# Transforming AI Data Pipelines with Advanced SSD Technology

## MonTitan™ PCIe Gen5 SSD Development Platform

### SMI's Enterprise PCIe Gen5 SSD Development Platform from Core to Edge for AI Era

Substantial amounts of data for artificial intelligence (AI) workloads are being collected from various sources at the edge, including Internet of Things (IoT) devices, consumer smartphones, and autonomous vehicles. This data is relayed to data centers, which must continuously advance their processing, memory, and storage capacities to meet these increasing demands efficiently.

Hyperscale data centers and high-performance computing must store large volumes of data efficiently and access it quickly for AI data pipelines. Key factors for AI data storage include performance, power consumption, and total cost of ownership (TCO), not just capacity.
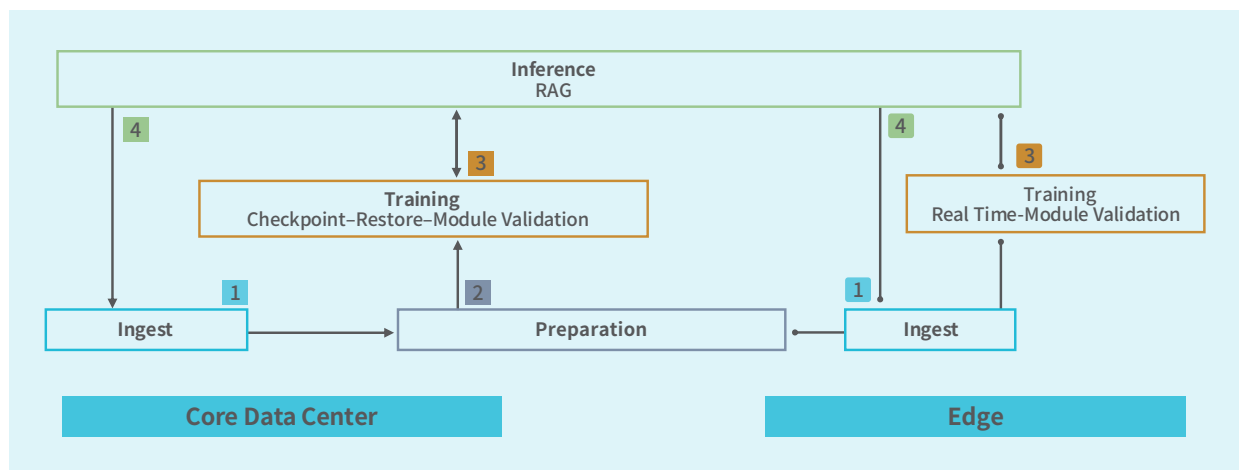
Hard drives can store large volumes of data collected from the edge for AI training, inference, and validation while it is at rest, but SSDs offer benefits that support the complete pipeline, including storing data and delivering it quickly and efficiently to numerous GPUs and other processors and accelerators. Among them are security features that protect the integrity of data as well as capabilities that optimize data placement and performance shaping that not only accelerate data transformation but also improve energy efficiency.

### Where is AI Data, How Is It Being Used and Trends Affecting Storage

The AI data pipelines relies heavily on substantial quantities of data, which are processed in data centers for ingestion, preparation, training, and inference. This data can be in various forms, such as text and video, resulting in extensive data sets that must be efficiently stored to facilitate manipulation and transformation.

Machine learning algorithms look for inter-dependencies and patterns with data sets and apply those learnings to any new data that is ingested. More data and higher quality data is how these algorithms are refined and improved over time.

This data isn't uniform, however, which means it must be prepared for training after it is ingested. The preparation and subsequent training and inference means data isn't sitting idle. It is being read, moved around and transformed.

## AI Data Pipeline and Its Demand on Storage

The AI data pipeline puts new pressures on storage. It is not enough to be able to store copious amounts of data that can be accessed quickly – high capacity is not the only requirement. Depending on the stage of the AI data pipelines, storage must manage both high performance sequential and mixed workloads.

A typical AI data pipeline for a large language model (LLM) requires storage that maximizes data efficiency and minimizes total system power while reducing training completion times.

The first stage of the pipeline ingests data in high volumes and velocity, requiring high throughput sequential write with append operations. It retrieves raw data from various external sources, including IoT devices, web scraping, databases, and autonomous vehicles.

This unstructured data is stored in its raw form; preparation data is then modified, which involves many different, complicated workloads. These workloads influence how sequential read and write performance must be handled.

The data preparation stage of an AI data pipelines is read intensive, sequential, and lower latency with sequential writes. Data must be cleaned because the collection methods are unstructured – corrupt or duplicate data, or simple "dummy data" that is not helpful for machine learning purposes must be sifted out. Because unstructured data is not categorized, formatted, or stored in a structured manner necessary for proper processing, it must be preprocessed, which involves automating classification and storage for use before processing can occur.

The actual training state of an AI data pipelines involves churning of data through many different processes. Random read performance is critical – there is a high degree of read and write throughout the training stages which requires low latency and extremely high IOPs. The training process usually has checkpoints and restore points in the event something goes wrong.

Checkpointing is more than writing data; it is more nuanced and complex than the auto save feature in Microsoft Word document. And if there is a problem that requires the training process to go back to a checkpoint, the storage media must support burst sequential read to enable recovery – a system restore of a large amount of data.

Inference requires both sequential read and write at low latency. Reference models for online shopping use inference to display to a customer what customers like have looked at or bought, which involves a great deal of random read and write performance from the storage media.

At the end of the pipeline there is data archiving, which requires high-capacity storage and sequential writing, but that data may get read again and even re-ingested.

All these stages together constitute an AI workload, and it is not a one-off process – it is a never-ending cycle of data being collected and transformed. These looping functions mean the storage environment must be tailored at any given point to meet performance requirements while achieving better return on investment (ROI).

### Massive Data Growth, Storage Options and Their Requirements

Massive data growth is not new. It started in the enterprise with the emergence of business intelligence applications and has been compounded by user generated content via smartphones and exponential growth of data collection being done at the intelligent edge.

In 2024, 402.89 million terabytes of data are created, captured, copied, or consumed every day, according to Statista. This adds up to 147 zettabytes of data per year. It's projected to grow to 181 zettabytes by 2025.[1] With all this data being ingested by AI data pipelines to create sizeable training models, the demand for storage is unending, and it must be high capacity, secure and efficient from both a data and power perspective.

Data security and privacy is more important than ever, not only because of legislation and regulations that govern how data is handled, but also because AI workloads represent valuable intellectual property.

Security is essential maintain the integrity of the data, which is critical for the AI data pipelines if the training models and ultimate inference are to be trusted. Data accuracy and validity must be ensured throughout the AI data pipeline from the point of ingestion to the final inference. In an AI data pipelines, security must ensure both data protection and platform protection by using data encryption and/or obfuscation techniques to prevent unauthorized access to data and secure boot with root of trust and prevent unauthorized mutable code access with authentication.

SSDs provide software encryption, hardware encryption, and advanced technology attachment (ATA) to secure data and the device from tampering. Software encrypts data on a logical volume using various software programs. AES 256-Bit (Advanced Encryption Standard) is a hardware-based method that uses a symmetric encryption algorithm to divide data into 128-bit blocks and encrypt them with a 256-bit key. Hardware-encrypted SSDs are designed to integrate with the rest of the drive without affecting performance.

Other available data security and privacy features include the Trusted Computing Group (TCG) Opal 2.0 protocol that can initialize, authenticate, and manage encrypted SSDs through usage of independent software vendors. Caliptra, meanwhile, defines core Root of Trust (RoT) capabilities that must be implemented in the System on Chip (SoC) or ASIC of any device in a cloud platform.

With the emergence of quantum computing, post quantum cryptography will become necessary to protect AI data pipelines as traditional cryptography methods will become ineffective.

The shift from core to edge computing supports more real-time inference. Despite the core being larger, the edge is expanding 3-5x faster, increasing storage needs outside data centers to meet AI demands. AI servers typically require two to four times more storage capacity than regular servers. Hard drives are good for storing substantial amounts of static data, while SSDs are ideal for handling data transformation in AI data pipelines.

AI workloads demand increasing memory bandwidth and capacity, as well as quick access to exabytes of data. Storage faces challenges due to intermittent data surges from AI data pipelines and straggler data causing tail latency. Additionally, storage must function efficiently over extended periods.

AI storage is fueling an exponential growth in data that not only requires high-capacity storage but also addresses the cost per terabyte, which is a leading metric for AI SSDs supporting ingest inference and checkpointing stages.

## Industry standards enable efficient AI data pipelines

The mature Non-Volatile Memory Express (NVMe) protocol plays a critical role in supporting the AI data pipelines. NVMe SSDs offer the ultra-fast read and write speeds necessary for each step in the data transformation process, as well as parallel data access by using PCIe lanes to enable concurrent read and write operations and optimize data transfer.

NVMe SSDs can handle large data sets in high-performance computing environments where AI and ML models require rapid data access and fast read/write capabilities as well as reduced latency for real-time analytics, whether it's the data center or the edge.

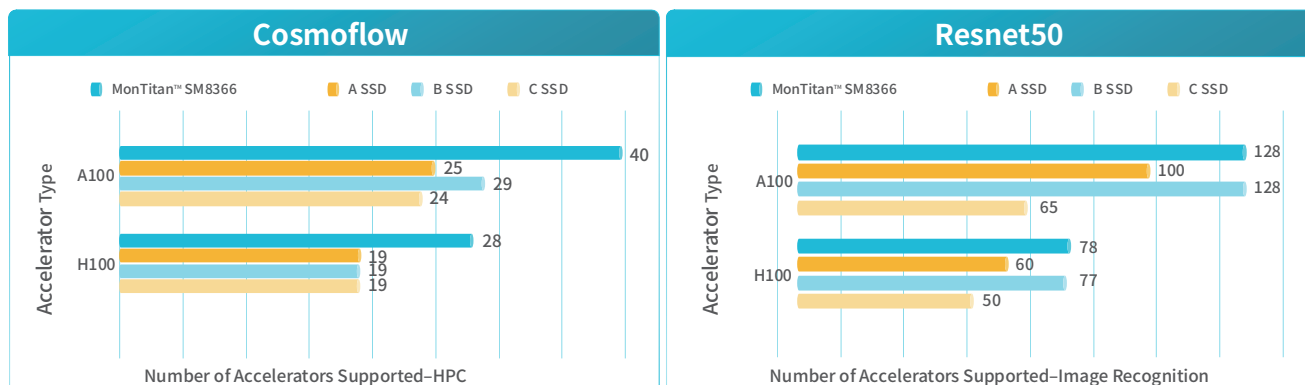## Smart SSD management improves performance, power efficiency

Because SSDs have such a critical role to play in the AI ecosystem, there must be a focus on how flash is managed. Controllers, not just the SSDs themselves, play a role in supporting the AI data pipelines, and are being optimized to enhance storage performance, while the adoption of QLC NAND in NVMe SSDs increases storage capacity at lower costs.

Capacity and fast data movement are not the only requirements AI data pipelines put on storage. Performance per watt has become a key benchmark for storage, especially in data center AI, as are data privacy and security.

Because the AI data pipelines must be fed data quickly, data efficiency is critical, as AI consumes an unprecedented amount of power. From an ASIC perspective, Silicon Motion provide power islands, frequency scaling, fast retry and exit, and low power modes, while optimizing data movement efficiency through industry standard techniques as well as proprietary technology.
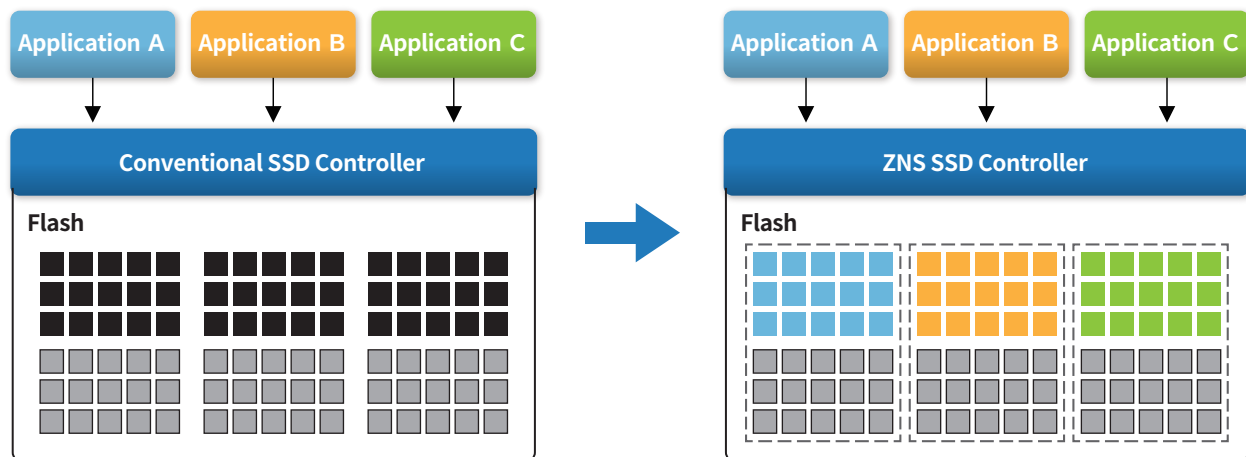
## AI Data Pipeline for Data Placement and Performance Shaping

Optimizing data placement is crucial for AI data pipelines efficiency. NVME Zone Named Spaces (ZNS) and Flexible Data Placement (FDP) reduce latency, boost performance, and enhance endurance for AI data access. Silicon Motion's MonTitan technology layers these technologies, providing additional capabilities in large-capacity storage.
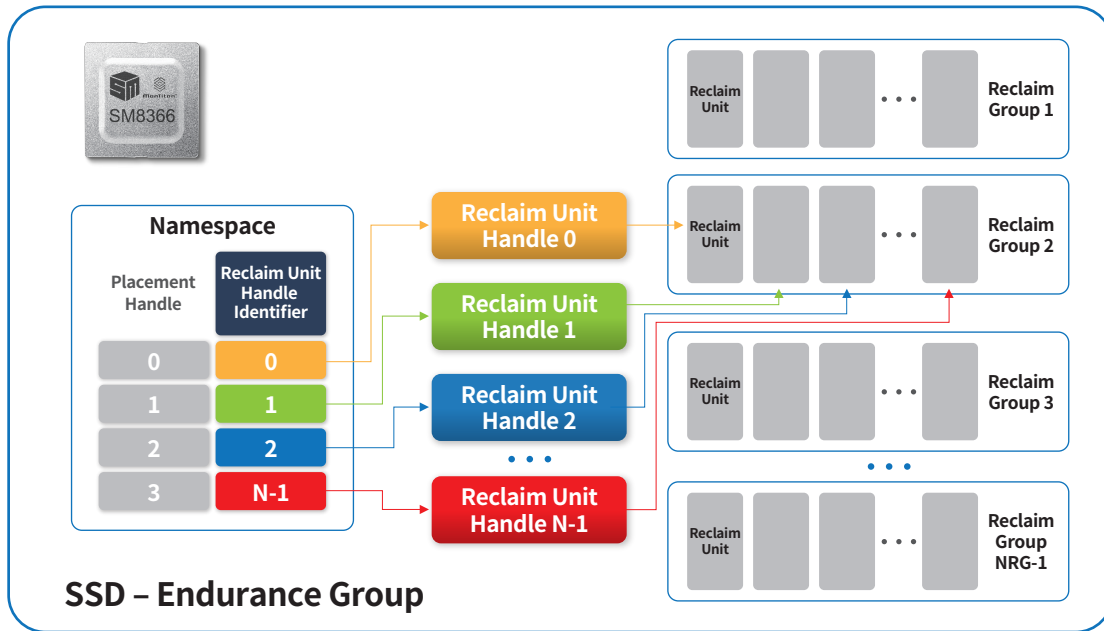


ZNS and FDP ensure proper placement of data at the ingest stage – by being able to accurately place the data where it needs to be, when it needs to be there, there's much less round tripping from a memory perspective, and the GPU are kept as busy as possible.

ZNS-enabled SSDs separate data into zones for more efficient placement and retrieval, reducing internal data movement. ZNS enhances SSD efficiency, longevity, lowers latency, and increases throughput, making them ideal for high-performance AI tasks.



FDP addresses the write-amplification problem in SSDs by allowing the host to have a simplified view and moderate awareness of media topology. An FDP SSD retains control over logical-2-physical mapping, garbage collection and the bad block management of the NAND media. For mixed random and sequential access, FDP optimization reduces write amplification near to one to improve write performance and endurance.

# Flexible Data Placement (FDP) Model



High-capacity SSDs in AI data centers face the IO blender effect due to multi-tenancy or the same storage being accessed by multiple stages of the AI data pipelines simultaneously, where different stages of the AI data pipelines try to access the same data. Silicon Motion's proprietary performance shaping technology reduces these conflicts by dynamically configuring Quality of Service (QoS) sets using a dual stage shaping algorithm tailored for specific workloads in the AI data pipeline.

## Silicon Motion's AI Strategy

Silicon Motion's comprehensive AI strategy spans the entire spectrum from enterprise core to edge applications with tailored SSD controllers that support the entire AI data pipelines across diverse markets, including data centers, enterprises, AI consumer edge devices, the artificial intelligence of things (AIoT), and automotive.

## AI storage for the core and the enterprise

For the data center, Silicon Motion offers both large capacity QLC SSDs and high performance, compute focused TLC SSDs to manage high volumes of structured and unstructured data generated for AI data pipelines, minimizing IO blender effect that results from multi-tenancy and simultaneous SSD access of different pipeline stages.

Silicon Motion's comprehensive portfolio to address AI data growth goes beyond SSDs and controllers, however. Silicon Motion's MonTitan Enterprise SSD Development platform provides purpose-built firmware, hardware and ASIC technologies that work synergistically to address customer needs, including robust data integrity and privacy capabilities, as well as performance / watt, which has become a key benchmark, especially in data center AI.

## The MonTitan Enterprise SSD Development Platform

MonTitan is a high-performance, user-programmable PCIe Gen5 development platform, and is available with Silicon Motion's production-ready SM8366 Flash controller ASIC, Turnkey and Layered Enterprise

firmware, and SSD Reference Design Kits to enable customers' rapid time-to-market design providing the best TCO.



## Controller innovation

The SM8366 is MonTitan platform's high performance, dual ported Enterprise and Data Center PCIe Gen5 x4 NVMe controller containing 16 channels and supporting up to 2400MT/s. The SM8366 provides industry leading fast sequential (>14.2 GB/s) and Random (>3.5M IOPS) SSD performance and contains a scalable Single / Dual Channel 40bit DDR4-3200 / DDR5-4800 DRAM interface.

The SM8366 is provided in a 21mm x21mm FCBGA package enabling a single chip and FW stack to be used across all new EDSFF form factors. MonTitan supports the NVM Express 2.0b, and OCP Data Center NVMe SSD 2.5 specifications with firmware optimized for power and performance in standard form factors including E1.S (9.5/15/25 mm), E3.x, and U.2/3. Its SSD Reference Design Kit (RDK) supports up to 128TB with QLC NAND. The platform also includes enterprise firmware to optimize data placement, including Silicon Motion's patented PerformaShape™ and NANDCommand™ technologies.

## Proprietary algorithms

PerformaShape is a multi-stage shaping algorithm configured in firmware that optimizes SSD performance on a per user defined QoS set bases. Combined with using true HW isolation technology, it enables the SM8366 controller ensure maximum bandwidth performance while maximizing user defined individual performance elements (QOS, Latency, RR/RW, power).

The NANDCommand capabilities of the SM8366 enhance enterprise NAND functionality through its fifth generation LDPC engine, machine learning, real-time media scan, data retention, and endurance extension algorithms. These algorithms are combined with high performance soft-decoding, multi-pass and multi-plane handling, buffer management, and read/write partitioning for effective command of next generation low latency TLC/QLC NAND.

With 16 channels running up to 2,400 MT/sand having DDR5 DRAM interfaces, the MonTitan platform and the SM8366 can deliver 3.5 million IOPS of random performance, which is 20% to 30% better than other off-the-shelf suppliers' specifications.

MonTitan addresses the IO blender effect with its focus on control plane efficiency and coherency with key hardware accelerators and processing elements, which makes it perfect for AI data pipelines and for data centers, as it isolates guaranteed traffic in different stages of the AI data pipelines which is critical for maximizing GPU utilization.



## Data protection

Aside from addressing key AI data pipelines performance requirements, the MonTitan platform also includes capabilities to ensure data integrity, including dual-engine inspection of encryption or scramblers, while all on-chip memory and DRAM interfaces are protected by SECDED ECC. Other data security features include support for NIST/FIPS (USA) standards and AES-XTS/SM4 Data at Rest protection.

At the platform level, MonTitan employs Caliptra's RoT capabilities, as well as supporting the SHA2/3 and SM3-based hash functions using the RSA/ECDSA/PQC/SM2 asymmetric algorithm to enable post-quantum encryption.

Silicon Motion's MonTitan Enterprise SSD Development platform is designed to address the storage challenges in the data center and emerging HPC, edge computing, and AI applications so that SSDs can meet key metrics organizations are looking for at every stage of the AI data pipelines.

## Intelligent AI storage from the core to the edge

As AI continues to generate massive amounts of data, efficient and high-capacity storage solutions are crucial. As a market leader in NAND storage and controllers, Silicon Motion offers a comprehensive AI strategy that spans the entire spectrum—from enterprise core to edge applications. Our solutions are tailored for diverse markets, including data centers, enterprises, AI consumer edge devices, AIoT, and automotive, ensuring fast, secure, and energy-efficient data management.

## References

1.Volume of data/information created, captured, copied, and consumed worldwide
from 2010 to 2023, with forecasts from 2024 to 2028
https://www.statista.com/statistics/871513/worldwide-data-created/

For more information about MonTitan family, please go to
**www.siliconmotion.com** or send email to **service@siliconmotion.com**

SiliconMotion