

## Case Study: Enhancing Data Querying Efficiency with Llama 2-7B Model

### *Executive Summary:*

In the ever-evolving landscape of data engineering and data analysis, organizations are continually seeking advanced tools to streamline their data-related tasks. This case study examines the implementation of the **Llama 2-7B model, fine-tuned on a specific dataset tailored for the fintech domain using Next AI console**, and compares its performance to ChatGPT, a widely-used language model. The results demonstrate that Llama 2-7B outperforms ChatGPT, offers enhanced security when hosted on an in-house GPU or trusted cloud, and is significantly more cost-effective, making it an invaluable asset for data teams.

### **Introduction:**

In this case study, we delve into how the Llama 2-7B model, fine-tuned on a custom dataset tailored for data engineering and data analysis within the fintech sector, surpassed ChatGPT in terms of query efficiency, security, and cost-effectiveness.

### **Methodology:**

- **Model Selection:** Our data team decided to assess the Llama 2-7B model, a cutting-edge NLP model, and fine-tuned it using a fintech-specific dataset for a direct comparison with ChatGPT fine-tuned for the same task.
- **Data Preparation:** We meticulously curated a high-quality dataset specific to SQL queries and data analysis tasks within the fintech industry to fine-tune both models, ensuring they were trained on relevant data.
- **Benchmarking:** The models were benchmarked on several key metrics crucial to data professionals, including query accuracy, response quality, task completion time, and cost-effectiveness, utilizing a standardized evaluation framework.

### **Results:**

#### **Performance Improvement:**

- Llama 2-7B consistently outperformed ChatGPT across all benchmarked metrics, displaying a significantly better understanding of fintech-specific data queries.
- Data team members reported that Llama 2-7B generated more accurate and contextually relevant SQL queries, leading to a substantial improvement in query efficiency and productivity.

**Security:**

- Hosting Llama 2-7B on an in-house GPU or a trusted cloud infrastructure provided an elevated level of security. The model's parameters and data were kept under the organization's direct control, reducing the risk of data breaches and ensuring compliance with strict privacy regulations, critical for fintech applications.
- ChatGPT, while secure to a certain extent, had limitations in terms of data control and customization, making it less adaptable to the stringent security requirements of the fintech sector.

**Cost-effectiveness:**

- The cost analysis revealed that using Llama 2-7B was approximately one-third the cost of relying on ChatGPT for similar data engineering and data analysis workloads. This substantial cost reduction allowed the organization to allocate resources more efficiently and invest in further data-driven research and development.
- The open-source nature of Llama 2-7B and the availability of pretrained checkpoints with quantization and optimized fine-tuning techniques further reduced costs associated with model development and maintenance, making it an economically sound choice.

**Conclusion:**

This case study underscores the advantages of adopting the Llama 2-7B model, fine-tuned on a dataset tailored for data engineering and data analysis in the fintech domain, over using ChatGPT in terms of query efficiency, security, and cost-effectiveness. Llama 2-7B's superior performance, coupled with its robust security measures and lower operational costs, positions it as an indispensable tool for data teams in the fintech industry.

By deploying Llama 2-7B on an in-house GPU or a trusted cloud infrastructure, organizations can fully leverage the potential of this model while maintaining control over their data, ensuring compliance, and realizing significant cost savings. As the field of data engineering and analysis continues to evolve, harnessing models like Llama 2-7B with customized fine-tuning will play a pivotal role in driving innovation and efficiency in fintech applications.

### Performance Metrics:

<b>Model</b>	<b>Query Accuracy</b>	<b>Response Quality</b>	<b>Task Completion Time</b>	<b>Cost</b>
Llama 2-7B with fine-tuning using Next AI	High	Accurate	Fast	1/3 of ChatGPT
ChatGPT	Moderate	Moderate	Moderate	Expensive