

The Lenovo logo is displayed in white text on a black rectangular background.

Reference Architecture: Lenovo ThinkSystem Edge Compute and Storage Solution for AI Inference Workloads

Last update: **09 August 2021**

Version 1.0

**Reference Architecture for
Edge AI Inference Workloads**

**Contains performance data for AI
inference using shared storage**

**Designed for small and
medium edge locations**

**Contains Bill of Materials for
compute servers and storage**

Miro Hodak, Lenovo
Satish Thyagarajan, NetApp



Table of Contents

1	Introduction	1
1.1	Target Audience	1
1.2	Solution Architecture	1
1.3	How to Use this Architecture	3
1.4	Solution Areas	3
1.4.1	Automobiles: Autonomous Vehicles	3
1.4.2	Healthcare: Patient Monitoring	3
1.4.3	Retail: Cashier-less Payment	4
1.4.4	Financial Services: Human Safety at kiosks and fraud prevention	4
1.4.5	Manufacturing: Industry 4.0	4
2	Technology Overview	5
2.1	Lenovo ThinkSystem DM Series System	5
2.2	DM Series Software	5
2.2.1	Simplify Data Management	5
2.2.2	Accelerate and Protect Data	6
2.2.3	Future-Proof Infrastructure	6
2.3	Lenovo ThinkSystem Server Portfolio	6
2.4	Lenovo ThinkSystem SE350 Edge Server	7
2.5	MLPerf	8
3	Test Overview	9
4	Test Configuration	11
5	Test Procedure	13
5.1	OS Setup	13
5.2	AI Inference Setup	13
5.3	AI Inference Runs	13
6	Test Results	14
6.1	AI Inference in Offline Scenario	14
6.2	AI Inference in Single Stream Scenario	14

6.3	AI Inference in Multi Stream Scenario.....	15
7	Architecture Adjustments	17
7.1	Compute Server Adjustment	17
7.2	Storage Capacity Increase.....	17
7.3	Other DM Series Storage Options	17
8	Conclusion	18
9	Appendix: Lenovo Bill of materials.....	19
9.1	BOM for compute servers	19
9.2	BOM for Storage Server	20
9.3	BOM for Network Switch.....	22
	Resources	23

1 Introduction

Companies are increasingly generating massive volumes of data at the network edge. To gain maximum business value from smart sensors and IoT data, organizations are looking for real-time event-streaming solutions that enable edge computing. Computationally demanding jobs are increasingly performed at the edge, outside of data centers. Artificial Intelligence (AI) inference is one of the drivers of this trend. Edge servers provide sufficient computational power for these workloads especially when using accelerators, but limited storage is often an issue especially in multi-server environment. Here we show how shared storage can be deployed in the edge environment and how it benefits AI inference workloads without imposing a performance penalty.

This document describes a reference architecture for Artificial Intelligence (AI) inference at the edge. It combines multiple Lenovo ThinkSystem edge servers with Lenovo ThinkSystem DM-series storage to create an easy to deploy and manage solution. It is meant as a base guide for practical deployments in various situations such as factory floor or point-of-sale systems.

The document covers testing and validation of a compute and storage configuration consisting of Lenovo ThinkSystem SE350 Edge Server and a Lenovo ThinkSystem DM5100F storage system. This architecture provides an efficient and cost-effective solution for AI deployments while also providing comprehensive data services, integrated data protection, seamless scalability and cloud connected data storage.

1.1 Target Audience

This document is intended for the following audiences:

- Business leaders and enterprise architects who want to productize AI at the edge.
- Data scientists, data engineers, architects, and developers of AI systems
- Enterprise architects who design solutions for the development of AI models and software
- Data scientists and data engineers looking for efficient ways to deploy deep learning (DL) and machine learning (ML) models.
- Edge device managers & administrators responsible for deployment & management of edge inferencing models.

1.2 Solution Architecture

This Lenovo ThinkSystem server and storage solution is designed to handle AI inference on large datasets using the processing power of GPUs alongside traditional CPUs. This validation demonstrates high performance and optimal data management with an architecture that uses either a single or multiple Lenovo SE350 edge servers interconnected with a single Lenovo DM5100F storage system.

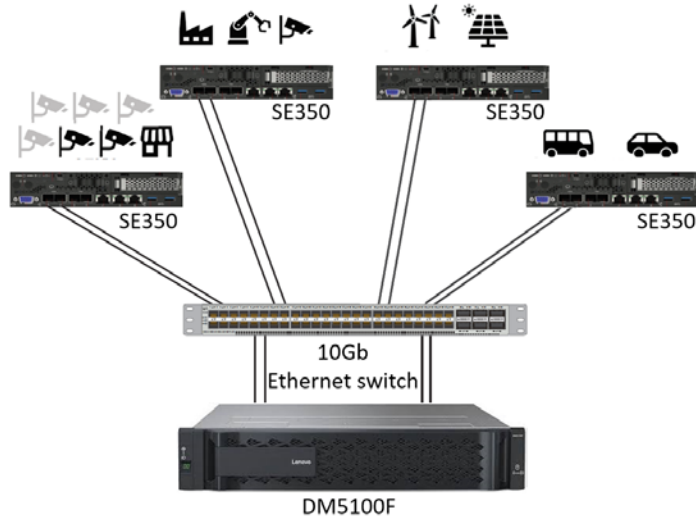


Figure 1: Physical Architecture Overview

The logical architecture in Figure 2 shows the roles of the compute and storage elements in this architecture. Specifically:

- Edge compute devices performing inference on the data they receive from cameras, sensors, etc.
- Shared storage element that serves multiple purposes:
 1. Central location for inference models and other data needed to perform the inference. Compute servers access the storage directly and use inference models across the network without the need to copy them locally.
 2. Updated models are pushed here.
 3. Archiving of input data that edge servers receive for later analysis. For example, if the edge devices are connected to cameras, the storage element would keep the videos captured by the cameras.

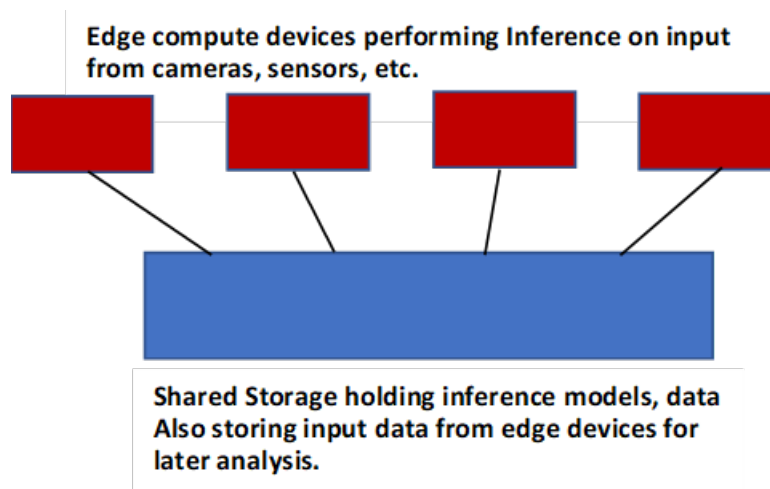


Figure 2: Logical Architecture Overview

This Lenovo solution offers the following key benefits:

- GPU-accelerated computing at the edge
- Deployment of multiple edge servers backed and managed from shared storage
- Robust data protection to meet low recovery point objectives (RPOs) and recovery time objectives (RTOs) with no data loss.
- Optimized data management with snapshots and clones to streamline development workflows

1.3 How to Use this Architecture

This document has validated the design and performance of the proposed Architecture. However, certain software-level pieces such as container/workload/model management have not been tested as they are specific to the particular deployment scenario. Here, multiple choices exist. At the container management level, Kubernetes container management is a good choice and is well supported in either fully upstream version (Canonical) or in a modified version suitable for enterprise deployments (Red Hat). In particular, Kubeflow, the machine learning toolkit for Kubernetes, provides additional AI capabilities along with support for model serving on several platforms such as TensorFlow Serving or NVIDIA Triton Inference Server. Another good option is NVIDIA EGX platform, which provides workload management along with access to a catalogue of GPU-enabled AI inference containers. However, these options may require significant effort and expertise to put them into production and thus a customized solution from an ISV specializing in the area of interest may be used instead.

1.4 Solution Areas

The key benefit of AI inferencing and edge computing is the ability of devices to compute, process and analyze data with high level of quality without latency. There are many examples of edge computing use-cases but here are a few leading ones below:

1.4.1 Automobiles: Autonomous Vehicles

The classic edge computing illustration is in the autonomous vehicle. The AI in driverless cars must rapidly process a lot of data from cameras and sensors to be a successful safe driver. Taking too long to interpret between an object and a human could mean life or death, therefore being able to process that data as close to the vehicle as possible is crucial.

In this case, one or more edge compute servers would handle the input from cameras, radar, and other sensors while shared storage would hold inference models and store input data from sensors.

1.4.2 Healthcare: Patient Monitoring

One of the greatest impacts of AI and edge computing is its ability to enhance continuous monitoring of patients for chronic diseases both in at-home care and Intensive Care Units (ICUs). Data from devices that monitor insulin levels, respiration, neurological activity, cardiac rhythm, and

gastrointestinal functions require instantaneous analysis of data which must be acted upon immediately because there is limited time to act.

1.4.3 Retail: Cashier-less Payment

Edge computing can power AI and Machine Learning to help retailers reduce checkout time and increase foot traffic. Cashier-less systems support various components like *authentication and access* by connecting the physical shopper to a validated account and permitting access to the retail space; *tracking systems*, a computer vision technology which associates a shopper with a product selection by virtue of proximity and *inventory monitoring* which uses sensors, RFID tags and vision systems that help confirm the selection or de-selection of items by shoppers.

Here, each of the edge servers would handle each checkout counter and the shared storage server would serve as a central synchronization point.

1.4.4 Financial Services: Human Safety at kiosks and fraud prevention

Banking organizations are using Artificial Intelligence and Edge computing to innovate and create personalized experiences. Interactive kiosks using real-time data analytics and AI inferencing, now enable ATMs with not only helping customers in withdrawing money but to proactively monitor through the images captured from cameras to identify risk to human safety or fraudulent behavior. In this scenario, edge compute servers and shared storage would be connected to interactive kiosks and cameras to help banks collect and process data with AI inference models.

1.4.5 Manufacturing: Industry 4.0

The fourth industrial revolution (Industry 4.0) has begun. To prepare for a data-dominated future large-scale machine-to-machine (M2M) communication and Internet of Things (IoT) are integrated for increased automation without the need for human intervention through some of the emerging trends like Smart Factory, 3D printing and Smart Sensors. Manufacturing is already highly automated and adding AI features is a natural continuation of the long-term trend. AI enables automating operations that could have been automated so far with the help of computer vision and other AI capabilities. Quality control or tasks that rely on human vision or decision-making can be automated and perform much faster with higher reliability to help factories meet with ISO standards for quality management. Here, each compute edge server would be connected to an array of sensors monitoring the manufacturing process and updated inference models would be pushed to the shared storage as needed.

2 Technology Overview

This Section describes key technologies leveraged in this solution.

2.1 Lenovo ThinkSystem DM Series

State-of-the-art Lenovo DM Series storage systems enable IT departments to meet enterprise storage requirements with industry-leading performance, superior flexibility, cloud integration, and best-in-class data management. Designed specifically for flash, DM Series storage systems help accelerate, manage, and protect business-critical data.

This Reference architecture is based on ThinkSystem DM5100F storage. Lenovo ThinkSystem DM5100F is an NVMe unified flash storage system that offers the following key features and benefits:

- Scalable up to 737 TB of raw storage capacity per system and up to 8.8 PB per NAS scale-out cluster of up to 12 systems
- Delivers up to 440K random read IOPS (8 KB blocks) per system.
- Supports 10/25 GbE NAS and iSCSI, 8/16/32 Gb Fibre Channel, and 32 Gb NVMe over Fibre Channel (NVMe/FC) host connectivity..A complete suite of data protection and replication features for industry-leading data management
- Designed for 99.9999% availability

Lenovo offers other storage systems, such as the DM7000F and DM7100F, that offer higher performance and scalability for larger-scale deployments.

2.2 DM Series Software

The Lenovo DM Series uses the latest generation of storage management software that enables businesses to modernize infrastructure and transition to a cloud-ready data center. Leveraging industry-leading data management capabilities, the DM Series enables the management and protection of data with a single set of tools, regardless of where that data resides. Data can also be moved freely to wherever it's needed—the edge, the core, or the cloud. The DM Series includes numerous features that simplify data management, accelerate and protect critical data, and future-proof infrastructure across hybrid cloud architectures.

2.2.1 Simplify Data Management

Data management is crucial to enterprise IT operations so that appropriate resources are used for applications and datasets. The DM Series includes the following features to streamline and simplify operations and reduce the total cost of operation:

- **Inline data compaction and expanded deduplication.** Data compaction reduces wasted space inside storage blocks, and deduplication significantly increases effective capacity. This applies to data stored locally and data tiered to the cloud.
- **Minimum, maximum, and adaptive quality of service (QoS).** Granular QoS controls help maintain performance levels for critical applications in highly shared environments.

- **FabricPool.** This feature provides automatic tiering of cold data to public and private cloud storage options, including Amazon Web Services (AWS) and Azure.

2.2.2 Accelerate and Protect Data

The Lenovo DM Series delivers superior levels of performance and data protection and extends these capabilities as follows:

- **Performance and lower latency.** The DM Series offers the highest possible throughput at the lowest possible latency.
- **Data protection.** The DM Series provides built-in data protection capabilities with common management across all platforms.
- **Volume Encryption.** The DM Series offers native volume-level encryption with both onboard and external key management support.

2.2.3 Future-Proof Infrastructure

The Lenovo DM Series helps meet demanding and constantly changing business needs:

- **Seamless scaling and nondisruptive operations.** The DM Series supports the nondisruptive addition of capacity to existing controllers as well as to scale-out clusters. Customers can upgrade to the latest technologies such as NVMe and 32Gb FC without costly data migrations or outages.
- **Cloud connection.** The DM Series is the most cloud-connected storage management software, with options for software-defined storage and cloud-native instances (Cloud Volumes Service) in all public clouds.
- **Integration with emerging applications.** The DM Series offers enterprise-grade data services for next-generation platforms and applications such as OpenStack, Hadoop, and MongoDB by using the same infrastructure that supports existing enterprise apps.

2.3 Lenovo ThinkSystem Server Portfolio

Lenovo ThinkSystem servers feature innovative hardware, software, and services that solve customers' challenges today and deliver an evolutionary, fit-for-purpose, modular design approach to address tomorrow's challenges. These servers capitalize on best-in-class, industry-standard technologies coupled with differentiated Lenovo innovations to provide the greatest possible flexibility in x86 servers.

Key advantages of deploying Lenovo ThinkSystem servers include:

- Highly scalable, modular designs to grow with your business
- Industry-leading resilience to save hours of costly unscheduled downtime
- Fast flash technologies for lower latencies, quicker response times, and smarter data management in real time

In the AI area, Lenovo is taking a practical approach to helping enterprises understand and adopt the benefits of ML and AI for their workloads. Lenovo customers can explore and evaluate Lenovo AI offerings in Lenovo AI Innovation Centers to fully understand the value for their particular use case. To improve time to value, this customer-centric approach gives customers proofs of concept for solution development platforms that are ready to use and optimized for AI.

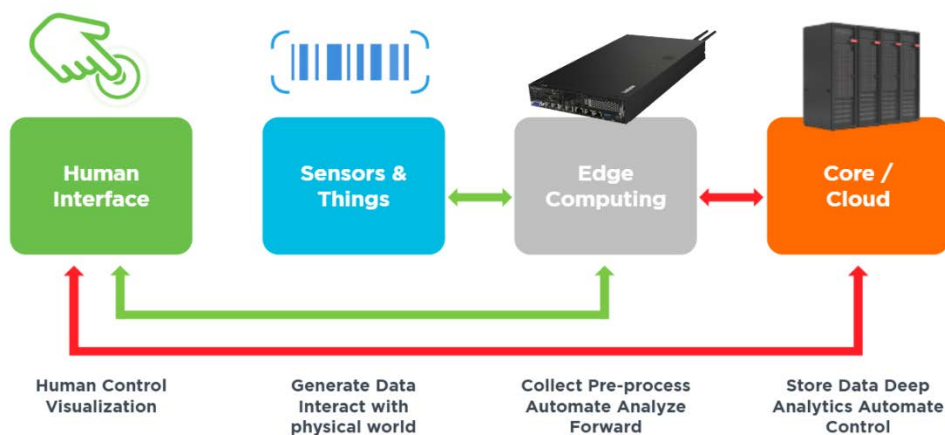
2.4 Lenovo ThinkSystem SE350 Edge Server

Edge computing allows data from internet of things devices to be analyzed at the edge of network before being sent to data center or cloud. The Lenovo ThinkSystem SE350 is designed for the unique requirements for deployment at the edge, with a focus on flexibility, connectivity, security and remote manageability in a compact ruggedized and environmentally hardened form factor.

Featuring the Intel® Xeon® D processor with the flexibility to support acceleration for edge AI workloads, the SE350 is purpose-built for addressing the challenge of server deployments in a variety of environments outside the data center.



Figure 2 ThinkSystem SE350



2.5 MLPerf

MLPerf is the industry-leading benchmark suite for evaluating AI performance. It covers many areas of applied AI including image classification, object detection, medical imaging, and natural language processing. In this validation, we used Inference v0.7 workloads, which is the latest iteration of the MLPerf Inference as of completion of this work.

MLPerf Inference results and code are publicly available and released under Apache license. MLPerf Inference has an Edge division, which supports the following scenarios:

- **Single Stream:** This scenario mimics systems where responsiveness is a critical factor, such as offline AI queries performed on smartphones. Individual queries are sent to the system and response times are recorded. 90th percentile latency of all the responses is reported as the result.
- **Multi Stream:** This benchmark is for systems that process input from multiple sensors. During the test, queries are sent at a fixed time interval. A Quality of Service constraint (maximum allowed latency) is imposed. The test reports the number of streams that the system can process while meeting the QoS constraint.
- **Offline:** This is the simplest scenario covering batch processing applications and the metric is throughput in samples per second. All data are available to the system and the benchmark measures time it takes to process all the samples.

Lenovo has published MLPerf Inference scores for SE350 with T4, the server used in this document. The results can be found at <https://mlperf.org/inference-results-0-7/> in the “Edge, Closed Division” section as the entry #0.7-145.

3 Test Overview

This document follows MLPerf Inference v0.7 [code](#) and [rules](#). We run benchmarks designed for inference at the edge as defined in the following Tables:

Area	Task	Model	Dataset	QSL Size	Quality	Multi-stream latency constraint
Vision	Image classification	Resnet50-v1.5	ImageNet (224x224)	1024	99% of FP32	50 ms
Vision	Object detection (large)	SSD-ResNet34	COCO (1200x1200)	64	99% of FP32	66 ms
Vision	Object detection (small)	SSD-MobileNets-v1	COCO (300x300)	256	99% of FP32	50 ms
Vision	Medical image segmentation	3D UNET	BraTS 2019 (224x224x160)	16	99% and 99.9% of FP32	N/A
Speech	Speech-to-text	RNNT	Librispeech dev-clean	2513	99% of FP32	N/A
Language	Language processing	BERT	SQuAD v1.1	10833	99% of FP32	N/A

Area	Task	Scenarios
Vision	Image classification	Single Stream, Offline, Multi-Stream
Vision	Object detection (large)	Single Stream, Offline, Multi-Stream
Vision	Object detection (small)	Single Stream, Offline, Multi-Stream

Vision	Medical image segmentation	Single Stream, Offline
Speech	Speech-to-text	Single Stream, Offline
Language	Language processing	Single Stream, Offline

We perform these benchmarks using the networked storage architecture developed in this work and compare results to those from local runs on the edge servers previously submitted to MLPerf. The comparison will be used to determine how much impact the shared storage has on inference performance.

4 Test Configuration

The test configuration is shown in Fig. 4. We have used Lenovo DM5100F storage, 2 Lenovo ThinkSystem SE350 servers (each one with 1 NVIDIA T4 accelerator). These components are connected through a 10GbE Network Switch.

The network storage holds validation/test datasets and pre-trained models. The servers provide computational capability. The storage is accessed over NFS protocol.

This section describes the tested configurations, the network infrastructure, the SE350 server, and the storage provisioning details.

We used the solution components listed in Table 2 for the validation.

Table 1) Base components for the solution architecture.

Solution Components	Details
Lenovo ThinkSystem Servers	<ul style="list-style-type: none">• 2 SR350 servers each with 1 NVIDIA T4 GPU cards• Each server contains 1 Intel Xeon D-2123IT CPU with 4 physical cores running at 2.20GHz and 128 GB RAM
Lenovo DM5100F storage system (HA pair)	<ul style="list-style-type: none">• DM Series software• 24 x 1.9TB SSDs• NFS protocol• 1 interface group (ifgrp) per controller, with 4 logical IP addresses for mount points

Figure 3) Network topology of tested configuration.

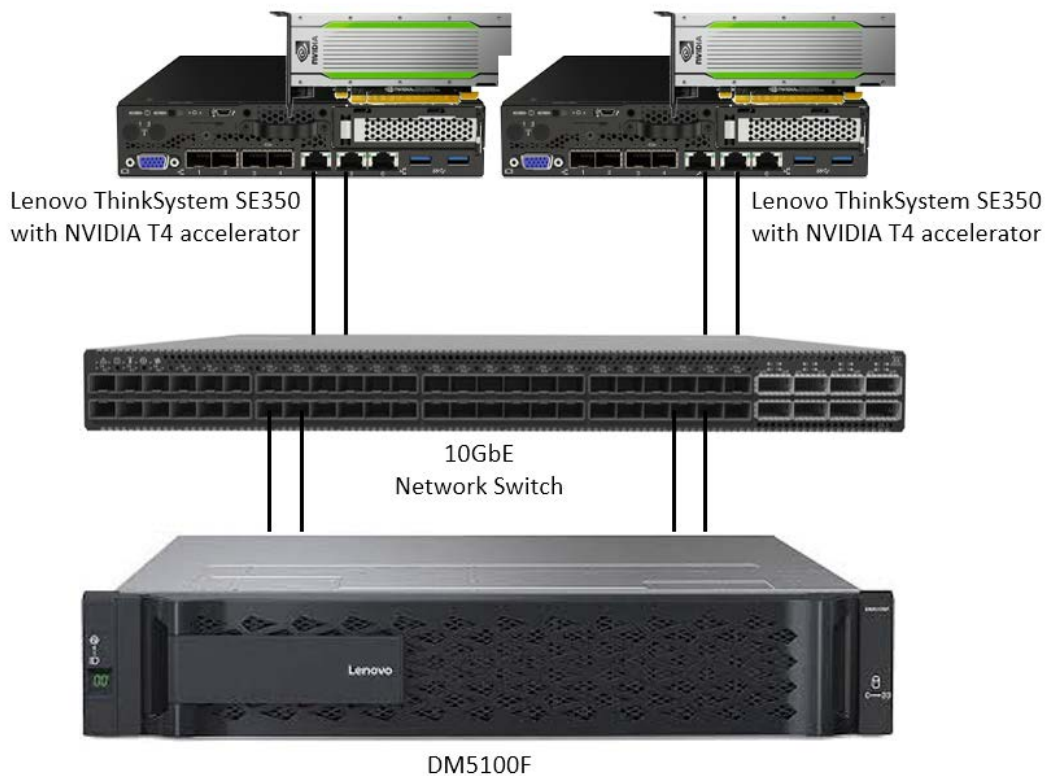


Table 2) Storage configuration.

Controller	Aggregate	FlexGroup Volume	Aggregate Size	Volume Size	Operating System Mount Point
Controller1	Aggr1	/lenovo_AI_fg	8.42TiB	15 TB	/lenovo_fg
Controller2	Aggr2		8.42TiB		

The /Lenovo_AI_fg folder contains the dataset used for validation.

5 Test Procedure

We used the following test procedure in this validation.

5.1 OS Setup

We used Ubuntu 18.04 with NVIDIA drivers and docker with support for NVIDIA GPUs.

5.2 AI Inference Setup

We used MLPerf [code](#) available as a part of Lenovo submission to MLPerf Inference v0.7. The setup requires the following steps:

1. Downloading datasets that require registration: ImageNet 2012 Validation set, Criteo Terabyte dataset, and BraTS 2019 Training set and unzip the files.
2. Create a working directory with at least 1TB and define environmental variable `MLPERF_SCRATCH_PATH` referring to the directory. This directory should be on the shared storage for the network storage use case, or local disk when testing with local data.
3. Run `make prebuild` command, which builds and launches docker container for the required inference tasks. The following commands are all executed from within the running docker container.
4. Run `make download_model` to download pretrained AI models for MLPerf Inference tasks.
5. Run `make download_data` to download additional datasets that are freely downloadable
6. Run `make preprocess_data` to pre-process the data.
7. Run `make build`
8. Run `make generate_engines` to build inference engines optimized for the GPU in compute servers
9. To run Inference workloads use the following (one command):

```
make run_harness RUN_ARGS="--benchmarks=<BENCHMARKS> --  
scenarios=<SCENARIOS>"
```

5.3 AI Inference Runs

Three types of runs are executed:

1. Single server AI Inference using local Storage
2. Single server AI Inference using network storage
3. Multi-server AI Inference using network storage

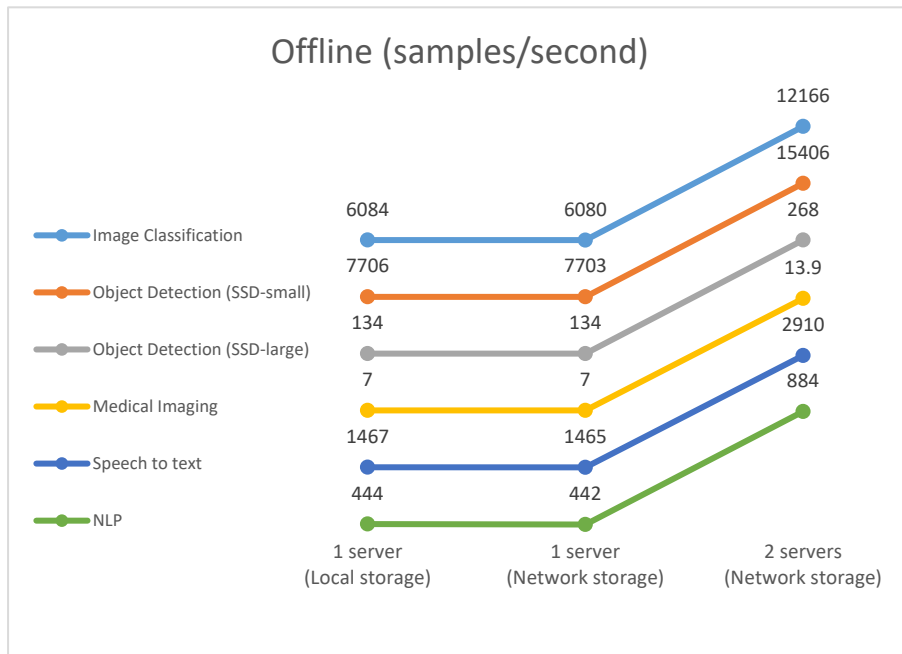
6 Test Results

A multitude of tests was run to evaluate performance of the proposed architecture. There are 6 different workloads (Image Classification, Object Detection (small), Object Detection (large), Medical Imaging, Speech-to-text, and Natural language Processing (NLP)), which can be run in 3 different scenarios: Offline, Single Stream, and Multi Stream. Note that the last scenario is implemented only for Image Classification and Object Detection. This gives 15 possible workloads, which were all tested under 3 different setups: (i) Single server/local storage, (ii) Single server/network storage, and (iii) multi-server/ network storage. The results are given below.

6.1 AI Inference in Offline Scenario

In this scenario, all the data are available to the server and the time it takes to process all of the samples is measured. Bandwidths in samples per second are reported as the results of the tests. When more than 1 compute server is used total bandwidth summed over all the servers is reported. The results for all three use cases are shown in the figure below.

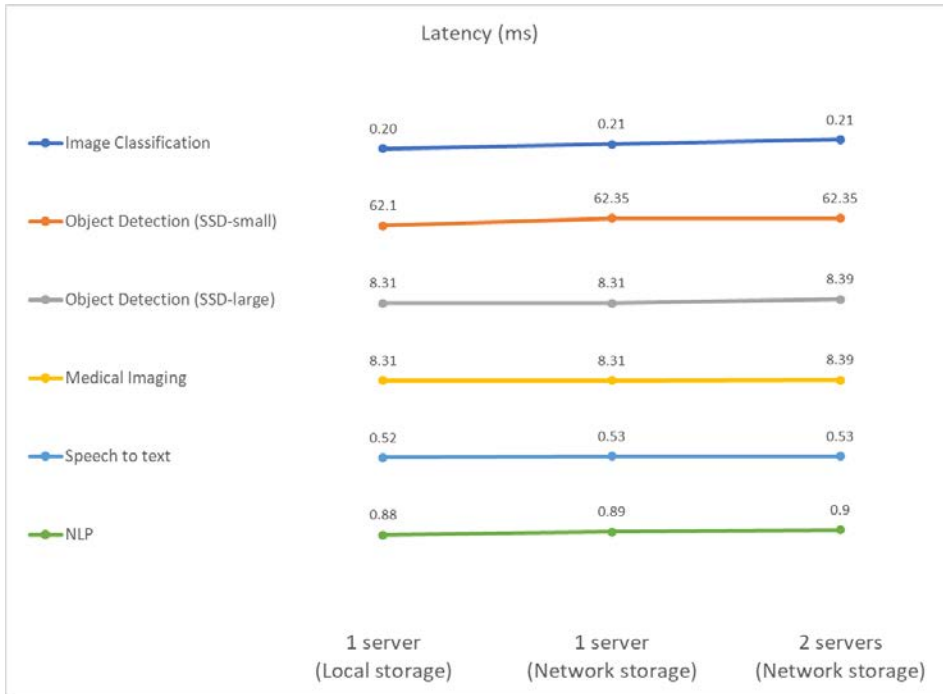
Figure 4) Bandwidth in samples per second for various inference tasks. For the 2 server case, combined bandwidth from both servers is reported.



6.2 AI Inference in Single Stream Scenario

This benchmark measures latency. For the multiple computational server case, the average latency is reported. The results for the suite of tasks are given in the Figure below.

Figure 5) Latency in milliseconds for various inference workloads. For the 2 server case, average latency from both servers is reported.

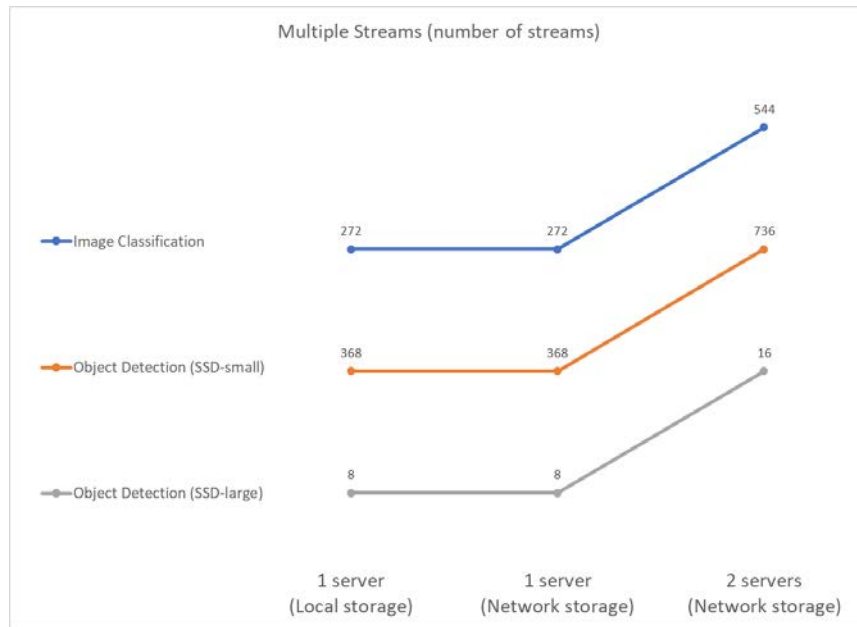


The results, again, show that the network storage is sufficient to handle the tasks. The difference between local and network storage in the one server case is minimal or none. Similarly, when two servers use the same storage, the latency on both servers stays the same or changes by a very small amount.

6.3 AI Inference in Multi Stream Scenario

In this case, the result is the number of streams that the system can handle while satisfying the quality-of-service constraint. Thus, the result is always an integer. For more than one server, we report total number of streams summed over all the servers. Not all the workloads support this scenario, but we have executed those that do. The results of our tests are summarized in the Figure below.

Figure 6) Number of streams in the Multi Stream scenario for various inference tasks. For the 2 server case, combined number of streams from both servers is reported.



The results show perfect performance of the setup – local and networking storage give the same results and adding the second server double the number of streams the proposed setup can handle.

7 Architecture Adjustments

The setup used for the validation can be adjusted to fit other use cases.

7.1 Compute Server Adjustment

We used an Intel Xeon D-2123IT CPU, which is the lowest level of CPU supported in SE350 with 4 physical cores and 60W TDP. While the server does not support replacing CPUs, it can be ordered with a more powerful CPU. The top CPU supported is Intel Xeon D-2183IT with 16 cores, 100W running at 2.20 GHz. This will increase CPU computational capability considerably. While CPU was not a bottleneck for running the inference workloads themselves, it will help data processing and other tasks related to inference.

At the present, Nvidia T4 is the only GPU available for accelerated edge use cases..

7.2 Storage Capacity Increase

DM5100F can be scaled out up to 12 High Availability pairs to 8.8 PB of total raw storage capacity.

7.3 Other DM Series Storage Options

This validation used the Lenovo Thinksystem DM5100F storage server, which is the entry-level storage option available in Lenovo's ThinkSystem storage portfolio. **Other DM Series options, DM7000F and DM7100F, provide higher performance and more features and thus it is expected that they would provide the same or better performance as the one obtained in this work.**

8 Conclusion

With edge computing, data is processed much faster since it does not have to travel to and from a data center. Lower latency and increased speed can be beneficial when businesses must make decisions in near-real time. Furthermore, when there is no need for a data center, cost associated with sending data back-and-forth to data centers or the cloud is also diminished.

This Lenovo solution is a flexible scale-out architecture that is ideal for enterprise AI inference deployments. Lenovo storage delivers the same or better performance as local SSD storage and offers the following benefits to data scientists, data engineers, and business or IT decision makers:

- Effortless sharing of data between AI systems, analytics, and other critical business systems. This data sharing reduces infrastructure overhead, improves performance, and streamlines data management across the enterprise.
- Independently scalable compute and storage to minimize costs and improve resource utilization.
- Streamlined development and deployment workflows using integrated snapshots and clones for instantaneous and space-efficient user workspaces, integrated version control, and automated deployment.
- Enterprise-grade data protection for disaster recovery and business continuity.

9 Appendix: Lenovo Bill of materials

This appendix contains the bill of materials (BOMs) for computational servers and a storage server.

The BOM lists in this appendix are not meant to be exhaustive and must always be double-checked with the configuration tools. Any discussion of pricing, support, and maintenance options is outside the scope of this document.

Within a specific BOM section, optional items are numbered with alternatives shown as lower-case letters. For example, a Fibre Channel adapter for a compute server is only needed for shared storage connected through a SAN.

9.1 BOM for compute servers

Part #	Description	Quantity
	SE350 with NVIDIA T4 AI Inference	2
7D1XCTOLWW	SE350 with T4 : ThinkSystem SE350 - 3yr Warranty - HPC&AI	1
B6EQ	ThinkSystem SE350 Edge Server Chassis	1
B6F4	ThinkSystem SE350 10GbE SFP+ 2-Port, 10/100/1GbE RJ45 2-Port Intel i350	1
B8ZR	Standard Shock & Vibration (15G & 0.21Grms)	1
B8ZT	Operational Temperature 0-45C	1
BFYE	Operating mode selection for: "Efficiency - Favoring Performance Mode"	1
B939	ThinkSystem SE350 Edge Server Intel Xeon D-2123IT 4C 60W 2.20 GHz	1
AUND	ThinkSystem 32GB TruDDR4 2666 MHz (2Rx4 1.2V) RDIMM	4
B6FF	ThinkSystem SE350 M.2 SATA/NVMe 4-bay Data Drive Enablement Kit	1
5977	Select Storage devices - no configured RAID required	1
B919	ThinkSystem M.2 5300 480GB SATA 6Gbps Non-Hot Swap SSD	4
B6FH	ThinkSystem SE350 M.2 Adapter SATA Cable	1
B88P	ThinkSystem SE350 M.2 Mirroring Enablement Kit	1
AUUV	ThinkSystem M.2 128GB SATA 6Gbps Non-Hot Swap SSD	2
B6FD	ThinkSystem SE350 PCIe Riser Cage	1
B4YB	ThinkSystem NVIDIA Tesla T4 16GB PCIe Passive GPU	1
B6FU	ThinkSystem SE350 - 12V PDM	1
B6FW	ThinkSystem SE350 240W AC Adapter	2
6311	2.8m, 10A/100-250V, C13 to IEC 320-C14 Rack Power Cable	2
B6KT	ThinkSystem SE350 - Mini USB to USB Type A (F) Console Cable	1

A1PJ	3m Passive DAC SFP+ Cable	2
B755	Desktop Mode	1
B6Q3	ThinkSystem SE350 Rubber Feet	1
B0MK	Enable TPM 2.0	1
B7XZ	Disable IPMI-over-LAN	1
BB98	Disable IPMI-over-KCS	1
A2N7	Planar Not Integrated With Chassis	1
B0ML	Feature Enable TPM on MB	1
B6EW	ThinkSystem SE350 Operator Panel	1
B6EY	ThinkSystem SE350 Pull Out Tag	1
B6Q1	ThinkSystem SE350 Node Handle	1
B756	ThinkSystem SE350 Tamper Intrusion Cable	1
B757	ThinkSystem SE350 SW GBM	1
B6F8	ThinkSystem SE350 Connector Filler Kit	1
B6G3	ThinkSystem SE350 Node WW LPK	1
B6G1	ThinkSystem SE350 Node WW Packaging	1
B7FA	ThinkSystem SE350 PCIe Riser Cage Label LI	1
B6FZ	ThinkSystem SE350 Node Label GBM	1
B6G6	ThinkSystem SE350 Node Top Cover	1
B6GA	ThinkSystem SE350 - Node Label SSL_LI	1
B6F7	ThinkSystem SE350 Passthrough Bezel	1
5PS7A34998	Premier Essential - 3Yr 24x7 4Hr Resp + YDYD SE350	1
BF94	AI & HPC - ThinkSystem Hardware	1

9.2 BOM for Storage Server

Part #	Description	Quantity
	DM5100F	2
7D3KCTOLWW	DM5100F : Lenovo ThinkSystem DM5100F All Flash Array - HPC&AI	1
BF3C	Lenovo ThinkSystem DM Series 2U NVMe Chassis	1
B5RJ	DM Series Premium Offering	1
BEVS	Lenovo ThinkSystem DM5100F NVMe Controller, 64GB	2

BEW2	Lenovo ThinkSystem 12TB (6x 1.9TB NVMe Non-SED) Drive Pack for DM5100F	2
BEVQ	Lenovo ThinkSystem DM Series HIC, 10/25Gb iSCSI,4-ports	2
AVFW	0.75m Green Cat6 Cable	2
A51N	1.5m Passive DAC SFP+ Cable	4
AV1W	Lenovo 1m Passive 25G SFP28 DAC Cable	2
B4BP	Lenovo ThinkSystem Storage USB Cable, Micro-USB	1
6400	2.8m, 13A/100-250V, C13 to C14 Jumper Cord	2
BF5Q	Lenovo ThinkSystem DM Series ONTAP 9.8 Software Encryption	1
B0W1	3 Years	1
B48J	Premier with Essential	1
B6Y6	Lenovo ThinkSystem NVMe Rail Kit 4 post	1
B4SF	DM Series CIFS Protocol License	2
B4SG	DM Series NFS Protocol License	2
B4SH	DM Series iSCSI Protocol License	2
B4SJ	DM Series FCP Protocol License	2
B4SK	DM Series SnapMirror License	2
B4SL	DM Series SnapRestore License	2
B4SM	DM Series FlexClone License	2
B4SN	DM Series Software Encryption License	2
B4SP	DM Series SnapManager License	2
B4SU	TPM	2
B5AZ	DM Series SnapVault License	2
B7AQ	SnapMirror Synchronous	2
B7N1	NVMe over FC Protocol	2
BH56	S3 Protocol License	2
BF3B	Configured with Lenovo ThinkSystem DM5100F	1
BEVT	Lenovo ThinkSystem DM5100F ShipKit-WW	1
BEVY	Lenovo Thinksystem DM5100F Product Label	1
B738	Lenovo ThinkSystem NVMe Accessory	1
BF4C	Lenovo ThinkSystem DM Series 2U24 Bezel	1

BEVX	Lenovo ThinkSystem DM5100F System Labels	1
B6Y5	Lenovo ThinkSystem NVMe SFF Filler	12
BFZ2	Lenovo ThinkSystem DM5100F Packaging	1
BHLW	DM5100F HIC faceplate ASM	2
5PS7A92253	Premier Essential- 3Y 24x7x4+YDYD ThinkSystem DM5100F AFA	1
BF94	AI & HPC - ThinkSystem Hardware	1

9.3 BOM for Network Switch

Part #	Description	Quantity
	10GbE network switch	1
7D5FCTO5WW-HPC	10GbE Switch : Mellanox SN2410B 10GbE Managed Switch with Cumulus (PSE)	1
BE2N	Mellanox SN2410B 10GbE Managed Switch with Cumulus (PSE)	1
BEGK	Mellanox SN24xx 1U Enterprise RMK w/Air Duct	1
A51P	2m Passive DAC SFP+ Cable	8
3802	1.5m Blue Cat5e Cable	1
6311	2.8m, 10A/100-250V, C13 to IEC 320-C14 Rack Power Cable	2
5WS7A88345	Premier Essential - 3Yr 24x7 4Hr Resp MLNX SN2410B PSE	1
BF94	AI & HPC - ThinkSystem Hardware	1

Resources

- [ThinkSystem SE350 Edge Server product web page](#)
- [ThinkSystem DM Series Storage product web page](#)
- [Lenovo/NetApp Compute Storage Technical Report](#)
- [MLPerf.org](#)
- [TensorFlow benchmarks](#)

Change history

Version 1.0	August 2021	Initial version validated on SE350 and DM5100F
-------------	-------------	--

Trademarks and special notices

© Copyright Lenovo 2021.

References in this document to Lenovo products or services do not imply that Lenovo intends to make them available in every country.

Lenovo and the Lenovo logo are trademarks or registered trademarks of Lenovo in the United States, other countries, or both. A current list of Lenovo trademarks is available on the Web at <https://www.lenovo.com/us/en/legal/copytrade/>.

The following terms are trademarks of Lenovo in the United States, other countries, or both:

Lenovo®
ThinkSystem
TruDDR4

The following terms are trademarks of other companies:

Intel® and Xeon® are trademarks of Intel Corporation or its subsidiaries.

Azure® is a trademark of Microsoft Corporation in the United States, other countries, or both.

Other company, product, or service names may be trademarks or service marks of others. Information is provided "AS IS" without warranty of any kind.

All customer examples described are presented as illustrations of how those customers have used Lenovo products and the results they may have achieved. Actual environmental costs and performance characteristics may vary by customer.

Information concerning non-Lenovo products was obtained from a supplier of these products, published announcement material, or other publicly available sources and does not constitute an endorsement of such products by Lenovo. Sources for non-Lenovo list prices and performance numbers are taken from publicly available information, including vendor announcements and vendor worldwide homepages. Lenovo has not tested these products and cannot confirm the accuracy of performance, capability, or any other claims related to non-Lenovo products. Questions on the capability of non-Lenovo products should be addressed to the supplier of those products.

All statements regarding Lenovo future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only. Contact your local Lenovo office or Lenovo authorized reseller for the full text of the specific Statement of Direction.

Some information addresses anticipated future capabilities. Such information is not intended as a definitive statement of a commitment to specific levels of performance, function or delivery schedules with respect to any future products. Such commitments are only made in Lenovo product announcements. The information is presented here to communicate Lenovo's current investment and development activities as a good faith effort to help with our customers' future planning.

Performance is based on measurements and projections using standard Lenovo benchmarks in a controlled environment. The actual throughput or performance that any user will experience will vary depending upon considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve throughput or performance improvements equivalent to the ratios stated here.

Photographs shown are of engineering prototypes. Changes may be incorporated in production models.

Any references in this information to non-Lenovo websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this Lenovo product and use of those websites is at your own risk.