

LARGE LANGUAGE MODELS WITH GENERATIVE AI

Ritu Sulam, Apoorva Mallikarjuna Aradhya, Kalyani Biradar, Saisri Vishwanath
Under Prof. Lu Xiao

OVERVIEW

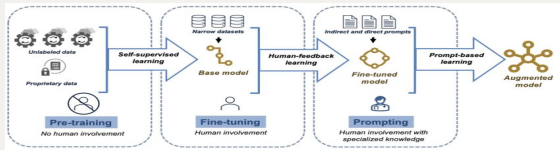
Large-scale language models are sophisticated AI systems engineered to interpret and produce language that closely resembles human communication. Generative AI specializes in the creation of novel content, data, or outputs, as opposed to merely analyzing existing information or making decisions based on set criteria. The integration of generative AI with large-scale language models yields a powerful tool capable of comprehending and crafting language in a diverse array of applications.

LANGUAGE MODELS & GENERATIVE AI

Generative Pre-trained Transformers (GPTs) represent a type of neural network model architecture employed in natural language processing (NLP) tasks. The Transformer model architecture is frequently utilized in large language models (LLMs). Notable LLMs include Google's PaLM (Pathways Language Model), BERT also from Google (Bidirectional Encoder Representations from Transformers), and OpenAI's GPT series. In generative artificial intelligence, models employ probabilistic algorithms to create outputs such as text, images, or audio.

LLM TRAINING STAGES

Large Language Models (LLMs) are crucial components in the fields of Natural Language Processing (NLP) and Natural Language Generation (NLG). These models are developed through a comprehensive training methodology, which includes initial pre-training on extensive datasets followed by the application of methods such as in-context learning, zero-shot, one-shot, and few-shot learning, as well as fine-tuning processes. The performance of LLMs is significantly dependent on the caliber of the pre-training data they utilize, necessitating a superior quality of data compared to smaller-scale language models.

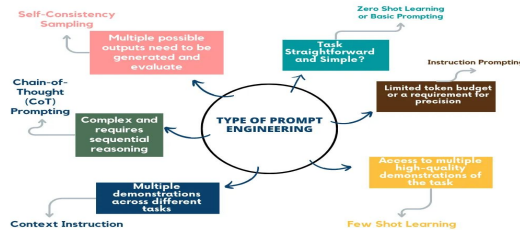


The stages of LLM training encompass several crucial steps:

- Data Collection and Preprocessing:** Gather diverse datasets from multiple sources like books, websites, and articles. Prepare and clean the data, performing tasks such as normalization, lowercasing, removing stopwords, and tokenization for subsequent analysis.

- Selection and Configuration of Model:** Choose a sophisticated model like OpenAI's GPT-3.5 or Google's BERT, primarily based on the Transformer architecture. Customize the model's structure by defining elements such as layer depth, attention mechanisms, loss functions, and hyperparameters, all tailored to the specific task and dataset. These configurations will directly influence the training duration and overall performance of the model.
- Model Training:** Utilize supervised learning to train the model in predicting subsequent words or sequences based on the provided text. Training these expansive models requires substantial resources and often involves techniques like model parallelism for efficiency. Alternatively, fine-tuning existing models proves more economical.
- Fine-tuning and Evaluation:** Evaluate the performance of the trained model using a distinct test dataset to measure its effectiveness. Improve the model's efficiency through fine-tuning, this could include making changes to hyperparameters, adjusting architectural components, or incorporating additional training based on insights obtained from the evaluation.

PROMPT ENGINEERING



The role of prompt engineering becomes vital, as it involves devising optimized input combinations for language models to improve the quality of their responses. This enhancement is achieved without having to resort to significant updates in parameters or extensive fine-tuning.

One strategy in prompt engineering is to integrate example answers within the prompts themselves. These examples are categorized into: Zero-Shot Learning (no examples provided), One-Shot Learning (one example given), and Few-Shot Learning (more than one example provided). Incorporating a diverse range of examples often leads to more accurate responses. These examples can be further broken down into seven distinct types based on their structure, functionality, and complexity, as shown in the figure.

In a given prompt, factors such as the structure, training instances, and their organization play a crucial role in determining the model's precision. Diverse decisions in these aspects result in different degrees of accuracy. When examining language models, biases emerge during few-shot learning, including majority label bias, recency bias, and common token bias, affecting the distribution of the model's outputs. To mitigate these biases, adjusting the output distribution entails gauging the model's bias using test inputs that lack specific content.

CHALLENGES

Generative AI, particularly Large Language Models (LLMs), grapples with challenges like high training costs, struggles with adapting to new data, and phenomena like "hallucination," leading to fabricated information. The use of LLM-generated content poses risks, demanding robust safeguards against harmful information dissemination, hate speech, and security threats. Privacy concerns, highlighted by inadvertent data leaks, emphasize the necessity for stringent privacy measures and responsible LLM usage to protect sensitive information.

CONCLUSION

In summary, LLMs with Generative AI are powerful tools with great potential, but they need to be used responsibly and ethically. Prompt engineering is a key technique for improving LLM outputs. LLMs can revolutionize the way we interact with computers and the world around us.

REFERENCES

- Feuerriegel, Stefan & Hartmann, Jochen & Janiesch, Christian & Zschech, Patrick. (2023). Generative AI.
- Brown, Tom B., et al. "Language Models are Few-Shot Learners." arXiv preprint arXiv:2005.14165, 2020.
- Gozalo-Brizuela, Roberto, and Eduardo C. Garrido-Merchan. "ChatGPT is not all you need. A State of the Art Review of large Generative AI models." arXiv preprint arXiv:2301.04655, 2023.
- Jeong, Cheonsu. "A Study on the Implementation of Generative AI Services Using an Enterprise Data-Based LLM Application Architecture." arXiv preprint arXiv:2309.01105, 2023.
- Yang, Jingfeng, et al. "Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond." arXiv preprint arXiv:2304.13712, 2023.
- Zhao, Tony Z., et al. (2021). Calibrate Before Use: Improving Few-Shot Performance of Language Models. arXiv preprint arXiv:2102.09690.