

Building Enterprise GenAI Infrastructure: A Practical Guide to Storage Requirements

Prabu Arjunan

Senior Technical Marketing Engineer

prabuarjunan@gmail.com

Abstract

In the paced world of business today as companies rush to embrace cutting edge Generative AI (GenAI) technologies such, as ChatGPT and custom language models they often encounter a yet underestimated challenge. Establishing a solid groundwork to support these advanced tools properly. This document delves into the fundamental storage considerations that business executives and IT decision makers need to grasp in order to ensure the integration of GenAI solutions into their operations effectively. Drawing insights from real life implementation experiences in Fortune 500 firms and latest research findings on large scale AI systems [1] I offer advice, for organizations at any phase of their AI exploration journey. When it comes to implementing intelligence, in enterprise storage and infrastructure planning, for data management and storage architecture design.

Keywords: Generative AI, Enterprise Storage, Infrastructure Planning, AI Implementation, Data Management, Storage Architecture

Introduction

Launching a GenAI project with the aim of revolutionizing customer service may lead to slow response times as the user base expands rapidly over time. Investments made in AI model development can quickly reveal challenges when it comes to storing and managing the increasing amount of data effectively. Such situations are common in businesses globally.

Typically stem from preparation, for the necessary infrastructure to sustain AI initiatives. Enterprises are facing hurdles when trying to expand their GenAI projects beyond the pilot phase due to a mistake of treating GenAI infrastructure needs like business applications.

Patterson et al.'s [3] recent study sheds light on the storage requirements of large AI models and explains why conventional infrastructure methods may not be sufficient.

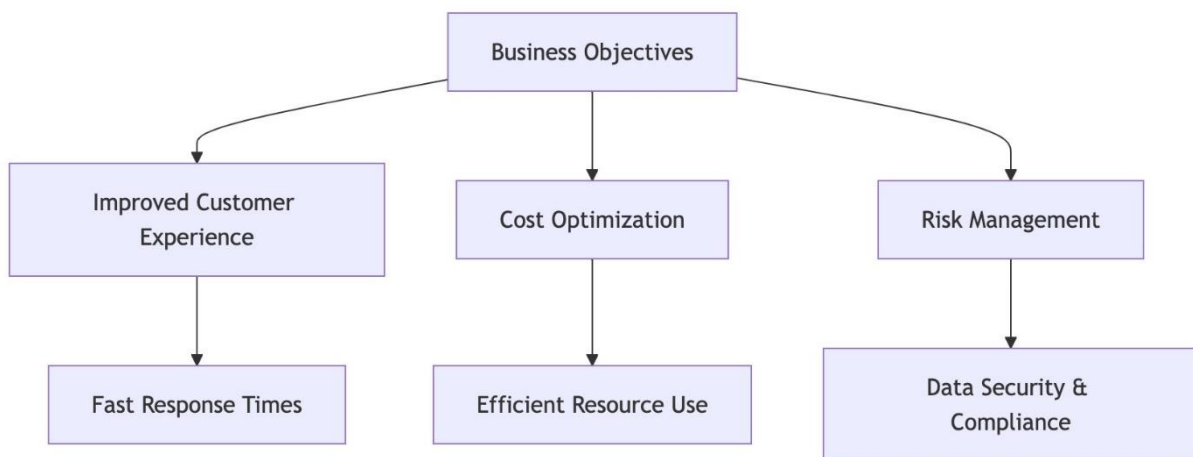
Understanding the Business Impact of Storage Decisions

Storing data systems may appear to be a tech matter at glance. However, its influence extends well beyond the realm of IT alone. The need for cutting edge infrastructure in today's AI models such as PaLM [4] is unparalleled in both size and intricacy. Imagine a institution that integrated a GenAI based customer support system recently. Their promising start with a trial initiative quickly spread throughout the company. Yet as usage expanded they faced obstacles that had a direct impact on their business results: Customers now have to wait a minutes for responses compared to the quick replies within seconds.

- The expenses for the system increased more than anticipated.
 - Keeping an eye on and reviewing AI interactions proved to be a challenge for compliance teams.
- The difficulties arose not from the AI technology per se, but from the infrastructural choices in the project's early phases, according to Raffel et al.'s research [5]. Their study emphasizes the importance of infrastructure planning to uphold AI systems' performance as they expand in scale.

The Business Case for Proper Storage Planning

Figure 1:



When constructing GenAI infrastructure, companies need to ensure that technical choices are inline with the goals of the business. This synchronization emphasizes three areas:

Customer Experience and Operational Efficiency

Crafted storage systems guarantee that GenAI apps remain responsive and reliable for any number of users—be it one or a thousand individuals being served at once. Take a corporation's AI-driven product recommendation platform as an illustration; it must deliver responses without delay both in everyday scenarios and during busy shopping periods.

Cost Management and Scalability

Cost control and adaptability are factors to consider for GenAI storage requirements as they can expand significantly over time, unlike traditional IT setups. For instance, a healthcare institution that integrated AI for analyzing records initially had 10 terabytes of data, but quickly scaled up to 100 terabytes in just half a year. To accommodate growth without the need for substantial initial expenditures, organizations must opt for storage solutions that can scale efficiently.

Risk Management and Compliance

GenAI applications deal with business information that requires a storage system meeting governance standards and ensuring security and compliance measures are in place for proper handling of data. This is especially vital in industries with regulations that mandate transparency and auditability of AI decisions. As organizations expand their use of GenAI technology, Microsoft's framework for assessing AI security risks underscores the significance of securing the lifecycle of AI systems[6], from data input to model training and deployment processes. Our observation aligns with this notion, emphasizing the

need for storage infrastructure to cater to the security and compliance demands throughout the AI pipeline.

A Practical Approach to Implementation

The path to successful GenAI infrastructure doesn't require massive upfront investments or complete system overhauls. Instead, organizations should follow a pragmatic, phased approach that aligns with business objectives and manages risks effectively.

Phase 1: Strategic Planning

Start by addressing inquiries related to the business sector:

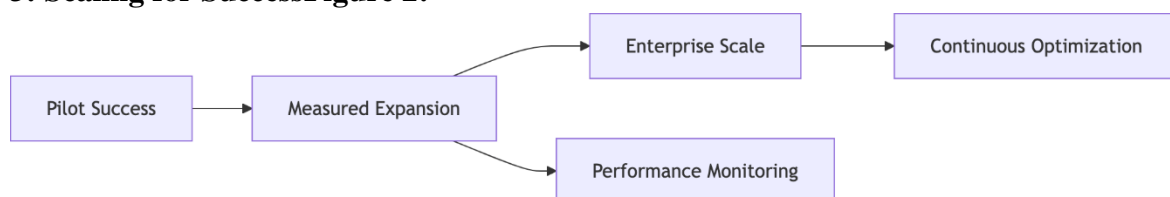
- What are the main ways you plan to use GenAI?
- How many individuals is the system designed to accommodate in the short and longrun?
- What kinds of data will your GenAI applications be handling?
- What performance standards and measures of success do you have in place?

Phase 2: Pilot Implementation

Begin by starting small but always keep growth opportunities in mind! A prominent manufacturing firm initiated their foray into AI technology with a targeted trial in their quality assurance division. They opted for a storage option that could begin on a small scale yet expand seamlessly when needed to accommodate their evolving needs:

- Confirm their method with minimal risk
- Collect real-world usage data
- Adjust their needs according to what they learned from real-world practice
- Develop expertise within the company over time

Phase 3: Scaling for Success



As your GenAI projects demonstrate success and growth opportunities arise, it's essential to expand in a way driven by metrics. A renowned consulting firm expanded its GenAI infrastructure successfully by following these steps:

- Keeping track of key metrics and performance indicators.
- Consistently reviewing and fine-tuning capacity plans.
- Ensuring that business needs are well-matched with the infrastructure capabilities.
- Implementing methods for allocating costs and charging back expenses.

Best Practices for Business Leaders

1. Align Infrastructure with Business Strategy

Prioritize achieving business objectives over fixating on technical details when setting up your storage system to align with your company's GenAI targets – be it enhancing customer satisfaction levels,

speeding up product innovation, or streamlining operational processes. As highlighted by Brown et al. [1], their study underscores the significance of infrastructure planning in facilitating seamless model implementation and expansion.

2. Build for Change

The field of AI is constantly advancing. Select options that allow for adjustments to meet evolving needs and technological advancements without starting from scratch. Rajbhandari et al. [2] provides insights on optimizing infrastructure for scalability.

3. Monitor and Measure

Track essential business metrics like application response times, user satisfaction levels, and total cost of ownership. Patterson et al. [3] offer frameworks for assessing and enhancing AI infrastructure efficiency.

Conclusion

Success in enterprise GenAI initiatives depends heavily on the foundational infrastructure decisions made early in the journey. By understanding the business implications of storage requirements and following a measured, practical approach to implementation, organizations can build scalable, cost-effective infrastructure that supports their AI ambitions. The key is not to over-engineer initial solutions but to build flexible foundations that can evolve with your business needs. Organizations that approach their GenAI infrastructure with this mindset are better positioned to achieve sustainable success in their AI initiatives.

References

1. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., . . . Amodei, D. (2020, May 28). Language Models are Few-Shot Learners. arXiv.org. <https://arxiv.org/abs/2005.14165>
2. ZeRO: Memory optimizations Toward Training Trillion Parameter Models. (2020, November 1). IEEE Conference Publication | IEEE Xplore. <https://ieeexplore.ieee.org/document/9355301>
3. Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L., Rothchild, D., So, D., Texier, M., & Dean, J. (2021, April 21). Carbon emissions and large neural network training. arXiv.org. <https://arxiv.org/abs/2104.10350>
4. Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., Schuh, P., Shi, K., Tsvyashchenko, S., Maynez, J., Rao, A., Barnes, P., Tay, Y., Shazeer, N., Prabhakaran, V., . . . Fiedel, N. (2022, April 5). PaLM: Scaling Language Modeling with Pathways. arXiv.org. <https://arxiv.org/abs/2204.02311>
5. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P. J. (2019, October 23). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. arXiv.org. <https://arxiv.org/abs/1910.10683>
6. Kumar, R. S. (2021, December 9). "Best practices for AI security risk management," Microsoft Security Blog. <https://www.microsoft.com/en-us/security/blog/2021/12/09/best-practices-for-ai-security-risk-management/>