

Generative AI Workloads

A Deep Dive into the Inner Workings of Autoregressive LLMs and How They Become AI Assistants

This file is meant for personal use by cwwdowns68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.
Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Agenda

- Large Language Models - Base vs Chat Models
- LLMs - The Training Procedure
- The LLM Inference Mechanism
- The Evolution of LLMs to Chat Assistants
- OpenAI's Assistant Training Recipe

This file is meant for personal use by cwdownloads68@gmail.com only.

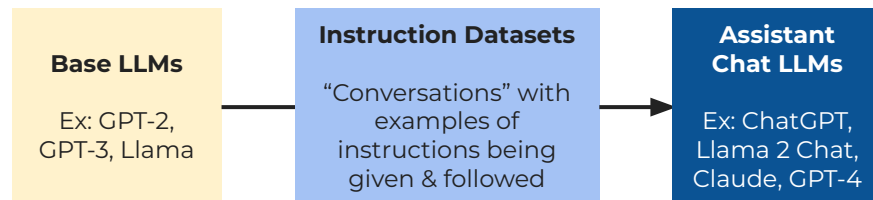
Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Large Language Models (LLMs)

Base Models vs Chat Models

- **LLMs are Deep-Learning based AI models built on the Transformer architecture.** They are trained on a humongous volume of text data to achieve **high-quality text understanding & text generation.**
- The training regimen of LLMs usually consists of a simple objective such as **predicting the missing (masked) word in large text corpuses** - just the process of learning how to predict this missing word very well has been shown to imbibe an impressive understanding of grammatical rules, good sentence construction, reasoning power and real-world knowledge.
- There are two variants of Large Language Models that we need to differentiate:
 - 1. Base LLMs** are purely trained on missing word prediction on generic internet text corpuses
 - 2. Assistant-style Instruction-tuned Chat LLMs** however, have additional training steps that involve missing word prediction on curated high-quality “instruction datasets”

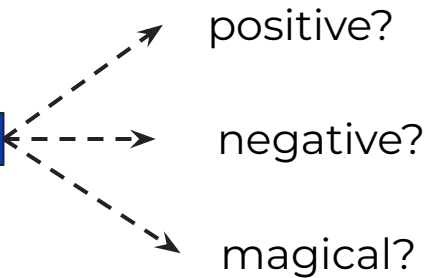


Large Language Models (LLMs)

LLMs are **train**ed using language modeling, that is, predicting the next word in a sequence.

Masked Sample

The movie was awesome. Overall, the experience was



This file is meant for personal use by cwwdowns68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

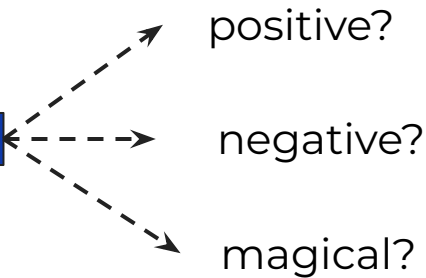
Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Large Language Models (LLMs)

LLMs are **train**ed using language modeling, that is, predicting the next word in a sequence.

Masked Sample

The movie was awesome. Overall, the experience was



positive

Ground Truth

Large Language Models (LLMs)

LLMs are **train**ed using language modeling, that is, predicting the next word in a sequence.

Masked Sample

The movie was **awesome**. Overall, the experience was **[]**

positive

negative

magical

Large Language Models (LLMs)

LLMs are **train**ed using language modeling, that is, predicting the next word in a sequence.

The movie is a visually stunning, action-packed, and emotionally resonant thrill ride that will leave you on the edge of the seat from the beginning to end. Overall, the experience was magical.

Vocabulary

positive $p = .03$

negative $p = .00001$

...

magical $p = .83$

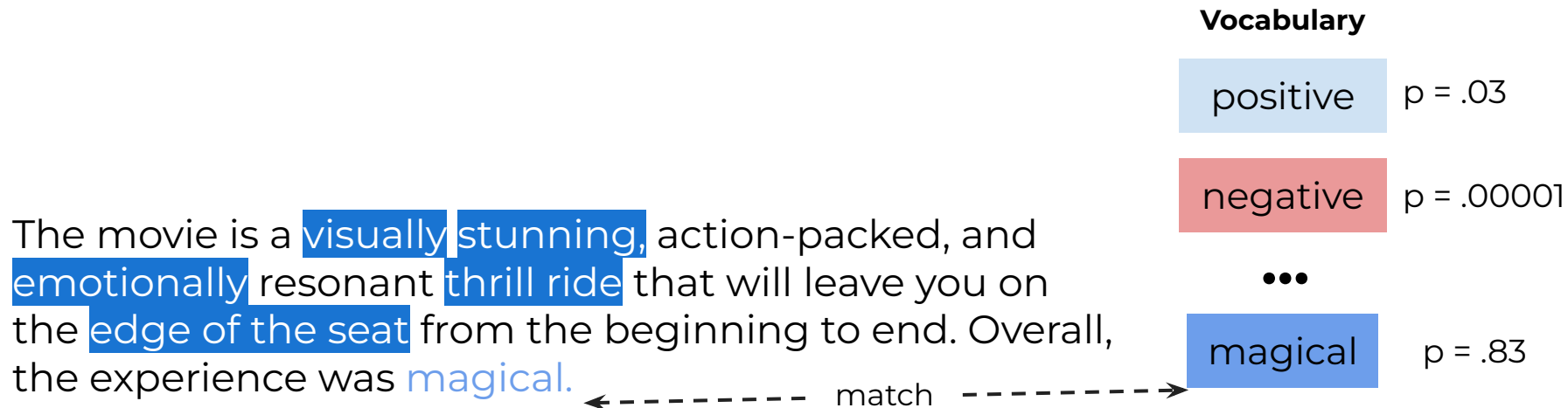
This file is meant for personal use by cwwdowns68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Large Language Models (LLMs)

LLMs are **train**ed using language modeling, that is, predicting the next word in a sequence.



This file is meant for personal use by cwwdowns68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Large Language Models (LLMs)

LLMs are **train**ed using language modeling, that is, predicting the next word in a sequence.

Original Sentence

The movie was awesome. Overall, the experience was positive.

Training Samples

The [REDACTED]

The movie [REDACTED]

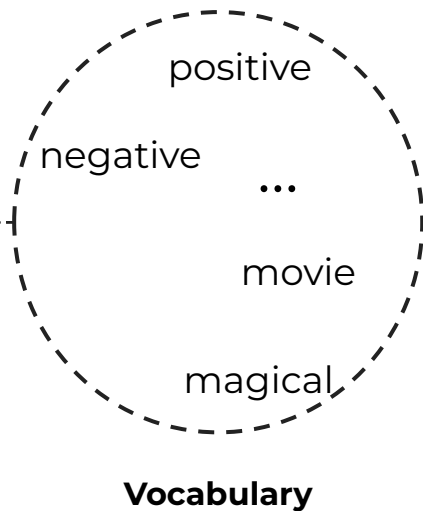
The movie was [REDACTED]

The movie was awesome. [REDACTED]

The movie was awesome. Overall, [REDACTED]

The movie was awesome. Overall, the [REDACTED]

⋮



Large Language Models (LLMs)

During ***inference***, the LLM predicts the next word in the input sequence.

Input Word = Prompt

The

Output, Word-by-Word

The []

The movie []

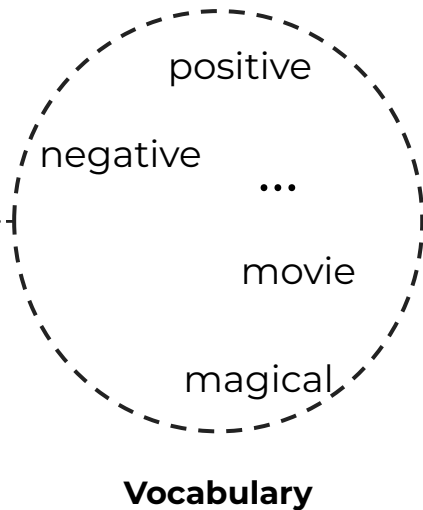
The movie was []

The movie was awesome. []

The movie was awesome. Overall, []

The movie was awesome. Overall, the []

⋮



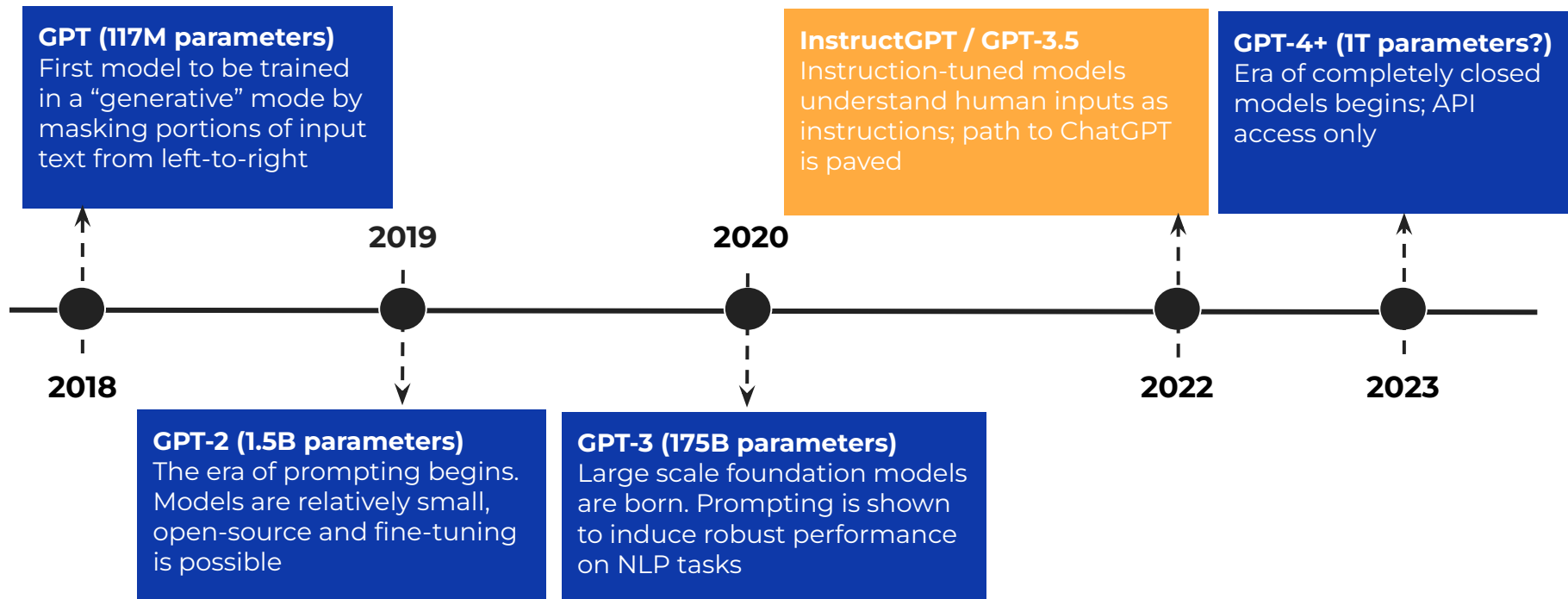
This file is meant for personal use by cwwdowns68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Large Language Models (LLMs)

Over the last 2 years, LLMs like ChatGPT have evolved to become **Conversational AI assistants**



This file is meant for personal use by cwwdowns68@gmail.com only.

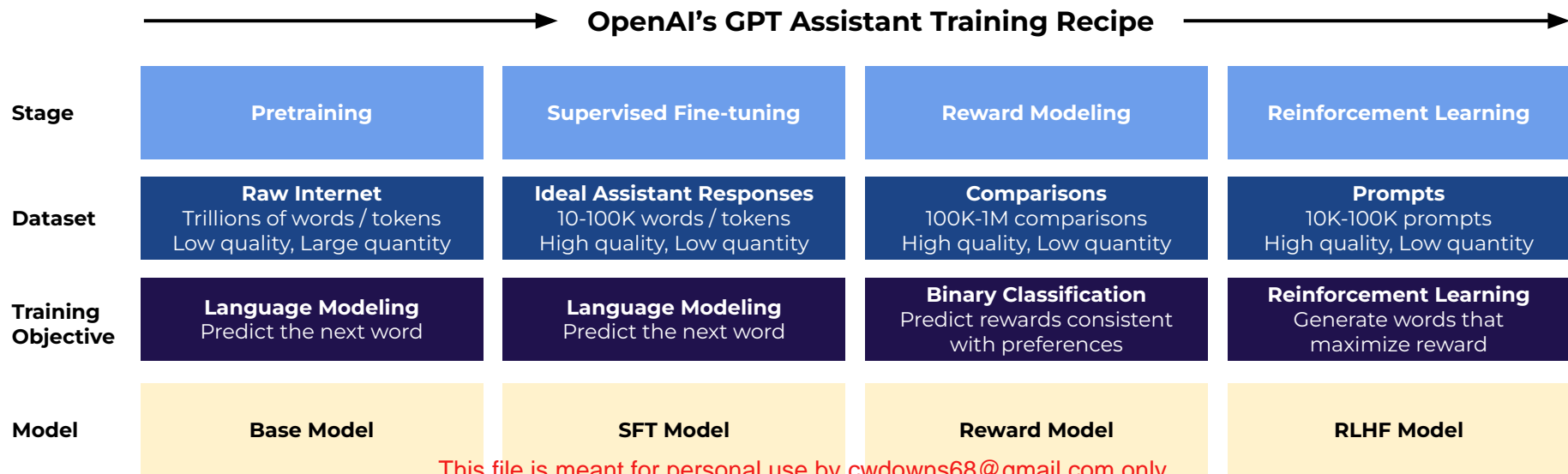
Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

The OpenAI GPT Assistant Training Recipe

High-Quality Assistant-style Chat Models

- **OpenAI's GPT-3.5** (behind the original ChatGPT of Nov 2022) and **GPT-4**, based on OpenAI's Generative Pre-trained Transformer (GPT) series, **first pioneered the multi-stage training procedure** for creating the **highest quality Instruction-tuned Chat Assistants**
 - honest, helpful and harmless by being aligned with human values.



This file is meant for personal use by cwdownloads68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

Summary

So in order to summarize:

- We got a glimpse into what the term Large Language Model (LLM) means, and we understood how to differentiate between Base LLMs and Chat LLMs
- We understood the training regimen and the next-word prediction objective behind Autoregressive LLMs, which governs how LLMs are trained to generate text one word at a time
- We also made the link from training to inference, and saw how the LLM inference process is merely the forward propagation step of the training process
- We also understood how LLMs (with this same Autoregressive Core Mechanism) have evolved from mere Text Completion Systems to Instruction-tuned Conversational AI Assistants
- Finally, we went through the modern multi-stage GPT Training Recipe pioneered by OpenAI's GPT series of LLMs, and understood each step of that multi-stage process in greater detail

This file is meant for personal use by cwwdowns68@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.