# Welcome!

My name is **Sinan Ozdemir** ( in/sinan-ozdemir + @prof_oz )

- Current **founder** of Loop Genius (using AI to help entrepreneurs get their first 100 customers)
- Current **lecturer** for O'Reilly and Pearson
- Founder of Kylie.ai (Funded by OpenAI Founder + Acquired)
- **Masters** in Theoretical Math from **Johns Hopkins**
- Former lecturer of Data Science at Johns Hopkins
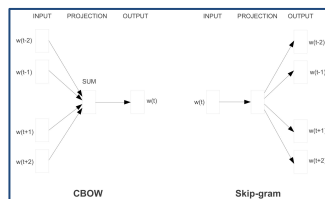
Author of ML textbooks and online series, including

- [The Principles of Data Science](#)
- [Introduction to Transformer Models for NLP](#)
- [Quick Start Guide to LLMs](#) (Top 10 in NLP on Amazon)

# Quick Recap

# Brief History of Modern NLP

**2001**
Neural Language Models

**2014 - 2017**
Seq2seq + Attention

**2013**
encoding semantic meaning with Word2vec

**2017 - Present**
Transformers + Large Language Models

Bengio et al. https://www.jmlr.org/papers/volume3/bengio03a/bengio03a.pdf
Mikolov et al. https://arxiv.org/abs/1301.3781
Xu et al. http://proceedings.mlr.press/v37/xuc15.pdf
https://papers.nips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

# Auto-__ Language Model Use Cases

## Auto-regressive Models

1. Predicting next word in a sentence (auto-complete)

2. Natural Language Generation (NLG)

3. GPT Family

## Auto-encoding Models

1. Comprehensive understanding and encoding of entire sequences of tokens

2. Natural Language Understanding (NLU)

3. BERT Family

# Using the ChatGPT Playground

Write an instruction to the LLM, and see the response

**Prompt**

Write a tweet talking about how great GPT-3 and ChatGPT are

"Just used #GPT3 and #ChatGPT to generate some amazing content! These AI language models are simply mind-blowing with their ability to comprehend, process and produce human-like language. Kudos to @OpenAI for creating such amazing tech! #AI #MachineLearning"

**LLM Response**

Source: ChatGPT Playground

Pearson

# Using LLMs

We can use LLMs in (generally) three ways:

1. **Encode** text into semantic vectors with little/no fine-tuning
   a. Eg. Creating an information retrieval system using BERT vectors

2. Fine-tune a pre-trained LLM to perform a very specific task using **Transfer Learning**
   a. Eg. Fine-tuning BERT to classify sequences with labels

3. Ask an LLM to solve a task it was pre-trained to solve or could intuit
   a. Eg. **Prompting** GPT3 to write a blog post
   b. Eg. **Prompting** T5 to perform language translation

Pearson

# Does the LLM know enough for my task?

A.  **Yes**, it has all knowledge encoded and it is ready to solve my task
    a.  May still need to format output to make it easier to work with


B.  **Mostly.** It knows the information but it lacks critical information (information is too new to be in the model or it knows a topic but not to the specifics that I need)
    a.  Create a secondary system to retrieve information on demand
    b.  Few-shots and chain of thought to help teach nuances/specifics


C.  **No,** not at all, I need to teach it pretty much everything from scratch
    a.  Just ask with explicit instructions
    b.  Few shot / chain of thought prompting
    c.  Fine-tuning for long term cost savings/speed

**Pearson**

# Identifying Patterns and Making Predictions

Pearson

# E.G., Classification

Still the most common ML application, **classification** assigns labels to a piece of data/text from a *finite set of labels*.

This generally requires fine-tuning to teach an LLM about the classes it needs to predict.

Examples: Assigning a news topic to an article, assigned parts of speech to words in a sentence

Pearson

# Classification with LLMs

We can perform classifications in many ways:

1. Prompting an LLM with a great set of instructions on how to classify
2. Using few shot to teach the model via *in-context* learning
3. Fine-tuning an LLM with labeled data (like we would before LLMs)
4. Using *0-shot classification*

Pearson

# 0-shot Classification

A variation of classical classification, **0-shot classification** is able to dynamically take in labels without fine-tuning and assign labels.

Models like BART (a variant of BERT from Meta) are great at this

facebook/**bart-large-mnli** ♡ like | 586

Pearson

# 0-shot Classification

⊞ Zero-Shot Classification          Examples  ⌄

> I have a problem with my iphone that needs to be resolved asap!!

Possible class names (comma-separated)

> urgent, not urgent, phone, tablet, computer

| | |
|---|---|
| urgent | 0.999 |
| phone | 0.995 |
| computer | 0.095 |
| not urgent | 0.001 |
| tablet | 0.000 |

⊶ facebook/**bart-large-mnli** ⧉          ♡ like  586

Pearson

# Transfer Learning with BERT



**Hugging Face**  Search models, datasets, use    Models    Datasets    Spaces    Docs    Soluti

**Tasks**

Fill-Mask    Question Answering
Summarization    Table Question Answering
Text Classification    Text Generation
Text2Text Generation    Token Classification
Translation    Zero-Shot Classification
Sentence Similarity    +16

**Libraries**

**Models**  35,367    Search Models

distilgpt2
Text Generation · Updated May 21, 2021 · ↓ 33.2M · ♡ 39

bert-base-uncased
Fill-Mask · Updated May 18, 2021 · ↓ 16.2M · ♡ 125

cross-encoder/ms-marco-MiniLM-L-12-v2
Text Classification · Updated Aug 5, 2021 · ↓ 11M · ♡ 5

**Selecting a source model**

Additional Task Layers

Pre-trained BERT

Training data for second task

**Reusing and training model**

# Fine-tuning OpenAI models



Finetuned ADA / DaVinci / ChatGPT Prices vs Days (assuming 1,000 classifications a day)
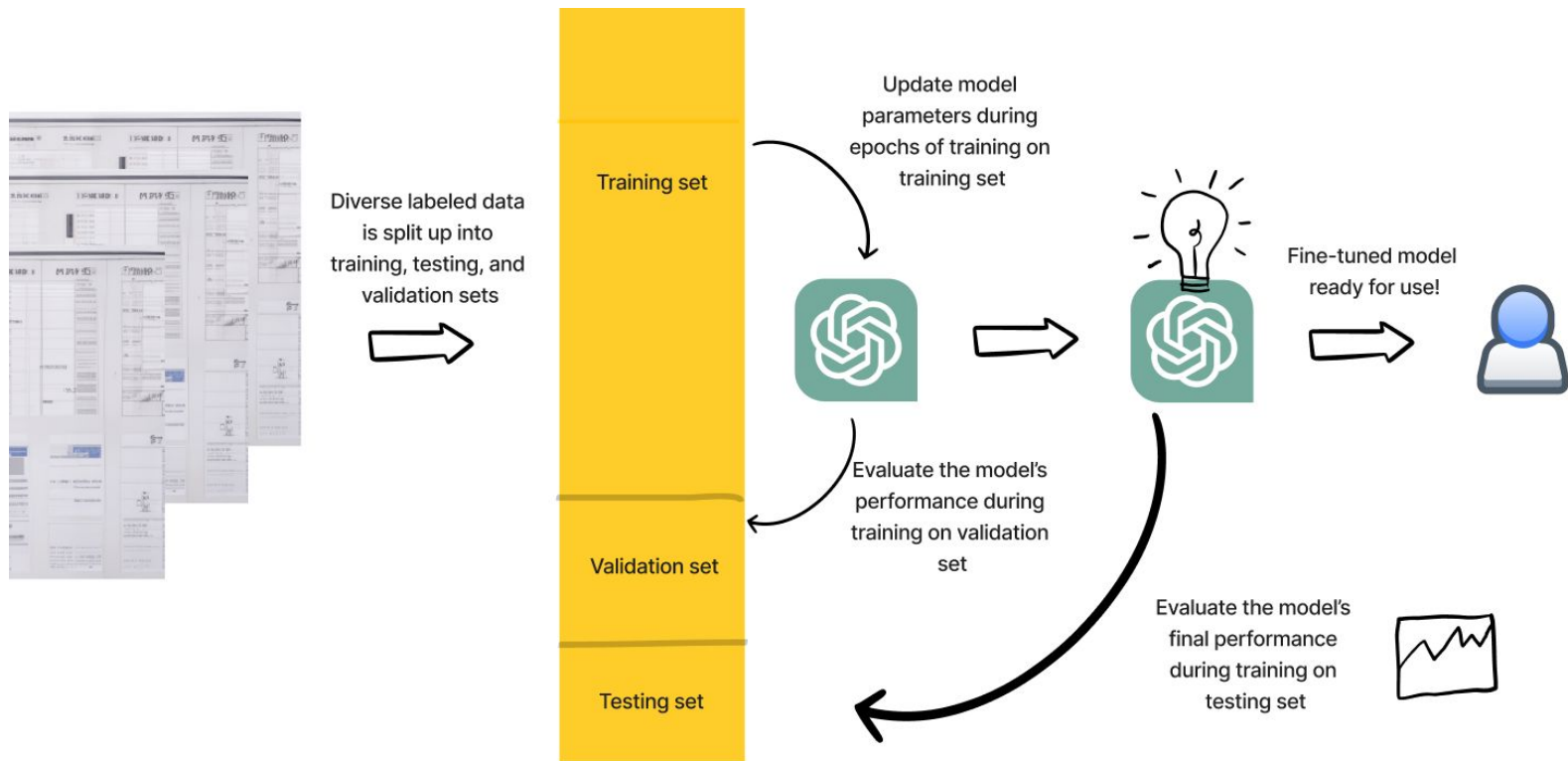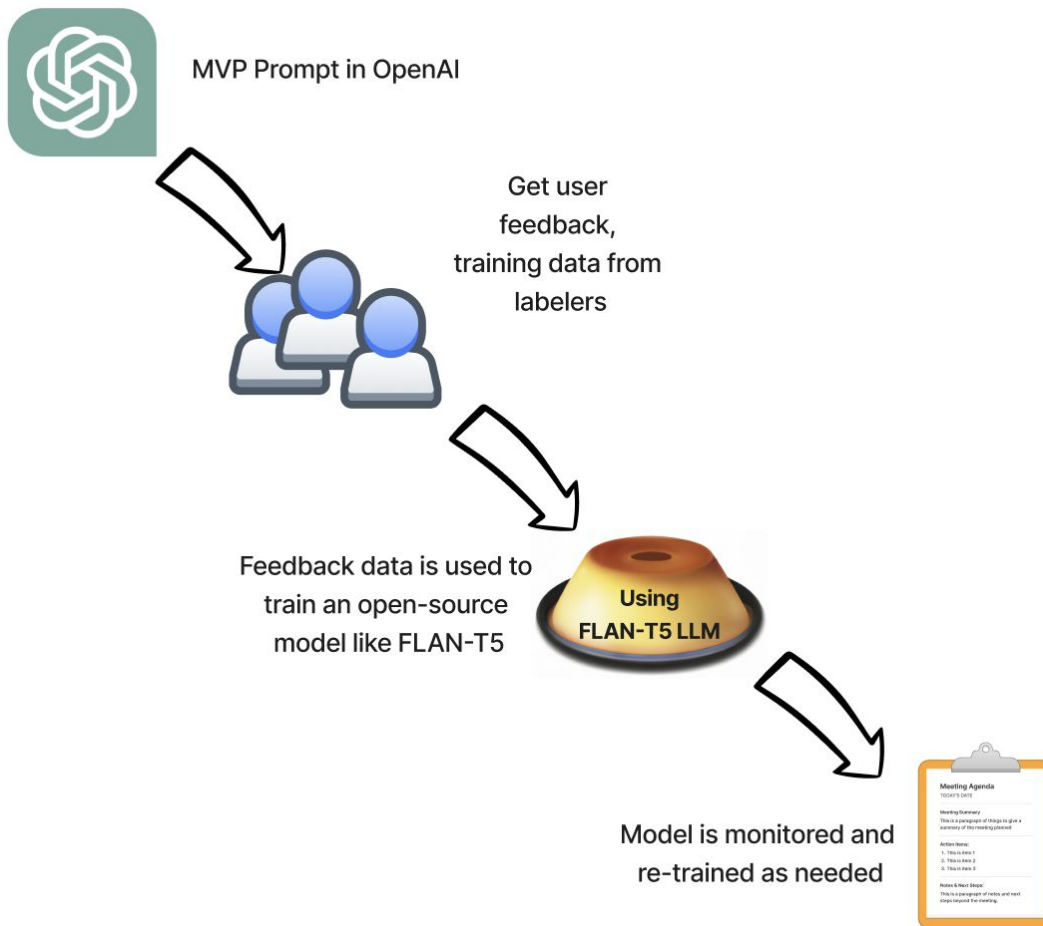
Assuming only 1,000 classifications a day and a relatively liberal prompt ratio (150 tokens (for few-shot examples, instructions, etc) for DaVinci or ChatGPT for every 40 tokens), **the cost of a fine-tuned model almost always wins the day overall cost-wise**
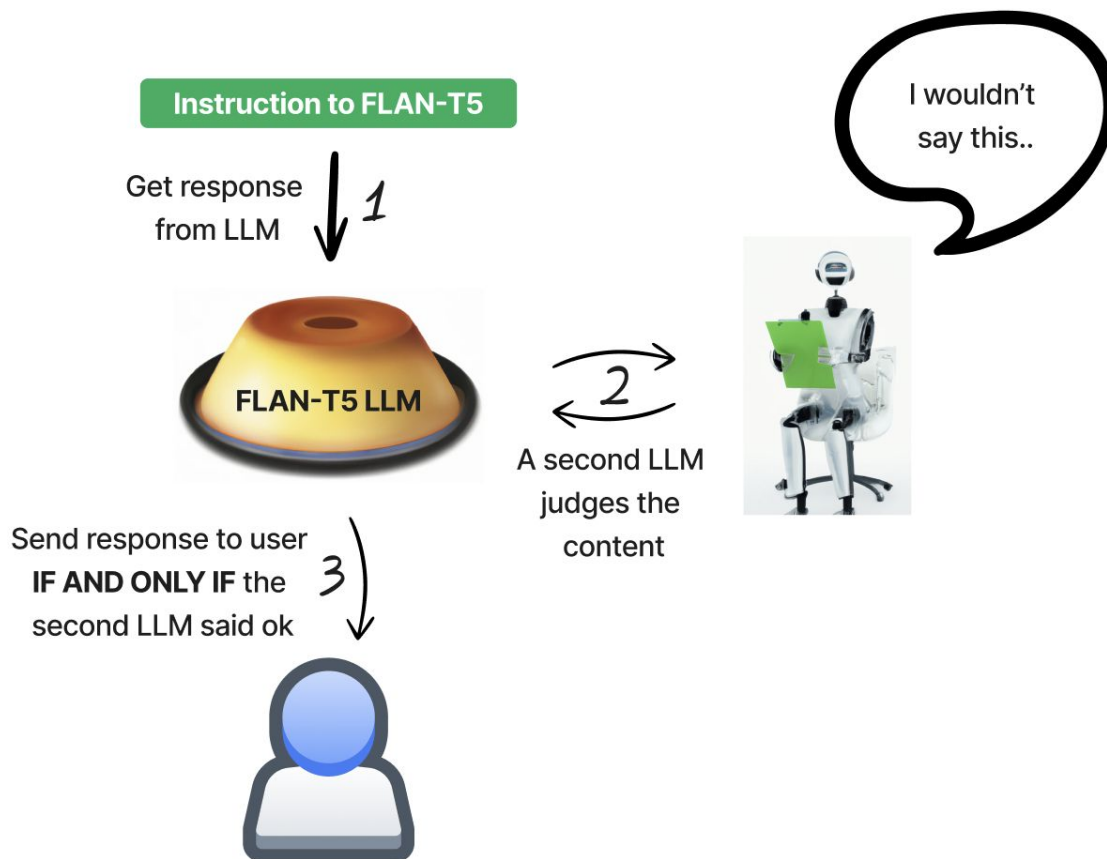
# Fine-tuning LLMs (e.g., OpenAI's Babbage)



Diverse labeled data is split up into training, testing, and validation sets

Training set

Validation set

Testing set

Update model parameters during epochs of training on training set

Evaluate the model's performance during training on validation set

Evaluate the model's final performance during training on testing set

Fine-tuned model ready for use!

# Moving from closed to open source

MVP Prompt in OpenAI

Get user feedback, training data from labelers

Feedback data is used to train an open-source model like FLAN-T5

Using FLAN-T5 LLM

Model is monitored and re-trained as needed

Pearson

- Think of LLMs as "**reasoning machines**" vs "thinking machines".

- LLMs excel at tasks that require **reasoning** - using context and input information in conjunction to produce a nuanced answer

Pearson

# Output Validation

# Output Structuring

Using an LLM's ability to reason comes in handy when we want it to output something in a more usable format.

Developers often prefer a *JSON* formatted output and we can ask an LLM to output things in a format if we ask it to or provide few-shot examples.

**USER**   Write a haiku as a JSON like this:

{"haiku": "(the haiku goes here)"}

**ASSISTANT**
```
{
 "haiku": "Autumn leaves falling\nWhispering secrets to earth\nNature's poetry"
}
```

Pearson

# Batch Prompting to save on latency/cost

## Standard Prompting

```
# K-shot in-context exemplars
Q: {question}
A: {answer}

Q: {question}
A: {answer}

...
# One sample to inference
Q: Ali had $21. Leila gave him half of her
   $100. How much does Ali have now?
---------------------------------------------
# Response
A: Leila gave 100/2=50 to Ali. Ali now has
   $21+$50 = $71. The answer is 71.
```

## Batch Prompting

```
# K-shot in-context exemplars in K/b batches
Q[1]: {question}
Q[2]: {question}           b(=2) samples
A[1]: {answer}             in one batch
A[2]: {answer}

...
# b samples in a batch to inference
Q[1]: Ali had $21. Leila gave him half of her
      $100. How much does Ali have now?
Q[2]: A robe takes 2 bolts of blue fiber and
      half that white fiber. How many bolts?
---------------------------------------------
# Responses to a batch
A[1]: Leila gave 100/2=50 to Ali. Ali now has
      $21+$50 = $71. The answer is 71.
A[2]: It takes 2/2=1 bolt of white fiber. The
      total amount is 2+1=3. The answer is 3.
```
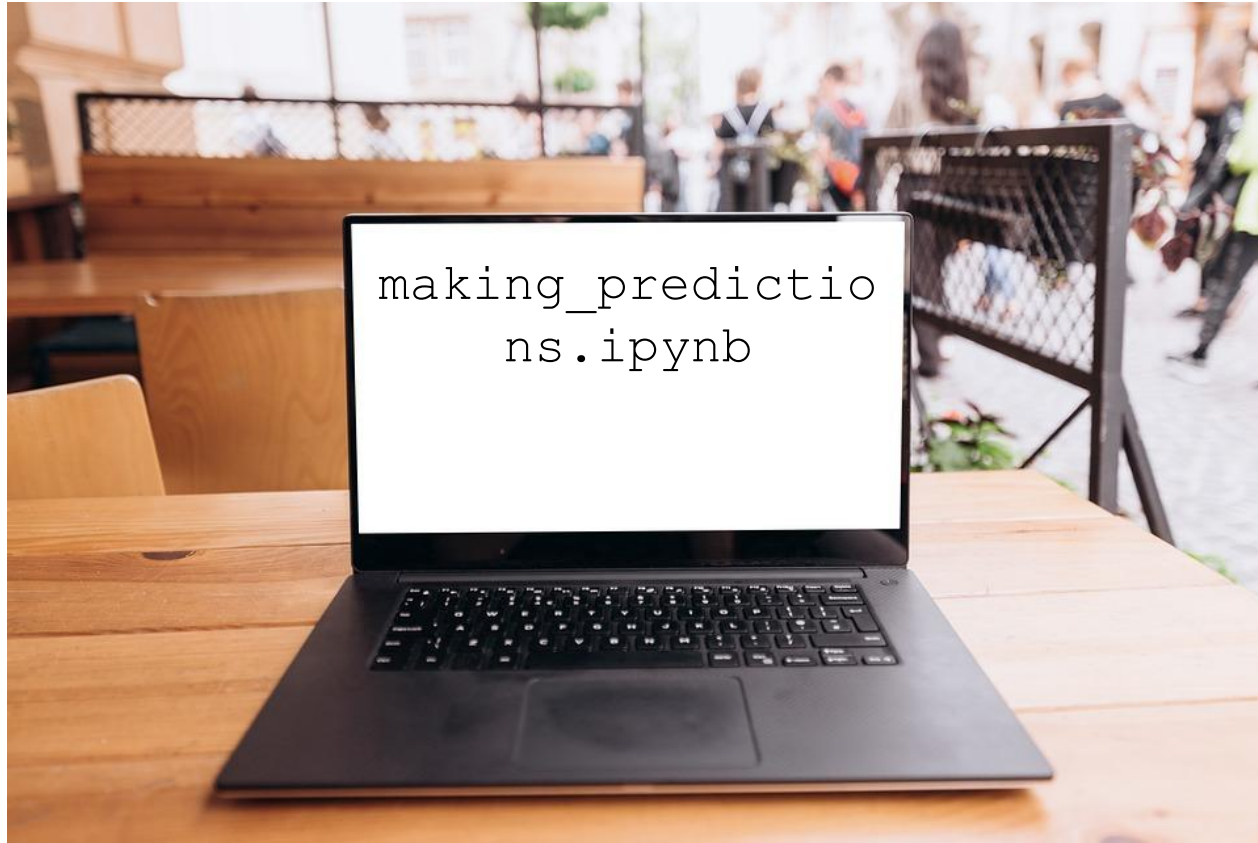
Pearson

# Code Time!



```
making_predictio
ns.ipynb
```

# Cost Projecting with LLMs and GPT

# Pricing with LLMs

Models like OpenAI GPT-3.5 Turbo (ChatGPT) charge per tokens inputted and tokens outputted.

Fine-tuned Classifiers have a cost for fine-tuning, inference (using them) and updating them

E.g., Hosting on HuggingFace for a small production-ready classifier starts at **$45/month**

Pearson

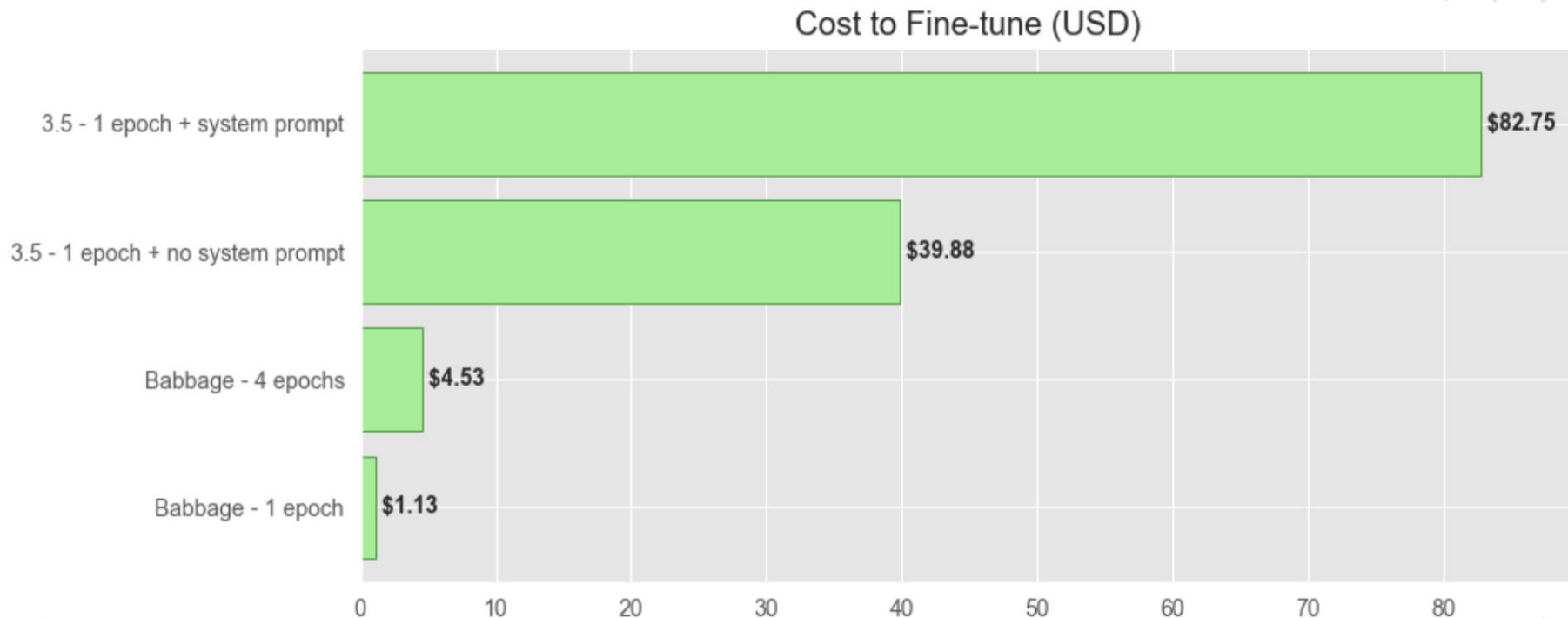# OpenAI Pricing

## GPT 3.5 Turbo (ChatGPT) (4K context window)

| Model | Input | Output |
|---|---|---|
| gpt-3.5-turbo-1106 | $0.0010 / 1K tokens | $0.0020 / 1K tokens |
| gpt-3.5-turbo-instruct | $0.0015 / 1K tokens | $0.0020 / 1K tokens |

## Fine-tuned Models

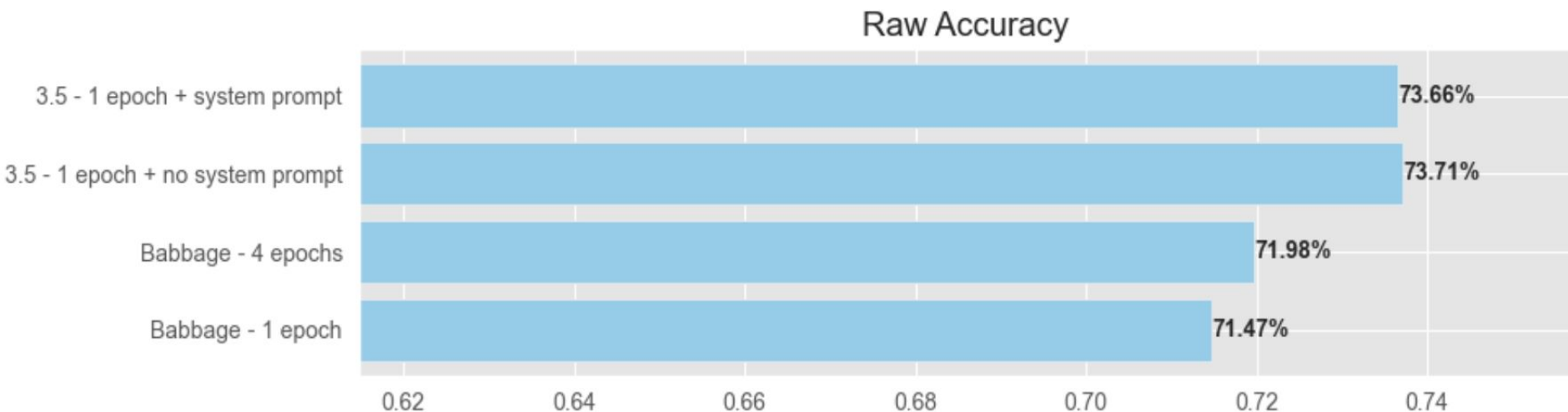| Model | Training | Input usage | Output usage |
|---|---|---|---|
| gpt-3.5-turbo | $0.0080 / 1K tokens | $0.0030 / 1K tokens | $0.0060 / 1K tokens |
| davinci-002 | $0.0060 / 1K tokens | $0.0120 / 1K tokens | $0.0120 / 1K tokens |
| babbage-002 | $0.0004 / 1K tokens | $0.0016 / 1K tokens | $0.0016 / 1K tokens |

# Cost of fine-tuning models

Classification with ~170K training examples



Cost to Fine-tune (USD)

| Model | Cost |
| --- | --- |
| 3.5 - 1 epoch + system prompt | $82.75 |
| 3.5 - 1 epoch + no system prompt | $39.88 |
| Babbage - 4 epochs | $4.53 |
| Babbage - 1 epoch | $1.13 |

# Cost of fine-tuning models

Accuracy bump probably not worth the cost

## Raw Accuracy

| Model | Accuracy |
|---|---|
| 3.5 - 1 epoch + system prompt | 73.66% |
| 3.5 - 1 epoch + no system prompt | 73.71% |
| Babbage - 4 epochs | 71.98% |
| Babbage - 1 epoch | 71.47% |

Pearson

# Cost Considerations for LLM applications

Open Source:

- **-** Data Collection
  - - Labelling, etc
- **-** Fine-tuning costs
  - - Machines, etc
- **-** Model Serving
  - - Machines, etc
- **-** Maintenance
  - - Future fine-tuning, etc

Closed Source

- **-** Number of tokens

Pearson

# Open vs Closed Source

Open Source:

- Pricing is under your control and generally cheaper

- Models are narrower but often more performant

Closed Source

- Easier to use, no need to think of hosting

- Often more expensive in the long term

**Pearson**

# Open vs Closed Source (deeper)

Open Source:

- Data privacy / security is controllable with on-premises systems

Closed Source

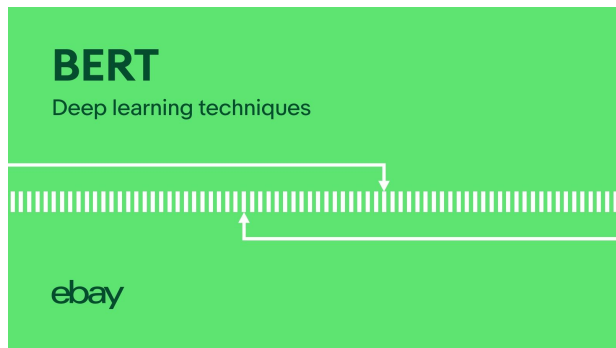- Companies control what parameters you can use (e.g. top-k is unavailable with OpenAI as are probabilities for tokens)

Pearson

# Code Time!



cost_projection.ipynb

Pearson

# Use Case Discussion

# Encoding Ebay's Recommendations with BERT

Ebay uses BERT to generate more relevant recommendations than traditional search techniques

BERT
Deep learning techniques

ebay

Pearson

# Visual Q/A with open source models



"What does the sign say?"

**Image Processor**
(Vision Transformer)

**Text Processor**
(DistilBERT)

GPT-2

STOP

Source: Quick start guide to LLMs

# Visual Q/A as a service

# Code Completion - Github's Copilot



```
File   Edit   Selection   View   Go   Run   Terminal   Help

{} en.json                    def sortByKey(key, array)  Untitled-1  ●

  1      def sortByKey(key, array) :
              for i in range(0, len(array)) :
                  for j in range(i, len(array)) :
                      if array[i][key] < array[j][key] :
                          array[i], array[j] = array[j], array[i]
              return array

  2
```

Source: Github Copilot

Pearson

# Legal - Klarity AI

Contract Amount ⚠

$1,000,000

...ber

...ber

New/Upsell/Renewal
**Upsell**

Governing Master Agreement
**MSA.docx**

Original Agreement Reviewed
**Yes**

**Service Credits**
- For any calendar month the Service Level is not met, if Customer has fulfilled all of its obligations under the Agreement and the Service Level Agreement, that month may be eligible for a Service Credit. The Service Credit will be calculated in accordance with the table below and must be used within tweClOe months of issuance. In no event will a refund be given.

- In the event that the System Availability falls below the percentages in the chart in this section in any given month, the Service Level for that month will be considered unmet.

| System Availability | Service Credit Eligibility |
|---|---|
| 99.95% or above | No Service Credit |
| 99.0% or above but below 99.95% | 10% of the pro-rated monthly fee paid |
| 95% or above but below 99.0% | 25% of the pro-rated monthly fee paid |

**Service Credits**
- For any calendar month the Service Level is not met, if Customer has fulfilled all of its obligations under the Agreement and the Service Level Agreement, that month may be eligible for a Service Credit. The Service Credit will be calculated in accordance with the table below and must be used within tweClOe months of issuance. In no event will a refund be given.

Source: Github Copilot

**Pearson**

# Code Time!



use_cases.ipynb

Pearson

# Week 2 Assignment

## For non-coders

Estimate how much it will cost to run a single instance of your task on the model (assuming you will use OpenAI 3.5 or 4)

1.  Think about how many tokens you are inputting and outputting
    a.   Use [OpenAI's tokenizer](#) to be exact
2.  Get a range of cost from what you might expect to be a shorter input vs a longer input (short vs long news article for a summarizer)

## For coders

Do the assignment for non-coders **AND** write a python function that performs the task and uses your MVP prompt. Your function will likely have at least one non-optional parameter that is your input and have hard-coded context (likely going in your system prompt). Your function will likely help you with the non-coder assignment

Bonus points if you do this for another closed source model

GPT-3   Codex

```
Write a haiku as a JSON like this:

{"haiku": "(the haiku goes here)"}
```
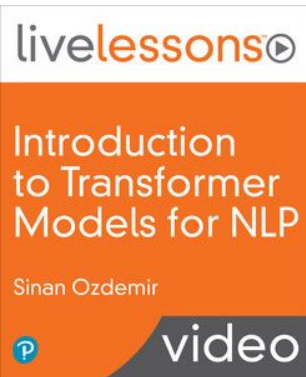
Clear    Show example

**Tokens**
24

**Characters**
70

Pearson

# Summary + Next Steps

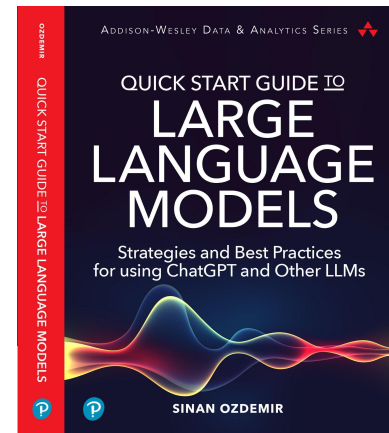A comprehensive introduction to LLMs + Transformers

https://learning.oreilly.com/videos/introduction-to-transformer/9780137923717

Check out my live trainings for more in depth content!

https://learning.oreilly.com/search/?q=Sinan%20Ozdemir&type=live-event-series

**New quick start guide to LLMs!**

Quick Start Guide to Large Language Models

# Large Language Models and ChatGPT in 3 Weeks

See you next week!

**Sinan Ozdemir**
Data Scientist, Entrepreneur,
Author, Lecturer