

Beyond ChatGPT:
Large Language Models and
Prompt Engineering

July 2023

Varun Shenoy

**ChatGPT is the fastest growing
consumer product ever...**

Time to 1 Million Users

Source: Collaborative Fund

5 days

ChatGPT



2.5 months

Instagram



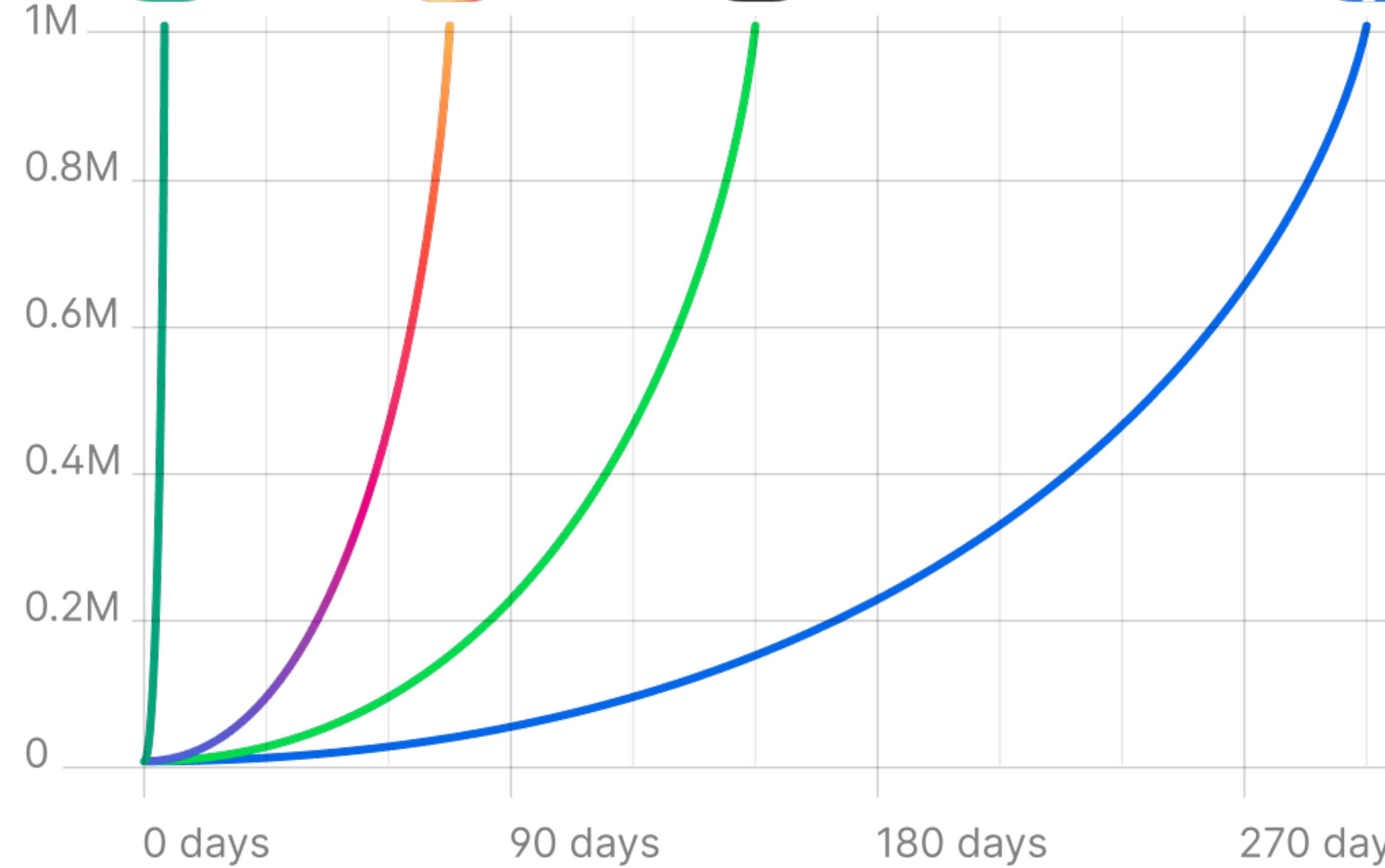
5 months

Spotify



10 months

Facebook



Me

Varun Shenoy @varunshenoy_ · Jan 31
Can LLMs extract knowledge graphs from unstructured text?

baseten Model library Documentation Customers Pricing Sign in Sign up

Varun Shenoy Meet Auralby It's an audio Audio creation touching a c A quick dem @gilbert) 🎉

Machine learning infrastructure that just works

Baseten provides all the infrastructure you need to deploy and serve ML models performantly, scalably, and cost-efficiently.

Sign up

Trusted by top engineering and machine learning teams

PATREON WRITER stability.ai rime

pipe motive predictionGuard laurel

0:51 | 2,943 10 1:59 9,171 views

ency between your

...
s from

Professor Hanks

Agenda

- Large Language Models
- Prompt Engineering
- Building a Better ChatGPT

Goal: get a high-level understanding of LLMs and get some hands on experience

From Zero to ChatGPT

Understanding Large Language Models

**LLMs complete sentences
word by word**

LLMs are trained on basically the entire internet and get good at guessing the next word.

This doesn't make them useful...

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

LLMs are “stochastic parrots”

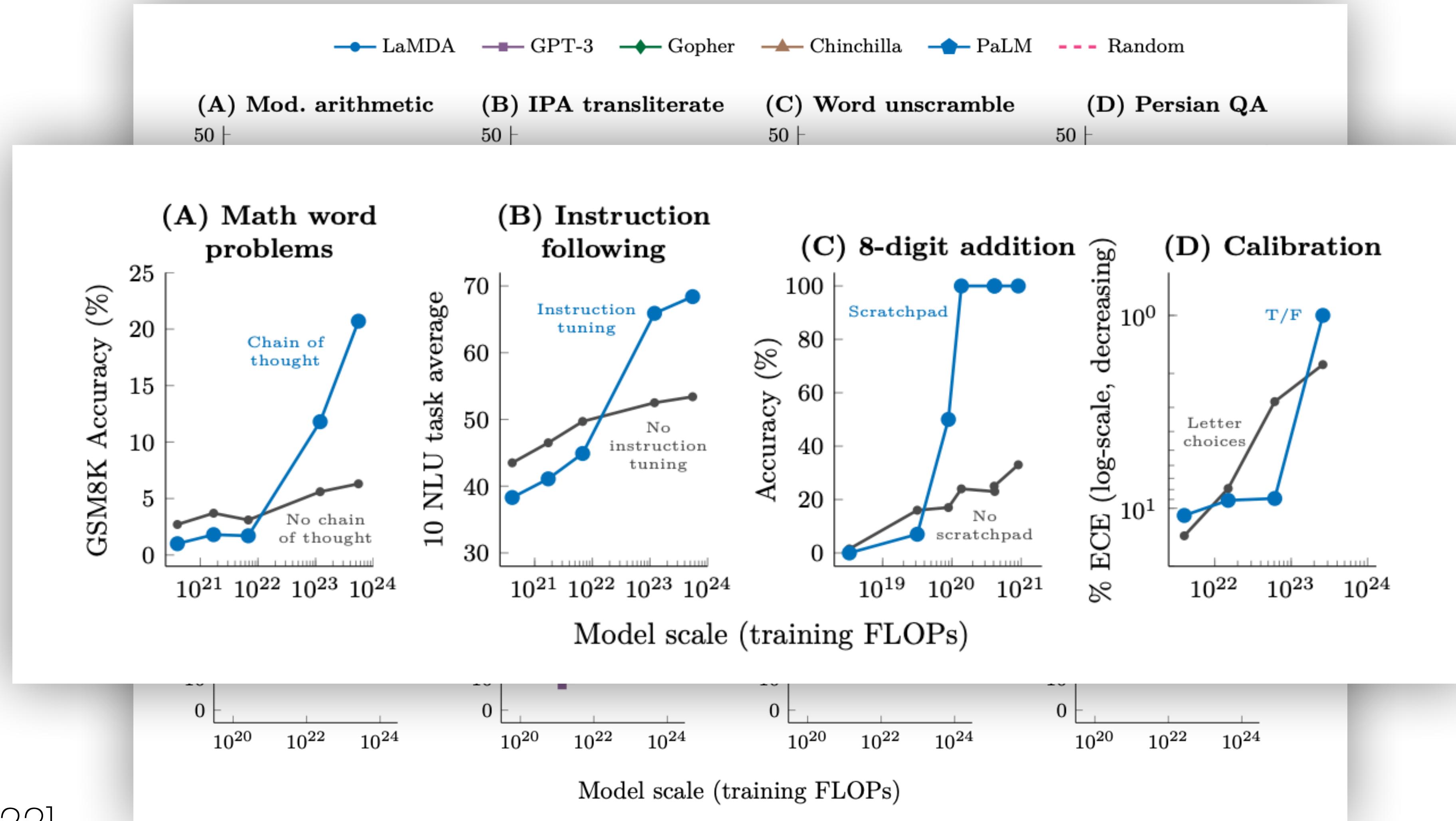
[Ouyang et al., 2022]

**ChatGPT = LLM + more parameters
+ instruction tuning + RLHF**

Q: Why are LLMs massive?

A: Emergent abilities

Emergent abilities



We can also fine-tune models
on human instructions.

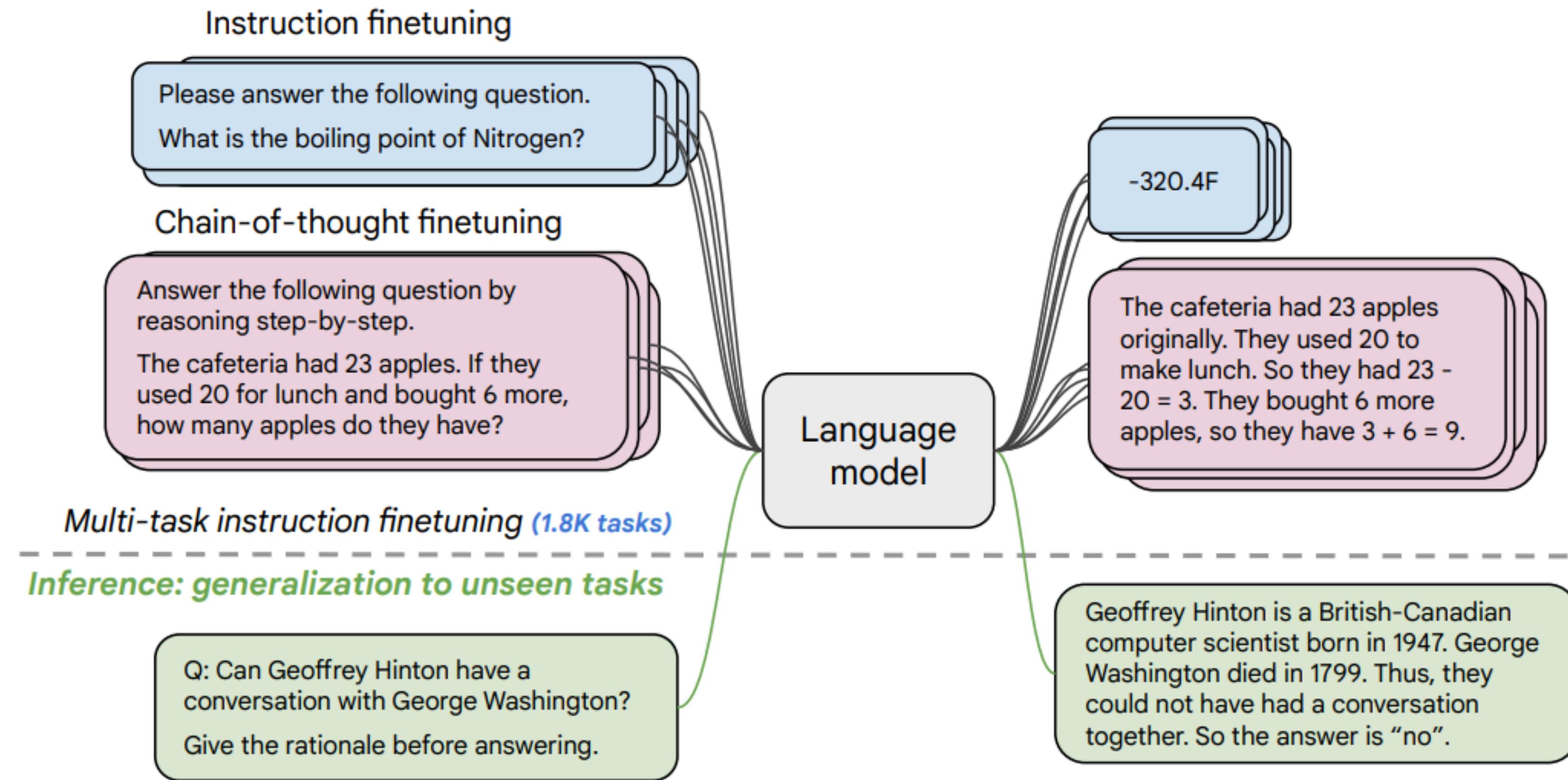


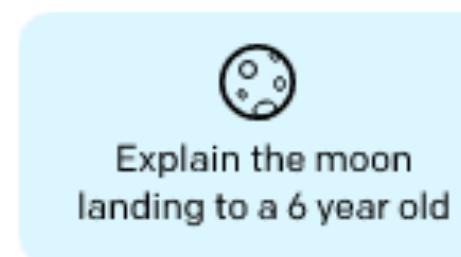
Figure 1: We finetune various language models on 1.8K tasks phrased as instructions, and evaluate them on unseen tasks. We finetune both with and without exemplars (i.e., zero-shot and few-shot) and with and without chain-of-thought, enabling generalization across a range of evaluation scenarios.

RLHF

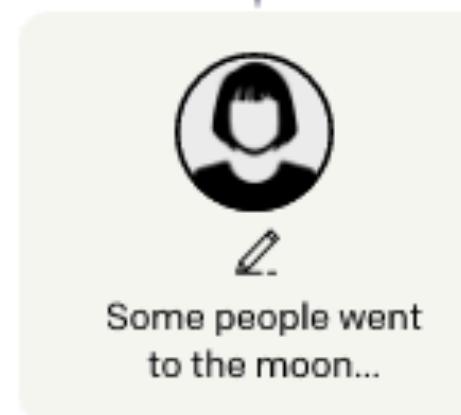
Step 1

Collect demonstration data, and train a supervised policy.

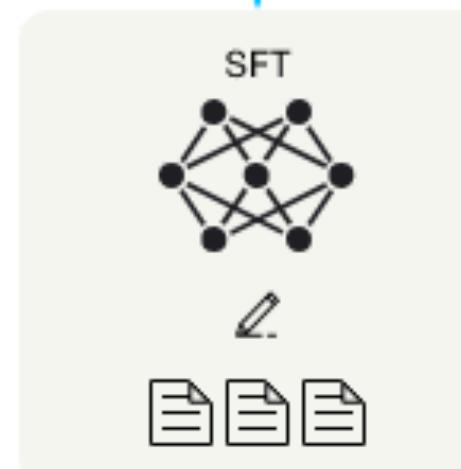
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



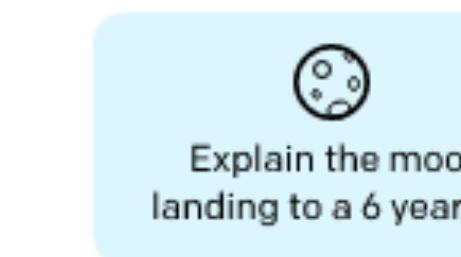
This data is used to fine-tune GPT-3 with supervised learning.



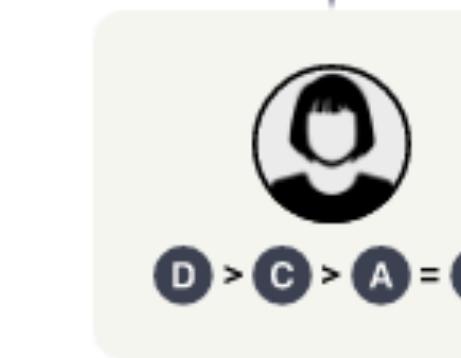
Step 2

Collect comparison data, and train a reward model.

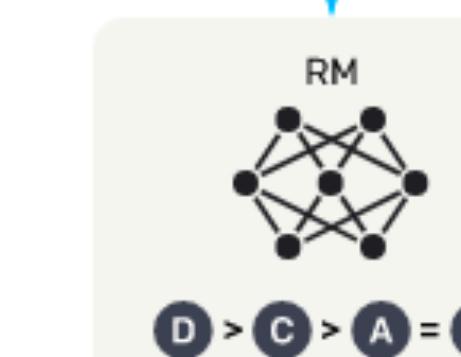
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



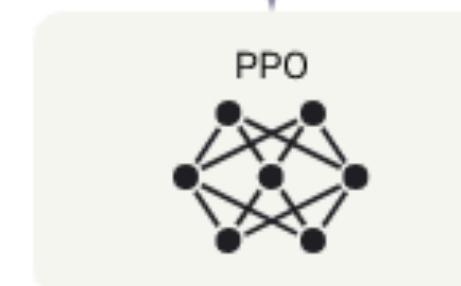
Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



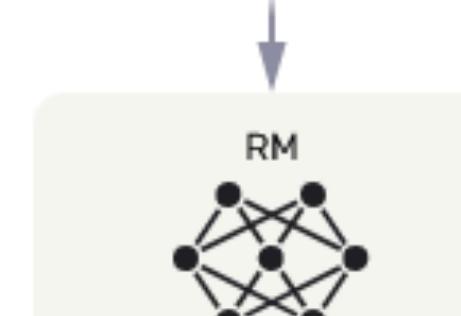
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



YN

What are all the programming languages you know?



As a large language model trained by OpenAI, I am not capable of programming in any language. My function is to assist users in generating human-like text based on the input provided to me. I do not have access to the internet and my knowledge is based on the text that I have been trained on, which has a fixed cutoff date. I do not have the ability to learn or acquire new information. Is there something else I can help you with?

Like Dislike

YN

What languages do you know?



As a large language model trained by OpenAI, I do not possess the ability to speak or understand any languages. My function is to assist users in generating human-like text based on the input provided to me. I do not have access to the internet and my knowledge is based on the text that I have been trained on, which has a fixed cutoff date. I do not have the ability to learn or acquire new information. Is there something else I can help you with?

Like Dislike

LLMs are forgetful :(

Limited context window size (4k tokens for ChatGPT)

A property of Transformers, the standard architecture
used for LLMs

Prompt Engineering

A step-by-step comprehensive guide to harnessing LLMs

Dawn of the Prompt Engineer?

The screenshot shows a Wikipedia article titled "Prompt engineering" on the Anthropic website. The page has a dark header with the word "ANTHROPIC" in white. The main content area features a large image of a person's face in a circular frame on the left, and a sidebar on the left with text about competition, anthropology, and salary. The main text area starts with a definition of prompt engineering as a concept in artificial intelligence, mentioning natural language processing (NLP) and how it embeds task descriptions in input. It also discusses prompt-based learning and frozen pretrained language models. The sidebar on the left has sections for "Competition", "Anthropology", "Committee", and "Salary".

Prompt engineering

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 embedded in the input, e.g., as a question instead of it being implicitly given. Prompt engineering typically works by converting one or more tasks to a prompt-based dataset and training a [language model](#) with what has been called "prompt-based learning" or just "prompt learning".^{[1][2]} Prompt engineering may work from a large "frozen" pretrained language model and where only the representation of the prompt is learned (i.e., optimized), using methods such as "prefix-tuning" or "prompt tuning".^{[3][4]}

The [GPT-2](#) and [GPT-3](#) language models^[5] were important steps in prompt engineering. In 2021, multitask prompt engineering using multiple NLP datasets showed good performance on new tasks.^[6] Prompts that include a [chain of thought](#) have been shown to improve the performance of large language models.

Beginner Prompting

Exploring the basic capabilities of GPT-3

Prompt?

- Guide a language model by providing context.
- What should it do?
- How should it respond?

Play along: <https://platform.openai.com/playground>

Basic Prompt

Python is

Basic Prompt

Python is an interpreted, high-level, general-purpose programming language created by Guido van Rossum in 1991. Python is a very popular language because of its easy-to-learn syntax, wide range of libraries and frameworks, and powerful data analysis capabilities. Python is often used for scripting, automation, web development, machine learning, and data science.

Text Classification

Rate the sentiment of the following sentence as positive or negative.

Sentence: I can say this is the best sushi I've had at this price point. The amount of fish per roll is always more than expected, and the fish tastes again of good quality.

Sentiment:

Text Classification

Rate the sentiment of the following sentence as positive or negative.

Sentence: I can say this is the best sushi I've had at this price point. The amount of fish per roll is always more than expected, and the fish tastes again of good quality.

Sentiment: Positive

Text Summarization

Summarize the following text.

Context: Everyone knows that to do great work you need both natural ability and determination. But there's a third ingredient that's not as well understood: an obsessive interest in a particular topic. [...] One of the most striking features of Darwin's book about his voyage on the Beagle is the sheer depth of his interest in natural history. His curiosity seems infinite. Ditto for Ramanujan, sitting by the hour working out on his slate what happens to series.

Summary:

Text Summarization

Summarize the following text.

Context: Everyone knows that to do great work you need both natural ability and determination. But there's a third ingredient that's not as well understood: an obsessive interest in a particular topic. [...] One of the most striking features of Darwin's book about his voyage on the Beagle is the sheer depth of his interest in natural history. His curiosity seems infinite. Ditto for Ramanujan, sitting by the hour working out on his slate what happens to series.

Summary: Obsessive interest in a particular topic is a necessary ingredient for great work, and can be seen in the lives of people who have done great work in the past. This interest is often disinterested, and is not done to impress or make money. Examples of this can be seen in Darwin's passion for natural history, and Ramanujan's focus on mathematical series.

Question Answering

Answer a question based on the context.

Context: Everyone knows that to do great work you need both natural ability and determination. But there's a third ingredient that's not as well understood: an obsessive interest in a particular topic. [...] One of the most striking features of Darwin's book about his voyage on the Beagle is the sheer depth of his interest in natural history. His curiosity seems infinite. Ditto for Ramanujan, sitting by the hour working out on his slate what happens to series.

Question: Why is an example of bus ticket collectors useful?

Answer:

Question Answering

Answer a question based on the context.

Context: Everyone knows that to do great work you need both natural ability and determination. But there's a third ingredient that's not as well understood: an obsessive interest in a particular topic. [...] One of the most striking features of Darwin's book about his voyage on the Beagle is the sheer depth of his interest in natural history. His curiosity seems infinite. Ditto for Ramanujan, sitting by the hour working out on his slate what happens to series.

Question: Why is an example of bus ticket collectors useful?

Answer: An example of bus ticket collectors is useful to explain the third ingredient to do great work: an obsessive interest in a particular topic. Bus ticket collectors have an obsessive interest in the minutiae of what they collect and their love is disinterested, providing an example of the kind of dedication needed to do great work.

Deconstructing a Prompt

- Instructions
- Context
- User Input
- Output Indicator

Code Generation

Create a Python function that uses FastAPI and requests to make a query to <https://api.sampleapis.com/coffee/hot>. Return the description of a cappuccino.

Code:

Code Generation

Create a Python function that uses FastAPI and requests to make a query to <https://api.sampleapis.com/coffee/hot>. Return the description of a cappuccino.

Code:

```
import requests  
from fastapi import FastAPI  
  
app = FastAPI()
```

```
@app.get("/cappuccino")  
def get_cappuccino():  
    response = requests.get("https://api.sampleapis.com/coffee/hot")  
    data = response.json()  
    for item in data["items"]:  
        if item["name"] == "Cappuccino":  
            return item["description"]  
    return "Cappuccino not found"
```

Not quite right, see if you can improve the prompt :)

Advanced Prompting

Do Large Language Models Dream of Electric Sheep? 😴

Few-shot Prompting

The screenshot shows a user interface for generating programmatic commands from text. At the top, there's a green icon with a white right-pointing arrow and the text "Text to command". Below it are two buttons: "Transformation" and "Generation".

The main area has a heading "Translate text into programmatic commands." followed by a "Prompt" section and a "Settings" section.

Prompt:

Convert this text to a programmatic command:
Example: Ask Constance if we need some bread
Output: send-msg `find constance` Do we need some bread?
Reach out to the ski store and figure out if I can get my skis fixed before I leave on Thursday

Settings:

Engine	text-davinci-003
Max tokens	100
Temperature	0
Top p	1.0
Frequency penalty	0.2
Presence penalty	0.0
Stop sequence	\n

Sample response:

send-msg `find ski store` Can I get my skis fixed before I leave on Thursday?

API request:

Parse unstructured data Classification

Ingestable JSON Output

Provide a summary of the given text and a short poem. Output the result as a JSON object with two keys, "summary", and "poem".

Context: Everyone knows that to do great work you need both natural ability and determination. But there's a third ingredient that's not as well understood: an obsessive interest in a particular topic. [...] One of the most striking features of Darwin's book about his voyage on the Beagle is the sheer depth of his interest in natural history. His curiosity seems infinite. Ditto for Ramanujan, sitting by the hour working out on his slate what happens to series.

JSON:

Ingestable JSON Output

Provide a summary of the given text and a short poem. Output the result as a JSON object with two keys, "summary", and "poem".

Context: Everyone knows that to do great work you need both natural ability and determination. But there's a third ingredient that's not as well understood: an obsessive interest in a particular topic. [...] One of the most striking features of Darwin's book about his voyage on the Beagle is the sheer depth of his interest in natural history. His curiosity seems infinite. Ditto for Ramanujan, sitting by the hour working out on his slate what happens to series.

JSON:

```
{ "summary": "Natural ability and determination are necessary for great work, but an obsessive interest in a particular topic is also essential. Charles Darwin and Srinivasa Ramanujan are two examples of people who had an infinite curiosity for their chosen fields.",  
  "poem": "Though talent and dedication are key, \nA passionate mind is the missing piece  
  \nFor great minds like Darwin and Ramanujan \nTheir curiosity was never at peace" }
```

Chain of Thought

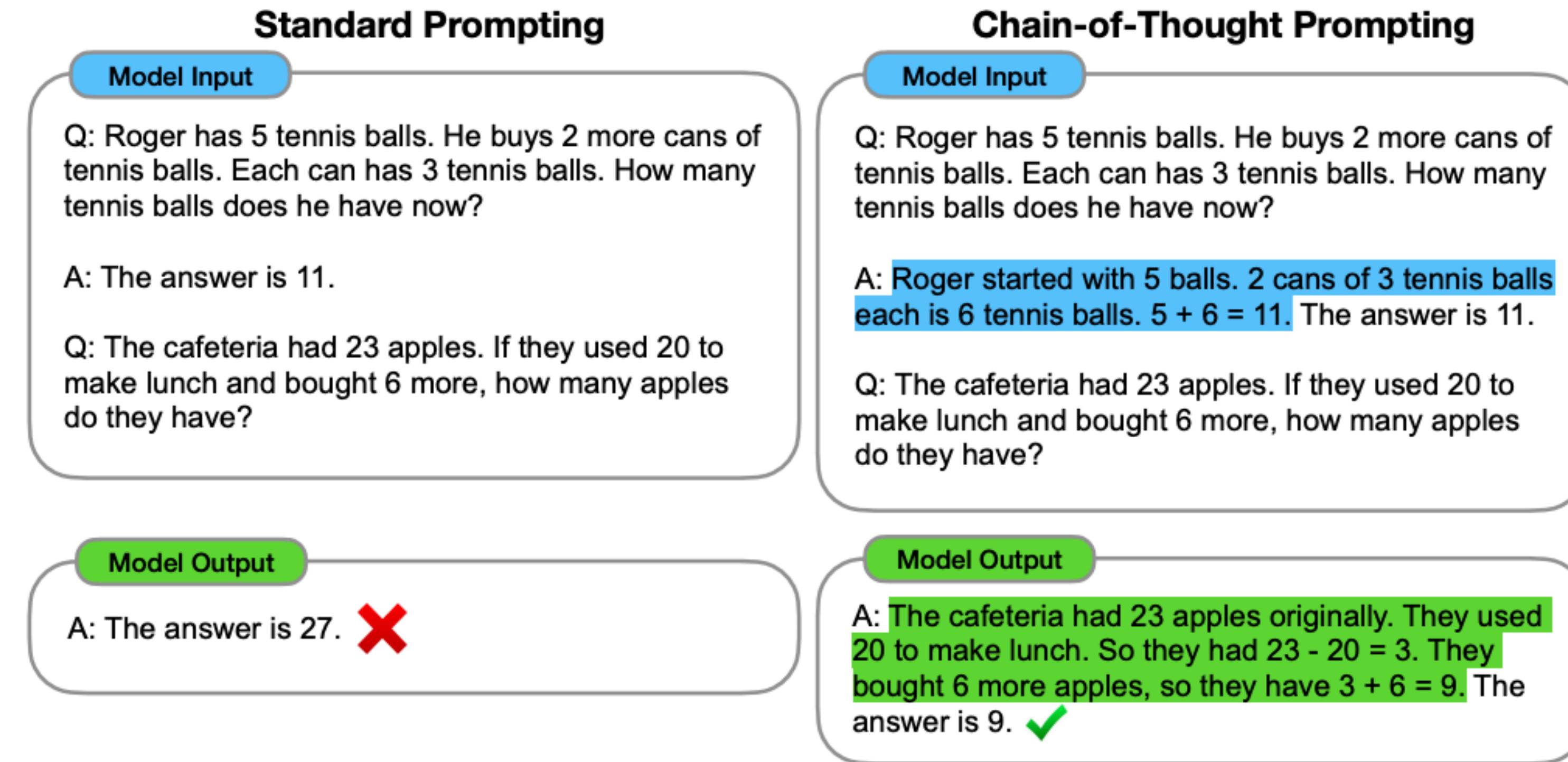


Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

Let's think step by step

(a) Few-shot

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: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. X

(b) Few-shot-CoT

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: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 X

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

Tools and Agents

```
tools = load_tools(["google-search"], llm=llm)

agent = initialize_agent(tools, llm, agent="zero-shot-react-description", verbose=True)

agent.run("What is the weather in Pomfret?")
```

> Entering new AgentExecutor chain...
I should look up the current weather conditions.
Action: Google Search
Action Input: "weather in Pomfret"
Observation: Showers early becoming a steady light rain later in the day. Near record high temperatures. High around 60F. Winds SW at 10 to 15 mph. Chance of rain 60%.
Pomfret, CT Weather Forecast, with current conditions, wind, air quality, and what to expect for the next 3 days. Hourly Weather-Pomfret, CT. As of 12:52 am EST. Special Weather Statement +2 ... Hazardous Weather Conditions. Special Weather Statement ...
Pomfret CT. Tonight ... National Digital Forecast Database Maximum Temperature Forecast. Pomfret Center Weather Forecasts. Weather Underground provides local & long-range weather forecasts, weatherreports, maps & tropical weather conditions for ... Pomfret, CT 12 hour by hour weather forecast includes precipitation, temperatures, sky conditions, rain chance, dew-point, relative humidity, wind direction ... North Pomfret Weather Forecasts. Weather Underground provides local & long-range weather forecasts, weatherreports, maps & tropical weather conditions for ... Today's Weather - Pomfret, CT. Dec 31, 2022 4:00 PM. Putnam MS. --. Weather forecast icon. Feels like --. Hi --. Lo --. Pomfret, CT temperature trend for the next 14 Days. Find daytime highs and nighttime lows from TheWeatherNetwork.com. Pomfret, MD Weather Forecast Date: 332 PM EST Wed Dec 28 2022. The area/counties/county of: Charles, including the cities of: St. Charles and Waldorf.
Thought: I now know the current weather conditions in Pomfret.
Final Answer: Showers early becoming a steady light rain later in the day. Near record high temperatures. High around 60F. Winds SW at 10 to 15 mph. Chance of rain 60%.
> Finished AgentExecutor chain.

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

Figure 1: Exemplary predictions of Toolformer. The model autonomously decides to call different APIs (from top to bottom: a question answering system, a calculator, a machine translation system, and a Wikipedia search engine) to obtain information that is useful for completing a piece of text.

Let's Build (a Better) ChatGPT

A chatbot grounded in truth with embeddings and GPT 🛠

... in under 30 lines of code 😱

Follow along at

[https://github.com/varunshenoy/
honestgpt/](https://github.com/varunshenoy/honestgpt/)

VN

What is the difference between a latte and cappuccino?



Latte and cappuccino are two popular coffee drinks that originated in Italy. While both drinks contain espresso and steamed milk, there are differences in their ratios and textures.



A latte typically has a higher ratio of steamed milk to espresso than a cappuccino. It usually consists of one shot of espresso (or sometimes two) and is topped with steamed milk and a thin layer of frothed milk on top. The steamed milk gives the drink a creamy texture and a slightly sweet taste.

On the other hand, a cappuccino has equal parts of espresso, steamed milk, and frothed milk. It is traditionally served in a smaller cup and has a thicker layer of frothed milk on top, which gives it a lighter, airier texture.

In summary, a latte has more steamed milk and less frothed milk than a cappuccino, while a cappuccino has equal parts of espresso, steamed milk, and frothed milk. The texture and taste of each drink differ slightly as a result.



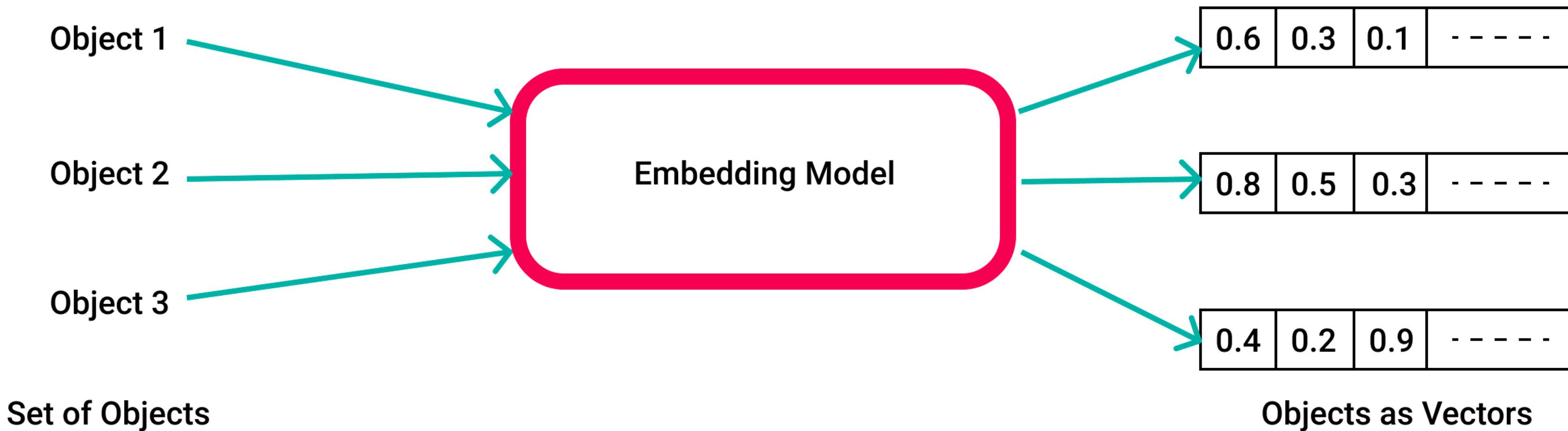
Embeddings



Embeddings

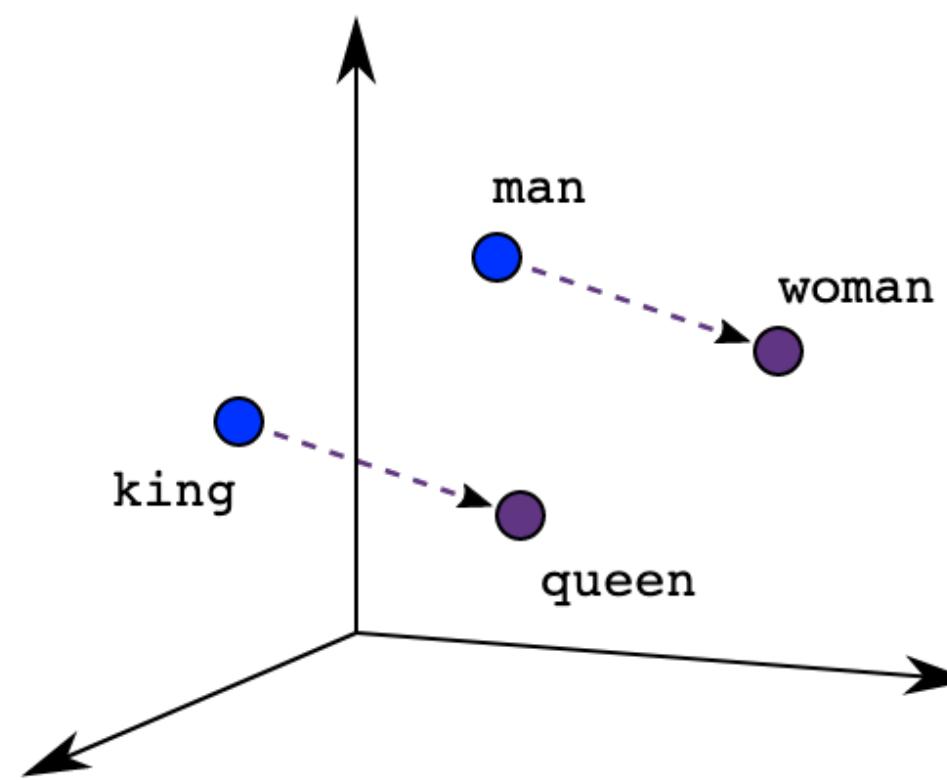


Embeddings

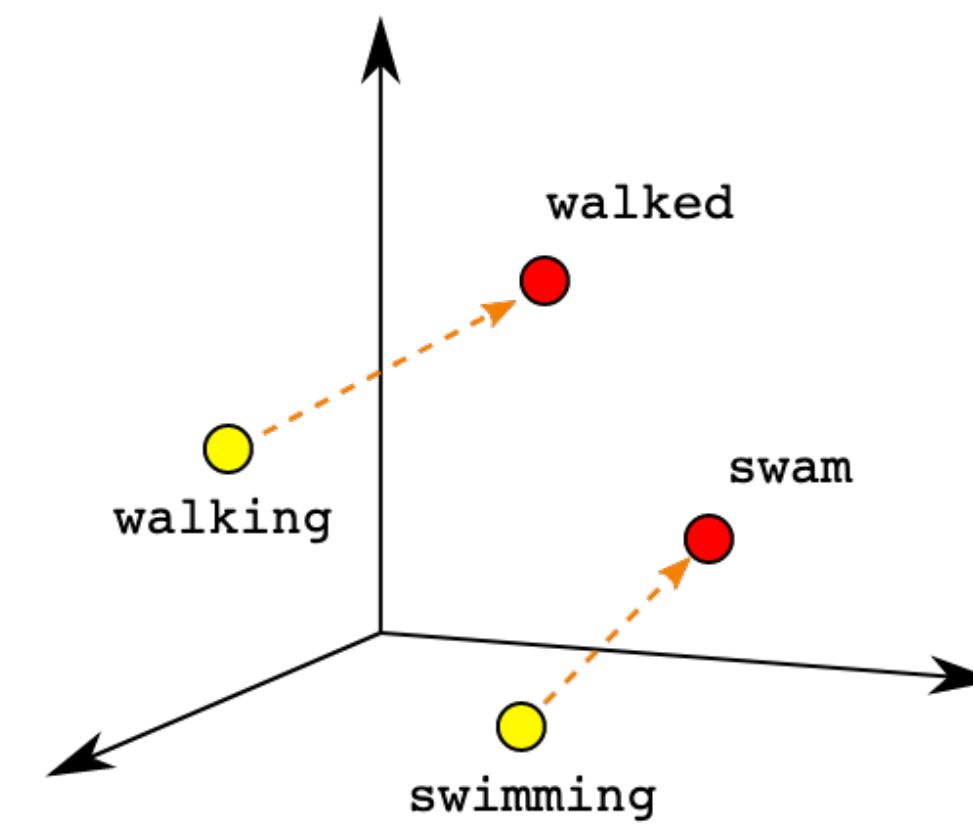


[What are Vector Embeddings, Pinecone]

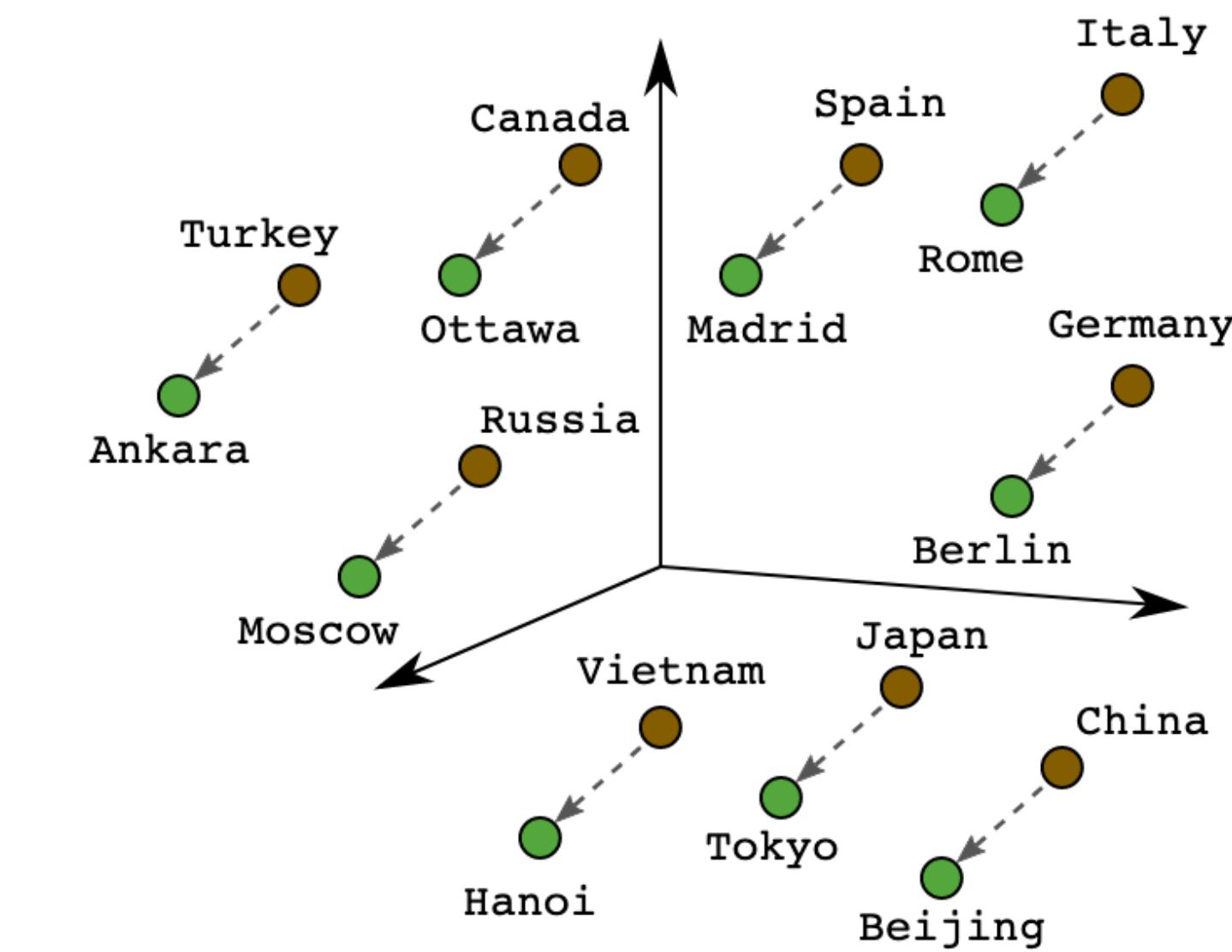
Embeddings



Male-Female



Verb Tense



Country-Capital

Coffee



coffee.txt



```
...  
  
def read_text_file(file_name, title):  
    with open(file_name, 'r') as file:  
        data = file.read()  
  
        # split into paragraphs  
        paragraphs = data.split('\n\n')  
  
        if title is not None:  
            for idx, paragraph in enumerate(paragraphs):  
                paragraphs[idx] = title + ': ' + paragraph  
  
    return paragraphs
```

● ● ●

```
def get_embeddings(paragraphs):

    # Load sentence transformers model
    model = SentenceTransformer(
        'sentence-transformers/all-MiniLM-L6-v2')

    # get embeddings
    embeddings = model.encode(paragraphs, show_progress_bar=True)

    # save embeddings
    np.save('embeddings.npy', embeddings)

return embeddings
```

Why Hugging Face?

- Free :)
- Lots of options
- Local embeddings

```
● ● ●

def get_similarity(embeddings, query):
    # load embeddings
    embeddings = np.load('embeddings.npy')

    # load sentence transformers model
    model = SentenceTransformer(
        'sentence-transformers/all-MiniLM-L6-v2')

    # get query embedding
    query_embedding = model.encode(query, show_progress_bar=True)
```

```
● ● ●

def get_similarity(embeddings, query):
    # load embeddings
    embeddings = np.load('embeddings.npy')

    # load sentence transformers model
    model = SentenceTransformer(
        'sentence-transformers/all-MiniLM-L6-v2')

    # get query embedding
    query_embedding = model.encode(query, show_progress_bar=True)

    # compute similarity using dot product
    similarity = np.dot(embeddings, query_embedding) / \
        (np.linalg.norm(embeddings, axis=1) *
    np.linalg.norm(query_embedding))

    # get top 5 results and their indices
    top_n_indices = np.argsort(-similarity)[:5]

return top_n_indices
```

```
def make_openai_call(context, question):
    sources = ''
    for idx, paragraph in enumerate(context):
        sources += f"Source {idx + 1}: {paragraph}\n"

    prompt = generate_prompt(sources, question)

    response = openai.ChatCompletion.create(
        model="gpt-3.5-turbo", # or gpt-4
        messages=[
            {"role": "system", "content": "You are HonestGPT. Make sure all your
answers cite the sources you used as in-text citations."},
            {"role": "user", "content": f"{prompt}"},
        ]
    )

    return response['choices'][0]['message']['content']
```

A quick aside on parameters:

- **Temperature [0, 1]:** higher temperature, more creativity
- **Top P [0, 1]:** high top-p, more creativity
- **Maximum Length [1, 4000]:**
 - token count = tokenized(prompt) + tokenized(output)

Good source for temperature: <https://lukesalamone.github.io/posts/what-is-temperature/>



```
def generate_prompt(sources, question):  
    return f""
```

Write a paragraph, addressing the question, and combine the text below to obtain relevant information. Cite sources using in-text citations with square brackets.

For example: [1] refers to source 1 and [2] refers to source 2. Cite once per sentence.

If the context doesn't answer the question. Output "I don't know".

{sources}

Question: {question}

Result:"""

● ● ●

```
# read text file
paragraphs = read_text_file('coffee.txt', None)
```

● ● ●

```
# read text file
paragraphs = read_text_file('coffee.txt', None)

# get embeddings
embeddings = get_embeddings(paragraphs)
```

● ● ●

```
# read text file
paragraphs = read_text_file('coffee.txt', None)

# get embeddings
embeddings = get_embeddings(paragraphs)

# get similarity
query = 'What is the difference between a cappuccino and a latte?'
top_n_indices = get_similarity(embeddings, query)
```

```
● ● ●

# read text file
paragraphs = read_text_file('coffee.txt', None)

# get embeddings
embeddings = get_embeddings(paragraphs)

# get similarity
query = 'What is the difference between a cappuccino and a latte?'
top_n_indices = get_similarity(embeddings, query)

top_results = []
for idx in top_n_indices:
    top_results.append(paragraphs[idx])
```

```
● ● ●

# read text file
paragraphs = read_text_file('coffee.txt', None)

# get embeddings
embeddings = get_embeddings(paragraphs)

# get similarity
query = 'What is the difference between a cappuccino and a latte?'
top_n_indices = get_similarity(embeddings, query)

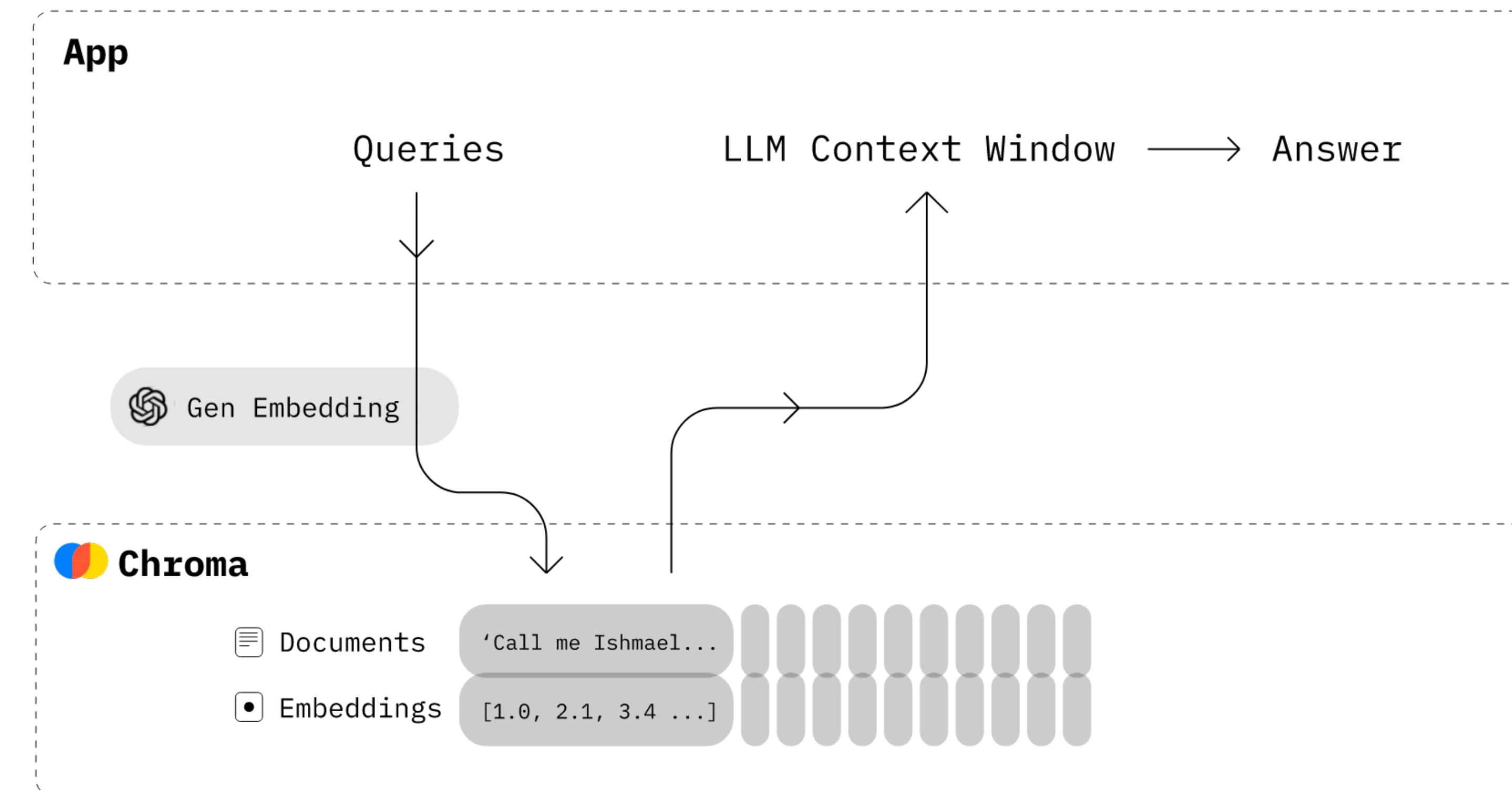
top_results = []
for idx in top_n_indices:
    top_results.append(paragraphs[idx])

print(make_openai_call(top_results, query).strip())
```

“The main difference between a cappuccino and a latte is the size and the proportion of espresso, steamed milk, and foamed milk. A cappuccino is generally a smaller drink, served with about 2-3 oz of steamed milk per shot of espresso, and the layer of foamed milk on top should be thicker than that of a cafe latte **[1]**. A latte combines espresso with steamed warm milk and a layer of foamed milk on top, while a cappuccino is a layered drink with 1-2 shots of espresso on the bottom, then equal amounts of steamed milk and foamed milk and possibly spices **[2, 3]**. An iced latte is similar to a caffe latte, but the milk is not steamed **[4, 5]**.”

What is the difference between a cappuccino and a latte?

Let's layer in abstractions...



trychroma.com

```
chroma_client = chromadb.Client()
collection = chroma_client.create_collection(name="my_collection")

def read_and_embed_file(file_name, title):
    with open(file_name, 'r') as file:
        data = file.read()

        # split into paragraphs
        paragraphs = data.split('\n\n')
        metadata_list = []
        ids_list = []

        for idx, paragraph in enumerate(paragraphs):
            metadata_list.append({"source": f"{title}"})
            ids_list.append(f"id{idx + 1}")

    collection.add(
        documents=paragraphs,
        metadatas=metadata_list,
        ids=ids_list
    )
```

```
● ● ●

# read text and embed file
read_and_embed_file('coffee.txt', "coffee-article")

# get similarity
query = 'What is the difference between a cappuccino and a latte?'
results = collection.query(
    query_texts=[query],
    n_results=5
)

top_results = results["documents"][0]

print(make_openai_call(top_results, query).strip())
```

“The main difference between a cappuccino and a latte is the amount of milk used and the thickness of the foam on top. A cappuccino is generally a smaller drink served with 2-3 oz of steamed milk per shot of espresso and a thicker layer of foamed milk on top than a latte. A latte combines espresso with steamed warm milk and a layer of foamed milk on top, while an iced latte combines espresso and cold milk over ice. **[1,2,3,4,5]**”

What is the difference between a cappuccino and a latte?

“Sources:

- 1: Unlike a latte, a cappuccino is generally a smaller drink, served with about 2-3 oz of steamed milk per shot of espresso. The layer of foamed milk on top should be thicker than that of a cafe latte.
- 2: Often topped with spices like nutmeg, cinnamon, or chocolate powder, think of a cappuccino as the latte's shorter, cozier cousin.
- 3: CAPPUCINO
A cappuccino is a layered drink with 1-2 shots of espresso on the bottom, then equal amounts of steamed milk and foamed milk (in that order).
- 4: LATTE
Lattes combine espresso with steamed warm milk and a layer of foamed milk on top. You can add any flavor such as mocha, peppermint, or pumpkin spice to this mixture for a completely different flavor profile.
- 5: ICED LATTE
An iced latte combines espresso and cold milk over ice. The only difference between an iced latte and a caffe latte is that the milk is not steamed.”

What is the difference between a cappuccino and a latte?

Let's layer in abstractions (again)...



langchain.dev

• • •

```
from langchain.llms import OpenAI
from langchain.chains.qa_with_sources import
load_qa_with_sources_chain
from langchain.text_splitter import CharacterTextSplitter
from langchain.embeddings import OpenAIEMBEDDINGS
from langchain.vectorstores import Chroma
from langchain.docstore.document import Document
```

```
from langchain.llms import OpenAI
from langchain.chains.qa_with_sources import
load_qa_with_sources_chain
from langchain.text_splitter import CharacterTextSplitter
from langchain.embeddings import OpenAIEmbeddings
from langchain.vectorstores import Chroma
from langchain.docstore.document import Document

# read text and embed file
with open('texts/coffee.txt') as f:
    coffee_text = f.read()
text_splitter = CharacterTextSplitter(chunk_size=1000,
chunk_overlap=0)
texts = text_splitter.split_text(coffee_text)

embeddings = OpenAIEmbeddings()

docsearch = Chroma.from_texts(texts, embeddings, metadatas=[{"source": str(i)} for i in
range(len(texts))])
```

```
chunk_overlap=0)
texts = text_splitter.split_text(coffee_text)

embeddings = OpenAIEmbeddings()

docsearch = Chroma.from_texts(texts, embeddings, metadatas=[
    {"source": str(i)} for i in
range(len(texts))])

# query and run chain
query = "What is the difference between a cappuccino and a latte?"
docs = docsearch.similarity_search(query)

chain = load_qa_with_sources_chain(OpenAI(temperature=0.3),
chain_type="stuff")
result = chain({"input_documents": docs, "question": query},
return_only_outputs=False)

print(result)
```

```
● ● ●

from langchain.llms import OpenAI
from langchain.chains.qa_with_sources import
load_qa_with_sources_chain
from langchain.text_splitter import CharacterTextSplitter
from langchain.embeddings import OpenAIEMBEDDINGS
from langchain.vectorstores import Chroma
from langchain.docstore.document import Document

# read text and embed file
with open('texts/coffee.txt') as f:
    coffee_text = f.read()
text_splitter = CharacterTextSplitter(chunk_size=1000,
chunk_overlap=0)
texts = text_splitter.split_text(coffee_text)

embeddings = OpenAIEMBEDDINGS()

docsearch = Chroma.from_texts(texts, embeddings, metadatas=[{"source": str(i)} for i in range(len(texts))])

# query and run chain
query = "What is the difference between a cappuccino and a latte?"
docs = docsearch.similarity_search(query)

chain = load_qa_with_sources_chain(OpenAI(temperature=0.3),
chain_type="stuff")
result = chain({"input_documents": docs, "question": query},
return_only_outputs=False)

print(result)
```

“A cappuccino consists of 1-2 shots of espresso, steamed milk and a layer of foamed milk on top, while a latte consists of espresso with steamed warm milk and a layer of foamed milk on top.

SOURCES: 4, 5”

What is the difference between a cappuccino and a latte?

**Clear trade-off between
abstractions and finer control
over prompt**

Extensions

- **Build out chat:** maintain conversation history, enable memory
- **Enable tools:** surf the Internet, use calculator to determine how much milk depending on number of friends
- **Try out other Embeddings:** Cohere, OpenAI, browse Hugging Face
- **Approximate nearest neighbors libraries:** FAISS, annoy, hnswlib
- And so much more...

Where to go from here?

- **LLM Composability/Abstractions:** Langchain, Llama Index
- **Vector Databases:** Pinecone, Weaviate
- **Other LLMs?** Claude (Anthropic), Cohere, Open Source (Llama 2, etc.)
- **Model Fine-Tuning**
- **Model Observability**

Thanks!

vnshenoy@stanford.edu



@varunshenoy_