

Adapting Large Language Model for JSON Extraction from Text Corpora

Van-Tuan Tran ^{1*}, Chin-Shiuh Shieh ¹, Ying-Chieh Chao ³, Casper Tsai ² and Mong-Fong Horng ¹

¹ National Kaohsiung University of Science and Technology, Taiwan

² Auray Technology Corp., Taiwan

³ ICP DAS CO., Ltd., Taipei, Taiwan

*Email: fl11169109@nkust.edu.tw

Abstract

This paper explores the adaptation of large language models (LLMs) for the task of extracting structured data in JSON format from extensive and unstructured text corpora. Specifically, we fine-tuned the Llama-2-7 billion model utilizing a combination of advanced techniques, including QLoRA (Quantized Low-Rank Adaptation), Fully Sharded Data Parallel (FSDP) training, and distributed training across multiple GPUs. Our approach addresses the challenges associated with processing large and complex datasets, which include HTML-tagged content, diverse textual paragraphs, and various formatting irregularities. By employing QLoRA, we achieve efficient low-rank adaptation, which helps in reducing the computational burden while retaining model performance. FSDP enables us to handle large-scale data by sharding model parameters and gradients, thus optimizing memory usage and speeding up the training process. Multi-GPU training further enhances scalability and accelerates the fine-tuning process, allowing us to manage and process extensive text corpora effectively. The fine-tuning process transforms raw textual information into structured JSON outputs, facilitating automated data extraction and processing. Our results demonstrate that the adapted model significantly improves the accuracy and efficiency of extracting relevant information from complex text sources, compared to traditional methods. This work not only showcases the effectiveness of combining these advanced techniques but also offers a scalable solution for applications such as web scraping, data parsing, and large-scale information retrieval.

Keywords: Fully Sharded Data Parallel (FSDP), Meta-Llama/Llama-2-7b, Meta-Llama/Llama-3-8b, Meta-Llama/Llama-3.1-8b, QLoRA, LoRA.