



Real-time, Managed Latency of Storage and ATTO's Specialized Solutions for AI Applications

By Tim Klein - President & CEO & Co-founder at ATTO Technology, Inc.

As artificial intelligence (AI) applications continue to evolve, the need for real-time, managed latency of storage becomes increasingly crucial. This article explores the impact of storage latency on AI applications and highlights the specialized solutions provided by ATTO Technology, a leading provider of high-performance storage and network connectivity solutions. I discuss the challenges posed by storage latency, the importance of real-time latency management, and how ATTO's [Host Bus Adapters \(HBAs\)](#), [Network Interface Cards \(NICs\)](#), and [Storage Appliances](#) optimize latency for AI workloads. Through intelligent features, software optimizations, and high-speed connectivity, ATTO empowers organizations to achieve optimal performance and responsiveness in AI systems.

1. Introduction

Artificial intelligence applications heavily rely on the availability and access speed of data stored in various storage mediums. Storage latency, the delay in data retrieval, directly impacts the performance and responsiveness of AI systems. This article explores the significance of real-time, managed latency and its implications for AI applications. We delve into the challenges posed by storage latency and introduce

ATTO Technology's specialized solutions, including HBAs, NICs, and Storage Appliances, designed to optimize latency and enhance AI application performance.

2. The Impact of Storage Latency on AI Applications

Performance Degradation: High storage latency results in performance degradation, affecting AI algorithms' execution, inference times, and overall system efficiency. Real-time interaction and decision-making applications, in particular, are highly sensitive to latency, as delays can significantly impact user experience and system responsiveness.

Data Access Bottlenecks: AI applications rely on large datasets stored in various storage mediums. Inefficient storage architectures or high latency between storage layers create data access bottlenecks, limiting the throughput and scalability of AI applications. Rapid data retrieval and processing become paramount as AI workloads become more demanding.

3. ATTO's Specialized Solutions for Real-time, Managed Latency

ATTO HBAs: ATTO HBAs are purpose-built to deliver high-speed, low-latency connectivity between servers and storage devices. They utilize industry-standard protocols, such as Fibre Channel, SAS, and NVMe, to ensure reliable and efficient data transfers with minimal latency. Intelligent features like Advanced Data Streaming (ADS™) technology optimize I/O performance and reduce latency by controlling data transfer and ensuring uninterrupted data flow.

ATTO NICs: ATTO NICs provide high-performance network connectivity for AI applications, facilitating fast and low-latency data communication between systems. Supporting Ethernet protocols like 10GbE, 25GbE, and 100GbE, ATTO NICs deliver reliable and high-bandwidth network connectivity necessary for real-time AI workloads. These NICs leverage advanced offload technologies and hardware accelerations to minimize latency and reduce CPU overhead, enabling efficient data transfer.

ATTO Storage Appliances: ATTO's Storage Appliances combine high-performance storage hardware with advanced software optimizations to deliver real-time, managed latency storage solutions. These appliances ensure low-latency access to data, facilitating fast and efficient data retrieval for AI applications. Technologies like NVMe over Fabrics (NVMe-oF) and RDMA (Remote Direct Memory Access) are leveraged to reduce storage latency, enabling high-throughput data access and improving AI application performance.

ATTO Software Ecosystem: ATTO offers a comprehensive software ecosystem that complements their hardware solutions, further enhancing real-time, managed latency. This ecosystem includes drivers, management tools, and performance tuning utilities designed to optimize and fine-tune storage and network connectivity for AI workloads. Seamlessly integrating ATTO's software solutions enables efficient monitoring, integration, and performance optimization of their hardware components, ensuring the best possible latency management.

4. Benefits and Impact on AI Applications

Improved AI Application Performance: ATTO's specialized solutions for real-time, managed latency significantly enhance AI application performance. Reduced storage and network latency enable faster data access, resulting in quicker model training, reduced inference times, and enhanced overall system responsiveness. This is particularly important for real-time decision-making, critical industries like healthcare, finance, and autonomous systems.

Increased Efficiency and Scalability: By optimizing latency, ATTO solutions overcome data access bottlenecks, increasing the efficiency and scalability of AI applications. Faster data retrieval and processing enhance the utilization of AI resources, allowing organizations to handle larger datasets and more complex AI workloads with improved performance and reduced time to insight.

Seamless Integration and Support: ATTO's specialized solutions offer seamless integration into existing AI infrastructures, providing organizations with comprehensive support and compatibility across different storage and network environments. This ensures a smooth adoption of ATTO's technology and simplifies the management and maintenance of latency-optimized AI systems.

6. ATTO SiliconDisk (unreleased): Empowering AI Applications with Ultra-Low Latency

In addition to ATTO's specialized solutions for storage and network connectivity, the soon to be released, ATTO SiliconDisk is a groundbreaking storage technology that brings ultra-low latency to AI applications. Designed with a focus on minimizing storage latency, the SiliconDisk provides significant benefits for AI workloads, enabling faster data access and enhancing overall system performance.

Ultra-Low Latency: The ATTO SiliconDisk leverages cutting-edge hardware and optimized software algorithms to achieve ultra-low storage latency. By employing advanced technologies such as integrated Ethernet cores into ATTO specialized performance ASIC, the SiliconDisk significantly reduces the time it takes to access

data, ensuring near-instantaneous retrieval for AI applications. This ultra-low latency allows AI systems to process data with minimal delay, leading to improved inference times and enhanced real-time decision-making capabilities.

Accelerated Data Processing: The SiliconDisk's ultra-low latency enables AI applications to efficiently access and process large volumes of data in real-time. This is particularly valuable for AI workloads that require fast and continuous data ingestion, such as real-time analytics, natural language processing, and video analysis. With accelerated data processing capabilities, the SiliconDisk ensures that AI algorithms can rapidly access the necessary data, enabling quick model training and efficient inferencing.

Enhanced Responsiveness and User Experience: AI applications that demand real-time responsiveness, such as interactive chatbots or autonomous vehicles, greatly benefit from the ultra-low latency offered by the SiliconDisk. The instant access to data enables AI systems to deliver rapid and accurate responses, enhancing user experience and system interactivity. This is especially critical in scenarios where immediate decisions or actions need to be taken based on AI-generated insights.

Scalability and Future-Proofing: As AI workloads grow in complexity and scale, the SiliconDisk's ultra-low latency ensures that storage performance can keep pace with the demanding requirements. By reducing latency bottlenecks, the SiliconDisk enables organizations to scale their AI infrastructure without compromising on performance. This future-proofing capability allows AI applications to handle larger datasets, accommodate higher user loads, and support more sophisticated AI models as the technology evolves.

The ATTO SiliconDisk with its ultra-low latency brings remarkable benefits to AI applications. By minimizing storage latency through advanced hardware and optimized software algorithms, the SiliconDisk enables faster data access, accelerated data processing, and enhanced system responsiveness. The ultra-low latency empowers AI systems to handle real-time workloads, deliver immediate insights, and provide a superior user experience. Additionally, the scalability and future-proofing capabilities of the SiliconDisk ensure that AI infrastructures can grow and adapt to the evolving demands of AI workloads. With the ATTO SiliconDisk, organizations can unlock the full potential of their AI applications by harnessing the power of ultra-low storage latency.

7. Conclusion

Real-time, managed latency of storage is paramount for optimal performance in AI applications. ATTO Technology's specialized solutions, including HBAs, NICs,

SiliconDisk and Storage Appliances, address the challenges of storage latency and deliver high-performance connectivity and low-latency data access. By leveraging intelligent features, software optimizations, and high-speed connectivity, ATTO empowers organizations to achieve real-time, managed latency for their AI workloads. The integration of ATTO solutions enhances AI application performance, increases efficiency and scalability, and provides seamless integration and support. With ATTO's focus on real-time, managed latency, organizations can confidently deploy their AI applications, confident in the performance and reliability of their storage and network connectivity.

Learn more at www.atto.com.