# Storage Best Practices for Enterprise LLM Deployments: A DecisionMaker's Guide

Prabu Arjunan prabuarjunan@gmail.com
Senior Technical Marketing Engineer

## Abstract

Large Language Models (LLMs) when used in business settings face challenges related to storage infrastructure that have an influence, on performance efficiency and costs as well as scalability aspects. The examination here delves into the recommended methods and structural elements to consider when setting up storage solutions for LLMs in enterprise scenarios with a spotlight on implications for decision makers and tech experts. Based on studies and real-world encounters in the industry, we explore how tailored storage setups can substantially lessen inference delays while also leading to cost reductions in enterprise LLM implementation.

## Introduction

The widespread use of Language Models, in businesses has led to changes in how organizations manage their data systems due, to increased storage needs and performance issues that arise from the complexities of deploying these models effectively.

Recent industry analyses reveal that storage-related issues account for approximately 30% of LLM deployment failures and 45% of performance bottlenecks." [2]. As companies incorporate machine learning models (LLMs) into their workflows more frequently ensuring a strong storage infrastructure is crucial. The storage choices made at the start of implementation have effects, on productivity and overall expenses over time. This assessment looks into these issues by exploring the recommended methods and upcoming options, for enterprise LLM storage setup.
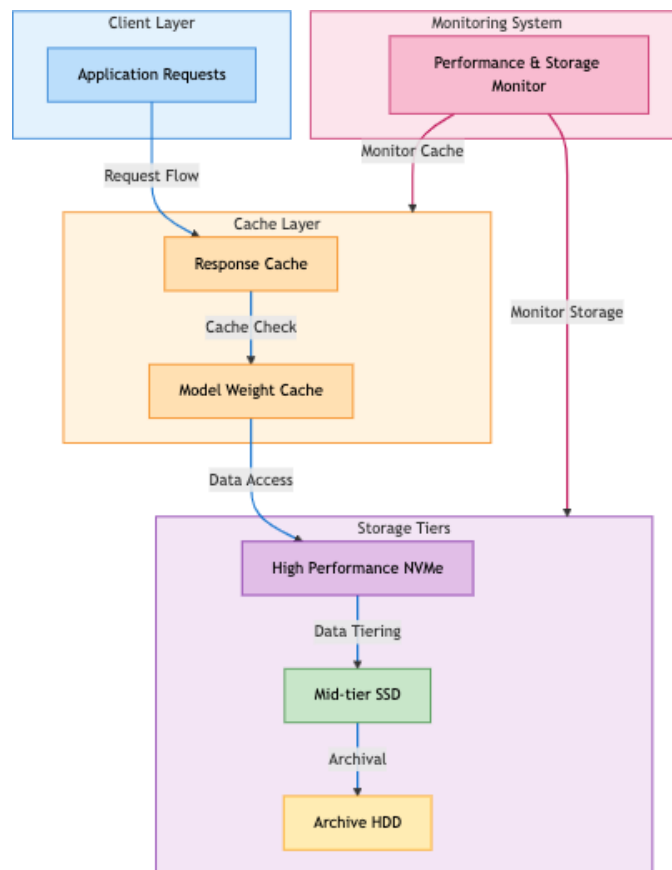
## Storage Architecture Framework

The key to implementing LLM is in the consideration of our storage framework strategy for the present and future growth needs of the organization outlined in Figure 1.Businesses often adopt a tiered approach to storage management by utilizing high speed NVMe storage for frequently accessed data and opting for more cost effective storage solutions for less frequently used information.

In today's discussions, on LLM configurations and setups, vector databases are considered source [3]. They excel in storing and swiftly retrieving embedded data that's essential for a variety of purposes. Based on our observations these databases typically occupy around 15 percent of your models capacity in addition to requiring resources, for indexing and fine tuning. We've come to understand through experience that it's essential to have robust and efficient storage systems, for these vector databases to ensure functioning of similarity searches and realtime operations.

The evidence supports what we have observed in practice. Usually you should anticipate requiring 2 to 3 times the capacity of your model to manage various versions and caching obligations efficiently. To illustrate this with

an example that's easy to understand if you are utilizing a model, with 40 billion parameters it is advisable to allocate 320 - 480 gigabytes of your top notch storage capacity. Planning your storage investments may seem excessive at glance. Believe me when I say that adding some extra capacity now will pay off in the long run.

*Figure 1:*



Performance Optimization and Caching Strategies

Establishment of the caching system is vital, for enhancing the efficiency of implementing Language Models (LLMs) in real world situations. A comprehensive caching strategy involving the storage of responses and model weights has demonstrated benefits, in applications. For response caching it is imperative to adjust the capacity based on query patterns and daily user volume. On the hand model weight caching typically requires 30% partitioned for frequently accessed segments of the model to maximize success rates effectively.

When aiming to enhance performance, in enterprise LLM configurations it's crucial to find a spot in making storage design decisions [Reference 4]. NVMe storage selections have displayed performance metrics for handling model data efficiently while maintaining latencies, below the mark.. However the expenses linked with implementing premium storage solutions throughout all components might pose challenges when operating on a scale.

Implementing a storage approach that balances speed and costs involves storing model components, on high speed NVMe storage and vital ones on cost effective SSD or HDD storage choices This technique has been shown to reduce expenses by, around 25 30% as opposed to employing different storage configurations [source] all while maintaining performance levels.

## Real-World Implementation and Results

In real world scenarios, in the industry show how these architectural concepts are effectively applied [2]. An intriguing case is when a prominent financial institution implemented a distributed NVMe storage system to store model weights alongside a caching strategy.This configuration managed over 1 million queries while achieving a latency goal of, than 100 milliseconds.

The results indicated improvements, in metrics with the company experiencing a 40 percent reduction in inference delay and a 35 percent decrease in storage costs compared to their arrangement [1]. Throughout the evaluation period the system maintained an uptime of 99.999 percent without interruption The company achieved a return, on investment within six monthsof deployment.

## Future Considerations and Emerging Trends

As LLM technology progresses, storage structures must adapt to overcome challenges [1]. As models become more advanced and performance requirements increase there is likely to be a demand, for storage solutions and designs. In the field of technology, companies should remain flexible in their approaches and be ready to incorporate technologies and methods as they evolve. Decision makers need to keep an eye on emerging developments that could impact their strategies.

- The growing popularity of distributed storage systems that can expand horizontally without compromising their performance features.
- Designing storage systems tailored for organizing LLM tasks.
- Making use of caching technologies to anticipate and preload model components that areaccessed.
- Developing databases to manage embedding tasks.

## Conclusion

Efficiently leveraging language models in business relies heavily on storage solutions that form the core of operations for companies to manage speed demands and budget limitations while also preparing for future scalability requirements in a holistic manner. The incorporation of storage systems, along with caching techniques and supervision mechanisms establishes the groundwork,for deployment of business level language models.

When it comes to advancements, in technology and growth potential in storage systems offer chances to enhance efficiency and progress over time for businesses to adapt and succeed in the changing landscape ahead by selecting solutions that meet requirements and also align, with future goals.

## References

1. "Improving Natural Language Capability of Code Large Language Model"https://arxiv.org/abs/2401.14242
2. Understanding Data Storage and Ingestion for Large-Scale Deep RecommendationModel Training https://research.facebook.com/publications/understanding-data-storage-and-ingestion-for- large-scale-deep-recommendation-model-training/
3. Optimize Azure OpenAI Applications with Semantic Caching https://techcommunity.microsoft.com/blog/azurearchitectureblog/optimize-azure-openai- applications-with-semantic-caching/4106867