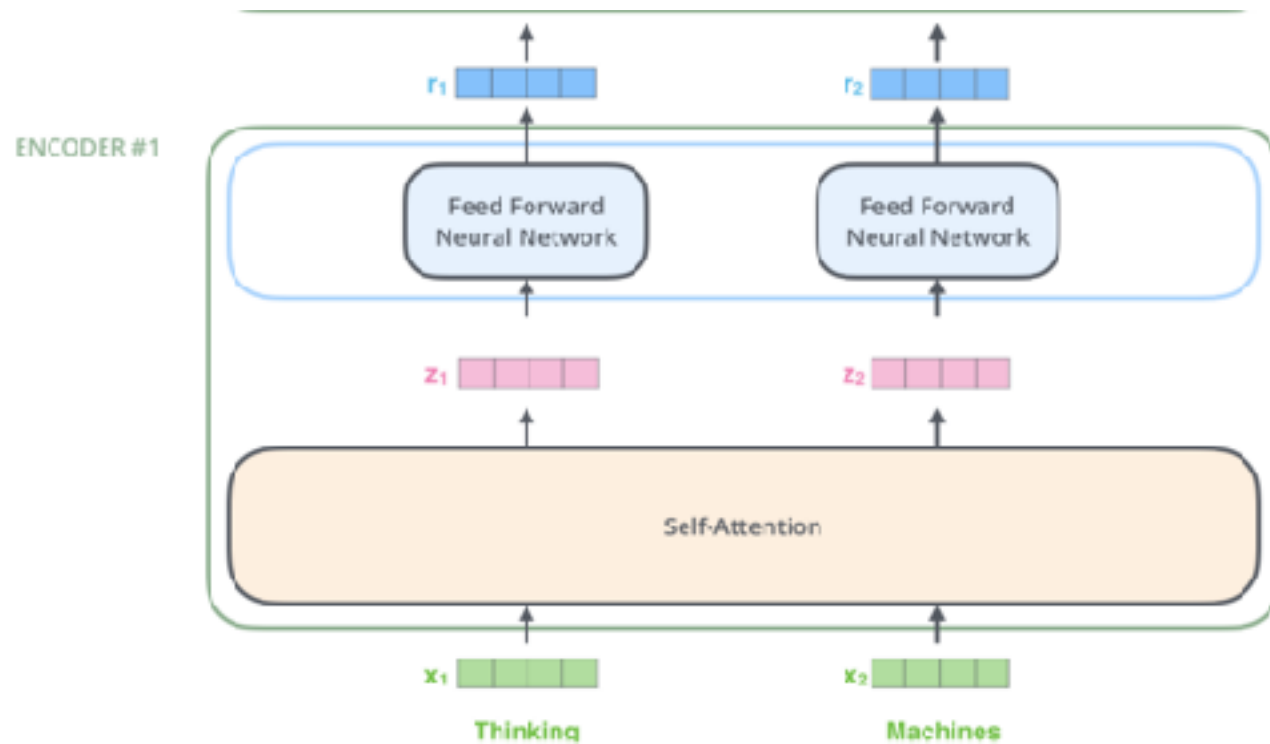Excellent resource: https://jalammar.github.io/illustrated-transformer/

# **Transformers**: Architecture

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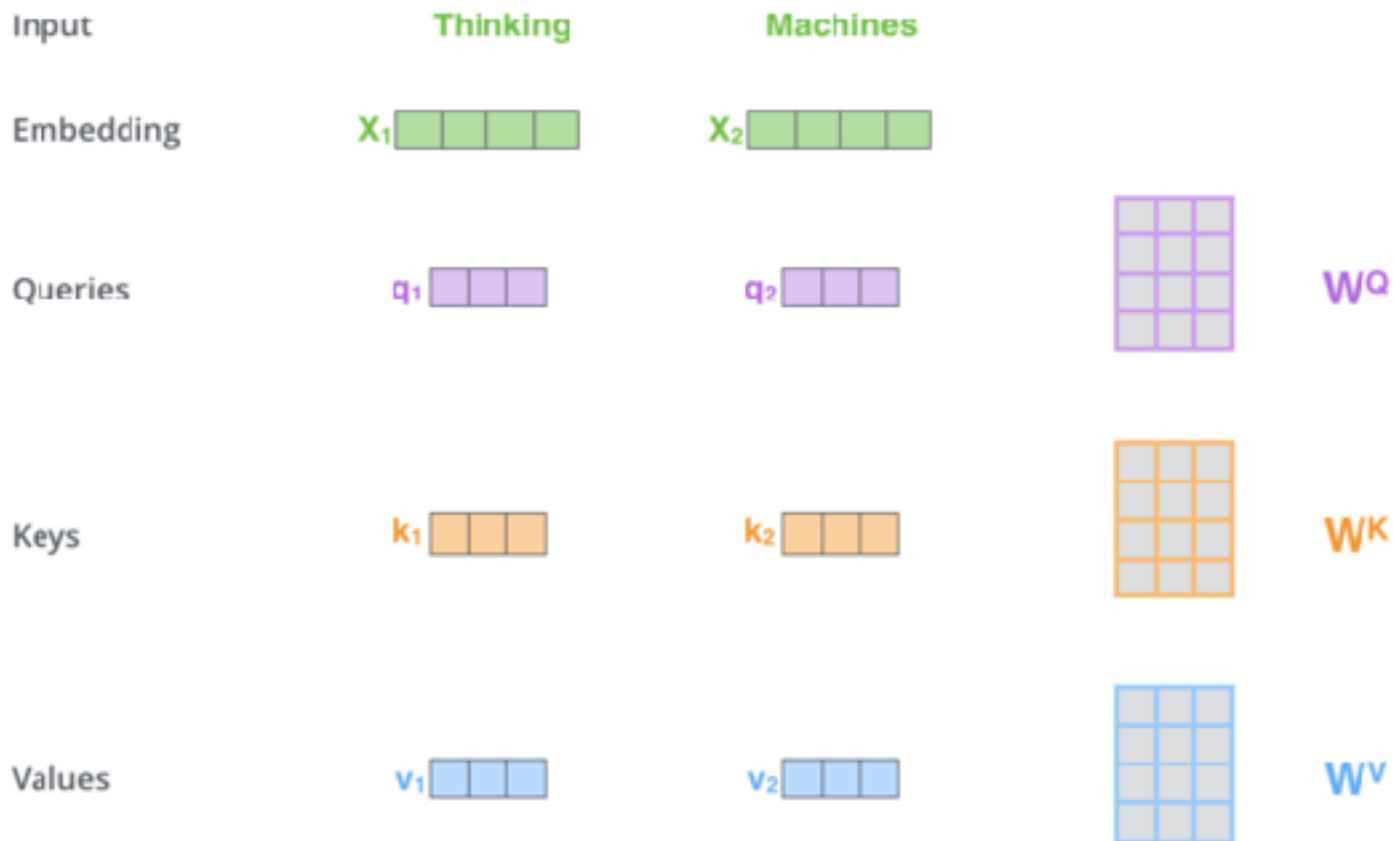  - What's inside each encoder/decoder unit?

# **Transformers**: Inside an Encoder

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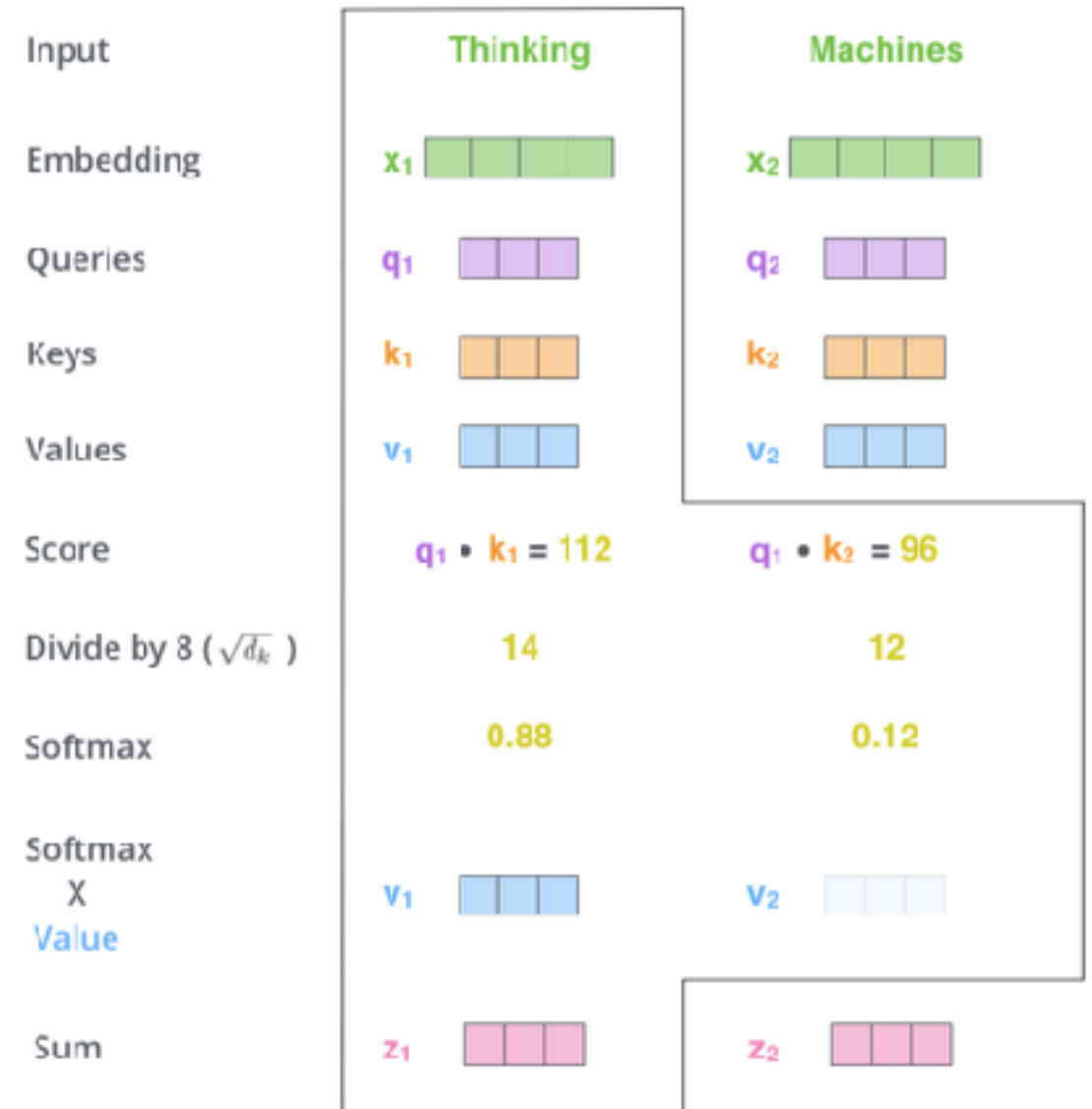  - 2. F**eedforward nets**

# **Transformers**: Self-Attention

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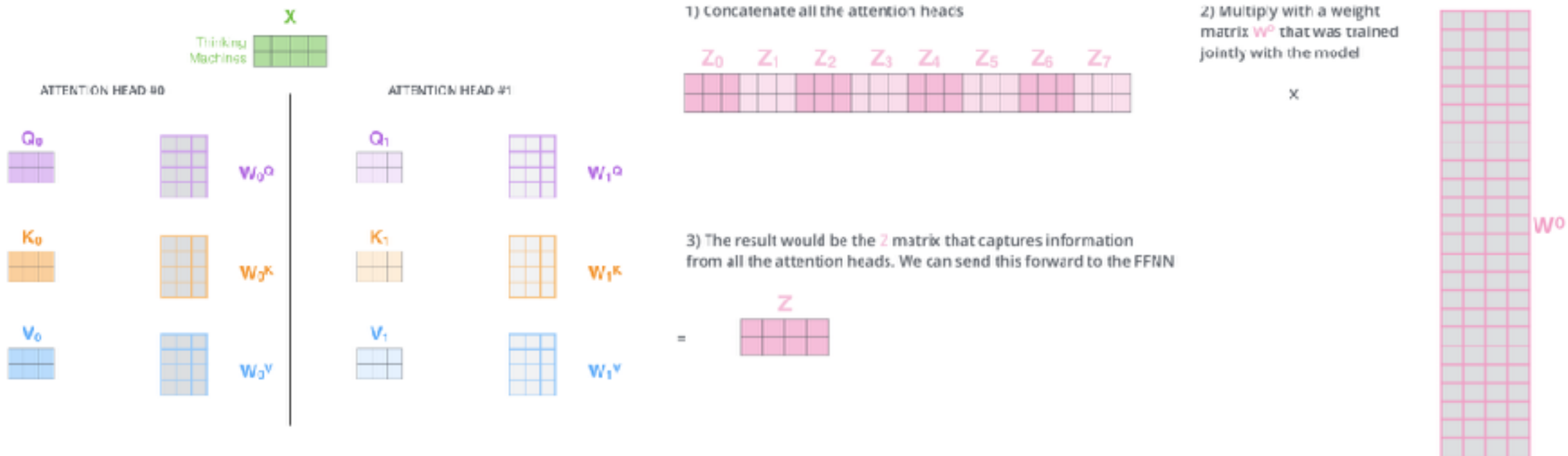  - Get 3 vectors for each embedding: **Query**, **Key**, **Value**

# **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

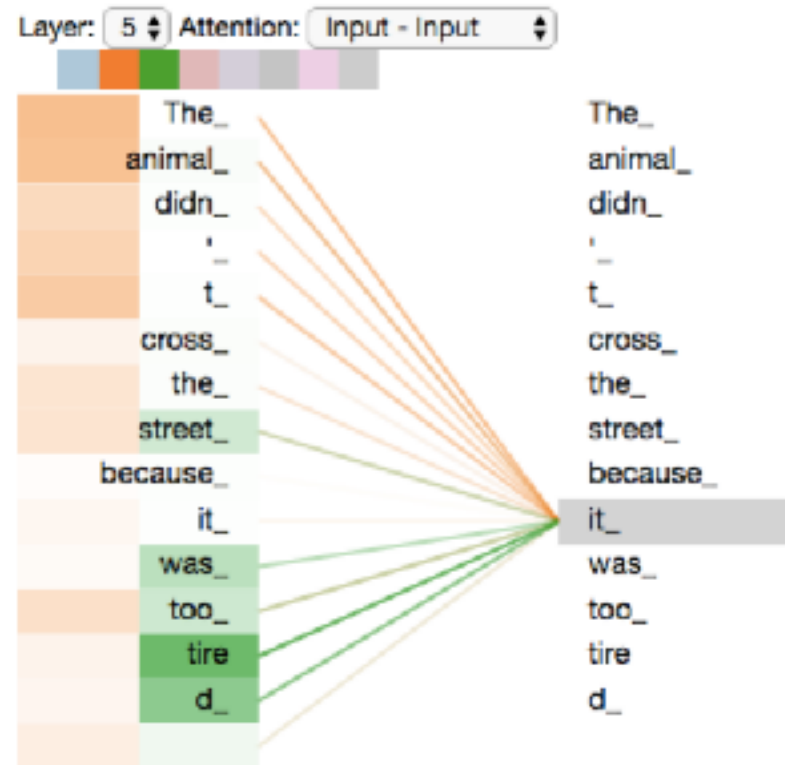  - Send these to feedforward NNs

- Highly parallelizable!



| | 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 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# **Transformers**: Multi-Headed Attention

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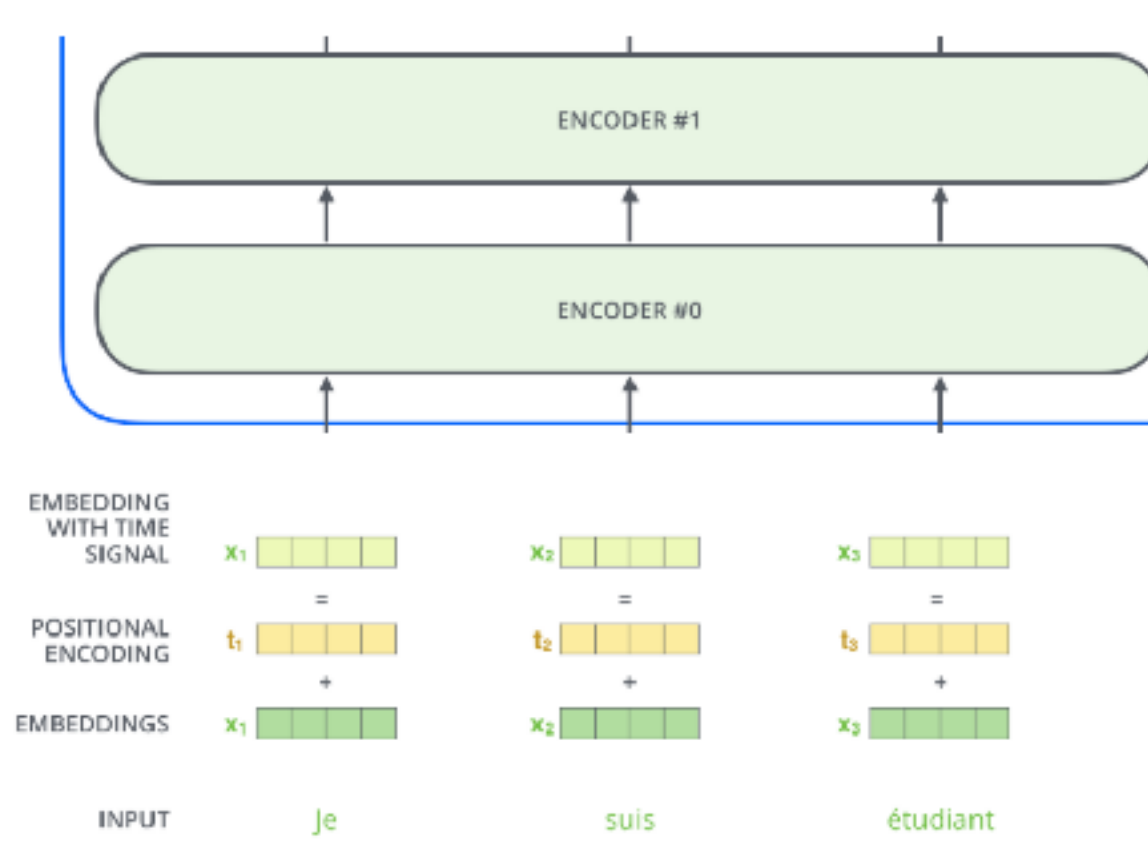  - Need to combine the resulting output sums

# **Transformers**: Attention Visualization

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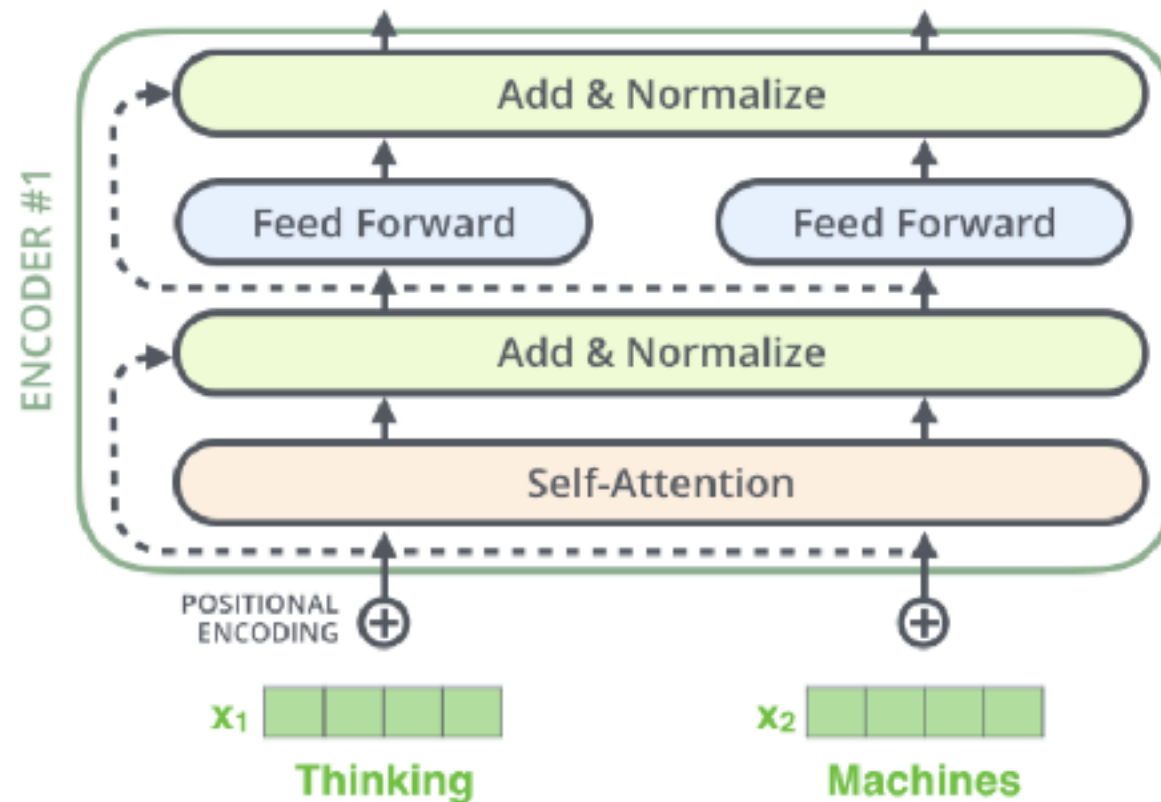  - Illustration for a sentence:

# **Transformers**: Positional Encodings

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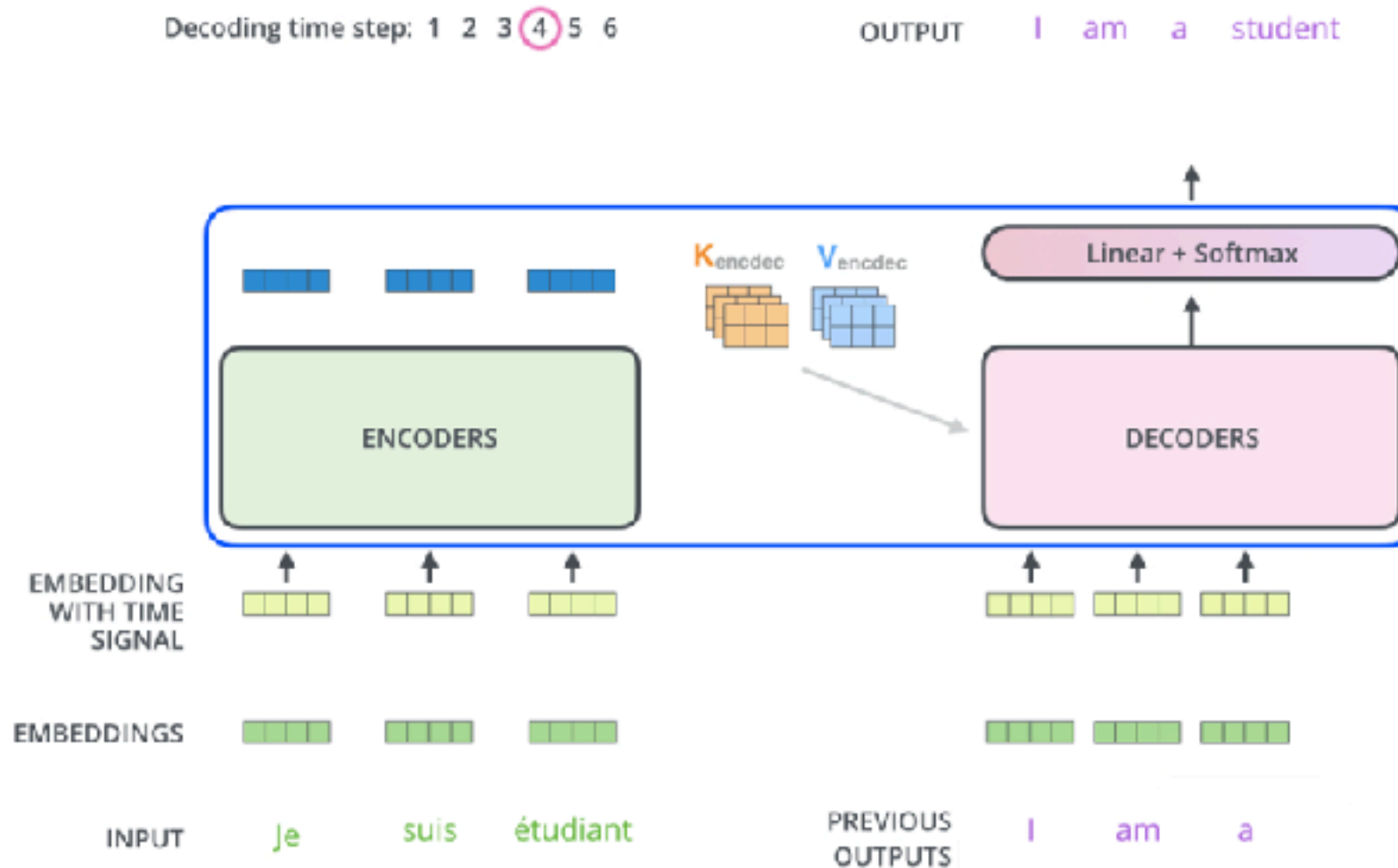  - Add a vector containing a special positional formula's embedding

# **Transformers**: More Tricks

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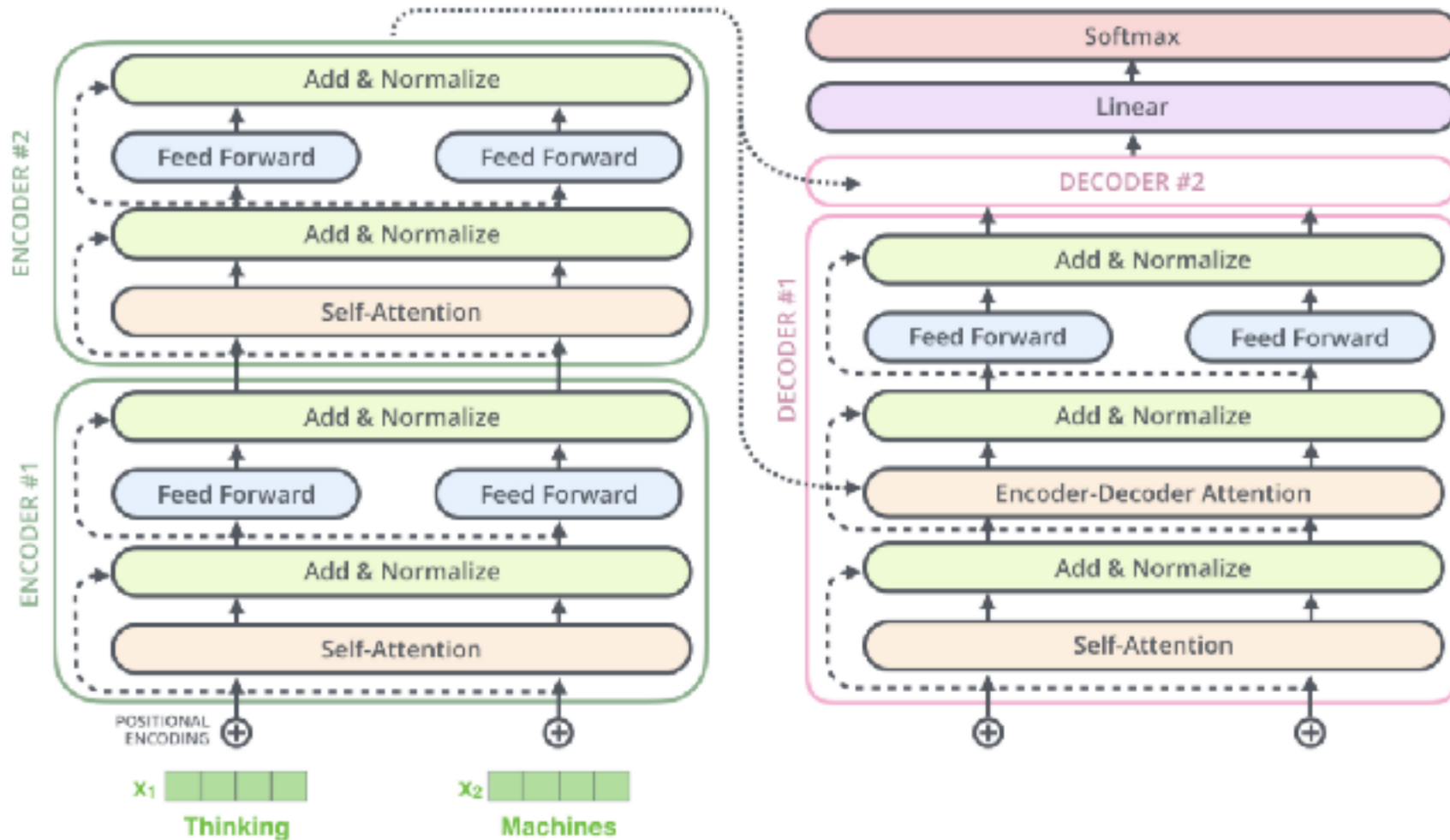  - Apply to our encoder layers

# **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
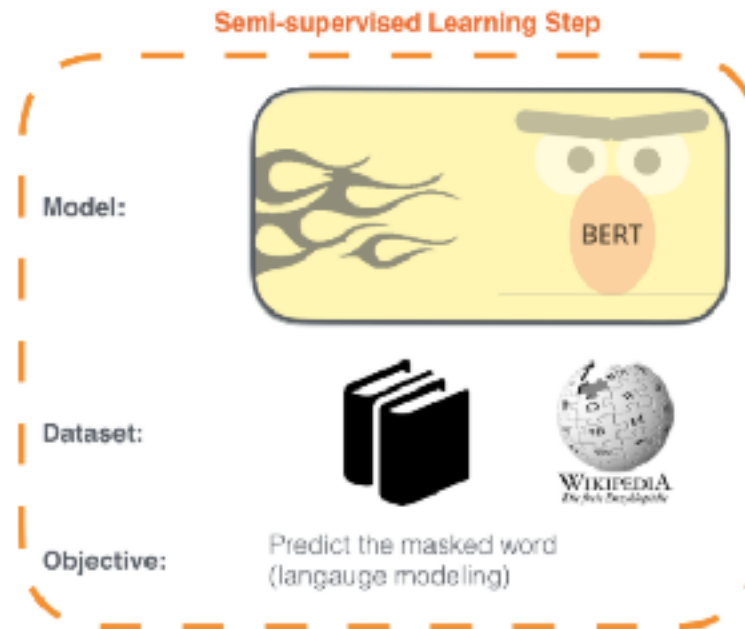- **Transformer-based Models**
  - BERT, GPTs, Foundation Models
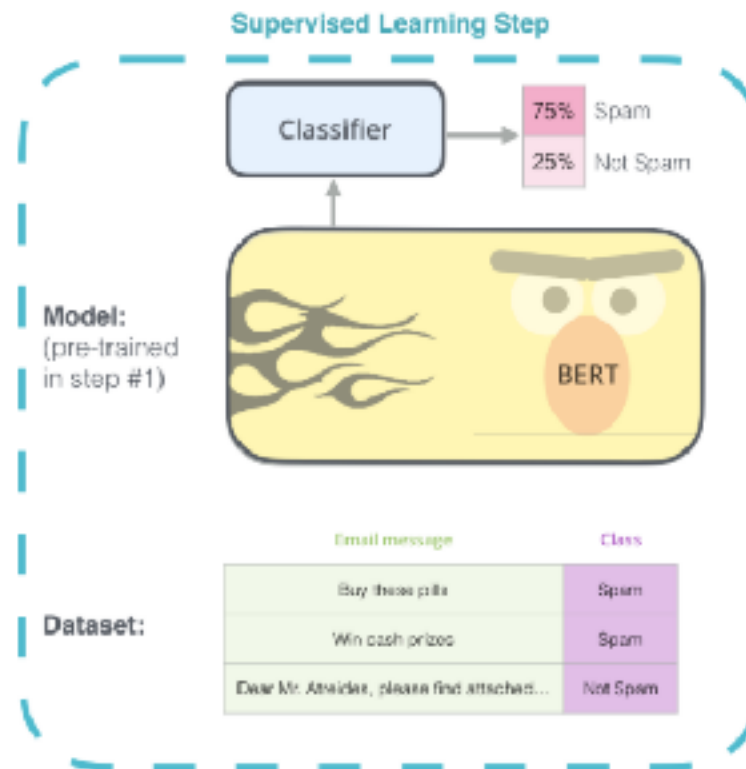
# Transformer-Based Models: **BERT**

- Semi-supervised learning + Transformers
  - Semi-supervised learning to learn embeddings in encoder
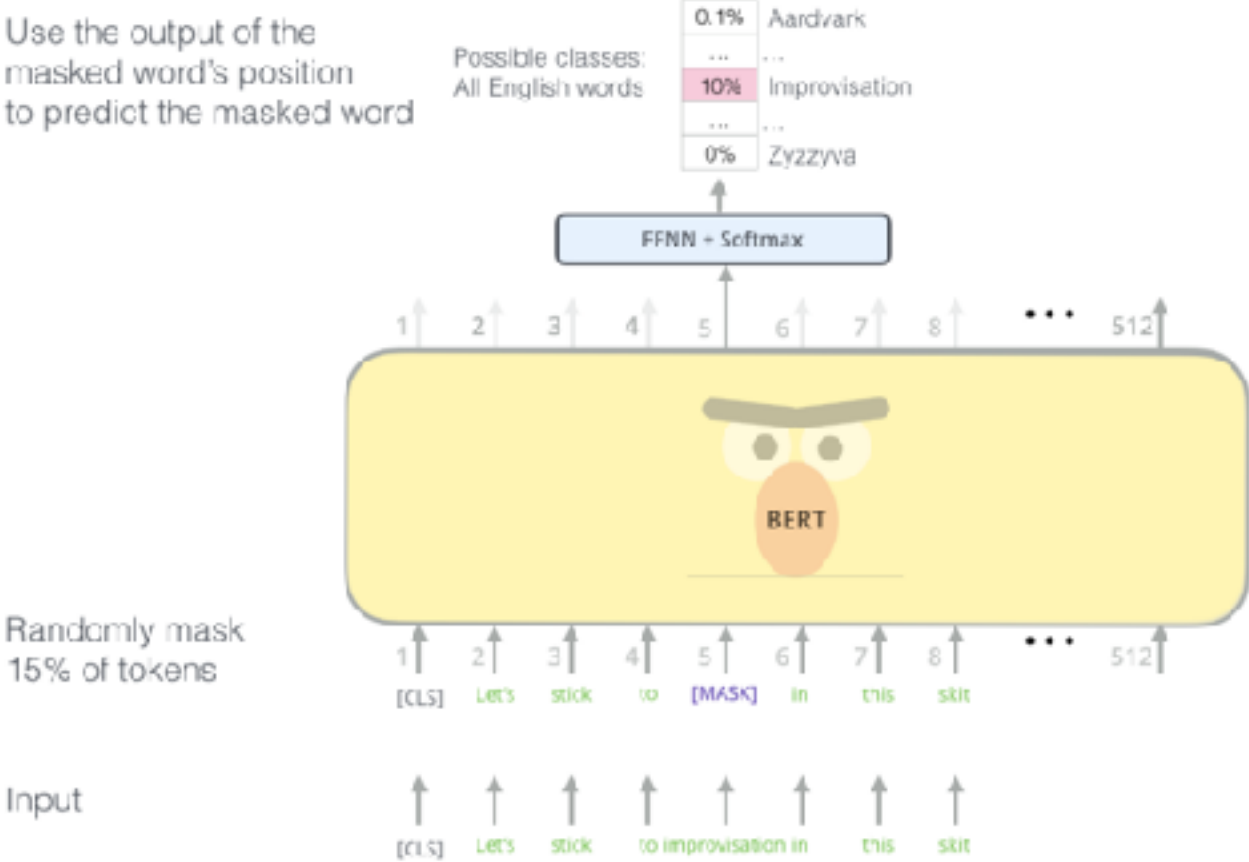
# **BERT**: Concepts

- 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) | QQP | QNLI | SST-2 | CoLA | STS-B | MRPC | RTE | Average |
|---|---|---|---|---|---|---|---|---|---|
|  | 392k | 363k | 108k | 67k | 8.5k | 5.7k | 3.5k | 2.5k | - |
| 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$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

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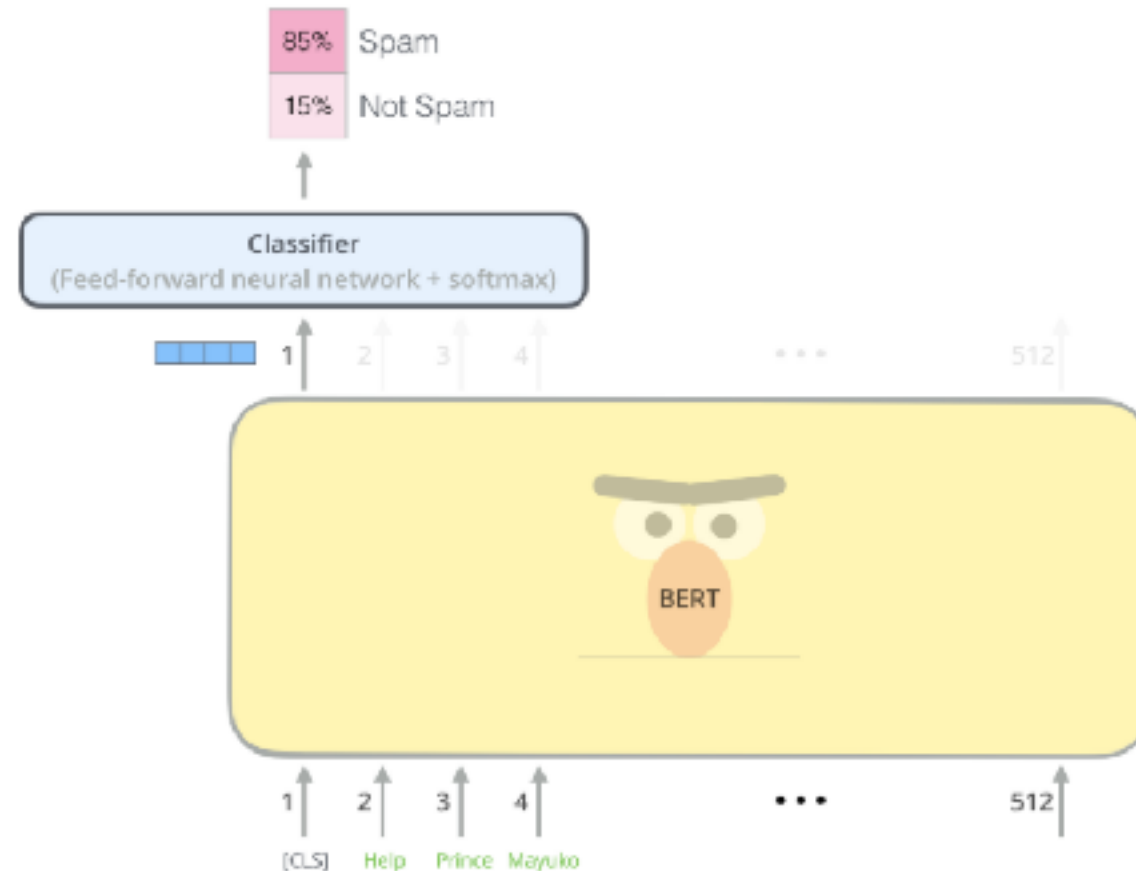
# **BERT**: Training

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

Use the output of the masked word's position to predict the masked word

Possible classes: All English words

| 0.1% | Aardvark |
|------|----------|
| ... | ... |
| 10% | Improvisation |
| ... | ... |
| 0% | Zyzzyva |

FFNN + Softmax

1 2 3 4 5 6 7 8 ••• 512

BERT

Randomly mask 15% of tokens

1 2 3 4 5 6 7 8 ••• 512

[CLS] Let's stick to [MASK] in this skit

Input

[CLS] Let's stick to improvisation in this skit

# BERT: Classification

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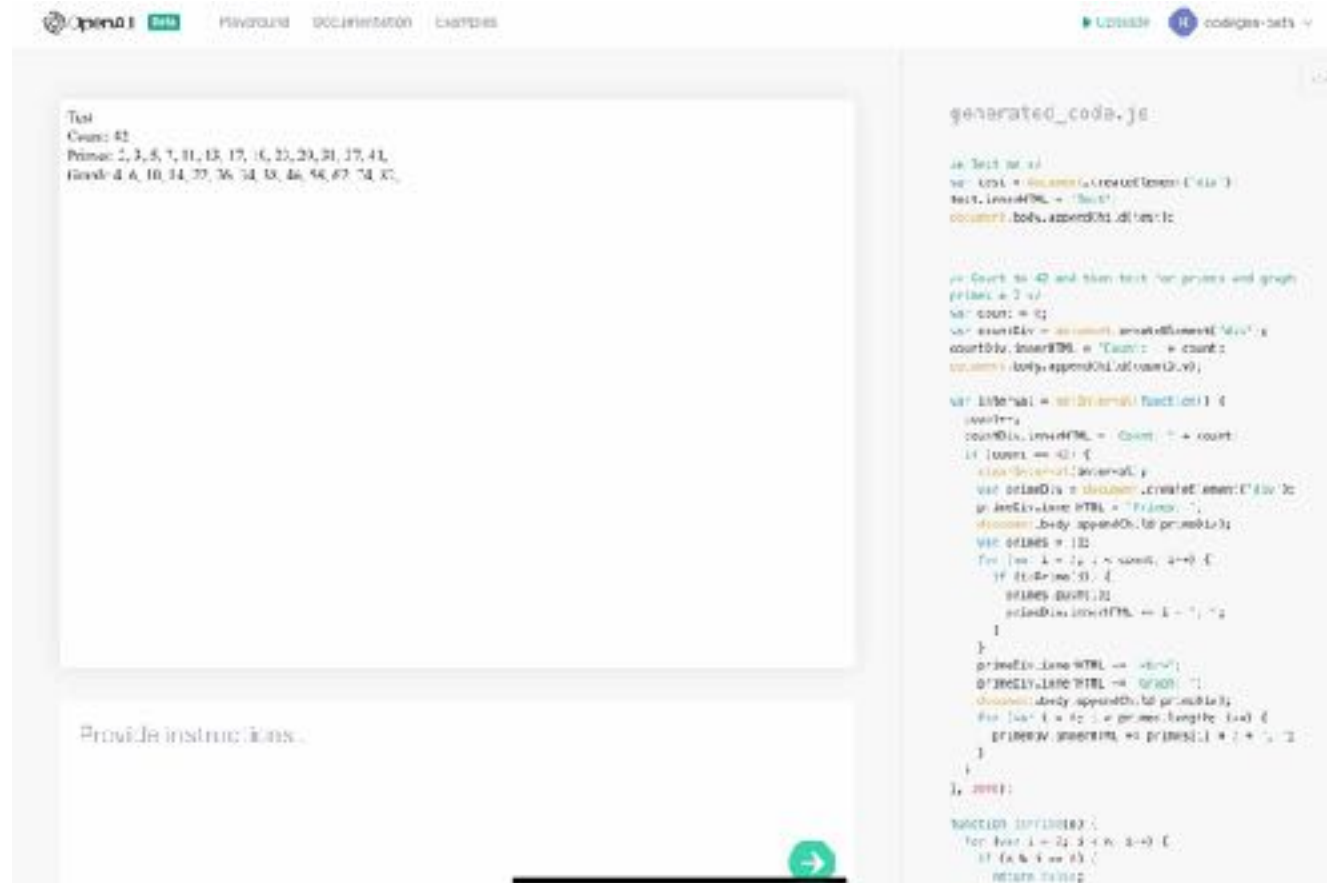  - Example: **binary classification**, spam VS not spam

# **GPT** Series of Models

- **GPT**: **G**enerative **P**re-trained **T**ransformer
  - 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



Russell Foltz-Smith

# DALL-E

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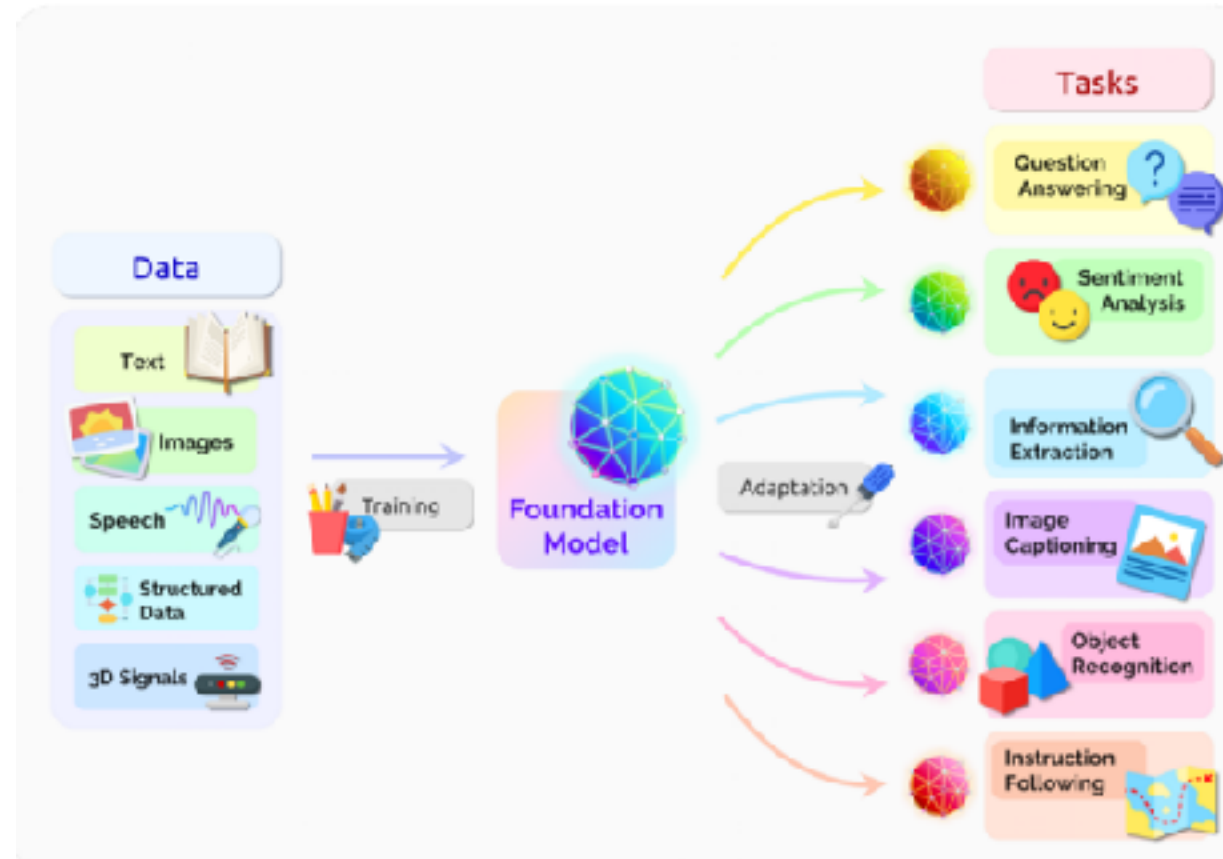


https://openai.com/blog/dall-e/

- Note: several online demos. Try it yourself!
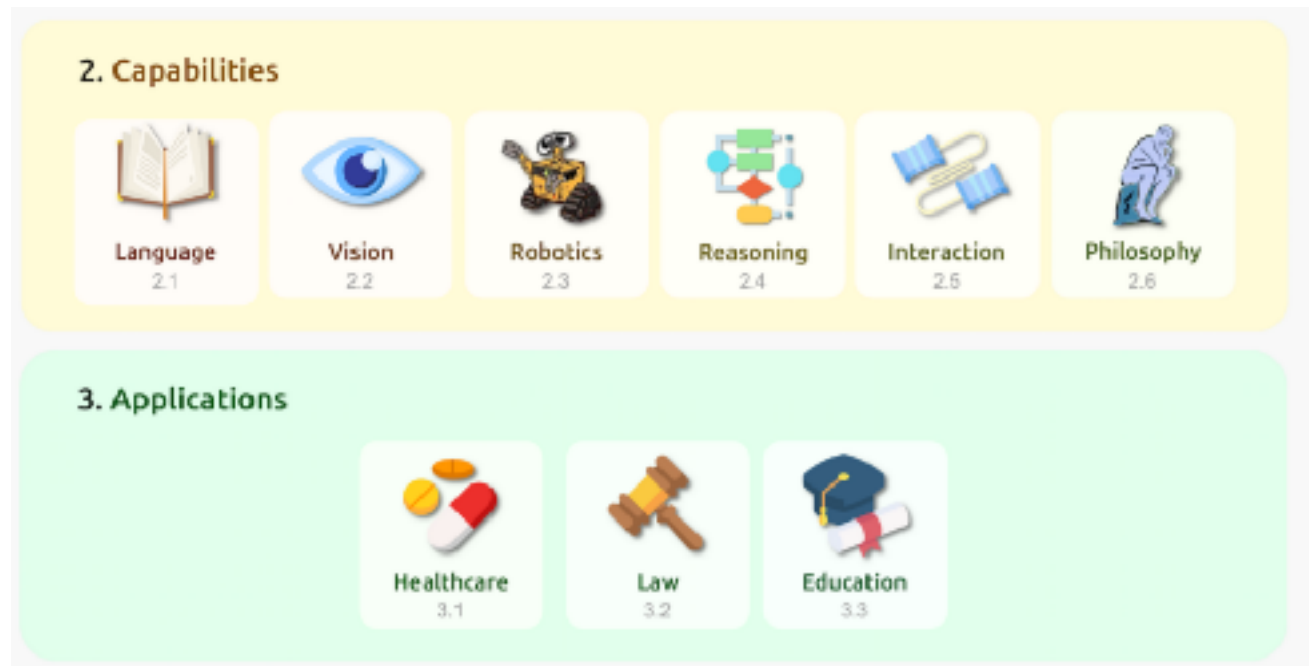
# Foundation Models

- 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, and Fred Sala