

PaLM Overview and Versions

Explain PaLM

PaLM, or the Pathways Language Model, is a series of advanced large language models developed by Google AI. The original PaLM, introduced in April 2022, is a dense decoder-only Transformer model with 540 billion parameters. It was trained on a diverse dataset of 780 billion tokens, including filtered web pages, books, Wikipedia articles, news articles, source code from GitHub, and social media conversations. This extensive training enabled PaLM to excel in various tasks such as commonsense reasoning, arithmetic reasoning, joke explanation, code generation, and translation. Notably, when combined with chain-of-thought prompting, PaLM demonstrated enhanced performance on multi-step reasoning tasks like word problems and logic-based questions.

ARXIV.ORG

In May 2023, Google announced PaLM 2, the next iteration of the model. PaLM 2 is reported to be a 340 billion-parameter model trained on 3.6 trillion tokens. It offers improved multilingual, reasoning, and coding capabilities, excelling in advanced reasoning tasks, including code and math problem-solving, classification, and question answering. PaLM 2 is built on a foundation of compute-optimal scaling, an improved dataset mixture, and refined model architecture, making it more efficient and versatile than its predecessor.

ARXIV.ORG

Google has also developed specialized versions of PaLM for specific domains. Med-PaLM, for instance, is fine-tuned on medical data and outperforms previous models on medical question-answering benchmarks. It was the first to obtain a passing score on U.S. medical licensing questions and can provide reasoning and self-evaluate its responses. Another variant, PaLM-E, integrates a vision transformer to create a state-of-the-art vision-language model used for robotic manipulation, enabling competitive performance in robotics tasks without the need for retraining or fine-tuning.

EN.WIKIPEDIA.ORG

PaLM's development underscores Google's commitment to advancing AI capabilities across various domains, ensuring responsible and ethical deployment of these powerful models.

Sources  W

The Pathways Language Model (PaLM) is a large-scale language model developed by Google AI, designed to advance natural language understanding and generation. Its architecture is based on the Transformer model, specifically utilizing a dense, decoder-only configuration. This design choice focuses on generating text by predicting subsequent tokens in a sequence, a method that has proven effective in various language modeling tasks.

Key Architectural Features:

- **Model Scale:** The original PaLM boasts 540 billion parameters, making it one of the largest language models to date. This extensive scale enables the model to capture intricate patterns and nuances in language.
- **Training Data:** PaLM was trained on a diverse dataset comprising 780 billion tokens. This dataset includes filtered web pages, books, Wikipedia articles, news articles, source code from GitHub, and social media conversations. The inclusion of social media data, which constitutes 50% of the corpus, enhances the model's conversational capabilities.
- **Training Infrastructure:** The training process utilized two TPU v4 Pods, each equipped with 3,072 TPU v4 chips, totaling 6,144 chips. This setup allowed for efficient large-scale training, achieving a hardware FLOPs utilization of 57.8%, a record for models of this magnitude.
- **Pathways System:** PaLM was trained using Google's Pathways system, which enables efficient training across multiple TPU Pods. This system facilitates the scaling of large models by optimizing resource utilization and streamlining the training process.

In May 2023, Google introduced PaLM 2, an enhanced version of the original model. PaLM 2 features 340 billion parameters and was trained on 3.6 trillion tokens. It offers improved multilingual, reasoning, and coding capabilities, demonstrating significant advancements over its predecessor.

The architecture of PaLM and its successor, PaLM 2, underscores the importance of scale, diverse training data, and efficient infrastructure in developing advanced language models capable of understanding and generating human-like text.