

LangChain for Generative AI

Bruno Gonçalves

graphs4sci.substack.com

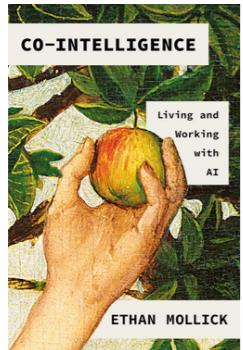
<https://github.com/DataForScience/LangChain>



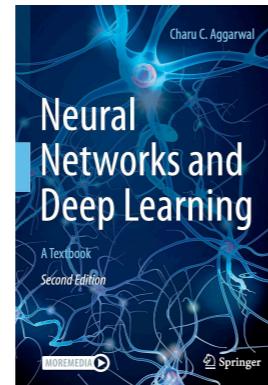
References

<https://github.com/sswamyn/LangChain-for-Gen-AT-pipeline>

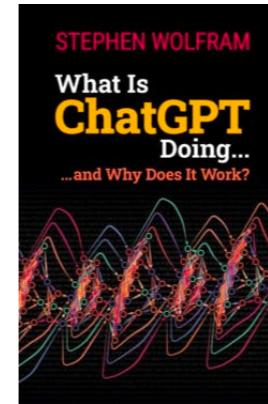
More business / application focused ...



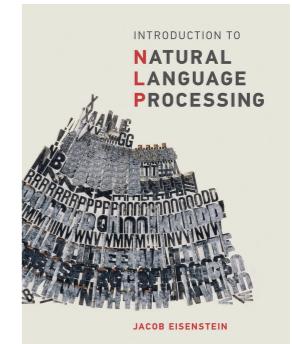
<https://amzn.to/3KcLlsf>



<https://amzn.to/48rZn9X>

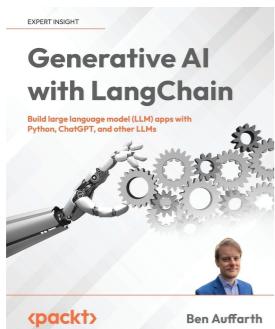


<https://amzn.to/3LBRdBY>

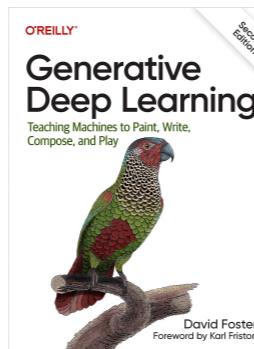


<https://amzn.to/3ZMnTih>

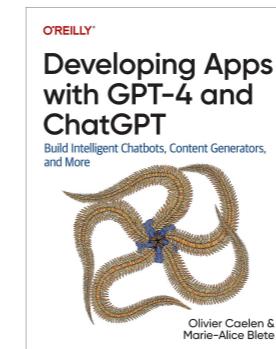
More technical ...



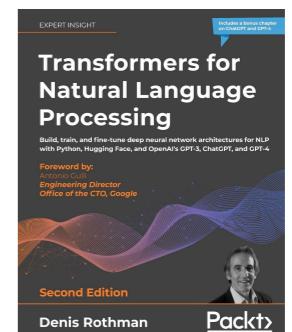
<https://amzn.to/3VquiPj>



<https://amzn.to/3t8PuxM>



<https://amzn.to/3RHRkQa>



<https://amzn.to/46kOo02>

<https://github.com/sswamyn/LangChain-for-Gen-AI-pipeline>

Table of Contents



1. Generative AI
2. LangChain
3. Information Processing
4. Chatbots
5. Prompt Engineering

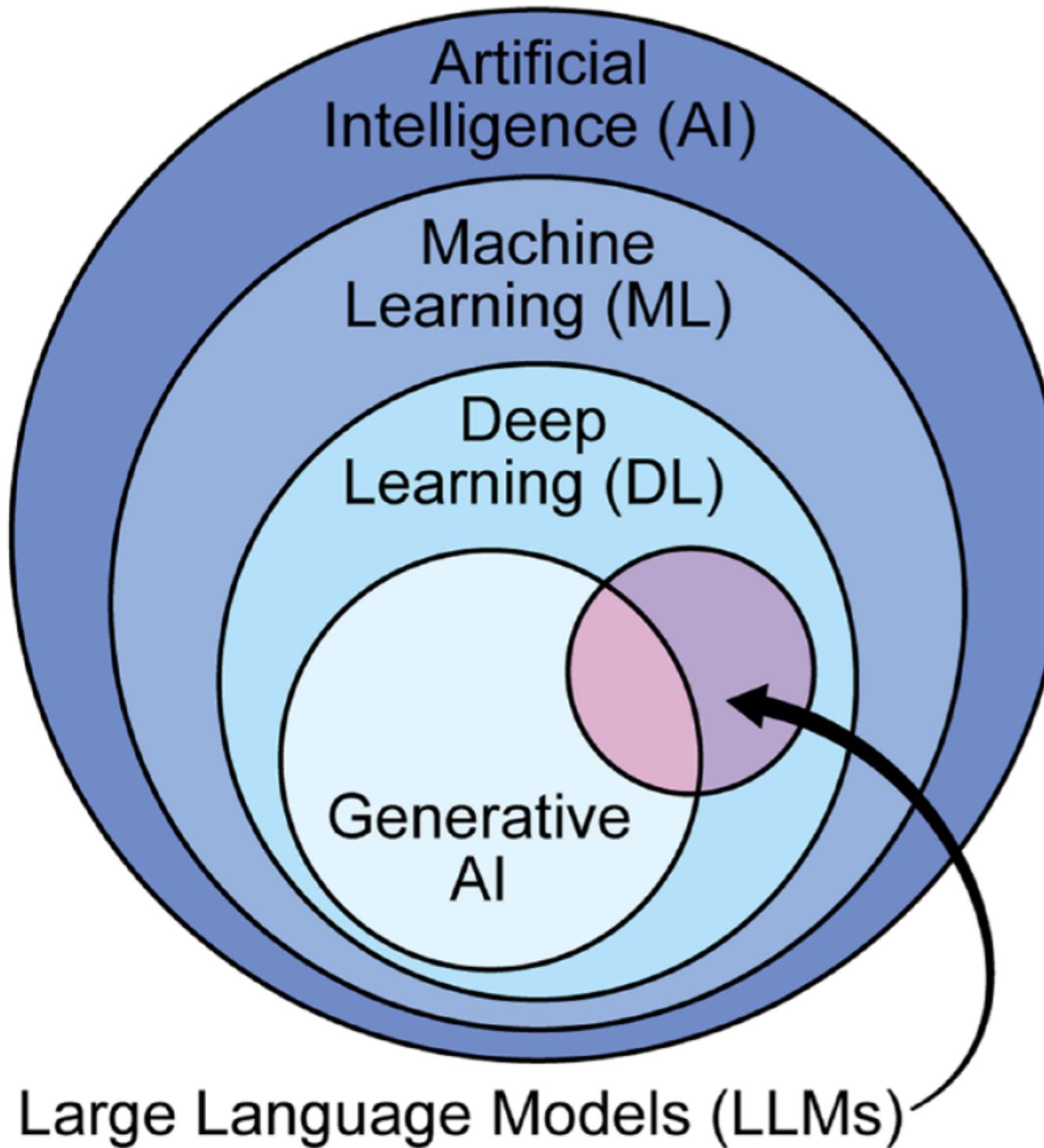


1. Generative AI

“The model that you are using today is the worst model that you will ever use.”

sidebar: OpenAI recently changed all Dictionaries to Tuples that will break older notebooks, but can be fixed with simple “find & replace”

Generative Models



- **Artificial Intelligence (AI):** A field of CS focused on creating intelligent agents.
- **Machine Learning (ML):** Subset of AI focused on developing algorithms that can learn from data.
- **Deep Learning (DL):** Neural networks with many layers
- **Generative Models:** Models that can generate new data based on learned patterns
- **Large Language Models (LLMs):** Statistical models used to predict words in natural language

Generative Models

- Come in many flavors, depending on what the input and outputs are
 - **Text-to-Text** - Dialog based systems, like ChatGPT
 - **Text-to-Image** - Capable of generating images based on prompts, like DALL-E
 - **Image-to-Text** - Capable of generating text that describes an image
 - **Text-to-Speech** - Generate audio from written text
 - **Speech-to-Text** - Transcribe audio into text
 - **Text-to-Code** - Generate programming code based on a problem descriptions, like GitHub CoPilot
 - **Image-to-Image** - Automatically image editing and manipulation

Language Models

- A foundational component of Natural Language Processing and AI.
- Software that predicts and generates human language based on patterns observed in training data.

- Applications to:

- Machine translation
- Summarization
- Sentiment analysis
- Chatbots and Virtual Assistants
- Content Generation
- etc...

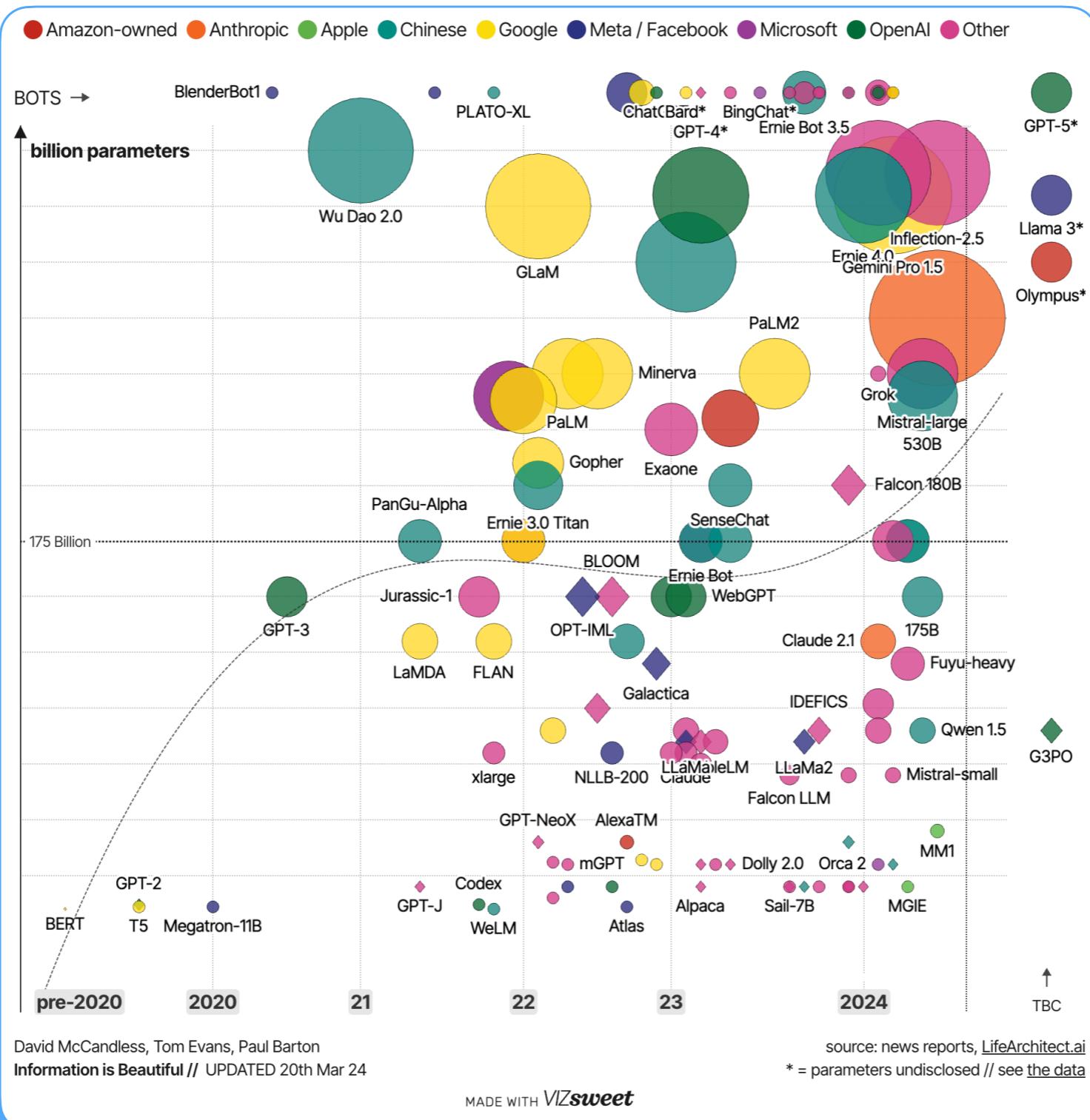
what is the difference between autoregressive and auto encoder models, what is chat GPT?

Autoregressive: The output is feed back as input. Ex: Chat GPT

Auto encoder: These models take the input and reduces it to a very small token. Then uses this small token and tries to reproduce the input. Do it as many times as needed to train the model to get the output accurate.

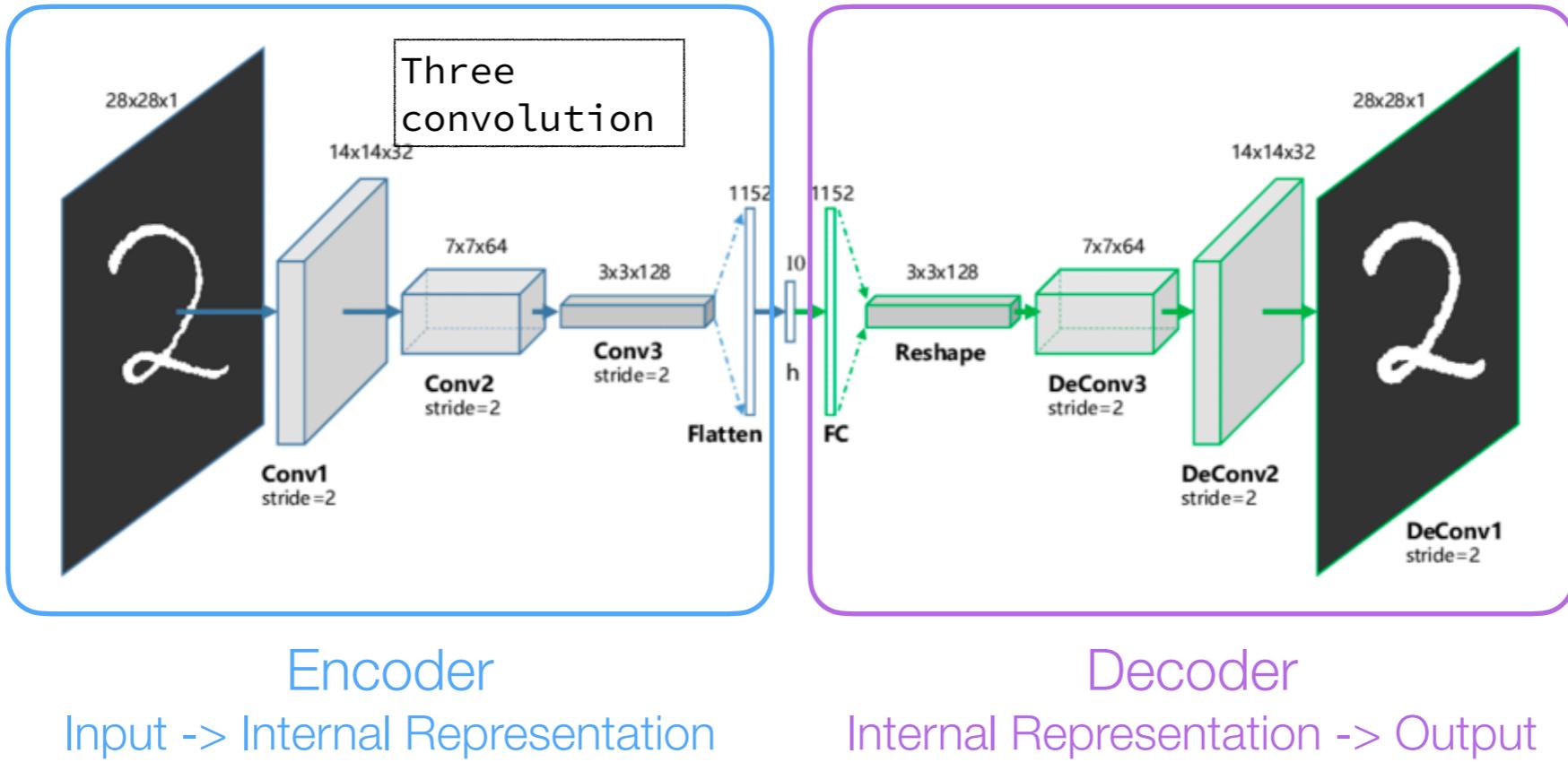
LLM Comparison

<https://informationisbeautiful.net/visualizations/the-rise-of-generative-ai-large-language-models-langs-like-chatgpt/>



size = no. of parameters
open-access

Encoder/Decoder Architecture



The magic that the Transformer model brought to the table.
In a big context window, Transformer added the capability to figure out what are the parts of the context are relevant to the most recent question being asked. This enabled the model to be accurate, instead of getting confused trying to use lots of information.

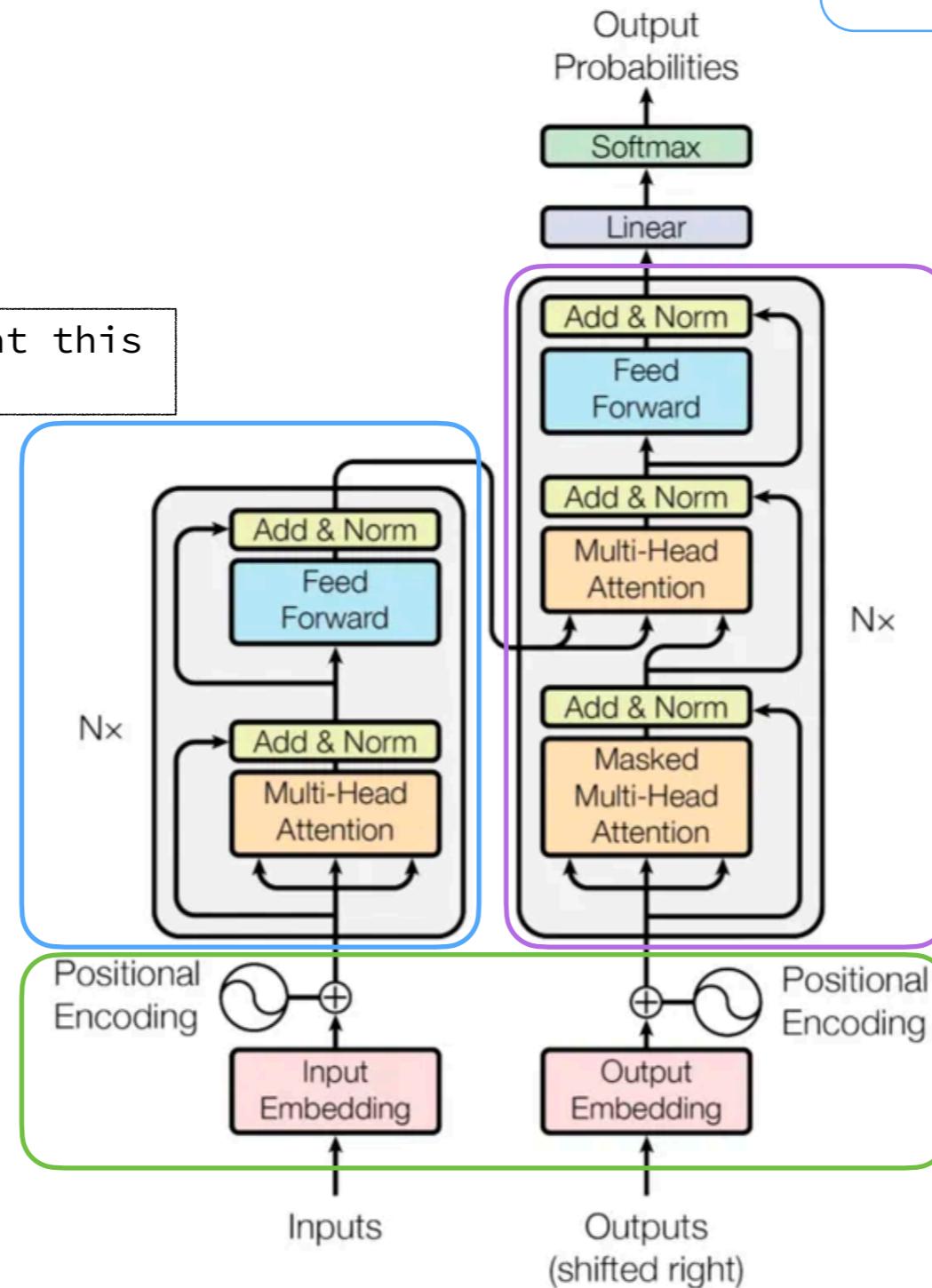
Pre-trnasformer - transfer learning was the big deal.

Transformers

arXiv:1706.03762

Nx - represents that this is repeated.

Encoder

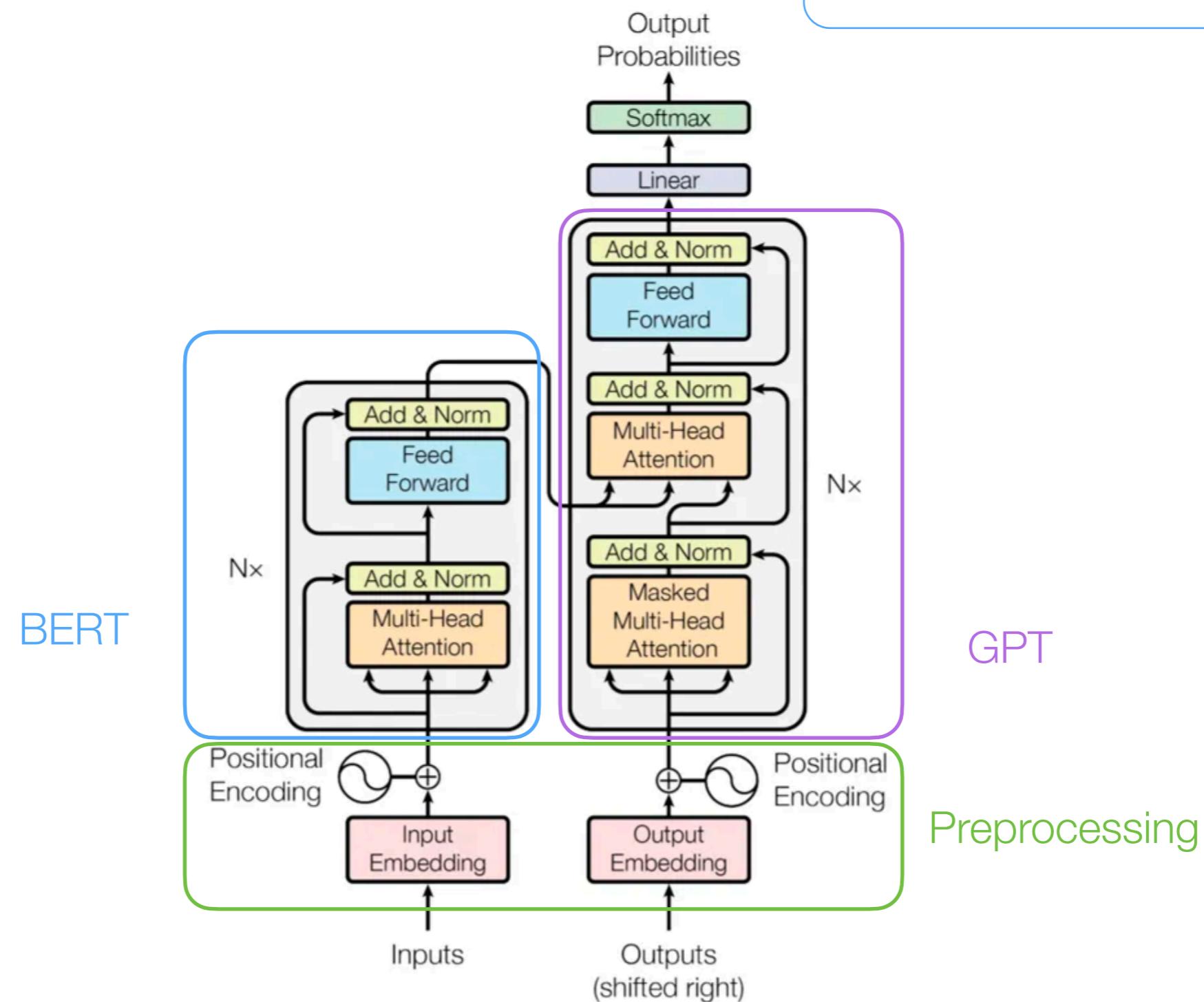


Decoder

Preprocessing

Transformers

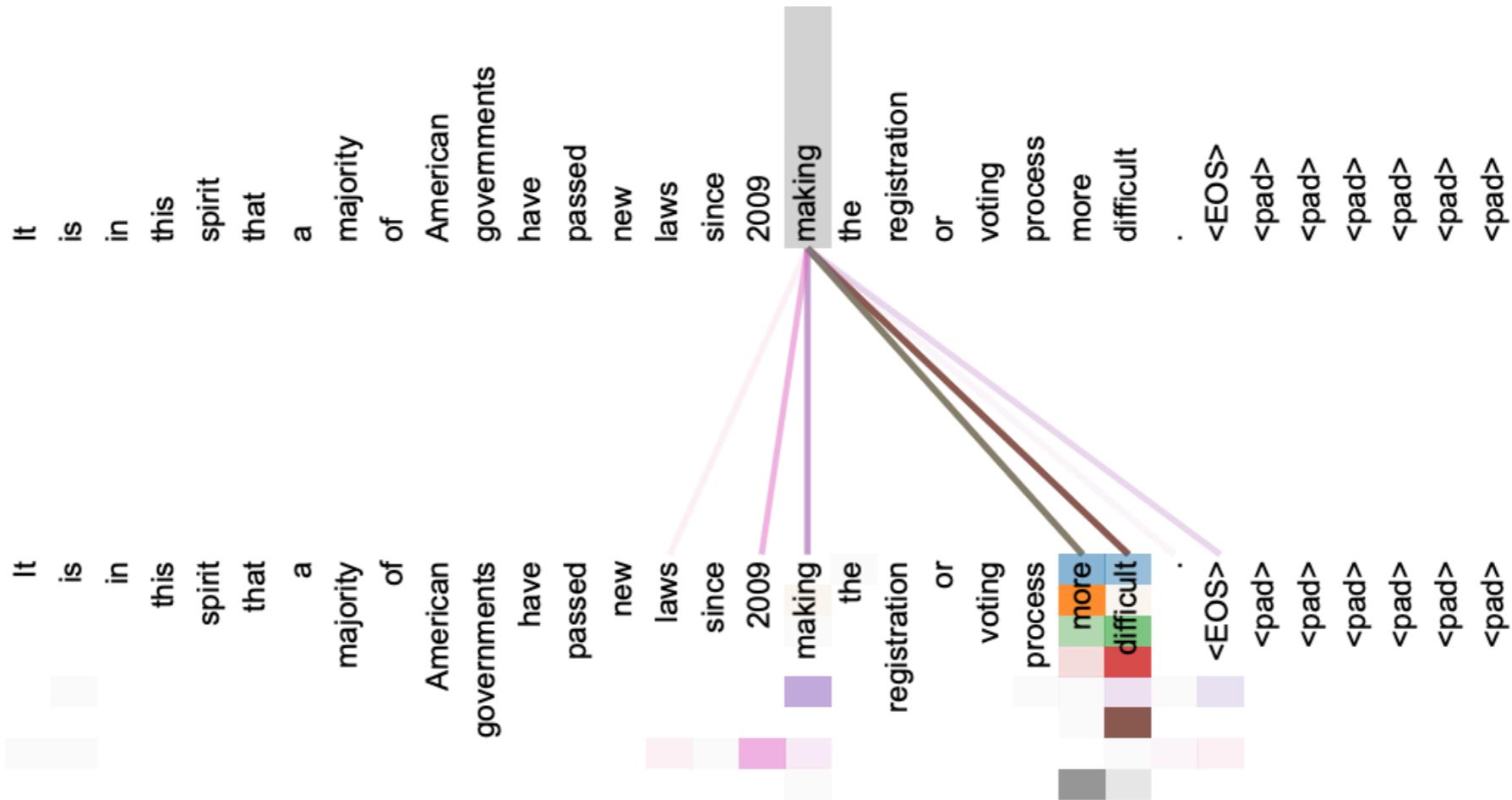
arXiv:1706.03762



Attention

arXiv:1706.03762

- “Simple” mechanism to allow each token to take the context it appears in into account
- Requires exponentially more weights to be computed (from each token to every other token)

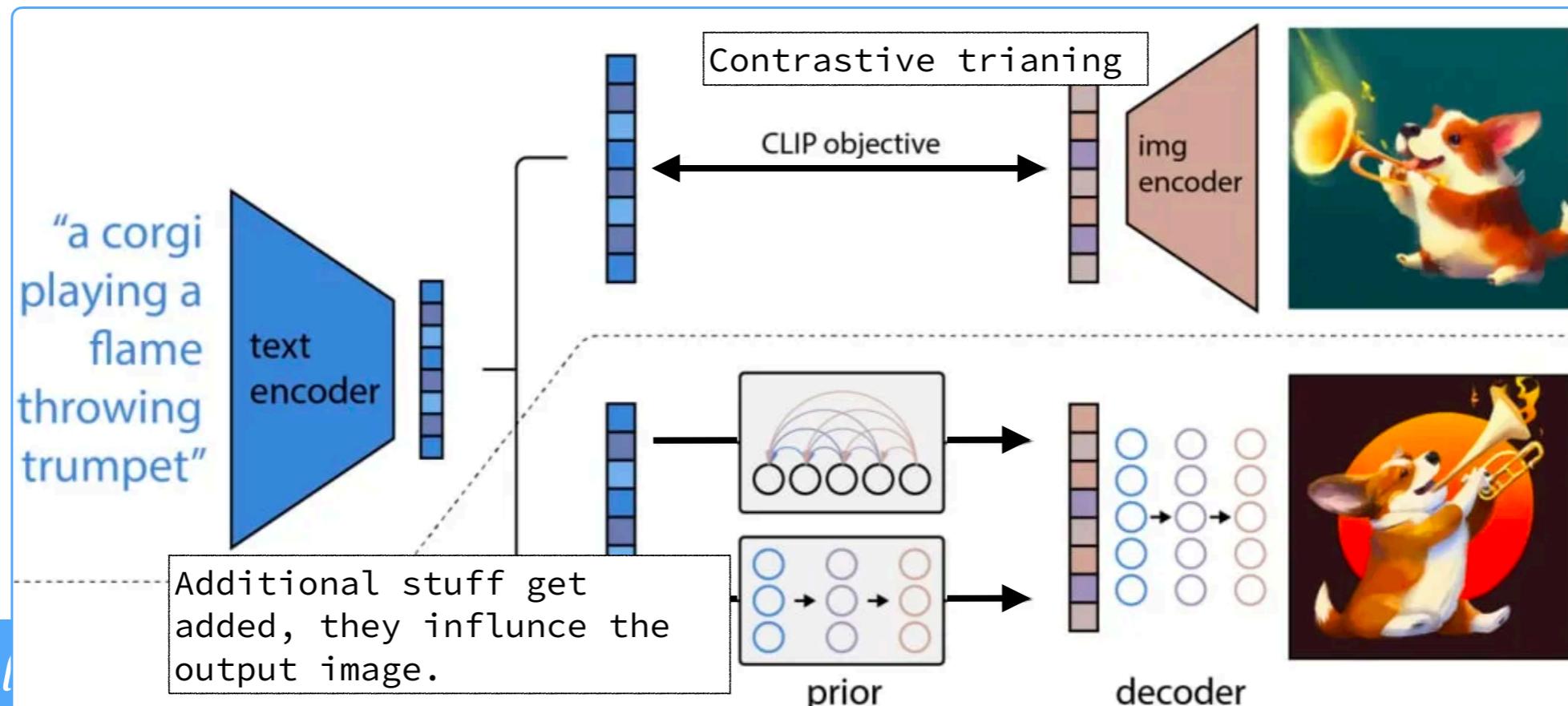


- Weights indicate the importance of each word relative to the current one

Text to Image Models

arXiv:2204.06125

- Image models like **DALL-E** build on top of the **LLM** approach outlined above
- Essentially, **DALL-E replaces the LLM decoder** with an image decoder
- The image decoder is trained using **Contrastive Training**:
 - A large list of images and textual descriptions is used to train the image decoder to generate an internal representation that matches the encoded description
 - Inverting the process produces a model that converts textual descriptions into images
- As a final improvement, we add a **prior** in between the output of the text encoder and the input of the image decoder.





Code - Generative AI

<https://github.com/DataForScience/LangChain>



2. LangChain

LangChain

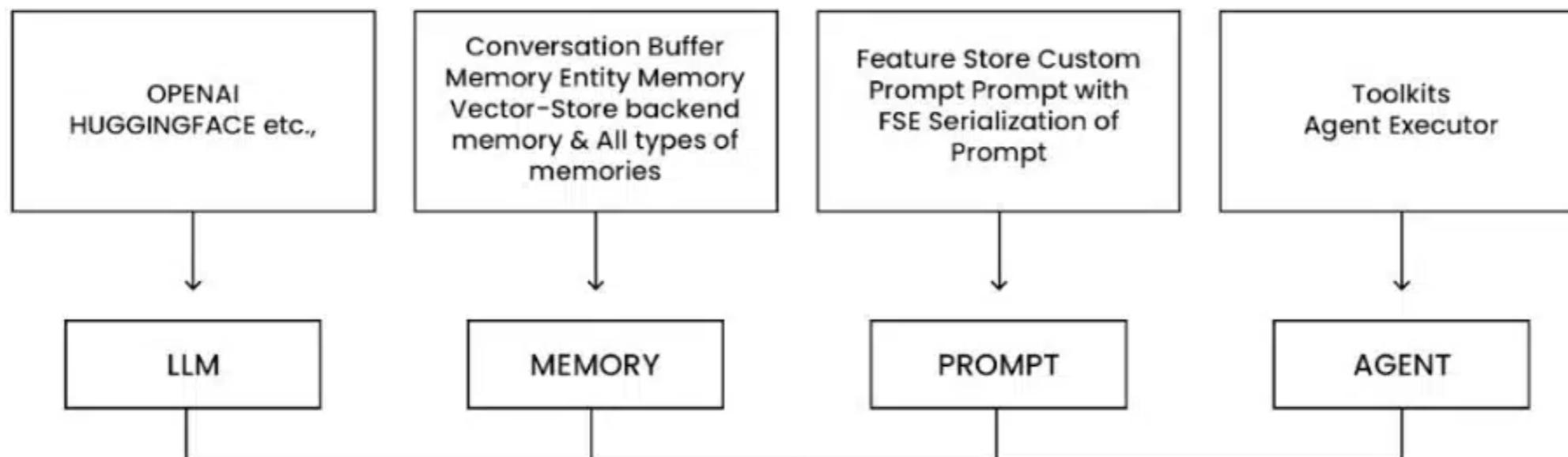
- LangChain is an open-source Python framework designed to simplify the development of applications using large language models (LLMs).
- Main features:
 - Model Abstraction: LangChain manages inputs and outputs seamlessly by abstracting away model details
 - Modular Structure: Allows developers to combine different prompts and even multiple LLMs within a single application.
 - Chains: Multiple components can be chained together to create complex applications
 - Data Integration: Allows LLMs to transform, store, and retrieve data from databases and online sources.
 - Memory Management: Provides utilities for adding memory to LLM systems, retaining conversation history or summaries for context.

LangChain Applications

- LangChain can be used to build a wide range of LLM-powered applications:
 - Chatbots
 - Coding assistants and code security analysis
 - Recommendation systems
 - Document analysis and summarization
 - Question-answering systems
 - Data augmentation
 - Text classification and summarization
 - Sentiment analysis
 - Machine translation

LangChain Structure

HIGH LEVEL STRUCTURE OF LANGCHAIN

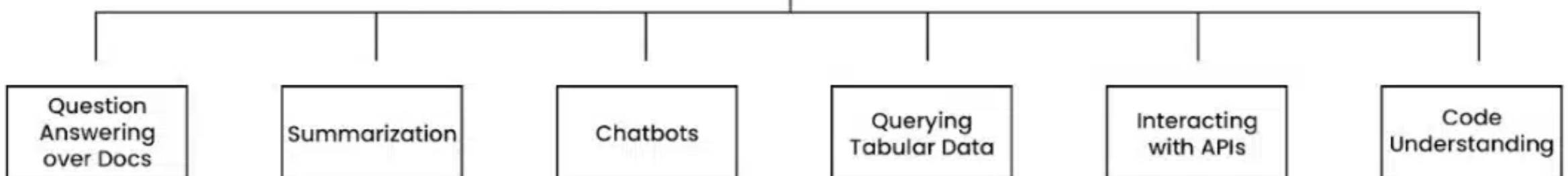


LangGraph is part of LangChain, that is used for

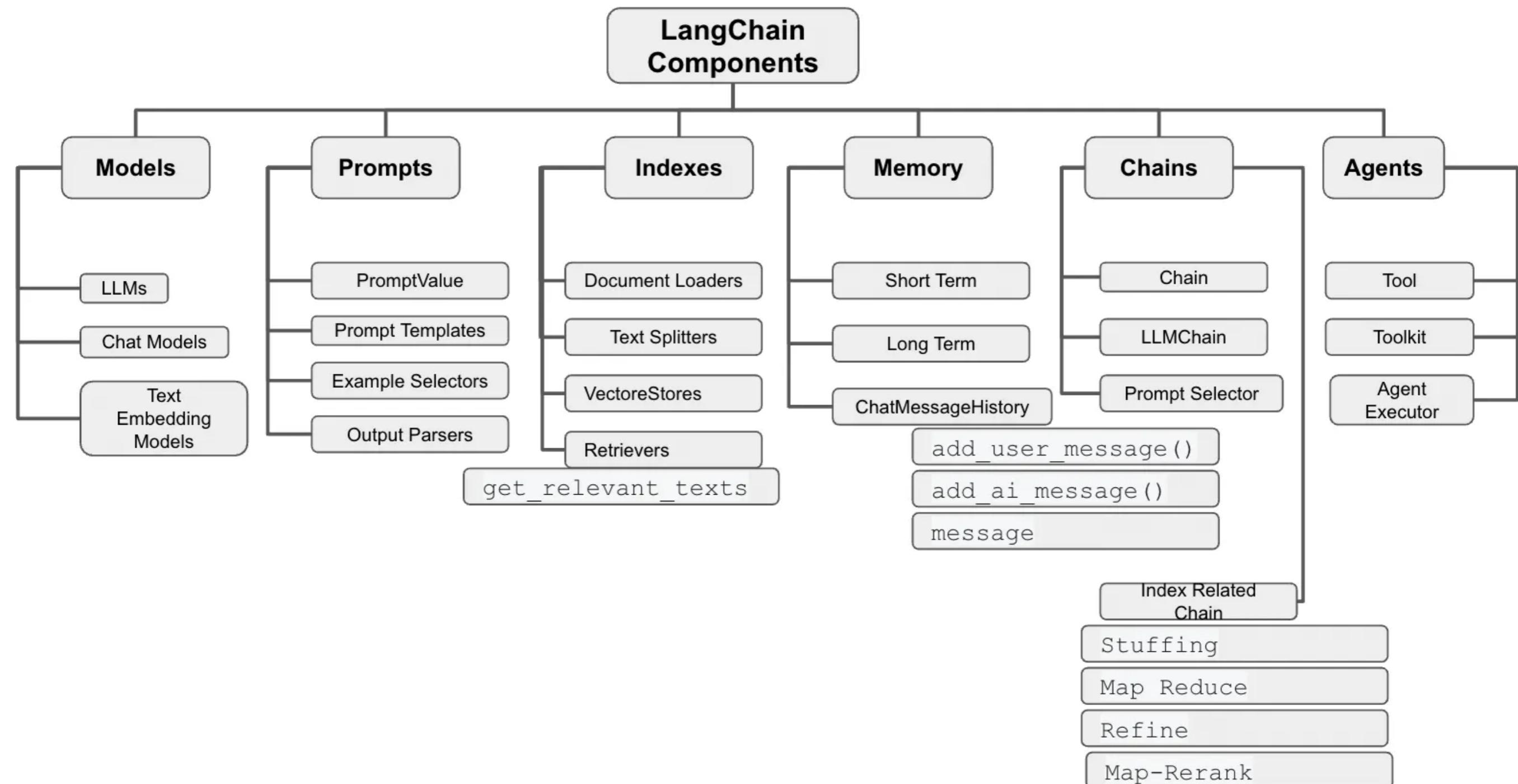
They are trying to push us into LG!

THere is a tool/agent for that!! :)

https://python.langchain.com/v0.2/docs/how_to/migrate_agent/



LangChain Structure



Prompts

- Prompts are the basic way to interact with an LLM model. LangChain provides different classes to manage and generate prompts
 - [PromptValue](#) - Represents the text input to a model
 - [PromptTemplate](#) - Jinja-like templating engine for a prompt. It takes input variables and returns a [PromptValue](#). Used for [HumanMessages](#), [SystemMessages](#), and [AIMessages](#).
 - [StrOutputParser](#) - Helps structure language model responses

Chains

- A Chain is a sequence of calls to components, including other chains.
- Conceptually similar to [sklearn Pipelines](#)
- Component Examples:
 - [LLMMath](#) - For math related queries
 - [SQLDatabaseChain](#) - To query databases
 - [OpenAIModerationChain](#) - Check content against moderation rules
 - [ConstitutionalChain](#) - Legal applications
 - [LLMCheckerChain](#) - Prevent hallucinations

Advantages of Chains

- **Modularity** - Each component is an isolated module
- **Composability** - Modules can be easily connected together
- **Readability** - Each step in the chain is clear
- **Maintainability** - The sequence of steps can be easily expanded or changed
- **Reusability** - Common functionality becomes chains
- **Tool Integration** - Integrate LLMs, databases, APIs, etc
- **Productivity** - Facilitates prototyping and development

Agents

- Agents are autonomous software that are capable of taking actions to accomplish specific goals and tasks
- Agents leverage chains to take goal-driven actions and can combine and orchestrate chains
- Advantages of agents:
 - **Goal-oriented execution:** Plan chains of logic targeting specific goals.
 - **Dynamic responses:** React and adapt to environment changes.
 - **Statefulness:** Maintain memory and context.
 - **Robustness:** Error handling.
 - **Composition:** Combine reusable component chains.

Memory

- A key limitation of agents and chains is their statelessness
- Memory refers to the persisting state between executions of a [Chain](#) or [Agent](#).
- Several versions:
 - [ConversationBufferMemory](#) - store all messages
 - [ConversationBufferWindowMemory](#) - retain only recent messages
 - [ConversationKGMemory](#) - Summarizes exchanges as a knowledge graph
 - [EntityMemory](#) - persists agent state and facts through a database
 - [SQLDatabase](#) - Database integration

Using tools

<https://python.langchain.com/v0.2/docs/integrations/tools/>

- **Tools** are interfaces that **Agents** can use to interact with the world.
- Examples:
 - [ShellTool](#) - To run commands on the command line
 - [DuckDuckGoSearchRun](#) - To search on DDG
 - [AlphaVantageAPIWrapper](#) - Financial data
 - [RequestsGetTool](#) - Access a specific URL
 - [WikipediaQueryRun](#) - Query Wikipedia
 - [YouTubeSearchTool](#) - Query YouTube
 - And many more...



Code - LangChain

<https://github.com/DataForScience/LangChain>



3. Information Processing

Information processing

- Information processing is perhaps the killer application for **LLMs**
- Large models ability to understand context and nuance lets them quickly and effectively extracting relevant information from various types of documents:
 - Unstructured text
 - Complex document formats (PDF, Word, JSON, etc)
 - Images with text (OCR)
 - mixed content formats.
- **LLMs** can sift through large volumes of data to identify and isolate critical details, making them invaluable for tasks like data categorization, entity extraction, sentiment analysis, etc
- **LLMs** are able to work across multiple languages
- And more importantly, **LLMs** can be easily automated

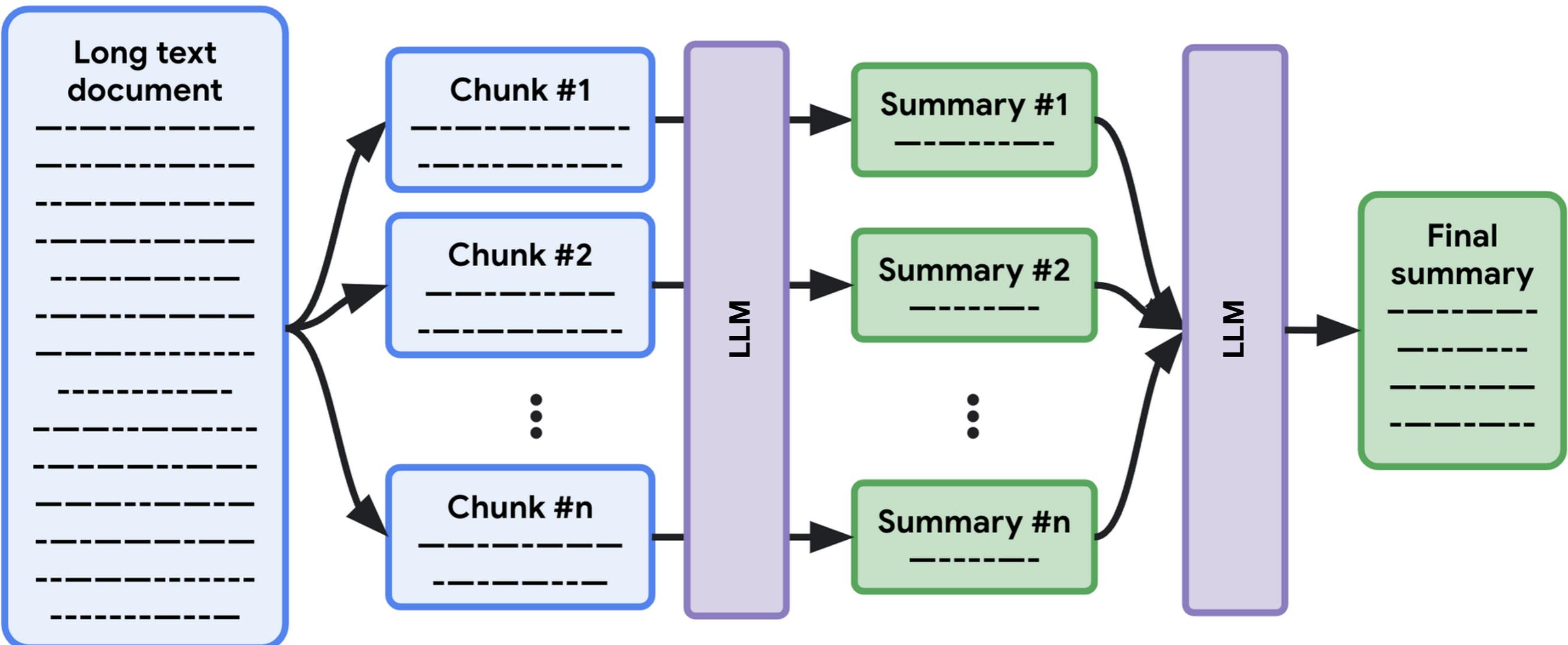
Text Summarization

- Text summarization is a classical NLP problem and a common application for LLMs
- Two main types of summarization approaches:
 - **Extractive Summarization**: Identify and extract the most important sentences from a text and combine them to form a summary.
 - **Abstractive Summarization**: Generate new sentences that capture the essence of the original text.
- Several approaches are possible based on text size:
 - Document fits in context window: A basic Prompt is enough - "**Write a summary of the following:**"
 - Multiple documents: **StuffDocumentsChain** combines all documents into a single object that can be processed by the LLM
 - Multiple Large Documents: **MapReduceDocumentsChain** to iteratively summarize each document (**Map**) and then combine all summaries into a single result (**Reduce**)

MapReduceDocumentsChain

<https://cloud.google.com/blog/products/ai-machine-learning/long-document-summarization-with-workflows-and-gemini-models>

- MapReduceDocumentsChain - uses multiple chains and prompts to obtain the final result
- Different prompting techniques can be used to improve the results



Information Extraction

- Information extraction is the process of automatically extracting structured data from unstructured or semi-structured sources such as webpages, PDFs, etc
- **LangChain** requires us to define the schema of the information we are trying to extract
- Once defined, the schema is passed to the `.with_structured_output()` function of the **LLM**

```
extraction_chain = prompt | llm.with_structured_output(schema=Data)
```



Code - Information Processing
<https://github.com/DataForScience/LangChain>



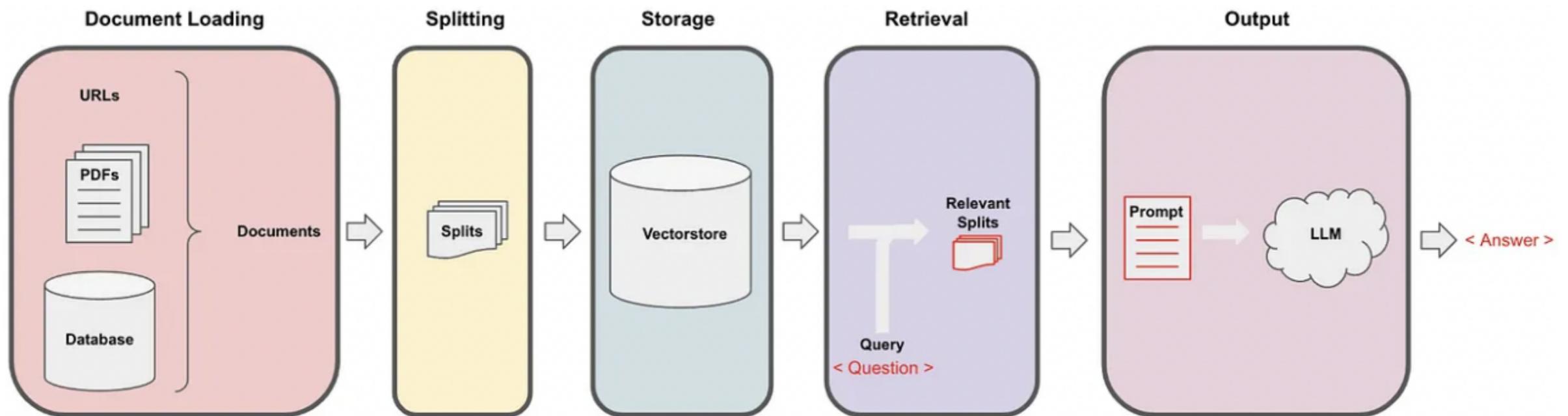
4. Chatbots

Retrieval Augmented Generation

- **Retrieval Augmented Generation (RAG)** is a technique to allow **LLMs** to retrieve and use relevant information from external sources
- Advantages of RAG
 - **Improved Accuracy:** **LLMs** generate more accurate and relevant responses with up-to-date information not included in their pre-training data.
 - **Verifiable Outputs:** **LLMs** can report the sources they used to generate their output.
 - **Reduced Hallucinations:** Grounds **LLMs** generation on factual information
 - **Adaptability:** **LLMs** can be adapted to new domains

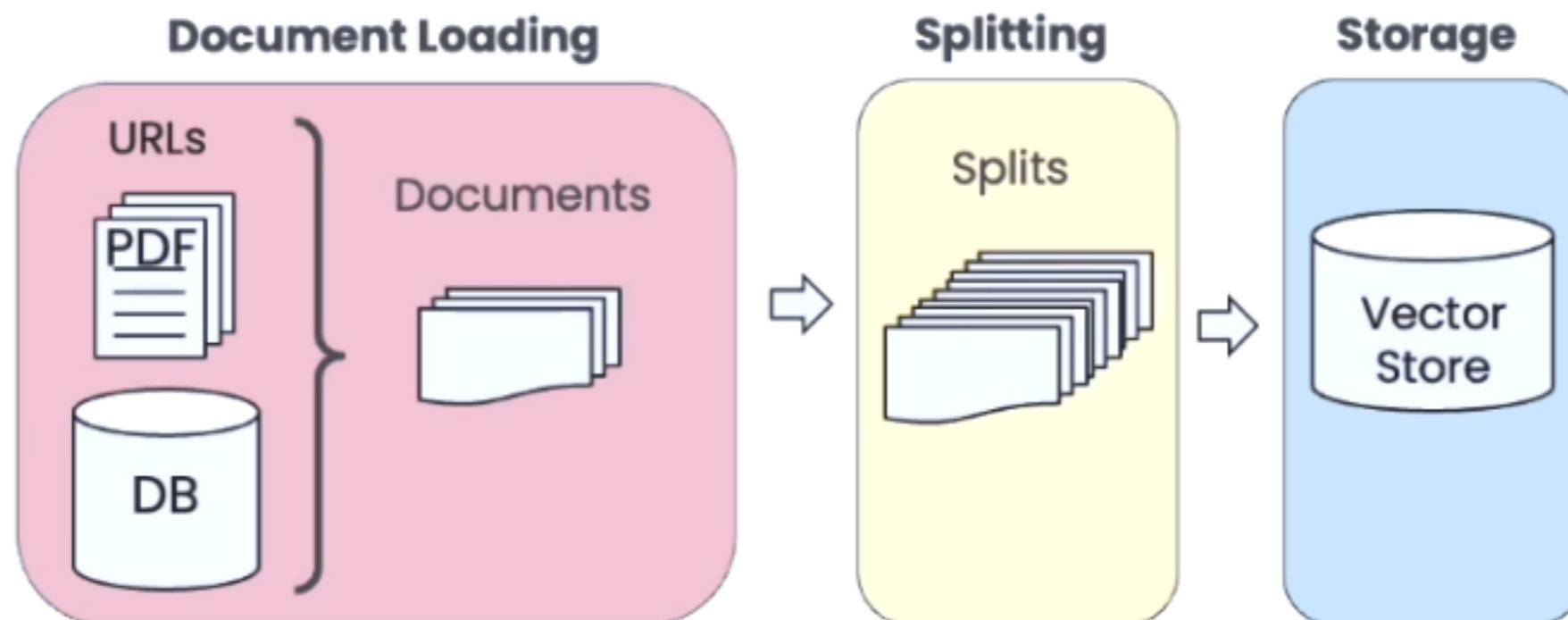
Retrieval Augmented Generation

- RAG is a multistep process



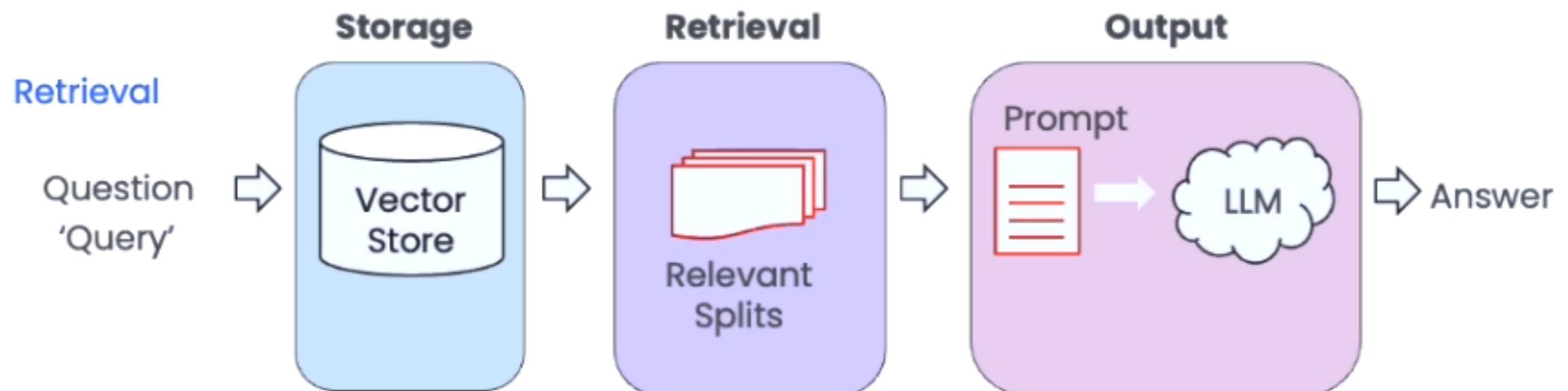
Vector Storage and Retrieval

- Vector embeddings are used to represent the semantic meaning of each document, chunk or query
- A vector database allows us to quickly find related documents to provide to the LLM



Vector Storage and Retrieval

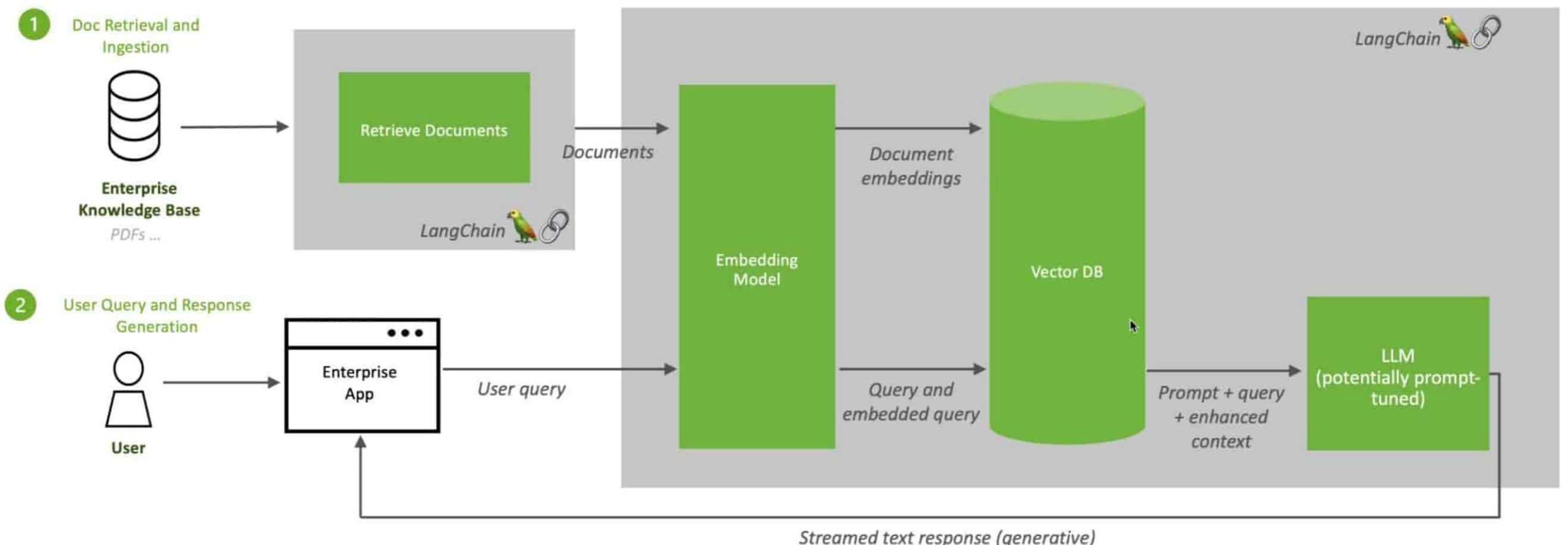
- Vector embeddings are used to represent the semantic meaning of each document, chunk or query
- A vector database allows us to quickly find related documents
- Retrieved documents can be added to the prompt context window to help the LLM provide better answers



Simple RAG ChatBot

<https://blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/>

Retrieval Augmented Generation (RAG) Sequence Diagram





Code - Chatbots

<https://github.com/DataForScience/LangChain>



5. Prompt Engineering

Prompts

<https://www.pinecone.io/learn/series/langchain/langchain-prompt-templates/>

- Prompts can have multiple parts, but now all prompts use all pieces

INSTRUCTIONS Tell the model what to do

"""Answer the question based on the context below. If the question cannot be answered using the information provided answer with "I don't know".

Context: Large Language Models (LLMs) are the latest models used in NLP. Their superior performance over smaller models has made them incredibly useful for developers building NLP enabled applications. These models can be accessed via Hugging Face's `transformers` library, via OpenAI using the `openai` library, and via Cohere using the `cohere` library.

Question: Which libraries and model providers offer LLMs?

Answer: """

PROMPTER INPUT

OUTPUT INDICATOR

- **PromptTemplates** allow us to quickly generate customized prompts for each use

LLM - is an hardworking intern. So you have to tell him/her what you want them to do very clearly - step-by-step.
As with any intern, you have to double check what they produce before using it! 😎

PromptTemplate

<https://www.pinecone.io/learn/series/langchain/langchain-prompt-templates/>

- Template parameters are identified by curly braces **{query}**

```
1 from langchain import PromptTemplate
2
3 template = """Answer the question based on the context below. If the
4 question cannot be answered using the information provided answer
5 with "I don't know".
6
7 Context: Large Language Models (LLMs) are the latest models used in NLP.
8 Their superior performance over smaller models has made them incredibly
9 useful for developers building NLP enabled applications. These models
10 can be accessed via Hugging Face's `transformers` library, via OpenAI
11 using the `openai` library, and via Cohere using the `cohere` library.
12
13 Question: {query}
14
15 Answer: """
16
17 prompt_template = PromptTemplate(
18     input_variables=["query"],
19     template=template
20 )
```

PromptTemplate

<https://www.pinecone.io/learn/series/langchain/langchain-prompt-templates/>

- Template parameters are identified by curly braces `{query}`
- To generate a new prompt, we must simply provide the value of all the parameters to the `format()` method

```
1 print(  
2     prompt_template.format(  
3         query="Which libraries and model providers offer LLMs?"  
4     )  
5 )
```

```
1 from langchain import PromptTemplate  
2  
3 template = """Answer the question based on the context below. If the  
4 question cannot be answered using the information provided answer  
5 with "I don't know".  
6  
7 Context: Large Language Models (LLMs) are the latest models used in NLP.  
8 Their superior performance over smaller models has made them incredibly  
9 useful for developers building NLP enabled applications. These models  
10 can be accessed via Hugging Face's `transformers` library, via OpenAI  
11 using the `openai` library, and via Cohere using the `cohere` library.  
12  
13 Question: {query}  
14  
15 Answer: """  
16  
17 prompt_template = PromptTemplate(  
18     input_variables=["query"],  
19     template=template  
20 )
```

Simple Prompting strategies

<https://lilianweng.github.io/posts/2023-03-15-prompt-engineering/>

- Prompting is how we interact with LLMs
- Choosing the correct type of prompt can have a major effect on the results we obtain
- Common prompting approaches:
 - **Zero-Shot** - Explain the task and ask for the result (most examples we've seen so far)
 - **Few-Shot** - Provide a few examples of both inputs and outputs that the model can generalize from
- The **FewShotPromptTemplate** makes it easy to generate custom prompts with a few examples, to test the effect of different number of examples, etc

```
39 few_shot_prompt_template = FewShotPromptTemplate(  
40     examples=examples,  
41     example_prompt=example_prompt,  
42     prefix=prefix,  
43     suffix=suffix,  
44     input_variables=["query"],  
45     example_separator="\n\n"  
46 )
```

A list of examples to be used
The template of each example
Instructions proceeding the examples
Instructions following the examples
variables for the few-shot prompt template
Characters to use to separate examples

Prompting Tips

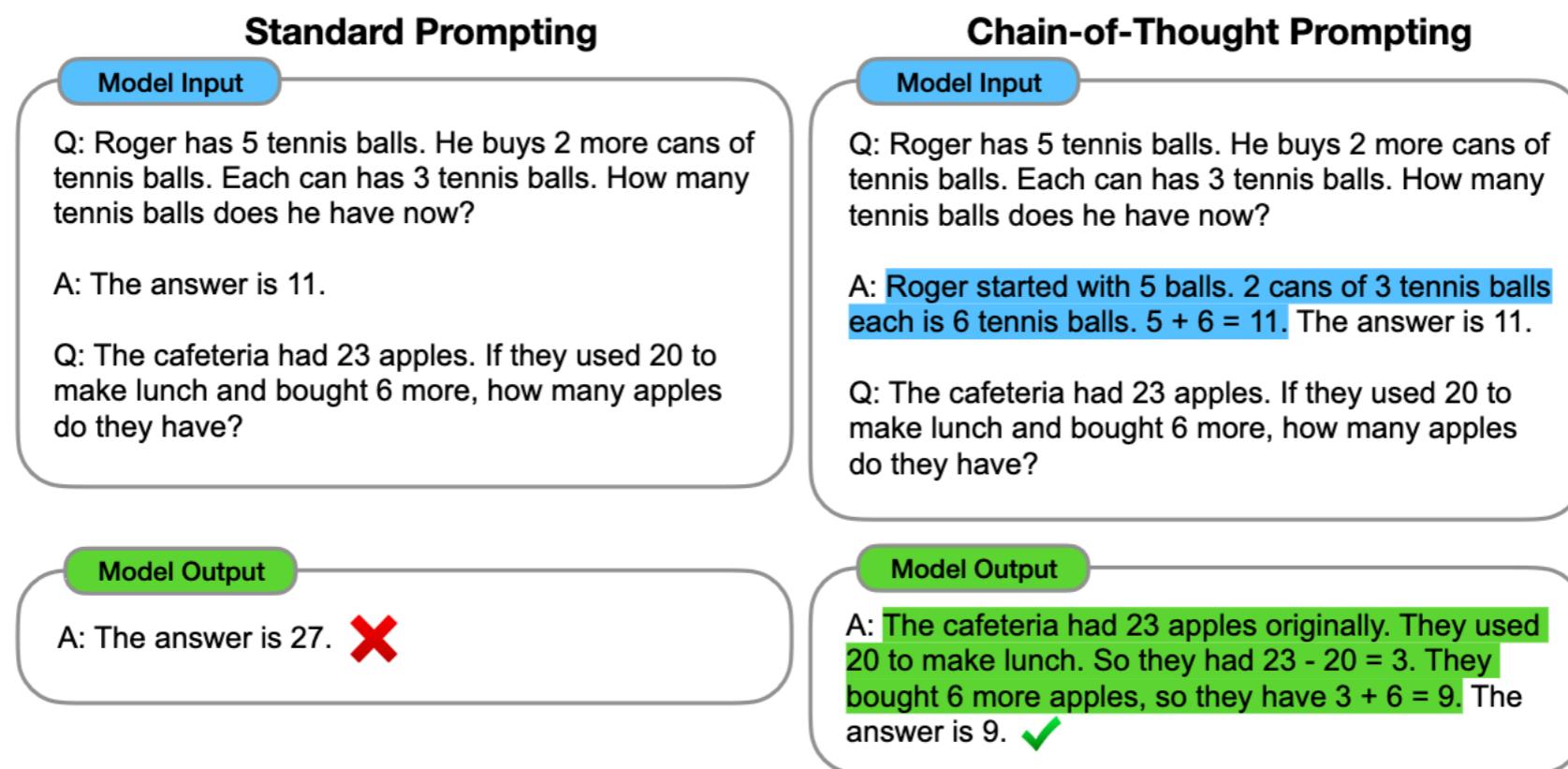
<https://www.promptingguide.ai/introduction/tips>

- Start simple and integrate by gradually adding the details necessary to improve results
- Break down complex tasks into smaller, more manageable, subtasks
- Use action commands ("Write", "Classify", "Summarize", "Translate", "Order", etc.) to indicate what the goal is, such as
- More descriptive and detailed the prompts produce better results
Chain of thought prompts

Chain-of-Thought prompts

arXiv: 2201.11903 (2022)

- **Chain of Thought (CoT)** prompts encourages LLMs to generate intermediate steps before providing the final solution to a problem.
- **CoT** breaks down complex problems into manageable, intermediate steps and mimic an human thought process when working through multi-step problems.
- Relatively small change to the prompt can have a profound effect



Chain-of-Thought prompts

arXiv: 2201.11903 (2022)

Math Word Problems (free response)

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.

Math Word Problems (multiple choice)

Q: How many keystrokes are needed to type the numbers from 1 to 500?
Answer Choices: (a) 1156 (b) 1392 (c) 1480
(d) 1562 (e) 1788

A: There are 9 one-digit numbers from 1 to 9. There are 90 two-digit numbers from 10 to 99. There are 401 three-digit numbers from 100 to 500. $9 + 90(2) + 401(3) = 1392$. The answer is (b).

CSQA (commonsense)

Q: Sammy wanted to go to where the people were. Where might he go?
Options: (a) race track (b) populated areas
(c) desert (d) apartment (e) roadblock

A: The answer must be a place with a lot of people. Race tracks, desert, apartments, and roadblocks don't have a lot of people, but populated areas do. So the answer is (b).

StrategyQA

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about 0.6 g/cm^3 , which is less than water. Thus, a pear would float. So the answer is no.

Date Understanding

Q: The concert was scheduled to be on 06/01/1943, but was delayed by one day to today. What is the date 10 days ago in MM/DD/YYYY?

A: One day after 06/01/1943 is 06/02/1943, so today is 06/02/1943. 10 days before today is 05/23/1943. So the answer is 05/23/1943.

Sports Understanding

Q: Is the following sentence plausible? "Joao Moutinho caught the screen pass in the NFC championship."

A: Joao Moutinho is a soccer player. The NFC championship is part of American football, not soccer. So the answer is no.

SayCan (Instructing a robot)

Human: How would you bring me something that isn't a fruit?

Explanation: the user wants something to eat that isn't a fruit. An energy bar is not a fruit, so I will bring the user an energy bar.

Plan: 1. find(energy bar) 2. pick(energy bar) 3. find(user) 4. put(energy bar) 5. done().

Last Letter Concatenation

Q: Take the last letters of the words in "Lady Gaga" and concatenate them.

A: The last letter of "Lady" is "y". The last letter of "Gaga" is "a". Concatenating them is "ya". So the answer is ya.

Coin Flip (state tracking)

Q: A coin is heads up. Maybelle flips the coin. Shalonda does not flip the coin. Is the coin still heads up?

A: The coin was flipped by Maybelle. So the coin was flipped 1 time, which is an odd number. The coin started heads up, so after an odd number of flips, it will be tails up. So the answer is no.

Chain-of-Thought - Zero Shot

arXiv: 2201.11903 (2022)

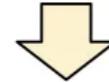
【1st prompt】 Reasoning Extraction

Q: On average Joe throws 25 punches per minute. A fight lasts 5 rounds of 3 minutes. How many punches did he throw?

A: Let's think step by step.



LLM



In one minute, Joe throws 25 punches.

In three minutes, Joe throws $3 * 25 = 75$ punches.

In five rounds, Joe throws $5 * 75 = 375$ punches.

Instruct the model to think step by step

【2nd prompt】 Answer Extraction

Q: On average Joe throws 25 punches per minute. A fight lasts 5 rounds of 3 ...
A: Let's think step by step.

In one minute, Joe throws 25 punches. ... In five rounds, Joe throws $5 * 75 = 375$ punches..

Therefore, the answer (arabic numerals) is



LLM



375.



Code - Prompt Engineering
<https://github.com/DataForScience/LangChain>

Question

- How was the technical level?
 - 1 — Too Low (too many details)
 - 2 — Low
 - 3 — Just Right
 - 4 — High
 - 5 — Too High (not enough details)

Question

- How was the level of Python code/explanations?
 - 1 — Too Low (too many details)
 - 2 — Low
 - 3 — Just Right
 - 4 — High
 - 5 — Too High (not enough details)

Events



graphs4sci.substack.com



Generative Artificial Intelligence with OpenAI

Jul 23, 2024 - 9am-1pm (PST)

Interactive Data Visualization with Python

Aug 7, 2024 - 10am-2pm (PST)

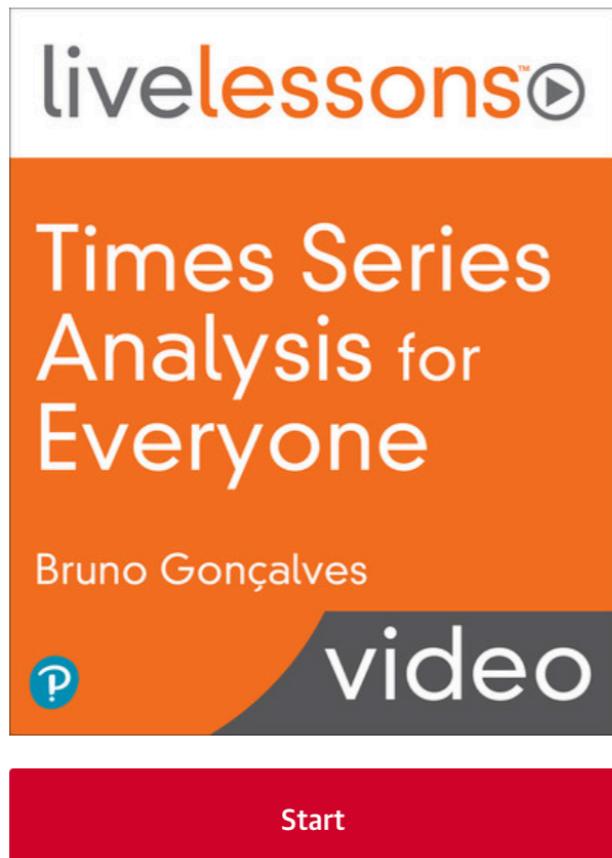
ChatGPT and Competing LLMs

Aug 13, 2024 - 10am-2pm (PST)

Times Series Analysis for Everyone

★★★★★ [1 review](#)

By [Bruno Gonçalves](#)



TIME TO COMPLETE:

6h

TOPICS:

[Time Series](#)

PUBLISHED BY:

[Pearson](#)

PUBLICATION DATE:

November 2021

https://bit.ly/Timeseries_LL

6 Hours of Video Instruction

The perfect introduction to time-based analytics

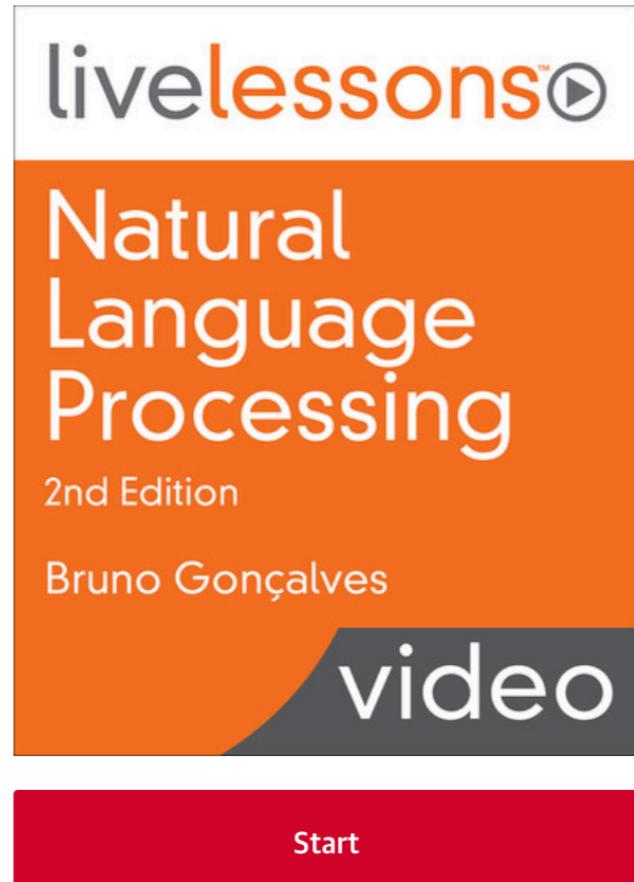
Overview

Times Series Analysis for Everyone LiveLessons covers the fundamental tools and techniques for the analysis of time series data. These lessons introduce you to the basic concepts, ideas, and algorithms necessary to develop your own time series applications in a step-by-step, intuitive fashion. The lessons follow a gradual progression, from the more specific to the more abstract, taking you from the very basics to some of the most recent and sophisticated algorithms by leveraging the statsmodels, arch, and Keras state-of-the-art models.

Natural Language Processing, 2nd Edition

Write the [first review](#)

By [Bruno Gonçalves](#)



TIME TO COMPLETE:

5h 23m

TOPICS:

[Natural Language Processing](#)

PUBLISHED BY:

[Addison-Wesley Professional](#)

PUBLICATION DATE:

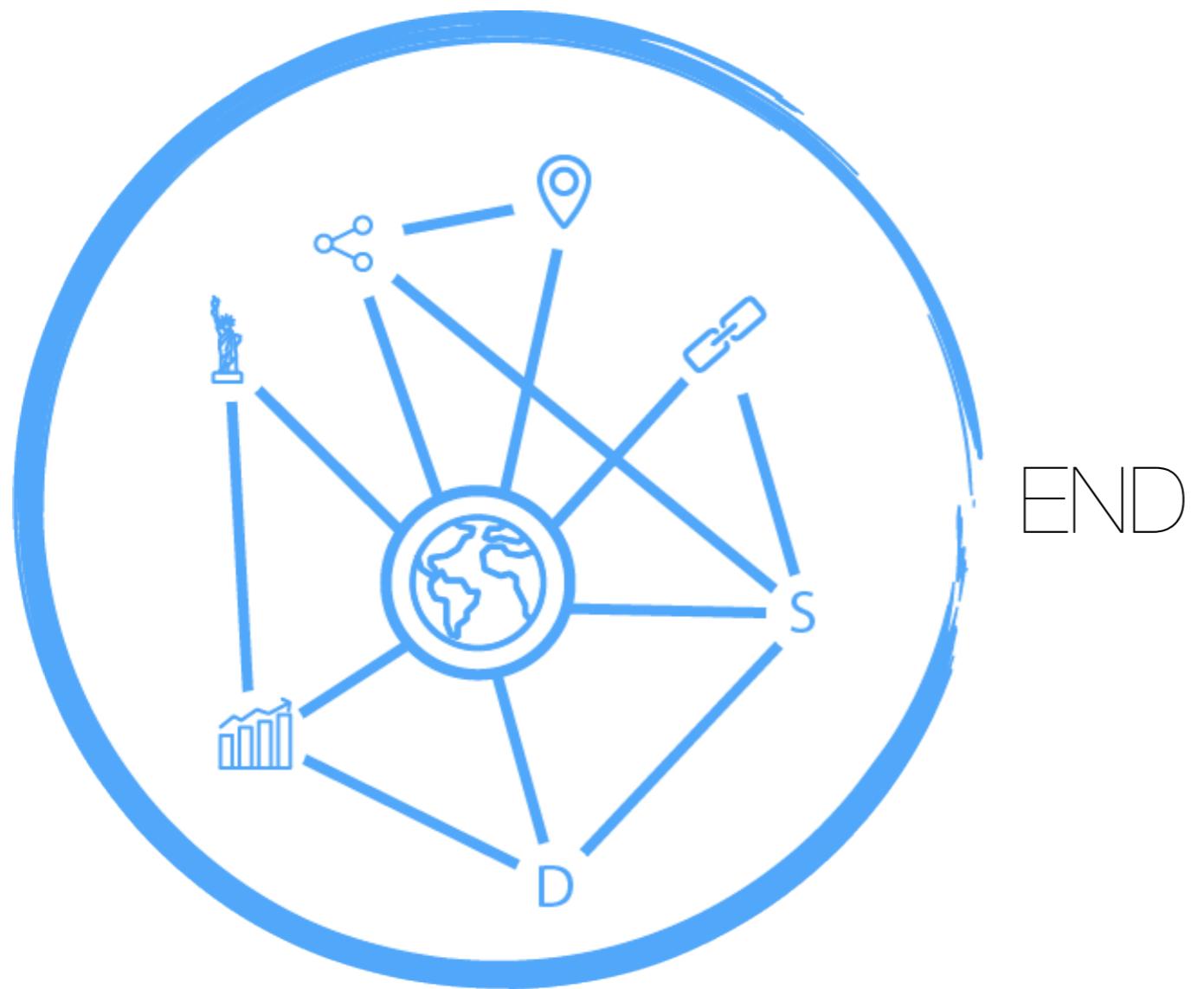
October 2021

https://bit.ly/NLP_LL

5 Hours of Video Instruction

Overview

Natural Language Processing LiveLessons covers the fundamentals of Natural Language Processing in a simple and intuitive way, empowering you to add NLP to your toolkit. Using the powerful NLTK package, it gradually moves from the basics of text representation, cleaning, topic detection, regular expressions, and sentiment analysis before moving on to the Keras deep learning framework to explore more advanced topics such as text classification and sequence-to-sequence models. After successfully completing these lessons you'll be equipped with a fundamental and practical understanding of state-of-the-art Natural Language Processing tools and algorithms.



END