# Large Language Models and How to Use Them

Ahmet Üstün
Research Scientist @C4AI @Cohere

May 11, 2023

✈ Cohere For AI

# Agenda

**01**

Preliminaries

**02**

Transformers

**03**

Foundation
Models
(or LLMs)

**04**

Instruction
Following
LLMs

⋌ Cohere For AI

# 01

Preliminaries

Cohere For AI

# Word Embeddings

Word embeddings are vectors that represent word meaning in high-dimensional space that are learned via their context.

> …government debt problems turning into **banking** crises as happened in 2009…
>
> …saying that Europe needs unified **banking** regulation to replace the hodgepodge…
>
> …India has just given its **banking** system a shot in the arm…

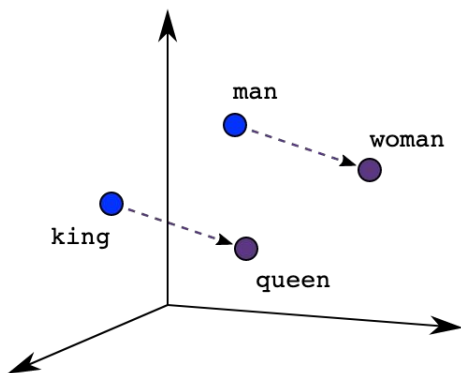These **context words** will represent **banking**

$$
banking =
\begin{pmatrix}
0.286 \\
0.792 \\
-0.177 \\
-0.107 \\
0.109 \\
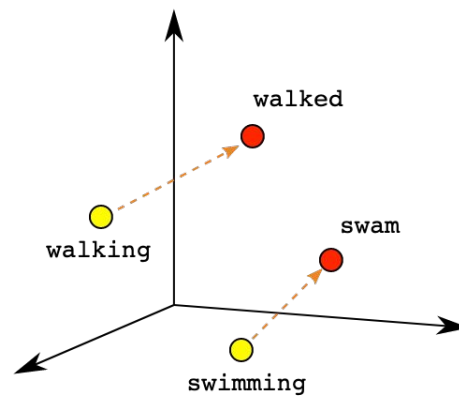-0.542 \\
0.349 \\
0.271
\end{pmatrix}
\qquad
monetary =
\begin{pmatrix}
0.413 \\
0.582 \\
-0.007 \\
0.247 \\
0.216 \\
-0.718 \\
0.147 \\
0.051
\end{pmatrix}
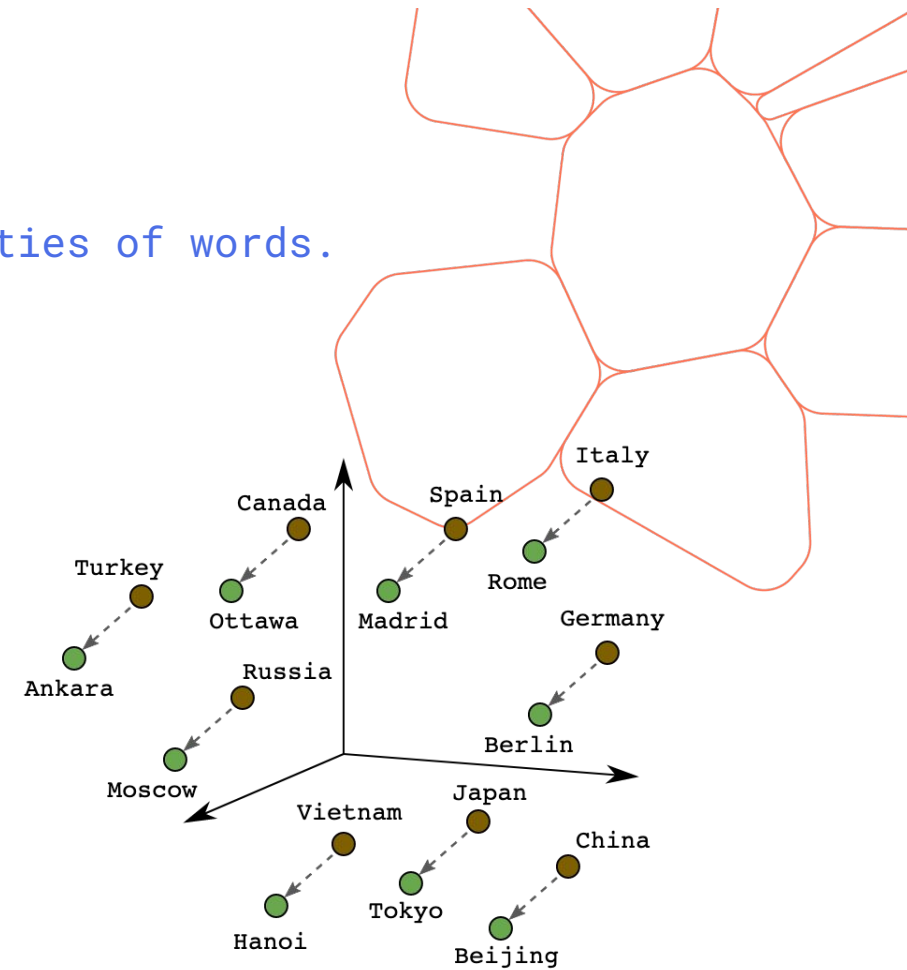$$

# Word Embeddings

Word embeddings can learn different properties of words.
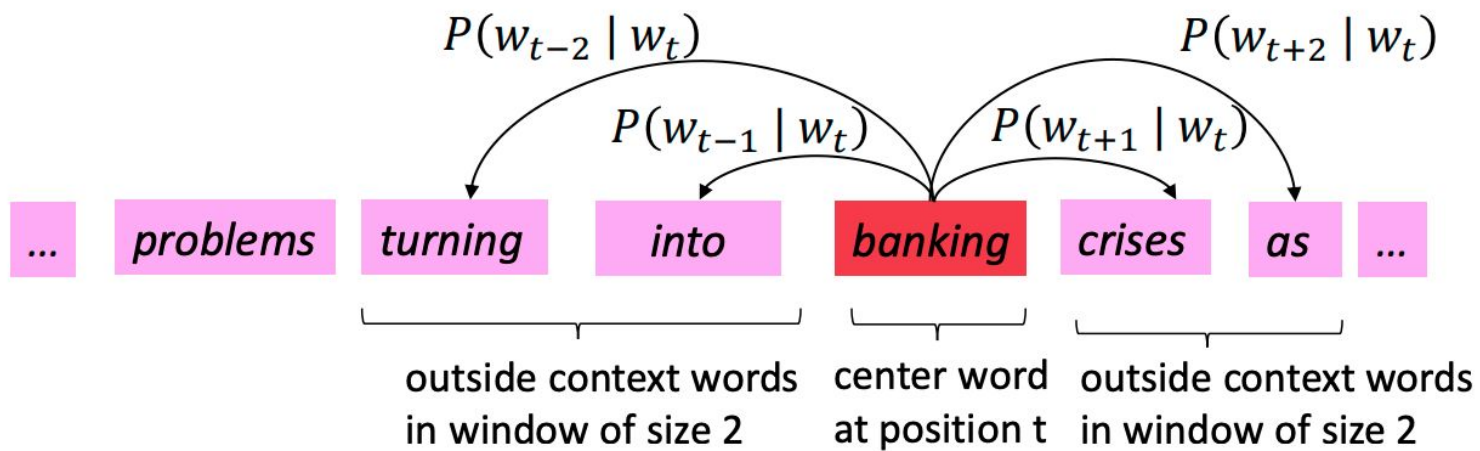


Male-Female

Verb Tense

Country-Capital

✈ Cohere For AI

# Word Embeddings

We learn word embeddings via language modeling.

$$P(w_{t-2} \mid w_t) \qquad P(w_{t+2} \mid w_t)$$

$$P(w_{t-1} \mid w_t) \qquad P(w_{t+1} \mid w_t)$$

| ... | problems | turning | into | banking | crises | as | ... |
|---|---|---|---|---|---|---|---|

outside context words in window of size 2    center word at position t    outside context words in window of size 2

✈ Cohere For AI

# Language Modelling

## Input
### Features

Thou →

shalt →

**Trained Language Model**

**Task:**
Predict the next word

## Output
### Prediction

| | |
|---|---|
| 0% | aardvark |
| 0% | aarhus |
| 0.1% | aaron |
| ... | |
| 40% | not |
| ... | |
| 0.01 | zyzzyva |

✈ Cohere For AI

# 02

---

Transformers

# Neural Language Models

**Recurrent neural networks**

Hidden-states

Outputs



$y^{<1>}$

$y^{<2>}$

$y^{<t>}$

$y^{<t+1>}$

$a^{<0>}$

$a^{<t-1>}$

$a^{<t>}$

$a^{<t+1>}$
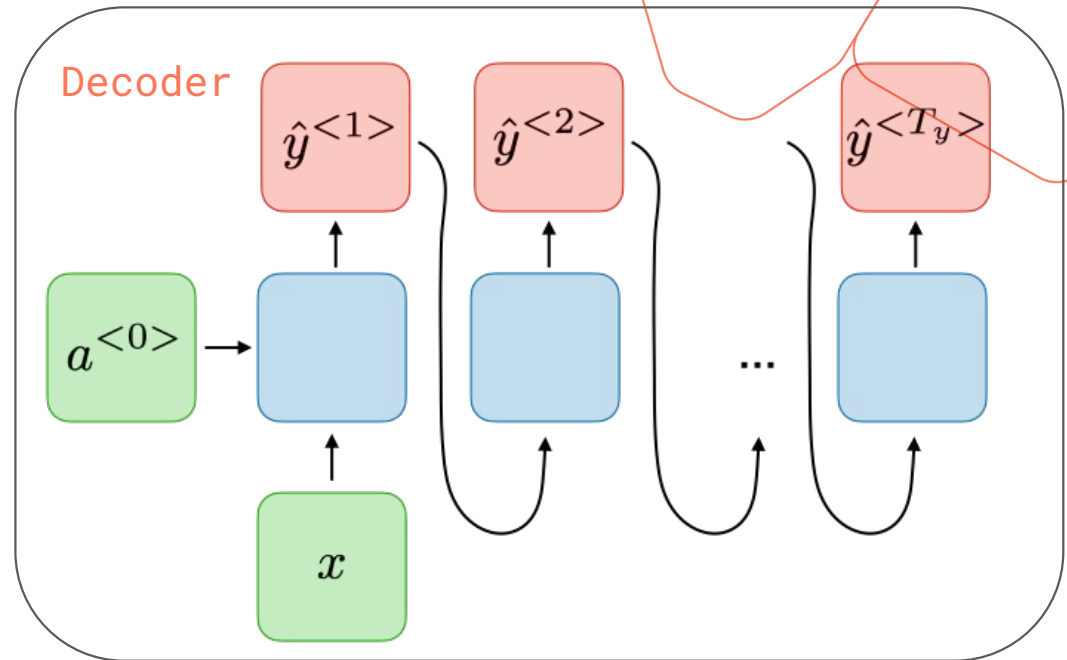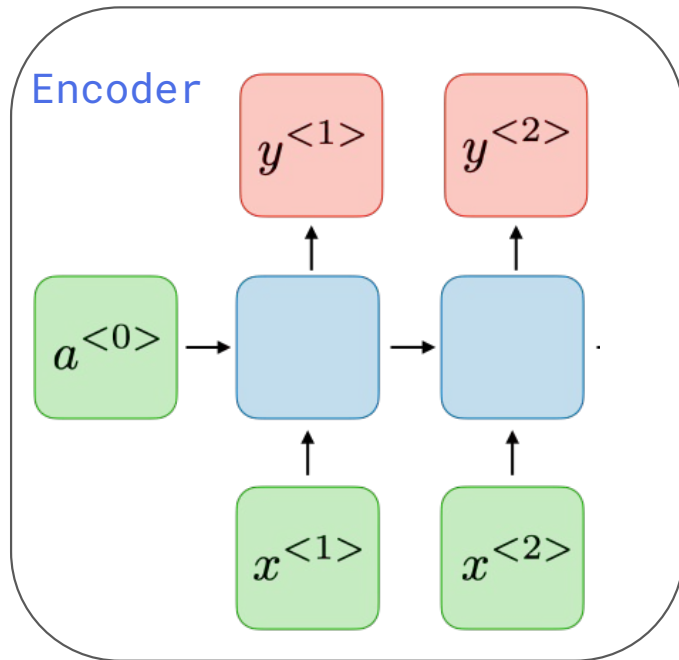
$\dots$

$x^{<1>}$

$x^{<2>}$

$x^{<t>}$

$x^{<t+1>}$

✈ Cohere For AI

# Neural Language Models

You can use RNNs for both encoding a sequence
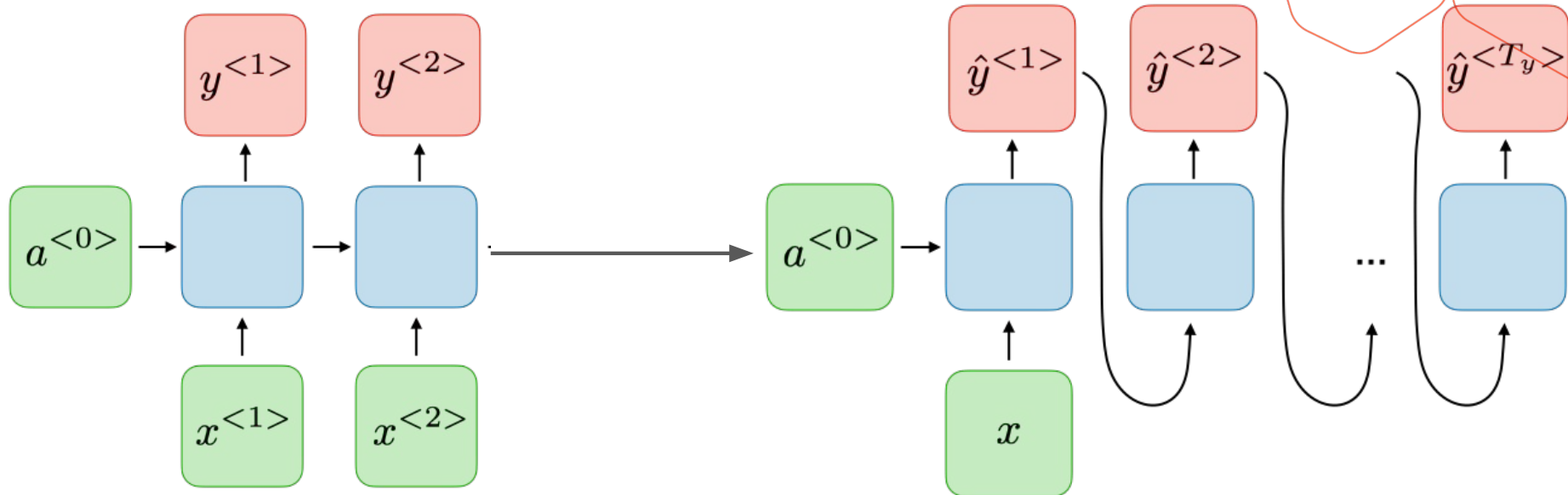
..also for decoding such as predicting next word

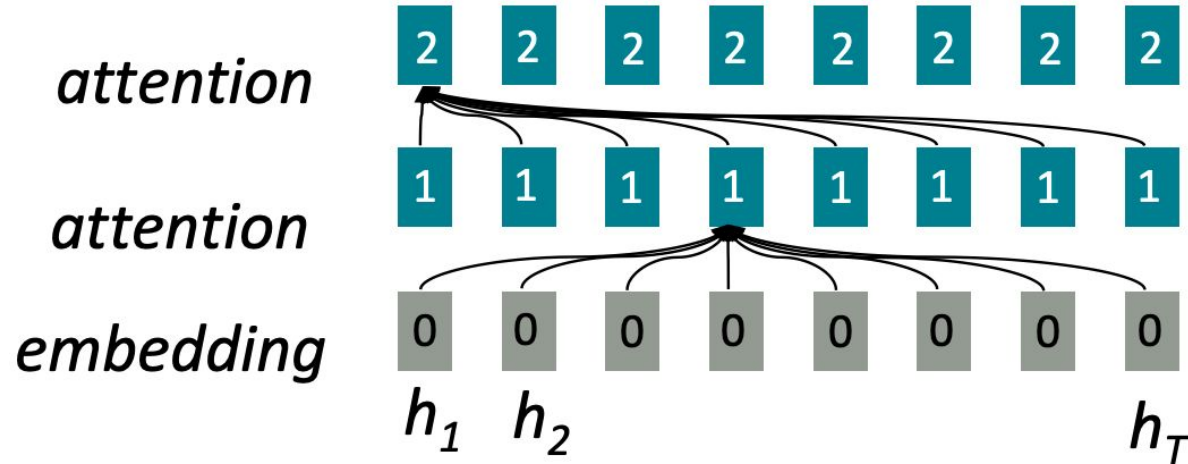# Neural Language Models

You can use RNNs for both encoding a sequence

..also for decoding such as predicting next word

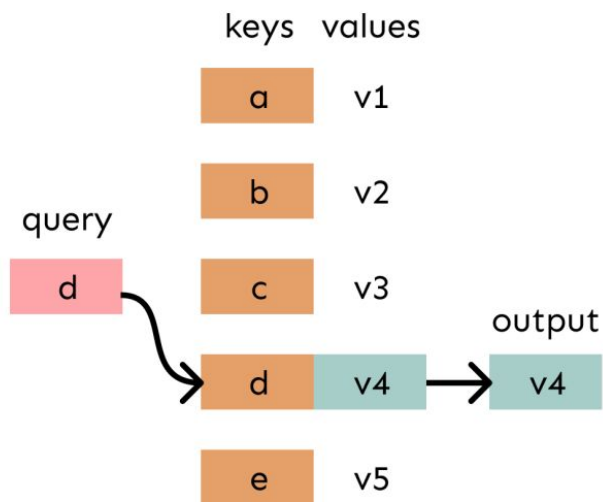## Encoder-Decoder

# Transformers and Self-Attention



All words attend to all words in previous layer; most arrows here are omitted

# Transformers and Self-Attention
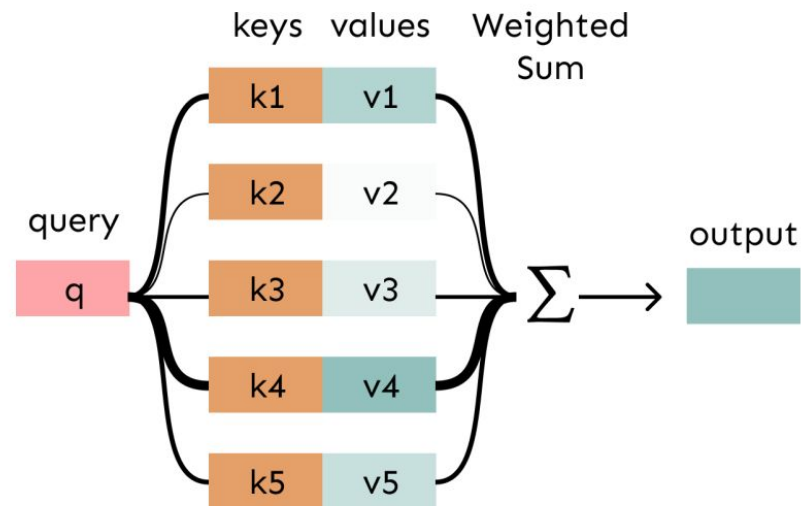
We can think of **attention** as performing fuzzy lookup in a key-value store.

In a **lookup table**, we have a table of **keys** that map to **values**. The **query** matches one of the keys, returning its value.

In **attention**, the **query** matches all **keys** *softly*, to a weight between 0 and 1. The keys' **values** are multiplied by the weights and summed.

✈ Cohere For AI

# Transformers and Self-Attention

Layer: 5 Attention: Input - Input

| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| ' | ' |
| _ | _ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

⊀ Cohere For AI

# What about the word order?



**Position embeddings**

⇕ Cohere For AI

# Why Transformers is so popular?

**The blessings of scale**

AI training runs, estimated computing resources used

Floating-point operations, selected systems, by type, log scale



● Drawing  ● Language
● Vision  ● Other

PaLM (540B)
GPT-3
LaMDA
GPT-2
DALL-E
BERT-Large
NPLM
NetTalk
Neocognitron
ADALINE
Theseus

$10^{24}$
$10^{20}$
$10^{16}$
$10^{12}$
$10^{8}$
$10^{4}$
1

1950  60  70  80  90  2000  10  22

Sources: "Compute trends across three eras of machine learning", by J. Sevilla et al., arXiv, 2022; Our World in Data

✈ Cohere For AI

# 03

---

Foundation Models (LLMs)

# Language modeling with encoder-decoder



**Targets**
<X> for inviting <Y> last <Z>

**Original text**
Thank you for inviting me to your party last week.

**Inputs**
Thank you <X> me to your party <Y> week.
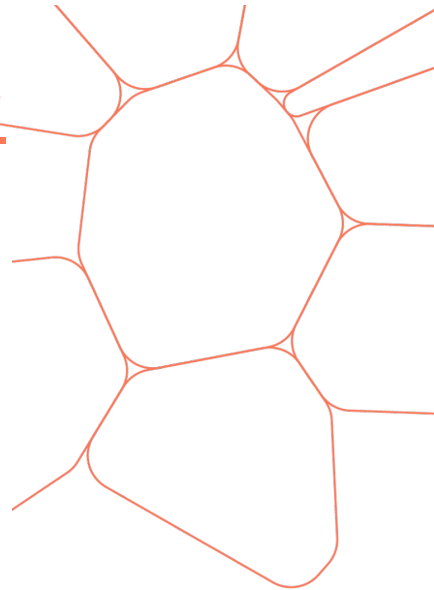
**T5: Text-to-Text Transformer**

⋊ Cohere For AI

# Decoder-Only LMs: Advantage of Scale



GPT-2 SMALL

| 12 | DECODER |
| ... | |
| 1 | DECODER |

Model Dimensionality: 768

GPT-2 MEDIUM

| 24 | DECODER |
| ... | |
| 2 | DECODER |
| 1 | DECODER |

Model Dimensionality: 1024

GPT-2 LARGE

| 36 | DECODER |
| ... | |
| 4 | DECODER |
| 3 | DECODER |
| 2 | DECODER |
| 1 | DECODER |

Model Dimensionality: 1280

GPT-2 EXTRA LARGE

| 48 | DECODER |
| ... | |
| 6 | DECODER |
| 5 | DECODER |
| 4 | DECODER |
| 3 | DECODER |
| 2 | DECODER |
| 1 | DECODER |

Model Dimensionality: 1600

✈ Cohere For AI

# Scale your data too!

200
Billion

30
Billion

<100
Million

3
Billion

1.4
Trillion

13 y.o.
Human

BERT
(2018)

RoBERTa
(2019)

GPT-3
(2020)

Chinchilla
(2022)

# tokens seen during training

# 04

**Few-shot Learning vs Instruction Tuning**

⟨ Cohere For AI

# Task-Specific Fine-tuning

Pretraining can improve NLP applications by serving as parameter initialization.



**Step 1: Pretrain (on language modeling)**
Lots of text; learn general things!

goes    to    make    tasty    tea    END

(Transformer, LSTM, ++ )

Iroh    goes    to    make    tasty    tea

**Step 2: Finetune (on your task)**
Not many labels; adapt to the task!

☺/☹

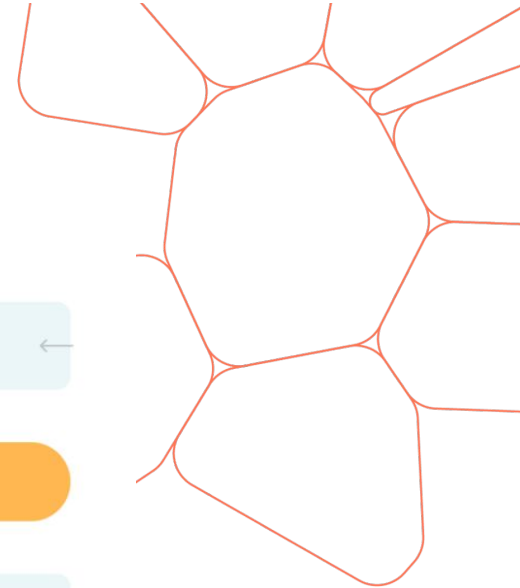(Transformer, LSTM, ++ )

*... the movie was ...*

# Few-shot Learning vs Fine-tuning

## New methods of "prompting" LMs

### Zero/few-shot prompting

```
1  Translate English to French:     ←
2  sea otter => loutre de mer        ←
3  peppermint => menthe poivrée      ←
4  plush girafe => girafe peluche    ←
5  cheese =>      ................   ←
```

### Traditional fine-tuning

```
1  sea otter => loutre de mer        ←
                  ↓
            gradient update
                  ↓
1  peppermint => menthe poivrée      ←
                  ↓
            gradient update
                  ↓
                 • • •
                  ↓
1  cheese =>      ...............    ←
```

✈ Cohere For AI

## Everyone will be a prompt engineer!

**WIKIPEDIA**
The Free Encyclopedia

Q    ...

# Prompt engineering

文A 5 languages ˅

Article    Talk                                                More ˅

From Wikipedia, the free encyclopedia

**Prompt engineering** is a concept in artificial intelligence, particularly natural language processing (NLP). In prompt engineering, the description of the task is

# Prompt Engineer and Librarian          APPLY FOR THIS JOB

SAN FRANCISCO, CA / PRODUCT / FULL-TIME / HYBRID

✈ Cohere For AI

# Chain-of-thought prompting

## Standard Prompting

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

## Chain-of-Thought Prompting

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?
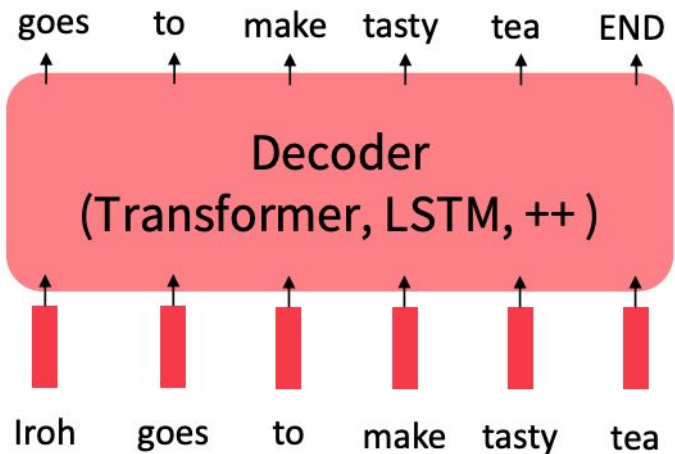
**Model Output**

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔️

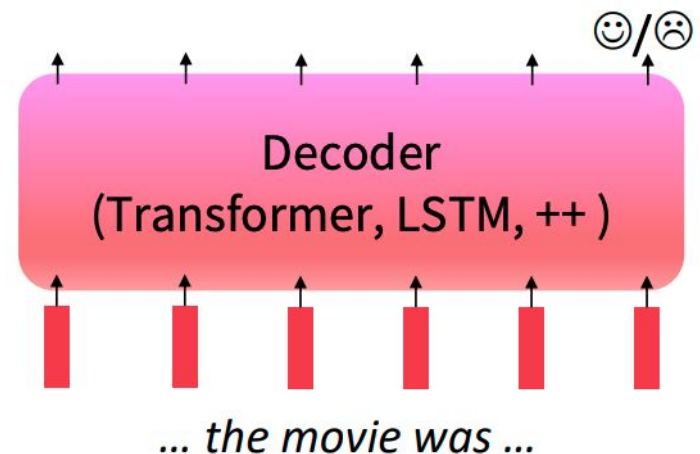[Wei et al., 2022; also see Nye et al., 2021]

# Instruction Tuning



**Step 1: Pretrain (on language modeling)**

Lots of text; learn general things!

goes    to    make    tasty    tea    END

Decoder
(Transformer, LSTM, ++ )

Iroh    goes    to    make    tasty    tea

**Step 2: Finetune (on many tasks)**

~~Not~~ many labels; adapt to the tasks!

☺/☹

Decoder
(Transformer, LSTM, ++ )

... the movie was ...

✈ Cohere For AI

# Instruction Tuning

- **Collect examples** of (instruction, output) pairs across many tasks and finetune an LM



- Evaluate on **unseen tasks**

[FLAN-T5; Chung et al., 2022]

✈ Cohere For AI

# Instruction Tuning

## Model input (Disambiguation QA)

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:
(A) They will discuss the reporter's favorite dishes
(B) They will discuss the chef's favorite dishes
(C) Ambiguous

A: Let's think step by step.

## Before instruction finetuning

The reporter and the chef will discuss their favorite dishes.
The reporter and the chef will discuss the reporter's favorite dishes.
The reporter and the chef will discuss the chef's favorite dishes.
The reporter and the chef will discuss the reporter's and the chef's favorite dishes.

❌ **(doesn't answer question)**

## After instruction finetuning

The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C). ✅

✈ Cohere For AI