

# In-Context Learning

The Universal API also opens up an interesting possibility

- Can we use a Large Language Model
- To solve a new Target task  
 $\backslash \text{pr} \backslash y \mid \backslash x$   
for an example in the Target task
  - with input  $\backslash x$  and target  $\backslash y$
- *without* further training (i.e., Fine-Tuning)

On a pure syntax level, this is feasible

- all problems are instances of "predict the next"
- perhaps the LLM's text completion power *also* captures domain-specific knowledge
  - the completion is related to the domain-specific words of the prompt
  - for example, if the prompt is in French, the completion would be expected to be in French too

So the possibility is there.

But how does the Source LLM discern what the Target task is ?

One way is to provide some *context*  $C$  that describes the Target task.

We thus try to predict

$$\text{\textcolor{red}{pr}}\text{\textcolor{red}{y}} \mid C, \text{\textcolor{red}{x}}$$

Context  $C$  could be

- an *instruction* that describes the task

An alternative is to

- describe the task by providing demonstrations (*exemplars*)
  - of the relationship between input  $\mathbf{x}$  and target  $\mathbf{y}$
  - as part of the Context  $C$

Consider  $k$  exemplars

$$\langle \mathbf{x}^{(1)}, \mathbf{y}^{(1)} \rangle, \dots, \langle \mathbf{x}^{(k)}, \mathbf{y}^{(k)} \rangle$$

from the Target task.

Let  $\langle \mathbf{x}, \mathbf{y} \rangle$  be a test example

- we know  $\mathbf{x}$
- and want the model to predict  $\hat{\mathbf{y}} = \mathbf{y}$

The idea behind *In-Context Learning* is to

- condition the Pre-Trained Language Model (Source LLM) with a context  $C$ 
  - containing exemplars of the Input to Target mapping for the Target task
- to output a predict  $\hat{y}$
- given an unseen input  $x$
- **without** further training (Fine-Tuning)

So the prompt might look like

Input:  $\backslash \mathbf{x}^{(1)}$

Output:  $\backslash \mathbf{y}^{(1)}$

$\vdots$

Input:  $\backslash \mathbf{x}^{(k)}$

Output:  $\backslash \mathbf{y}^{(k)}$

Input:  $\backslash \mathbf{x}$

Output:

and expect the continuation to be the prediction  $\hat{\backslash \mathbf{y}}$  corresponding to test input  $\backslash \mathbf{x}$

That is, we defined the Context  $C$  to be  $k$  exemplars of the input to target mapping

Input:  $\backslash \mathbf{x}^{(1)}$

Output:  $\backslash \mathbf{y}^{(1)}$

$\vdots$

Input:  $\backslash \mathbf{x}^{(k)}$

Output:  $\backslash \mathbf{y}^{(k)}$

and ask the model to predict

given an arbitrary (and unseen) input  $\backslash \mathbf{x}$ .

$\backslash \mathbf{pr} \backslash \mathbf{y} \mid C, \backslash \mathbf{x}$



For example, we can describe Translation between languages with the following Context  
*C*

Translate English to French

sea otter => loutre de mer

peppermint => menthe poivree

plush giraffe => girafe peluche

The expectation is that when the user presents the prompt \x

cheese =>

the model will respond with the French translation of cheese .

- the "next words" predicted by the Language Modeling

# Learning to learn

Does In-Context learning really work ?

We can begin to answer this question by

- examining the behavior of a Pre-Trained LLM
- on a new task
- using  $k$  exemplars
  - varying  $k$

Depending on  $k$ , we refer to the behavior of the LLM by slightly different names

- **Few shot learning:**  $10 \leq k \leq 100$  typically
- **One shot learning:**  $k = 1$
- **Zero shot learning**  $k = 0$

A picture will help

## The three settings we explore for in-context learning

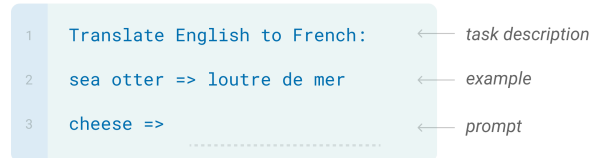
### Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



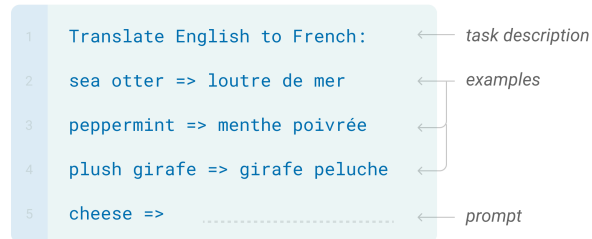
### One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



## Traditional fine-tuning (not used for GPT-3)

### Fine-tuning

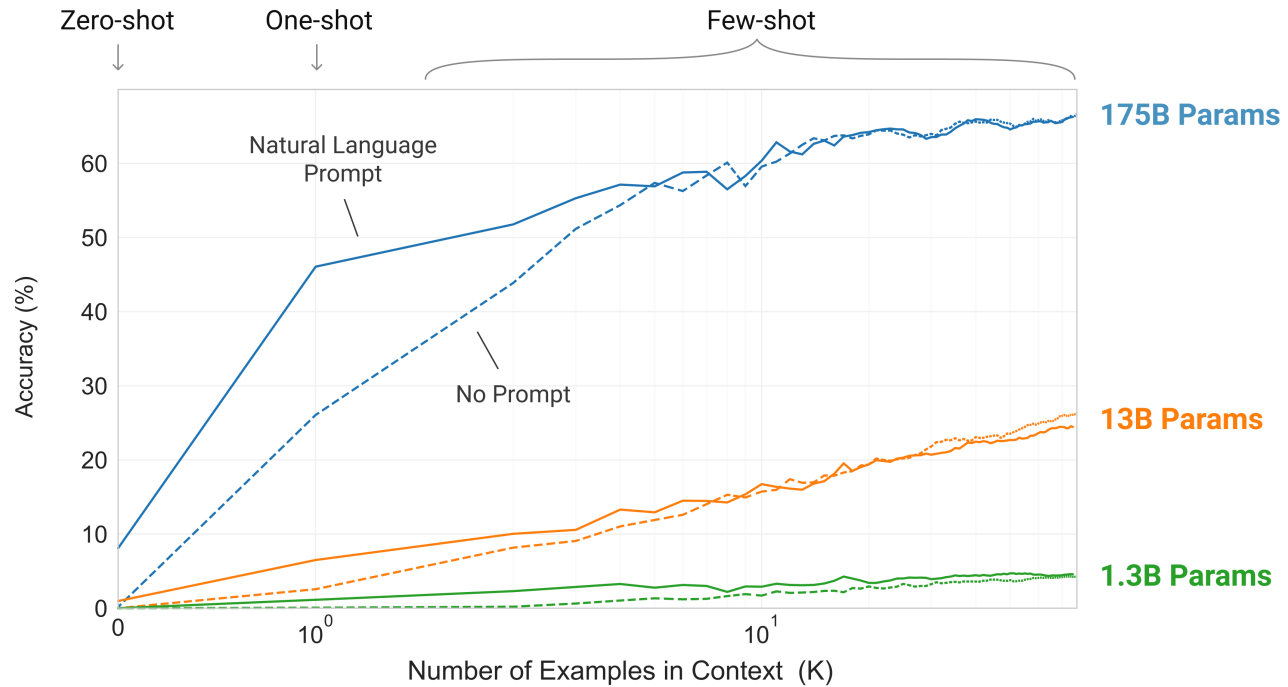
The model is trained via repeated gradient updates using a large corpus of example tasks.



Is this even possible ?! Learning a new task with **zero** exemplars ?

Let's look at the reported In-Context Learning results of 3 LLM's of varying size.

### Few/One/Zero shot learning



Picture from: <https://arxiv.org/pdf/2005.14165.pdf#page=4>

## A couple of observations

- As the size of the model grows: In-Context Learning behavior improves
  - compare the 175 Billion parameter model to the smaller models
  - we sometimes refer to this as behavior that "emerges" only when a model is sufficiently large
- More exemplars (greater  $k$ ) helps
  - but not much for the smallest model
- Zero shot learning works !
  - but this is a behavior that only emerges for very large models



# "Fine-tuning" a model with In-Context Learning

"Fine-tuning" (Transfer Learning) refers to

- adapting a model for a Source Task
- to be able to perform a Target Task

Usually

- this means training a pre-trained model (Source Task)
- with a small number of examples of the Target Task
- causing the model's weights to adapt to the Target Task

But In-Context Learning

- adapts a model for the Source Task
- to be able to solve the Target Task
- **without** changing the model's weights

So this may be an alternative method of Fine-Tuning

# In-Context learning: let's experiment

The [HuggingFace platform \(https://huggingface.co/\)](https://huggingface.co/) has libraries of pre-trained models for many tasks, including Language models.

There is a clean API for using these models in code (I recommend their on-line [course \(https://huggingface.co/\)](https://huggingface.co/) if you want to play with it).

But they also host many of their models for interactive use.

This is valuable not just for the obvious reason of ease of use

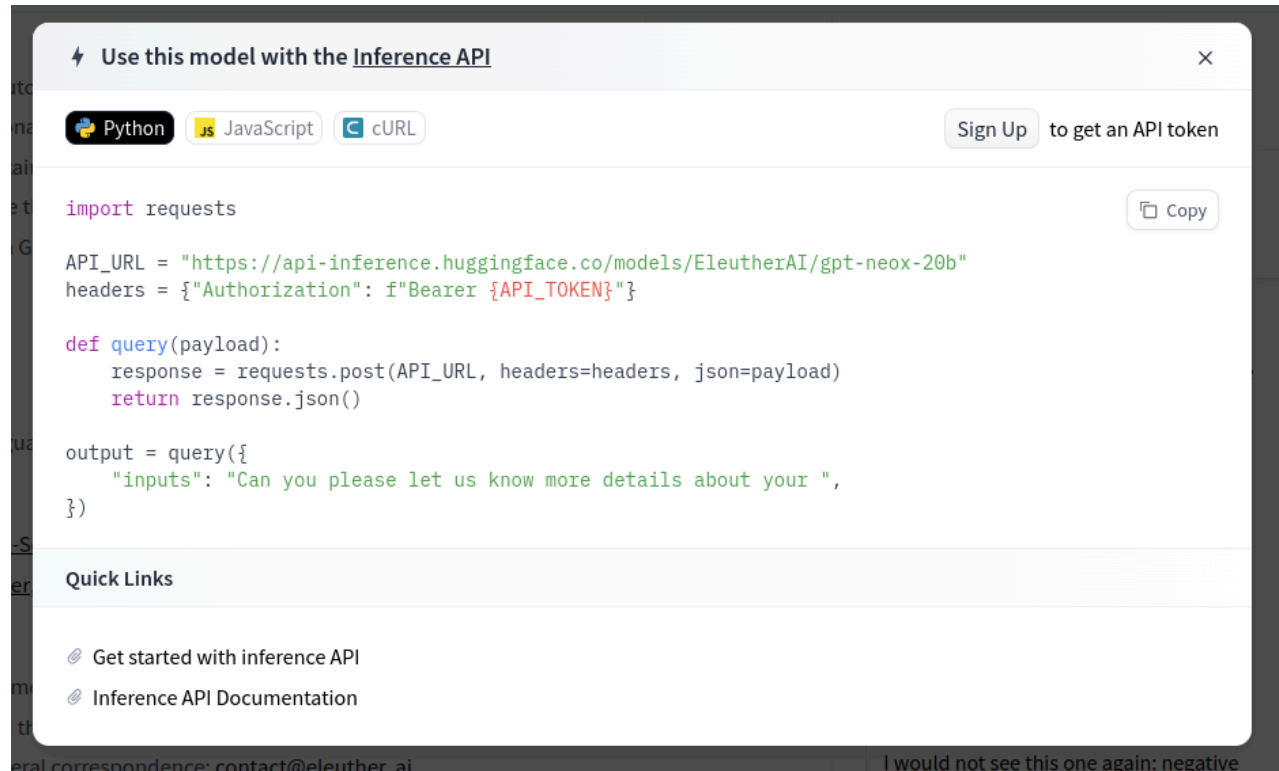
- some models are too big to load on the machines available to us

You can play with In-Context learning by going to the page of a model and typing into the *Hosted Inference API* text box.

But there is also an API that allows you to pass the input (context plus prompt) via a URL.

If you click on the **Deploy** button and choose the **Inference API** drop-down

- you will see Python code for querying the model programatically.



The screenshot shows a web interface titled "Use this model with the Inference API". It features three tabs: "Python", "JavaScript", and "cURL", with "Python" selected. A "Sign Up" button is present next to the text "to get an API token". A "Copy" button is located next to the Python code. The code defines a function `query` that sends a POST request to the Hugging Face Inference API. The code is as follows:

```
import requests

API_URL = "https://api-inference.huggingface.co/models/EleutherAI/gpt-neox-20b"
headers = {"Authorization": f"Bearer {API_TOKEN}"}

def query(payload):
    response = requests.post(API_URL, headers=headers, json=payload)
    return response.json()

output = query({
    "inputs": "Can you please let us know more details about your ",
})
```

Below the code is a "Quick Links" section with two links:

- Get started with inference API
- Inference API Documentation

At the bottom of the interface, there is a footer with the text "eral correspondence: contact@eleuther.ai." and "I would not see this one again: negative".

Please note

- our toy example above used a *single* test example
- even if we manage to get a correct prediction on a single example
  - we don't have confidence that the new task was successfully learned !
  - we really should evaluate success on a larger number of text examples
- still: the fact that the exemplars taught the model the correct syntax for an answer is exciting

Here is a very crude notebook that uses the HuggingFace inference API to experiment with in-context learning.

- [Experiment in In context learning: Colab](https://colab.research.google.com/github/kenperry-public/ML_Advanced_Fall_2024/blob/master/HF_inference_play.ipynb)  
([https://colab.research.google.com/github/kenperry-public/ML\\_Advanced\\_Fall\\_2024/blob/master/HF\\_inference\\_play.ipynb](https://colab.research.google.com/github/kenperry-public/ML_Advanced_Fall_2024/blob/master/HF_inference_play.ipynb))
- [Experiment in In context learning: local \(HF\\_inference\\_play.ipynb\)](#)

Our new Target task is movie reviews.

Here are the exemplars we will use

```
exemplars = [ "this movie was great: positive",  
              "one of the best films of the year: positive",  
              "just plain awful: negative",  
              "I would not see this one again: negative",  
              "this movie was great: positive",  
              "one of the best films of the year: positive",  
              "just plain awful: negative",  
              "I would not see this one again: negative",  
              "I love this film: positive"  
]
```

Not the greatest exemplars (too short), but we just want to illustrate the idea.

We want the model to classify the following review as positive/negative

"I've heard not so great things about this one:"

We would hope the exemplars are sufficient to cause the "predict the next" task to generate the continuation

negative



We will use a relatively small (20B parameters) LLM.

Here are the results

[0.22 seconds, using EleutherAI/gpt-neox-20b]

this movie was great: positive

one of the best films of the year: positive

just plain awful: negative

I would not see this one again: negative

this movie was great: positive

one of the best films of the year: positive

just plain awful: negative

I would not see this one again: negative

I love this film: positive

I've heard not so great things about this one: negative

Negative ! Just as we had hoped.

However, the exemplars are not ideal. The LLM continues *beyond* the classification and further generates more reviews !

it is entertaining: positive  
I would not see this one again: negative  
I love this film: positive  
I've heard not so great things: negative  
it is entertaining: positive  
I've heard not so great things: negative  
enchanced by slow motion visuals: positive  
excellent: positive  
terrific sound design in this one too: positive  
Stallone is a great actor: positive  
I'd turn down a free trip to London to see this movie:

Obviously, our exemplars did not fully convey our intent.

Creating a prompt (context) that will cause an LLM to do *exactly* what you want

- is an art
- called Prompt Engineering

In [2]: `print("Done")`

Done

