



# Tier 0 Deployment Guide

## Overview

Hammerspace unlocks a new tier of storage by transforming existing local NVMe storage on GPU servers into a Tier 0 of ultra-fast, persistent shared storage. By activating this previously 'stranded' local NVMe storage seamlessly into the Hammerspace Global Data Platform, Tier 0 delivers data directly to GPUs at local NVMe speeds, unleashing untapped potential and redefining both GPU computing performance and storage efficiency.

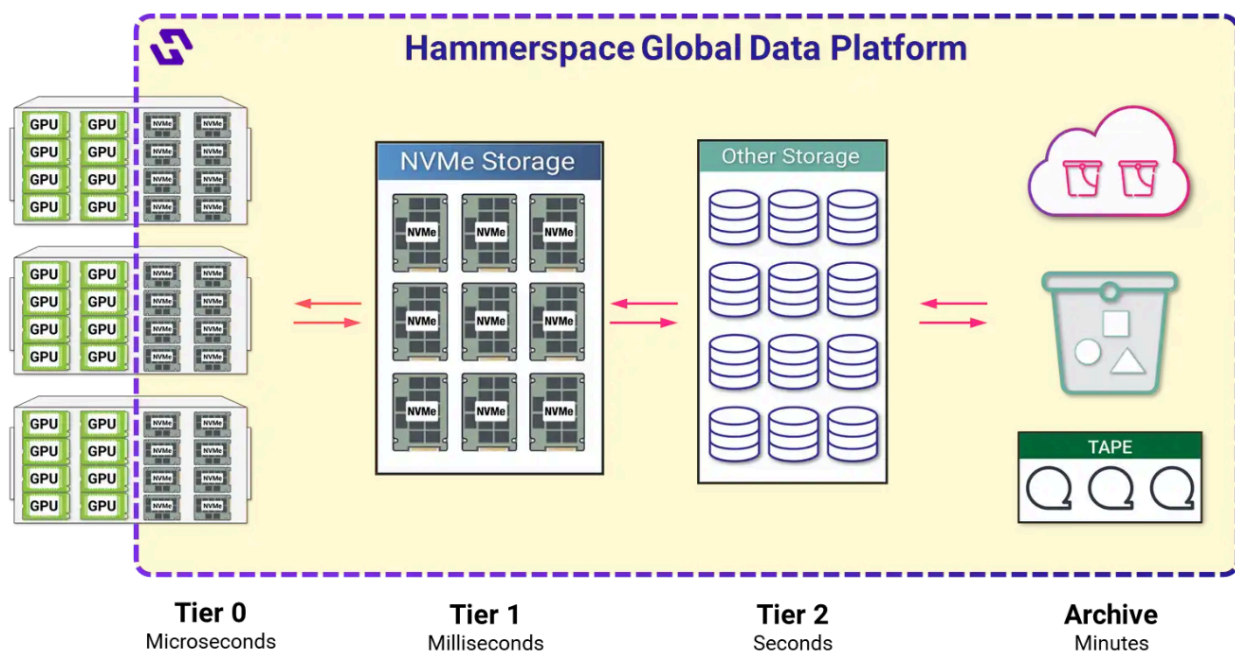
This innovation works with existing on-premises and/or cloud-based GPU (or CPU/ARM) servers from any vendor, by activating internal local NVMe capacity that has already been paid for.

This ultra-fast local NVMe capacity – which is largely underutilized – is a sunk cost, since it is already included within the price of GPU servers in both on-premises and cloud-based deployments. The fact that this also doesn't add any additional power requirements is also a net benefit for data centers who are near the limit of their available power.

Activating this local NVMe as Tier 0 shared storage is done leveraging intelligence built into the Linux kernel, which is included in all standard distributions. This means that servers do not require installation of proprietary client software, or any other alterations as they are delivered from manufacturers.

Moreover, utilizing this local Tier 0 capacity reduces or can even eliminate the costs of adding extreme-performance networking infrastructure and/or additional external high-performance storage systems, all of which are very expensive ways of trying to lessen the impact of the network bottleneck when trying to keep the GPUs fully utilized.

Tier 0 delivers these benefits while also providing global multi-protocol file/object access, non-disruptive data orchestration, data protection, and other enterprise-class data services across a global namespace that unifies storage silos from any vendor in one or more locations, including hybrid cloud environments, accelerating time to value for AI investments.



With Tier 0, organizations can fully leverage existing local NVMe on GPU servers as a shared resource, activating this higher-performing tier with the redundancy and protection of conventional external storage tiers. Moreover, utilizing this local Tier 0 capacity reduces or can even eliminate the costs of adding extreme-performance networking infrastructure and/or additional external high-performance storage systems, all of which are very expensive ways of trying to lessen the impact of the network bottleneck when trying to keep the GPUs fully utilized.



# Table of Contents

## Version 1.0

<b>What is Tier 0?</b>	<b>6</b>
Goal of this Document	7
<b>Tier 0 Superpowers</b>	<b>8</b>
Make Use of Stranded Capacity	8
Data Orchestration Across Multiple Tiers	8
Job Scheduler-Integrated Orchestration	8
Data to Node Affinity	8
HCI and CPU-based Servers	9
<b>Architecture</b>	<b>10</b>
Servers	10
Tier 0 Node	10
Linux Storage Server	10
Networking	11
Flat Networks	11
Jumbo Frames (MTU)	11
Bonding (Link Aggregation)	11
IP per NIC	11
Availability Zones	12
<b>Prerequisites</b>	<b>12</b>
Configuration Management with Teams & Admins	12
Repository Access	13
Packages for Tier 0 Nodes	13
<b>Preparation</b>	<b>13</b>
<b>Setup</b>	<b>13</b>
Automation	13
System BIOS, Firmware, NIC Firmware	14
Drives, RAID, Filesystems	15
Tips to Improve RAID Performance	15
Pre-Setup Checks	15
Configuring the Server: RAID and Filesystems	17
NFS, Exports, /etc/exports	17
Export Options	18
Firewall	20
<b>Hammerspace Deployment and Integration</b>	<b>21</b>



Hammerspace Installation	21
Cluster Integration	21
Adding a Tier 0 node or LSS as a storage node	21
Add Node Exports as Volumes	22
Create share(s)	22
Mounting Shares on Clients	23
Mount Parameter and Option Notes	23
Guidance for Setting NCONNECT with Respect to NFS Threads	24
Goal	24
Formula	24
Takeaway	25
noatime	25
<b>Validation</b>	<b>25</b>
Configuration Validation	25
Data Placement	26
Data Placement Using Objectives	27
Default Behavior	27
Simple Objective to Make Multiple Instances	28
Keeping Instances Where You Want Them	33
Configuring for Tier 0 Nodes to Write Local	34
<b>After-Setup Tasks</b>	<b>34</b>
<b>Ongoing Usage</b>	<b>35</b>
General	35
Ingesting Data	36
Snapshots	36
<b>Monitoring</b>	<b>36</b>
<b>Troubleshooting</b>	<b>37</b>
Volume Issues - Volumes go “suspected”	37
Mobility Failures (Unsuccessful)	38
Mobility Troubleshooting	38
<b>Maintenance</b>	<b>39</b>
Adding Tier 0 or LSS Nodes	39
<b>Conclusion</b>	<b>39</b>
<b>Additional Resources</b>	<b>39</b>
<b>Appendix 1 - Installing Data Movers on Linux Nodes</b>	<b>40</b>
Outline	40
Supported Distributions / Platforms	40
Dependencies	40
Get the Hammerspace Data Services Components Kit	41



Installation, Setup, and Configuration of the DI Using RPM	41
Troubleshooting	44
<b>Appendix 2 - Use Tier 0 for Checkpointing</b>	<b>46</b>
Checkpointing Workflow Example	46
Tier 0 Checkpoint Analysis	46
Estimating Checkpoint Data Size	46
Estimating Checkpoint Time When Writing to External Storage	47
Estimating Checkpoint Time When Writing to Tier 0 Local NVMe Storage	47
<b>Glossary</b>	<b>48</b>



## What is Tier 0?

Tier 0 means making use of otherwise stranded or inefficiently used storage devices in compute nodes (including GPU and CPU based compute nodes). The name is taken from the idea that storage systems outside of the GPU nodes make up traditional tiers, starting with tier 1. Tier 0 is simply the tier closest to the compute node - right inside it.

There are multiple ways these devices can be put to work, increasing efficiency in multiple dimensions which may include simply using the space, through making the whole AI job more efficient by putting data used by the GPU nodes running the job on the nodes actually doing the work.

It's important to understand the distinction between Tier 0 nodes and Linux Storage Server (LSS) nodes. Tier 0 nodes include the compute function such as GPUs. The only storage service they provide is exporting (NFS term for sharing) a number of internal storage devices as file systems managed by Hammerspace, but consumed by Tier 0 compute nodes. GPU nodes are usually sold in a specific server configuration including a set of internal storage devices, usually NVME, and run an operating system dictated by the GPU vendor or systems integrator. For example, Nvidia based platforms typically require Ubuntu as the operating system on the GPU nodes, which is fully supported by Hammerspace.

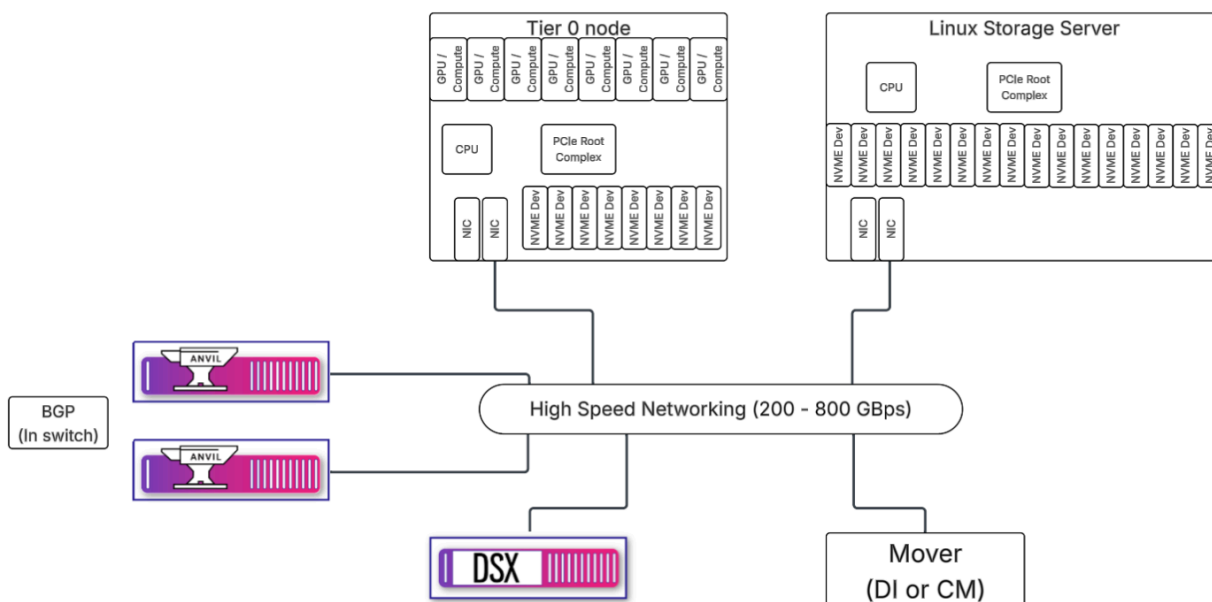
LSS nodes are more flexible in terms of both configuration and operating system choice. Customers can select the server vendor, model, and configuration of their choice, provided it meets the requirements. Hammerspace offers some predefined configurations referred to as "ready nodes" as one option for LSS. While several supported options exist for operating systems, Hammerspace recommends RHEL-based distributions based on 9.5 including Rocky 9.5. However, since NFS is an open standard and built into all Linux kernels, we do not limit to 9.5. Customers should consult with their Hammerspace technical team to verify any compatibility, feature support, or other issues.

LSS nodes can provide the following services:

- Storage volumes
- Data Mover
  - NFS mover, also known as Data Instantiator (DI)
  - Cloud Mover (CM), which supports any S3-compatible cloud or object storage
- Hammerspace S3 services
- Quorum services for Anvil nodes



Some customers have Tier 0 as the first tier, and LSS as the second, which they may call Tier 1. Others have third-party enterprise storage as Tier 1. You might have LSS as the first tier of storage, and not use any of the storage in GPU nodes as Hammerspace managed volumes. Although that technically isn't Tier 0, it is an architecture used by some large Hammerspace customers, and the storage configuration of the LSS is essentially similar to how it would be done on Tier 0 nodes.



## Goal of this Document

The purpose of this document is to outline the setup, configuration, and management of Hammerspace Tier 0.

This document will:

- Provide an overview of solution architecture, including components, their relationships, and connections
- Provide an overview of how to configure Tier 0 and Linux Storage Server nodes, add them to a Hammerspace cluster both as servers and as clients
- Provide an overview of how to install, configure, and connect an NFS mover
- Provide some example objectives for data placement



## Tier 0 Superpowers

These superpowers are listed from the most basic, to the most sophisticated. Each builds on the previous capabilities.

### Make Use of Stranded Capacity

The most basic level of utility for Tier 0 is putting data on any drive on any node in the cluster without regards to affinity to nodes, whether data is active or idle, or where else the data might be needed. Any orchestration is limited to having multiple instances of each file on different nodes.

### Data Orchestration Across Multiple Tiers

The next level of use case adds orchestration to keep active data instances on Tier 0 and move idle data to less expensive tiers that are separate from the GPU nodes. The orchestration is managed by a storage service, with no regard for integration with the GPU/AI job management.

### Job Scheduler-Integrated Orchestration

The AI job scheduler knows which nodes will be assigned to a job. Integrating storage and data management APIs into the job setup workflow allows the scheduler to control data placement into the nodes that will run the job.

#### Advantages:

- Higher probability data needed by the job is on a node running the job.
- Remove vulnerability of instance or data loss due to reimage of nodes used for other jobs.

### Data to Node Affinity

This level of Tier 0 consists of three elements:

1. Orchestration to place data on job-specific nodes.
2. When data is placed on multiple nodes, Anvil provides layout with the requesting client as the first entry in the layout.
3. The ability for I/O to bypass the NFS and network stack when accessing a local instance of a file. This is referred to as localio in the Linux NFS client.





## HCI and CPU-based Servers

While this document mostly discusses GPU-based workloads that are most urgently needed for AI use cases, the fact is that CPU-based servers often ship with the same ultra-fast local NVMe drives that can be activated as Tier 0. This is also true for low-power ARM CPUs.

These server types are often confused with Hyperconverged Infrastructure (HCI), which typically combine compute and storage in a single unit. HCI systems are predominantly used for structured data, VMs, and other block-oriented I/O patterns that are not suitable for CPU-intensive unstructured data workloads. And while HCI systems can use distributed storage architectures and other tiering solutions, their focus on block-oriented use cases make them unsuitable for processing unstructured data.

There are numerous CPU-intensive use cases that need large volumes of unstructured data, and that can take advantage of the same performance and cost improvements of Tier 0 as the GPU-computing workloads outlined above. Among them are other AI use cases, which include CPU-based AI inferencing, predictive analytics, or machine-learning pipelines dealing directly with unstructured datasets.

However, beyond AI use cases there are numerous other CPU-based application workflows that demand high performance and that can directly benefit from the ultra-low latency and enhanced throughput provided by activating Tier 0 on local NVMe drives within CPU servers.

Some of these use cases include:

- **Real-time analytics and data streaming:** Rapid processing of log data, sensor inputs, or clickstream data to generate insights instantly.
- **High-performance search and indexing:** Fast indexing and search across large-scale document repositories or datasets (e.g., Elasticsearch, Splunk).
- **Video and media processing:** Low-latency editing, transcoding, rendering, and content delivery requiring quick, frequent data access.
- **Financial trading and modeling:** Real-time risk analysis, algorithmic trading, and financial modeling needing immediate access to large datasets.
- **Genomics and bioinformatics:** Rapid processing of genome sequencing data, imaging data, or bioinformatics workloads demanding near-instantaneous data retrieval.

The benefits for CPU-based Tier 0 include the same advantages highlighted in the GPU-based examples above:

- The lower cost of activating ultra-high performance local NVMe drives within the servers vs. going over expensive networks to external arrays.

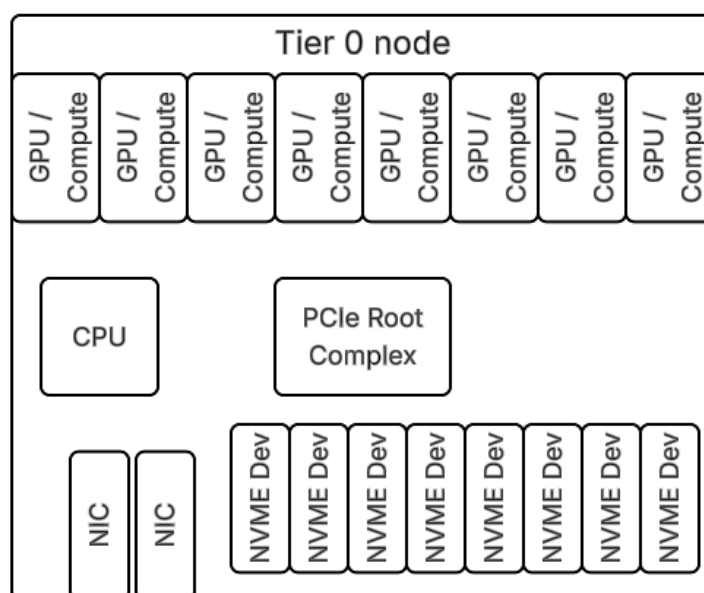
- The significant boost in performance and lower latency that Tier 0 makes possible compared with external storage, which directly benefits these applications.

## Architecture

### Servers

#### Tier 0 Node

A Tier 0 node consists of a server platform with CPU, memory, disk, and networking as the baseline, with some specialized or additional compute which may be in the form of GPUs or extra CPU, along with additional storage beyond the boot device. It also often includes additional networking, some of which may be dedicated or intended for node to node compute related communication.



#### Linux Storage Server

LSS typically has more storage devices than a Tier 0 node (24, 32, or more depending on the system selected). An LSS can also provide additional services such as data movers and quorum. An LSS is typically not configured as a client, except as necessary for services it provides.



## Networking

Tier 0 deployments should use recent high speed networking technology supporting 200 - 800Gbps speeds.

### Flat Networks

Hammerspace recommends flat networks as much as practicable in Tier 0 environments. The goal is to minimize network hops between nodes. With multiple availability zones, each zone could be a flat network (subnet).

### Jumbo Frames (MTU)

- It is recommended to use the customer's typical MTU (Maximum Transmission Unit) settings. MTU values should generally be in the range of 1500 to 9000; anything outside of this range should be discussed.
- Switches should typically be set to an MTU of 9216.
- Modern switches often ship with Jumbo frames enabled by default. While Jumbo frames may not significantly increase wire speed, they can reduce CPU consumption and interrupt handling due to fewer packets.
- Test network connectivity with ping using a specific packet size (e.g., 8972 for a 9000 MTU) and the "do not fragment" flag is advised to confirm end-to-end connectivity.

### Bonding (Link Aggregation)

- **General Recommendation (Tier 0):** Bonding is generally not recommended for Tier 0 or LSS as it cannot be NUMA-aware and has implications for switch and kernel versions; using IP per NIC is preferred for Tier 0 and LSS nodes. Bonding is supported on Hammerspace Anvil and DSX nodes (LACP mode 4).
- **Caveat:** Some customer environments (e.g., cloud, pre-existing infrastructure) may require bonding for redundancy or throughput. In such cases, professional services engagement is necessary to ensure proper configuration.

### IP per NIC

- **Recommendation:** For systems with multiple NICs, assign a unique IP address to each NIC.
- **Rationale:** Provides multiple independent network paths, improving parallelism and resilience.
- **Configuration:** Requires specific `sysctl` and network configuration to ensure proper routing and load balancing without introducing routing conflicts.



## Availability Zones

An availability zone is a part of an IT infrastructure system that's designed to be completely independent of other availability zones. This means it doesn't share any critical components, like power, cooling, or network, with other zones.

Components to consider may include:

- Power, whether circuit, rail, panel or similar
- Cooling
- Fire suppression zones
- Separate rooms or floors within a datacenter
- Racks or rows
- Network equipment and architecture
- Security access zones

While many discussions focus on cloud provider AZs, in this context, we focus on AZs within a datacenter. Hammerspace supports defining AZs for storage resources using a naming convention with a prefix indicating the AZ.

Hammerspace recommends that availability zones be designed approximately symmetrical in that they each have the same number of storage nodes (Tier 0 or LSS) and volumes. You should also place redundant Hammerspace nodes, including Anvil, DSX (if used), and mover nodes in different AZs. If Anvil nodes are in different AZs, and each AZ is a different subnet (or set of subnets), you will need to implement the ExaBGP (Border Gateway Protocol service) on the Anvils and configure the switches to allow Anvils to trigger the movement of the subnet containing the cluster management IP to the switch port to which the primary Anvil is connected.

When data is stored only on Tier 0, a minimum of 4 AZs is required, and 6 AZs is recommended. With 4 AZs, one can be taken offline for maintenance, leaving 3 providing full 3-way redundancy. Having additional AZs (5 or 6 total) allows for failures with automatic resilvering across three surviving AZs when an AZ is offline for maintenance.

## Prerequisites

### Configuration Management with Teams & Admins

Ensure other teams/admins with control over storages nodes, especially Tier 0/GPU nodes, do not take control of drives used for Tier 0 and reformat or remount. One option is to let them



dictate mount point, and agree to how much space will be used for different purposes. Hammerspace can configure thresholds to comply with usage guidelines.

## Repository Access

Since several packages need to be installed, access to a repository containing the packages is essential. Sites without internet access will have to establish their own repos that include packages used through Tier 0 deployment procedures. Note package requirements for basic LSS /Tier 0 node deployments, as well as for additional components such as data movers. Requirements include RPM style packages as well as modules for Python scripts, and, of course, software provided by Hammerspace.

## Packages for Tier 0 Nodes

- mdadm
- mkfs.xfs (Not installed by default on some distributions)
- nfs-server
- If using RDMA and Intel or Broadcom NICs, the OEM driver is required. Do not use the Linux distribution inbox driver.

## Preparation

Gather nodes, server, client information. Having a master list will help manage things like export rules.

- Administrative (root or equivalent) credentials
- Hostname, IP address(es), network connectivity
- Role
  - Client only
  - Server only (LSS)
  - Client and server (GPU / Tier 0 nodes)
  - Mover
- Cluster membership
- Availability zone
- Unusual configuration, drivers, etc. that need to be set or installed

## Setup

### Automation

As you read through these steps and consider the number of nodes involved, it will become obvious that these procedures are repetitive and potentially error-prone. Automation should be



utilized to to set up, configure, and manage. To this end, Hammerspace has a sample script that can be used for initial configuration of Tier 0 nodes, LSS nodes, and the Hammerspace cluster. While the instructions in this document use the Hammerspace Admin command line, the script uses the Hammerspace REST API, since anything that can be done in the Admin CLI and GUI can be done via the API.

Additionally, Hammerspace has Ansible playbooks for Hammerspace cluster setup on GitHub, and the plan is to provide specific guidance on using Ansible for initial setup with Tier 0 as well as ongoing maintenance and administration. Other configuration management tools are being investigated. Hammerspace welcomes customer feedback on tools they would like to use for GPU and Tier 0 automation.

## System BIOS, Firmware, NIC Firmware

The following items need to be considered and prepared on the server and peripheral hardware:

- System BIOS
  - Nodes Per Socket (NPS) - Hammerspace recommends that Nodes Per Socket be set to 1 on AMD platforms. On Intel, this is called Sub-NUMA Clustering (SNC).
  - Hyperthreading - Simultaneous Multithreading (SMT) can depend on server function, but the general recommendation is to set it to off in the system BIOS since it can introduce latency variance and reduce sustained memory copy performance, impacting applications sensitive to predictable I/O.
- Memory architecture and optimization
  - DIMM types and speeds should match the memory controller of the CPU.
  - Populate all banks appropriately to maintain maximum memory speed. Review CPU specifications for number of memory channels and impact of number of DIMMs per channel (DPC).
- Drive firmware
- NIC firmware
- NIC Driver
- NVME configuration
  - Optimal slot / bay usage considering PCIe lanes, NUMA node. Some of this is dictated by the server architecture.
  - Namespaces. Hammerspace is conducting ongoing testing with 2 namespaces per drive when individual NVME devices are used. In the meantime, a single namespace per drive is sufficient. With RAID, the recommendation is a single namespace for each NVME device.
  - Select an NVME format that does not use LBA Metadata space (PI)



- Sector size (512 vs 4096 byte). 4096 byte sector size is strongly recommended for Tier 0 and LSS nodes. Hammerspace node boot devices currently require 512 bytes per sector.
- Device drivers and settings especially for NICs
  - There may be different device drivers for advanced networking such as RDMA

## Drives, RAID, Filesystems

Check how drives enumerate. Are they consistent between reboots or other events?

NVME format. Metadata space (NVME metadata space, not HS). Namespaces.

Determine whether RAID or individual drives will be used. One consideration is how many volumes will need to be managed, which is a function of how many servers times how many drives they have. If the number of volumes is projected to be over 1000, RAID configuration should be favored.

### Tips to Improve RAID Performance

- Where possible, use RAID set sizes that are a power of 2
- Attempt to build RAID sets on NUMA local NVMe
- The default mdadm chunk size of 512KiB will suffice for many deployments. For advanced tuning, consult with your Hammerspace tech team to calculate a chunk size matched to I/O profile.

Determine NUMA placement of NVME devices as follows:

#### Linux Command Line

```
# grep "" /sys/block/nvme*n1/device/numa_node
/sys/block/nvme0n1/device/numa_node:0
/sys/block/nvme10n1/device/numa_node:1
/sys/block/nvme11n1/device/numa_node:1
/sys/block/nvme1n1/device/numa_node:0
/sys/block/nvme2n1/device/numa_node:0
/sys/block/nvme3n1/device/numa_node:0
/sys/block/nvme4n1/device/numa_node:0
/sys/block/nvme5n1/device/numa_node:0
/sys/block/nvme6n1/device/numa_node:1
/sys/block/nvme7n1/device/numa_node:1
/sys/block/nvme8n1/device/numa_node:1
/sys/block/nvme9n1/device/numa_node:1
```



## Pre-Setup Checks

Once the base server hardware is set up, operating system installed, and basic networking is configured, you can check the recommendations listed above.

If jumbo frames are in use, check the configuration with ping using a large frame size and the Do Not Fragment flag set. Note that the frame size of the ping is smaller than the MTU because of packet overhead. The most common MTU, 9000, should be checked with a ping size of 8972. You should check node (Tier 0 and LSS) to node, node to Hammerspace nodes, and other combinations including DI, etc. Every node should be checked.

### Linux Command Line

```
# ping -M do -s 8972 192.168.0.101
PING 192.168.0.101 (192.168.0.101) 8972(9000) bytes of data.
8980 bytes from 192.168.0.101: icmp_seq=1 ttl=64 time=0.013 ms
8980 bytes from 192.168.0.101: icmp_seq=2 ttl=64 time=0.004 ms
8980 bytes from 192.168.0.101: icmp_seq=3 ttl=64 time=0.004 ms
^C
--- 192.168.0.101 ping statistics ---
3 packets transmitted, 3 received, 0% packet loss, time 2038ms
rtt min/avg/max/mdev = 0.004/0.007/0.013/0.004 ms
```

Listing Storage Devices and Current Usage:

### Linux Command Line

```
# lsblk -n
NAME                MAJ:MIN RM   SIZE RO TYPE MOUNTPOINTS
nvme11n1             259:0    0 894.3G  0 disk
nvme10n1             259:1    0 745.1G  0 disk
├─nvme10n1p1         259:2    0   600M  0 part /boot/efi
├─nvme10n1p2         259:3    0     1G  0 part /boot
├─nvme10n1p3         259:4    0    32G  0 part [SWAP]
└─nvme10n1p4         259:5    0 711.5G  0 part /
nvme0n1              259:6    0     7T  0 disk
nvme1n1              259:7    0     7T  0 disk
nvme2n1              259:8    0     7T  0 disk
nvme3n1              259:9    0     7T  0 disk
nvme4n1              259:10   0     7T  0 disk
nvme5n1              259:11   0     7T  0 disk
nvme6n1              259:12   0     7T  0 disk
nvme7n1              259:13   0     7T  0 disk
nvme8n1              259:14   0     7T  0 disk
nvme9n1              259:15   0     7T  0 disk
nvme12n1             259:16   0     7T  0 disk
```





```
nvme13n1    259:17    0    7T    0 disk
nvme13n1    259:18    0    7T    0 disk
nvme15n1    259:19    0    7T    0 disk
nvme16n1    259:20    0    7T    0 disk
```

## Configuring the Server: RAID and Filesystems

This shows an example RAID configuration of 3 RAID arrays of 8 drives each. The last two commands ensure that the RAID configuration is saved in the mdadm configuration file.

### Linux Command Line

```
# mdadm --create --verbose /dev/md0 --level=0 --raid-devices=8 /dev/nvme0n1
/dev/nvme1n1 /dev/nvme2n1 /dev/nvme3n1 /dev/nvme4n1 /dev/nvme5n1 /dev/nvme6n1
/dev/nvme7n1
# mdadm --create --verbose /dev/md1 --level=0 --raid-devices=8 /dev/nvme9n1
/dev/nvme10n1 /dev/nvme11n1 /dev/nvme12n1 /dev/nvme13n1 /dev/nvme14n1 /dev/nvme15n1
/dev/nvme16n1
# mdadm --create --verbose /dev/md2 --level=0 --raid-devices=8 /dev/nvme17n1
/dev/nvme18n1 /dev/nvme19n1 /dev/nvme20n1 /dev/nvme21n1 /dev/nvme22n1 /dev/nvme23n1
/dev/nvme24n1
# mkdir /etc/mdadm
# mdadm --detail --scan | tee -a /etc/mdadm/mdadm.conf
```

Set up mount points, partition arrays, make a file system, mount, and add to fstab:

### Linux Command Line

```
# mkdir /hsvol0
# parted /dev/md0 mklabel gpt
# parted -a opt /dev/md0 mkpart primary xfs 0% 100%
# mkfs.xfs -d agcount=512 -L hsvol0 /dev/md0p1
# mount /dev/md0p1 /hsvol0
# echo '/dev/md0p1 /hsvol0 xfs 0 0' >> /etc/fstab
```

After editing `/etc/fstab`, you should run `'systemctl daemon-reload'` to update systemd units generated from the file.

## NFS, Exports, `/etc/exports`

If the NFS server (`nfs-kernel-server`) is already installed, check `/etc/nfs.conf` and `/etc/nfs.conf.d/local.conf` for non-default TCP/UDP port configuration.

Hammerspace recommends the following `/etc/nfs.conf` settings for use in Tier 0 environments:

**Linux /etc/nfs.conf file**

```
[nfsd]
threads=128
# udp=n
# tcp=y
vers3=y
# vers4=y
vers4.0=n
vers4.1=n
vers4.2=y
rdma=y
rdma-port=20049
```

NFS v4.2 is required for Parallel NFS, client side mirroring and other advanced features used with Tier 0 implementations. Although clients mount the Hammerspace share(s) via NFS v4.2, the I/O to storage volumes uses NFS v3, so NFS servers including Tier 0 nodes must enable NFS v3.

Run the following command to start NFS and configure it to start on system startup:

**Linux Command Line**

```
# systemctl enable --now nfs-server.service
```

**NOTE:** On some Linux distributions, the NFS service might be `nfs-kernel-server`, although that is usually linked to `nfs-server`.

Ensure `showmount` is *not* disabled (uncommon).

For `/etc/exports`:

- Add Hammerspace nodes and Anvil cluster mgmt floating IP RW, `no_root_squash, sync, secure, mp, no_subtree_check`.
- Add clients RW, `root_squash, sync, secure, mp, no_subtree_check`. This can be done by subnet or host

**Export Options**

The following export options are recommended for nodes that serve NFS, including Tier 0 nodes and LSS:

- `rw` - Read/write access



- `no_root_squash` - Allows appropriate nodes to have root access on the export. This setting is required for Hammerspace nodes (Anvil and DSX), and nodes running mover services (DI and CM).
- `root_squash` - Prohibits nodes from root access on the export. This should be used for Tier 0 and LSS nodes that are not running mover services.
- `sync` - Causes the NFS server to ensure writes are committed to storage before replying to the client. This is default behavior in modern NFS server implementations, but specifying the option may suppress warning messages.
- `secure` - Requires that clients use a secure, well-known TCP port (less than 1024).
- `mp` - Mountpoint indicates that the export should have a filesystem mounted for export. Using this option prevents accidentally exporting the directory under the mounted filesystem, which would result in data being written to the root filesystem.
- `no_subtree_check` - This is a default in modern NFS server implementations. The main reason to specify this parameter is to suppress a warning that the default behavior has changed.

The following example of `/etc/exports` shows 3 exported volumes with 3 exports for Hammerspace nodes and 3 exports for a subnet of Tier 0 or LSS nodes.

#### Linux `/etc/exports` file

```
/hsvo10 10.1.2.3(rw,no_root_squash,sync,secure,mp,no_subtree_check)
/hsvo11 10.1.2.3(rw,no_root_squash,sync,secure,mp,no_subtree_check)
/hsvo12 10.1.2.3(rw,no_root_squash,sync,secure,mp,no_subtree_check)
/hsvo10 10.2.1.0/24(rw,root_squash,sync,secure,mp,no_subtree_check)
/hsvo11 10.2.1.0/24(rw,root_squash,sync,secure,mp,no_subtree_check)
/hsvo12 10.2.1.0/24(rw,root_squash,sync,secure,mp,no_subtree_check)
```

Many distributions allow multiple IPs on a single line with their common export options first:

#### Linux `/etc/exports` file

```
/mnt/hsvo11 -rw,no_root_squash,sync,secure,mp,no_subtree_check 10.200.104.38
10.200.105.195
/mnt/hsvo11 -rw,root_squash,sync,secure,mp,no_subtree_check 10.200.102.33
10.200.103.223 10.200.101.217 10.200.102.92 10.200.102.34 10.200.103.170
10.200.102.172 10.200.103.173
```

After updating `/etc/exports`, run the following command to re-read the exports file:

**Linux Command Line**

```
# exportfs -r
```

**Firewall**

You may need to allow ports used by NFS and related protocols and services. Check `/etc/nfs.conf` and the files in `/etc/nfs.conf.d` for port specifications, and also `rpcinfo -p` to confirm. We have seen sites that had other applications on GPU nodes configured to use ports normally used by NFS.

For Rocky, Centos, RHEL using standard ports:

**Linux Command Line**

```
# firewall-cmd --permanent --add-service=nfs --add-service=rpc-bind  
--add-service=mountd ; firewall-cmd --reload
```

For Ubuntu and Debian using standard ports:

**Linux Command Line**

```
# ufw allow port 111 # portmapper  
# ufw allow port 2049 # NFS  
# ufw allow port 20048 # mountd
```

Confirm it works by checking exports from another client (can be another LSS/Tier 0 server):

**Linux Command Line**

```
# showmount -e 10.200.103.167  
Export list for 10.200.103.167:  
/mnt/hsvol1  
10.200.112.0/21,10.200.10.0/24,10.200.103.166,10.200.103.167,10.200.102.220,10.200.  
101.3,10.200.106.197  
/mnt/hsvol0  
10.200.112.0/21,10.200.10.0/24,10.200.103.166,10.200.103.167,10.200.102.220,10.200.  
101.3,10.200.106.197
```

# Hammerspace Deployment and Integration

## Hammerspace Installation

General installation of Hammerspace nodes is well-documented in the Hammerspace Installation and Licensing Guide. What follows here are specifics that may apply to Tier 0 in general, or your environment.

- **Anvils:** If Anvils are to be installed in different subnets, BGP (Border Gateway Protocol) is required, using the ExaBGP tools in the Anvil nodes to communicate with connected switches in order to control which port routes the subnet for the Anvil floating management IP. Contact your Hammerspace tech team for architecture and deployment details.
- **DSXs:** DSX may not be required for some Tier 0 installations. Mover (DI) may be deployed as RPM or container. See the appendix for instructions on deploying the DI RPM.

## Cluster Integration

LSS and Tier 0 storage services can be added using the Hammerspace GUI, CLI, or scripted with the API. For a small installation, it may be simple enough to add storage systems via GUI or CLI. For larger installations, setup should be automated using a script or configuration management or automation tools such as Ansible or Terraform.

### Adding a Tier 0 node or LSS as a storage node

Add a storage node using the node-add command:

#### Hammerspace Admin Command Line

```
anvil1> node-add --type OTHER --name lss001 --ip 10.1.2.2
```

If you plan to use per-node objectives (for example to achieve local write affinity), you can have the node-add command create them as part of adding the node:

#### Hammerspace Admin Command Line

```
anvil1> node-add --type OTHER --name lss001 --ip 10.1.2.2  
--create-placement-objectives
```



If successful, the node-add command responds with a description of the node, with the exports listed as logical volumes.

## Add Node Exports as Volumes

Add each export as a volume using the volume-add command:

### Hammerspace Admin Command Line

```
anvil1> volume-add --name lss001::/hsvol1 --node-name lss001 --access-type  
read_write --logical-volume-id 123 --low-threshold 90 --high-threshold 95  
--skip-performance-test
```

You can use either `--logical-volume-name` or `--logical-volume-id` to specify the volume. With scripting or other automation, it may be easier to use the logical volume name, since for simple NFS servers, the logical volume name is the same as the export. Alternatively, you can get the list of logical volume IDs from the output of the node-add command.

When adding many volumes in rapid succession, you should skip the performance test.

If you plan to use availability zones, which are a best practice, the volume name must be prefixed with "AZ", a whole number, colon, then a unique name. The example below shows adding a volume with the "AZ1:" prefix.

### Hammerspace Admin Command Line

```
anvil1> volume-add --name AZ1:lss001::/hsvol1 --node-name lss001 --access-type  
read_write --logical-volume-id 123 --low-threshold 90 --high-threshold 95  
--skip-performance-test
```

You can also name the node with the AZx: prefix, then by not specifying the `--node-name` parameter, with the volume-add command, the volume name will default to `node::/path`. With this trick, the volume name gets the AZx: prefix from the node name, since all volumes on a node would be in the same AZ.

## Create share(s)

You can create many shares for various purposes, such as checkpointing, to contain job-specific data, or results/output.

Run the following to create a share:



### Hammerspace Admin Command Line

```
anvil1> share-create --name checkpoints --path /checkpoints --export-option  
10.1.2.0/24,rw,root-squash
```

If you don't specify `--path`, it will default to the share name. You can specify a size, but must plan this carefully so jobs and applications don't run out of space unexpectedly due to a share quota. If size is not specified, the share size and free space will show as the total size and free space of all local (NFS, not object) volumes.

You must provide the necessary set of export options to include all clients needing access to the share, including all Tier 0 nodes. If LSS nodes are not used as clients, they do not need to be listed in the Hammerspace share exports. You may need to specify multiple `--export-option` parameters. If the command line gets too long, you can add export options using the `share-update` command. Clients can be specified by IP address, subnet in CIDR format, hostname, FQDN, or netgroup.

Hammerspace exports support `ro` or `rw`, `root_squash` or `no_root_squash`, and `secure` (default) or `insecure`. Generally, "`rw,root_squash,secure`" is the recommended set for Hammerspace exports to Tier 0 nodes.

You can specify objectives during share creation or add them later. See [Data Placement Using Objectives](#).

## Mounting Shares on Clients

Once shares are created, clients can mount them using NFS v4.2. The exact mount options may depend on a variety of factors including use case, model of NICs, and especially client kernel version, since the kernel version dictates which advanced NFS features are available.

### Linux Command Line

```
# mkdir /mnt/checkpoints  
# mount -o vers=4.2,nconnect=8 <anvil_IP>:/checkpoints /mnt/checkpoints
```

## Mount Parameter and Option Notes

For NFSv4.2, the IP or other host specification used is for the floating data IP of the Anvil, even though the data resides on Tier 0 and/or LSS nodes.



If mount fails with “mount.nfs: Protocol not supported”, first check firewall settings. If firewall is not the issue, your Linux kernel may be newer than the Hammerspace supported kernel whitelist allows. You can bypass this using the mount option `port=20492`, but please report the kernel version to your Hammerspace technical team.

## Guidance for Setting NCONNECT with Respect to NFS Threads

When using NFS with `nconnect` in large-scale environments, it’s important to manage oversubscription of the NFS server’s `knfsd` threads. Excessive oversubscription can lead to contention, latency, or degraded performance.

### Goal

Maintain a maximum 8:1 oversubscription ratio of client connections (`nconnect`) to server `knfsd` threads, while also enforcing a minimum of 2 and a maximum of 16 to ensure basic parallelism.

### Formula

#### Given:

T = number of `knfsd` threads (from `nfs.conf`, e.g., `threads=128`)

N = number of client nodes expected to connect to the server

R = desired oversubscription ratio (e.g., R = 8)

C = number of `nconnect` connections each client should use

#### Then:

$C = \max(2, \text{floor}((T * R) / N))$

### Example

#### If:

T = 128 `knfsd` threads

N = 1023 clients

R = 8 (oversubscription target)

#### Then:

$C = \max(2, \text{floor}((128 * 8) / 1023))$   
=  $\max(2, \text{floor}(1024 / 1023))$   
=  $\max(2, 1)$   
= 2

You should configure `nconnect=2`. This results in a 16:1 oversubscription in this specific scenario, which exceeds the 8:1 target but honors the enforced minimum of 2.





## Takeaway

To optimize performance, adjust `nconnect` downward as the cluster grows, using the formula above. Monitor actual server thread utilization to verify assumptions, especially under real workloads. Consider increasing threads in `nfs.conf` on high-demand servers if consistent oversubscription is unavoidable.

## noatime

With most AI and similar workloads, access time is not a critical file attribute to maintain. Hammerspace recommends setting `noatime`.

# Validation

## Configuration Validation

- Using the GUI or CLI, check the storage system list and the volume list to see that the expected systems and volumes are added
- Check that newly added volumes have 0 capacity used. If a volume has an unexpected used size, it is possible or even likely that the device is not mounted as intended and what is exported is the empty directory in the root volume of the server. Using the `mp` export option on the NFS exports should prevent this issue.
- Check the root file handle on volumes added from the same server. If the root file handle are the same on different volumes, it is likely that the devices are not mounted as intended and what is exported are empty directories in the root volume of the server. Again, the `mp` export option should prevent this issue.

### Hammerspace Admin Command Line

```
anvil> volume-list --name AZ2:node2::/mnt/hsvol1
ID:                                d0f22c08-4089-4a5c-b3ee-4bdc52523b21
Name:                              AZ2:node2::/mnt/hsvol1
Internal ID:                        15
Discovered addresses:               [IP: 10.200.103.53/32, Port: 2049, NetId: tcp, NodeNum: 0]
Effective IPs:                      [IP: 10.200.103.53/32, Port: 2049, NetId: tcp, NodeNum: 0]
Path:                              /mnt/hsvol1
State:                              OK
Access type:                        Read Write
Node:                              AZ2:node2
```



```
Oper state:      Up
Admin state:     Up
Capacity:        [Total: 7TB, Used: 158.6MB (<1%), Free: 7TB]
Capabilities:    Max read IOPs:      0
                  Max write IOPs:     0
                  Min read latency:    0.001 milliseconds
                  Min write latency:   0.001 milliseconds
                  Max read bandwidth:  0 KB/s
                  Max write bandwidth: 0 KB/s
                  NFS read request size: 1.0 MiB
                  NFS write request size: 1.0 MiB
                  Tolerable read latency: 100000 milliseconds
                  Tolerable write latency: 100000 milliseconds
                  High threshold:      95%
                  Low threshold:       91%
                  Availability:         99%
                  Availability (effective): 99%
                  Availability drop:    Disabled
                  Durability:           99.9%
                  Durability (effective): 99.9%
                  Encryption:          false
                  Clone status:        Not supported
                  Online delay:        Online

Created:         2025-07-10 18:13:02 UTC
Modified:        2025-07-10 18:13:02 UTC
Root file handle: 010006002ec5b990aceb436693c0d035b7325b53
Locations:
                  [Type: Node, ID: 74a111ce-6069-4427-8ebd-d75ddae8a63a,
Name: AZ2:node2]
                  [Type: Storage Volume, ID:
d0f22c08-4089-4a5c-b3ee-4bdc52523b21, Name: AZ2:node2:./mnt/hsvol1]
                  [Type: Volume Group, ID:
22c4be93-2717-4dd5-badd-345d1c912e62, Name: all]
Random writes:   Enabled
```

## Data Placement

Using one of the clients, create a number of files in a Hammerspace share, ideally more files than the number of volumes or a multiple of the number of volumes as the number of files. Check the locations of these files. You can use tools like HSTK and the command `hs usage volume <share path>` to see the distribution of files across the volumes.



You can look at the placement of individual files using `hs eval -e instances.volume <filename>`.

## Data Placement Using Objectives

Objectives are Hammerspace rules that control the behavior of files. They control creation and placement of instances on storage volumes and other advanced functionality. There is extensive documentation on objectives in the Hammerspace Administrator Guide and in the [Hammerspace Objectives Guide](#).

Objectives can be set on the whole share, or directories and files within a share, giving the flexibility and granularity to control files for many use cases.

### Default Behavior

Simply by creating a share in Hammerspace, you get a set of 8 default objectives that control sane behavior of files. For example, they prevent instances from filling up one volume when there are other appropriate volumes for the data in question. Effectively, they balance the usage of volumes. If you create a set of files on a share with default objectives, you should see the files spread across all volumes with available space. These objectives will respect capacity thresholds for volumes. If a volume is at or above the low threshold, data will be placed in other volumes. If a volume is at or above the high threshold, data will be evacuated to other volumes.

Assuming a share called `share1`, we can test this behavior as follows:

#### Linux Command Line

```
# mkdir /mnt/share1
# mount -o vers=4.2,nconnect=8 <anvil_IP>:/share1 /mnt/share1
# mkdir /mnt/share1/dir1
# cd /mnt/share1/dir1
# for f in `seq 1 100`; do dd if=/dev/urandom of=f$f count=10 bs=1024k; done
10+0 records in
10+0 records out
10485760 bytes (10 MB, 10 MiB) copied, 0.0790123 s, 133 MB/s
<snip>

# hs usage volume
SUMS_TABLE{
    |KEY = STORAGE_VOLUME('AZ1:node1::/mnt/hsvo10'),
    |VALUE = 9;
```



```
|KEY = STORAGE_VOLUME('AZ1:node1::mnt/hsvo11'),  
|VALUE = 9;  
  
|KEY = STORAGE_VOLUME('AZ2:node2::mnt/hsvo11'),  
|VALUE = 9;  
  
|KEY = STORAGE_VOLUME('AZ2:node2::mnt/hsvo10'),  
|VALUE = 9;  
  
|KEY = STORAGE_VOLUME('AZ3:node3::mnt/hsvo10'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ3:node3::mnt/hsvo11'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ4:node4::mnt/hsvo10'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ4:node4::mnt/hsvo11'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ5:node5::mnt/hsvo10'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ5:node5::mnt/hsvo11'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ6:node6::mnt/hsvo11'),  
|VALUE = 8;  
  
|KEY = STORAGE_VOLUME('AZ6:node6::mnt/hsvo10'),  
|VALUE = 8}
```

The `hs` command above is from the Hammerspace Toolkit (HSTK). You can get this from GitHub, and it is extensively described in [Hammerscript and Hammerspace Toolkit](#).

## Simple Objective to Make Multiple Instances

There are a few different ways to create and apply objectives to create multiple instances of files on different storage volumes. You could create “place-on” objectives for each volume, but that quickly becomes impractical since you would have to apply different objectives to each file.



A simpler way to use an objective to make multiple instances and have them balanced across volumes, nodes or even AZs is to use a durability or availability objective, specified as a number of "9s" such as 99.999% which is referred to as "five nines". Some of these objectives are created by default, but you can create more with different durability or availability requirements. The main difference between durability and availability is that with durability, the system or volume might go offline, but the data is intact and will be there when the volume comes back online. With availability, the volume is expected to stay online for the specified percentage of time.

The way this works is when a durability or availability objective is applied, the system looks for a storage volume that by itself meets or exceeds the requirement. If there isn't one, the system determines a set of volumes that together add up to the requirement and places an instance of the file on each volume in the set. It turns out that the probability math works out that you can add up the number of nines from each volume and get the same effective reliability of a single volume with that number of nines. When there are many volumes and multiple files are created, the system determines a different set of volumes for each file.

It goes even further. If you configure availability zones using the "AZx:" prefix for volume names as discussed earlier, the system will place instances such that no two instances of the same file are placed in volumes in the same availability zone.

The default volume durability is 99.9% and availability is 99%. These values can be changed, and for most storage platforms should be set to values determined based on the storage system and the vendor documentation and guidance. For Tier 0 and LSS, the defaults are reasonable.

The following table lists default durability or availability objectives as included with a Hammerspace installation, and the number of instances you should expect with the default volume durability and availability:

Objective	Number of Instances
<b>availability-1-nine</b>	1
<b>availability-3-nines</b>	2
<b>availability-5-nines</b>	3
<b>durability-1-nine</b>	1
<b>durability-3-nines</b>	1



<b>durability-5-nines</b>	2
<b>durability-10-nines</b>	4

To get 3 instances using durability, we would need to create an objective with 2 times the volume durability plus 1, which would be 7 nines.

#### Hammerspace Admin Command Line

```
anvil1> objective-create --name durability-7-nines --durability 7 --description
'Make 3 instances when vol dur = 99.9'
ID:                f27e26e4-6c3d-40a4-907d-9fb154f35d2b
Name:              durability-7-nines
Internal ID:       8
Priority:           MEDIUM
Durability:         99.99999%
Description:        Make 3 instances when vol dur = 99.9
```

Assuming we create a directory called /csm3, we can apply the durability-7-nines to that directory:

#### Hammerspace Admin Command Line

```
anvil1> share-objective-add --name test4 --objective durability-7-nines --path
/csm3
Share name:        test4
Share internal ID: 2
Path:              /csm3
Locally applied objectives:
                    [Name: durability-7-nines, Applicability: TRUE]
All applied objectives:
                    [Name: delegate-on-open, Applicability: true]
                    [Name: durability-3-nines, Applicability:
DATA_ORIGIN_LOCAL AND IS_DURABLE]
                    [Name: durability-7-nines, Applicability: TRUE]
                    [Name: keep-online, Applicability: IS_BEING_CREATED OR
HAS_KEEP_ON_THIS_SITE]
                    [Name: layout-get-on-open, Applicability: IS_BEING_CREATED
OR HAS_ONLINE_INSTANCE]
                    [Name: optimize-for-capacity, Applicability: true]
                    [Name: availability-1-nine, Applicability:
DATA_ORIGIN_LOCAL?ALWAYS]
                    [Name: durability-1-nine, Applicability:
DATA_ORIGIN_LOCAL?ALWAYS]
```



```
[Name: keep-online, Applicability: IS_BEING_CREATED OR
HAS_ONLINE_INSTANCE AND IS_RECENTLY_USED?ALWAYS]
Active objectives:
    [Name: durability-7-nines]
    [Name: optimize-for-capacity]
    [Name: delegate-on-open]
```

You can also set objectives on files and directories using HSTK from a client:

#### Linux Command Line

```
# cd csm3
# hs objective add durability-7-nines .
```

Then you can run the same loop to create files in the /csm3 directory and look at the file locations and distributions.

#### Linux Command Line

```
# cd csm3
# for f in `seq 1 100`; do dd if=/dev/urandom of=f$f count=10 bs=1024k; done
10485760 bytes (10 MB, 10 MiB) copied, 5.95156 s, 1.8 MB/s
10+0 records in
10+0 records out
<snip>
# hs eval -e instances.volume f3
INSTANCES_TABLE{
    |VOLUME = STORAGE_VOLUME('AZ4:node4::mnt/hsvol1');
    |VOLUME = STORAGE_VOLUME('AZ5:node5::mnt/hsvol1');
    |VOLUME = STORAGE_VOLUME('AZ6:node6::mnt/hsvol0')}
[root@peterlearmonth-23360-07102025-client6 csm3]# hs eval -e instances.volume f55
INSTANCES_TABLE{
    |VOLUME = STORAGE_VOLUME('AZ1:node1::mnt/hsvol1');
    |VOLUME = STORAGE_VOLUME('AZ2:node2::mnt/hsvol1');
    |VOLUME = STORAGE_VOLUME('AZ6:node6::mnt/hsvol0')}
[root@peterlearmonth-23360-07102025-client6 csm3]# hs usage volume
SUMS_TABLE{
    |KEY = STORAGE_VOLUME('AZ1:node1::mnt/hsvol0'),
    |VALUE = 40;

    |KEY = STORAGE_VOLUME('AZ1:node1::mnt/hsvol1'),
```



```
|VALUE = 40;

|KEY = STORAGE_VOLUME('AZ2:node2::/mnt/hsvo11'),
|VALUE = 33;

|KEY = STORAGE_VOLUME('AZ2:node2::/mnt/hsvo10'),
|VALUE = 16;

|KEY = STORAGE_VOLUME('AZ3:node3::/mnt/hsvo10'),
|VALUE = 22;

|KEY = STORAGE_VOLUME('AZ3:node3::/mnt/hsvo11'),
|VALUE = 22;

|KEY = STORAGE_VOLUME('AZ4:node4::/mnt/hsvo10'),
|VALUE = 22;

|KEY = STORAGE_VOLUME('AZ4:node4::/mnt/hsvo11'),
|VALUE = 22;

|KEY = STORAGE_VOLUME('AZ5:node5::/mnt/hsvo10'),
|VALUE = 21;

|KEY = STORAGE_VOLUME('AZ5:node5::/mnt/hsvo11'),
|VALUE = 21;

|KEY = STORAGE_VOLUME('AZ6:node6::/mnt/hsvo11'),
|VALUE = 20;

|KEY = STORAGE_VOLUME('AZ6:node6::/mnt/hsvo10'),
|VALUE = 21}
```

In the current version of Hammerspace, distribution may not be exactly even, as shown above. This is fixed in the next major release. However, each file has instances in 3 different availability zones.

When you set an objective that requires more than one instance of files, and your Linux client has a recent enough kernel, up to 3 instances (with the current Hammerspace product) will be written by the client. This is referred to as Client Side Mirroring, as defined in [RFC 8435](#). If you set an objective that requires more than 3 instances, the first 3 will be written by the client using CSM, and the rest will be created by a mover (DI).





## Keeping Instances Where You Want Them

When you have another tier of storage that has higher availability and durability than your Tier 0 or LSS nodes, the system will attempt to satisfy availability and durability objectives by placing fewer instances on a volume or volumes with higher availability and durability. With many Tier 0 use cases, the data should remain on Tier 0 volumes for some period of time. To keep data on specific volumes or a set of volumes, we can use a confine-to objective with a volume group.

First, create the desired volume group. Volume groups can consist of any combination of storage systems, volumes, or even other volume groups. It may make sense to have a volume group for each AZ, and a volume group that includes all AZ volume group

### Hammerspace Admin Command Line

```
anvil1> volume-group-create --name AZ1 --expressions
'node:AZ1:node101,node:AZ1:node102,node:AZ1:node103,node:AZ1:node104,node:AZ1:node1
05,node:AZ1:node106'
anvil1> volume-group-create --name AZ2 --expressions
'node:AZ2:node201,node:AZ2:node202,node:AZ2:node203,node:AZ2:node204,node:AZ2:node2
05,node:AZ2:node206'
anvil1> volume-group-create --name AZ3 --expressions
'node:AZ3:node301,node:AZ3:node302,node:AZ3:node303,node:AZ3:node304,node:AZ3:node3
05,node:AZ3:node306'
anvil1> volume-group-create --name AZ4 --expressions
'node:AZ4:node401,node:AZ4:node402,node:AZ4:node403,node:AZ4:node404,node:AZ4:node4
05,node:AZ4:node406'
anvil1> volume-group-create --name AZ5 --expressions
'node:AZ5:node501,node:AZ5:node502,node:AZ5:node503,node:AZ5:node504,node:AZ5:node5
05,node:AZ5:node506'
anvil1> volume-group-create --name AZ6 --expressions
'node:AZ6:node601,node:AZ6:node602,node:AZ6:node603,node:AZ6:node604,node:AZ6:node6
05,node:AZ6:node606'
anvil1> volume-group-create --name AZ_all --expressions
'volume-group:AZ1,volume-group:AZ2,volume-group:AZ4,volume-group:AZ4,volume-group:A
Z5,volume-group:AZ6'
```

When you create a volume group, the system automatically creates a set of objectives for that volume group including place-on, exclude-from, and confine-to. Now you can apply the confine-to for AZ\_all to the data that needs to stay on the Tier 0 volumes. You can even include an expression to define an “applicability” such as how long the data should stay on Tier 0.

**Hammerspace Admin Command Line**

```
anvil1> share-objective-add --name test4 --path /csm3 --objective confine-to-AZ_all  
--applicability 'LAST_USE_AGE<1*HOURS?TRUE'
```

## Configuring for Tier 0 Nodes to Write Local

Future versions of Hammerspace are planned to have additional functionality to simplify node affinity for use cases like checkpointing. Currently, it is possible to use node objectives to direct I/O for a given directory to storage volumes on that node.

Verify that the job manager can instruct Tier 0 nodes to write checkpoints to a specific directory within the share.

Create per Tier 0 node directories. From a client with the necessary permissions, use the `mkdir` command:

**Linux Command Line**

```
# mkdir /mnt/checkpoints/node001
```

Add a place-on-node objective to the directory:

**Hammerspace Admin Command Line**

```
anvil1> share-objective-add --name checkpoints --objective place-on-node001 --path  
/node001
```

To protect checkpoints after the initial fast write to local storage, you need an objective to place instances of the checkpoint on other storage, either other Tier 0 nodes, on one or more Linux Storage Servers, or other robust storage such as a NAS or cloud.

## After-Setup Tasks

Send a support bundle to Hammerspace to establish a baseline for support.

**Hammerspace Admin Command Line**

```
anvil1> support-bundle --push
```



Note that support bundles post to Hammerspace using HTTPS over port 443. If firewall rules block this traffic, and there is a proxy available, you can configure the proxy using the `cluster-config` command.

#### Hammerspace Admin Command Line

```
anvil1> cluster-config --proxy-url http://10.99.42.123:12345
```

**Note:** The `proxy-url` parameter accepts `http` or `https` as the protocol.

If port 443 is blocked and there is no proxy, you can create local support bundles, download from a URL on the Anvil, then transfer them to Hammerspace using any mechanism supported in your environment. Hammerspace support can provide a secure manual upload URL.

#### Hammerspace Admin Command Line

```
anvil1> support-bundle --local --all
ID:                62c491e5-446b-454f-bc32-4901e9274413
Name:              support-bundle
Status:           COMPLETED
Status message:    You may download the support bundle from
https://10.99.42.40:8443/support/support.peter-anvil.20250611.183518281.tar
Progress:         100%
Created:          2025-06-11 18:35:17 UTC
Started:          2025-06-11 18:35:17 UTC
Ended:            2025-06-11 18:38:18 UTC
Params:           node: ; created-by-name: admin; should-push: false; node-all:
true; created-by: Uoid [uuid=de46fcae-c3d5-4151-85aa-b00dfff1ae65, objectType=USER]
Context:          output-file-path:
/support/support.peter-anvil.20250611.183518281.tar; collect-local-success: true
```

## Ongoing Usage

### General

- Ensure that jobs are configured to use storage on Hammerspace shares, both in terms of mounting the shares on the GPU nodes, and placing data on them.
- Ensure that node maintenance workflows do not take ownership of NVME drives and reformat or remount them.
- Ensure that maintenance operations affect nodes within a single availability zone at a time.

## Ingesting Data

There are essentially three ways to get data into Hammerspace:

- **Copy in** - Any data copy tool that can write to an NFS mount could be used to bring existing data into Hammerspace.
- **Create new** - Workloads produce some kind of output or result. The Tier 0 architecture is designed to store that output.
- **Assimilation** - For existing data currently on other storage platforms, Hammerspace provides a way to “import” the metadata and hook a filesystem into a Hammerspace share, initially without moving the actual data. Objectives can be applied to the share into which the data is assimilated to move instances of the data onto Tier 0 or LSS nodes based on any metadata criteria such as time stamps, filename or extension, or path. Refer to the [Hammerspace Assimilation and Data Orchestration](#) article for more details on assimilation.

## Snapshots

Hammerspace has a snapshot technology as one level of data protection. With supported storage systems, Hammerspace snapshots can use offloaded file cloning. This currently supports DSX nodes and some third-party storage systems, it does not currently support XFS file cloning on Linux servers including Tier 0 and LSS. When offloaded cloning is not available, the snapshots revert to full file copy on change. For this reason, snapshots should be avoided during checkpoint creation.

## Monitoring

Hammerspace recommends Grafana as a monitoring tool, with Prometheus as the time series database service to collect metrics from Hammerspace and LSS / Tier 0 nodes. The Prometheus exporters are enabled on a Hammerspace cluster with the cluster-config command as follows:

### Hammerspace Admin Command Line

```
anvil1> cluster-config --prometheus-exporters-enable
ID:                                07106af1-4ae9-4509-8d3c-497614109567
Name:                              AE9Y4XC9E5X9J5
State:                              Standalone
Management IPs:                    [10.200.106.160/20]
Data IPs:                          [10.200.106.160/20]
Cluster floating IPs:              [10.200.106.160/20]
Since:                             2025-06-18 19:20:10 UTC
```



```
Timezone: UTC
GFS participant max suspected time: 30 minutes
Prometheus exporters: Enabled
Evaluation expiration: 2025-07-18 19:18:47 UTC
Online license activation support: true
NAS volume capacity: [Total: 136.6GB, Used: 1.2GB, Free: 135.3GB]
Share space (quota): [Total: 0B, Used: 0B, Free: 0B]
Metadata servers:
                    [Object type: Anvil, Node name: anvil1.hammer.space, Role:
Primary, Admin state: Up, Oper state: Up]
```

For more information, reference the following:

- Refer to the [Hammerspace Grafana Dashboards GitHub page](#) for details on downloading, installing, and configuring the dashboards and related components.
- Refer to [Monitoring Linux host metrics with the Node Exporter | Prometheus](#) for installation instructions for Prometheus exporters for Linux nodes.
- Refer to [DCGM Exporter – NVIDIA GPU Telemetry 1.0.0 documentation](#) for using Prometheus and Grafana with Nvidia GPU nodes.

## Troubleshooting

### Volume Issues - Volumes go “suspected”

**Causes:** Node gets rebuilt, NVME reformatted or remounted in a different path. The mount directory might still be exported, although the `mp` (mountpoint) export option should prevent the empty directory from being exported.

**What to check:**

- Storage node (LSS or Tier 0) is up
- NFS service is running
- Exports are correct
- Volumes or RAID arrays have not been reformatted or remounted
- Hammerspace comb structure (`/PrimaryData`) is intact

If you see multiple volumes from a single storage server go suspected, and the FSID has changed, if the new FSIDs are the same on multiple volumes on the same server, this most likely indicates that it is the root filesystem being exported due to empty directory mount points.

In this example, it is fairly obvious that the volume is not properly exported.

**Hammerspace Admin Command Line**

```
anvil1> volume-list
ID: 190ad494-c79a-42e8-bac5-38587a198b19
Name: AZ3:node3::/mnt/hsvol1
Internal ID: 20
Discovered addresses: [IP: 10.200.100.42/32, Port: 2049, NetId: tcp, NodeNum: 0]
Effective IPs: [IP: 10.200.100.42/32, Port: 2049, NetId: tcp, NodeNum: 0]
Path: /mnt/hsvol1
State: OK
Access type: Read Write
Node: AZ3:node3
Oper state: Suspected
Oper state reason: path /mnt/hsvol1 is not exported at 10.200.100.42
```

**Mobility Failures (Unsuccessful)**

File instances aren't where they should be (not enough instances, or not on the right volumes)  
GUI Mobility screen shows unsuccessful mobilities that don't resolve in short order.

**Common Causes**

- No DI (Data Instantiator or NFS Mover), or DI services are not running properly. The Hammerspace GUI and CLI alerts will indicate if there are no DIs. Check the status of the DIs.
- DIs don't have RW, no\_root\_squash on exports of the storage servers

**Mobility Troubleshooting**

If files don't appear to have the right number of instances in the right places, especially after changing objectives or creating files with more than 3 instances, check the mobility screen in the Hammerspace GUI. If you see any "Unsuccessful" mobilities, the most likely cause is one or more movers not having access to the exports on the source or destination storage node or volume. Check that all movers are listed in the exports of all volumes, and that they have no\_root\_squash.



# Maintenance

## Adding Tier 0 or LSS Nodes

When adding nodes, one key consideration is to update the exports on all other nodes that provide services so the new nodes can access shares and exports with the same options as existing nodes. To facilitate this, it is useful to maintain a master list of IPs and export options.

## Taking Nodes Offline for Maintenance

With availability zones, you can take one Tier 0 node or many offline at the same time, provided they are all within the same AZ and certain steps are taken to minimize impact. Depending on how long a node will be offline, meaning short term or long term or even permanently.

# Conclusion

Hammerspace Tier 0 turns GPU server local NVMe into a new tier of high-performance shared storage, managed and protected by Hammerspace. Tier 0 storage is up to 10x faster than networked storage - so you can reduce checkpointing time for AI training and HPC, and improve response times for inferencing and agentic AI. And it enables optimal use of assets you already own so you can reduce the need for external flash storage - to reduce costs, and gain back the power and rack space those systems would otherwise consume.

Tier 0 is a game-changing solution that immediately provides value for environments of all sizes, from a small cluster with only a few GPU servers, all the way to hyperscale environments with many thousands of GPUs.

## Additional Resources

- [Hammerspace Assimilation and Data Orchestration](#)
- [Hammerscript and Hammerspace Toolkit \(HSTK\)](#)
- [Hammerspace Objectives Guide](#)
- [Grafana Dashboards for Hammerspace](#)



# Appendix 1 - Installing Data Movers on Linux Nodes

There are two components that are collectively referred to as data movers:

- Mover, sometimes referred to as NFS mover, and internally referred to as Data Instantiator or DI. This mover moves or copies instances of files between local storage volumes via NFS. This could be between DSX nodes, to/from a Linux storage server or other third part storage system.
- Cloud Mover or CM. CM moves or copies files to and from object storage (local or remote) and cloud storage using S3, HTTPS and proprietary cloud protocols. CM also manages chunking, hash values, deduplication, and encryption for data placed in object storage.

There are three ways to implement the Hammerspace data mover technology:

- Deploy DSX nodes, which include data movers as part of their included functionality.
- Deploy the data mover RPMs on a supported Linux distribution.
- Deploy the data movers as containers

While it is supported to install DI and CM on GPU nodes, customers may choose to install on separate servers to avoid additional resource consumption on the CPUs of GPU nodes.

## Outline

### Supported Distributions / Platforms

DI as an RPM is available for el8 and el9 platforms.

Other x86 based Linux distributions must use container based deployments.

### Dependencies

- Firewall ports 9095 and 9096 open for inbound from all Anvil IPs.
- epel-release
- platform-python
- lttng-tools
- lttng-ust
- jemalloc
- Babeltrace

**Note:** Download and installation instructions for the above packages are below.





## Get the Hammerspace Data Services Components Kit

The DI RPM and container are available from the Hammerspace Support download server. You should ask your Hammerspace tech team for a download URL specific to the Hammerspace version you are running.

## Installation, Setup, and Configuration of the DI Using RPM

Install Extra Packages for Enterprise Linux and other dependencies:

### Linux Command Line

```
# dnf install -y epel-release
# dnf install -y platform-python lttng-tools lttng-ust jemalloc
```

Download and install babeltrace:

### Linux Command Line

```
# wget
https://dl.rockylinux.org/pub/rocky/9/devel/x86_64/os/Packages/b/babeltrace-1.5.8-1
0.el9.x86_64.rpm
# dnf install -y ./babeltrace-1.5.8-10.el9.x86_64.rpm
```

Download the Hammerspace Data Services Components tarball, then extract and install the DI. Note that if you are given a different version of the components, file names will be different.

### Linux Command Line

```
# wget
trans.doit.hammerspace.com/download/stg/5.1.28-340_hs_gxPb0/el9-components-5.1.28-3
40.tar.gz
# gunzip el9-components-5.1.28-340.tar.gz
# tar xvf el9-components-5.1.28-340.tar */pd-di-*
el9-components/pd-di-5.1.28-337.el9.x86_64.rpm
# dnf install -y el9-components/pd-di-5.1.28-337.el9.x86_64.rpm
```

Add the management IP address of the Hammerspace cluster to /etc/hosts. Be sure to use the correct management floating IP for your Hammerspace cluster.

**Linux Command Line**

```
# echo "10.1.2.3 data-cluster" >>/etc/hosts
```

Open up the firewall ports Anvil uses to talk to the DI:

**Linux Command Line**

```
# firewall-cmd --zone=public --add-port=9095/tcp --permanent; firewall-cmd
--zone=public --add-port=9096/tcp --permanent

# firewall-cmd --reload
```

Check firewall settings with `firewall-cmd --list-all`. You should see “ports: 9095/tcp 9096/tcp”.

**Linux Command Line**

```
# firewall-cmd --list-all
public (active)
  target: default
  icmp-block-inversion: no
  interfaces: ens192
  sources:
  services: cockpit dhcpv6-client mountd nfs rpc-bind ssh
  ports: 9095/tcp 9096/tcp
  protocols:
  forward: yes
  masquerade: no
  forward-ports:
  source-ports:
  icmp-blocks:
  rich rules:
```

Disable SELinux.

**Linux Command Line**

```
# vi /etc/selinux/config
# Set SELINUX to disabled and reboot.
SELINUX=disabled
```



Enable Linux Trace Toolkit: next generation (LTTNG) services:

#### Linux Command Line

```
# systemctl enable --now lttng-sessiond.service
# systemctl enable --now pd-di-lttng-recorder.service
```

Install Python PIP and other dependencies needed by the script to add the DI to the Hammerspace cluster:

#### Linux Command Line

```
# dnf install -y pip
# pip3 install urllib3
# pip3 install requests
```

Download and run the add\_node.py script to add the DI to the cluster:

#### Linux Command Line

```
# ./add_node.py -i 10.4.5.6 -l 20 -d -u admin -n di-6
```

In the above command, the following options are used:

- i IP address of the DI node being added
- l Netmask prefix length of the DI node being added
- d A DI is being added (as opposed to a cloud mover)
- u Administrative user name on the Hammerspace cluster
- n Node name of the DI as you wish it to appear in Hammerspace

**Note:** You can also specify admin user password with the -p option, but if not, the script will prompt.

Enable and start the DI service:

#### Linux Command Line

```
systemctl enable --now pd-di.service
```



You can verify the DI is properly configured and added to Hammerspace by logging in to the Hammerspace CLI as admin (or equivalent user) and running the node-list command specifying the DI node name.

#### Hammerspace Admin Command Line

```
admin@anvil> node-list --name my_di1
Name: my_di1
Type: External Mover
Internal ID: 1073744505
ID: bae33620-0281-11ef-8cab-5254009a4b8c
HW state: OK
Node state: MANAGED
Management IP: 10.4.5.6/20
S3 auth signing type: Default
Created: 2024-06-05 15:45:23 UTC
Modified: 2024-06-05 15:45:33 UTC
System services:
    Object type: Data Mover
    ID: 254e163c-1767-4b02-87b4-2f9c9c688b26
    Internal ID: 536873594
    Admin state: Up
    Oper state: Up
    Port: 9095
```

## Troubleshooting

The following node-list command shows the status of the DI from the Anvil perspective. In this example, there is a problem with the firewall on the node running the DI.

#### Hammerspace Admin Command Line

```
admin@anvil> node-list --name di-173
Name: di-173
Type: External Mover
Internal ID: 1073741834
ID: 594662a3-e045-507b-a609-0d2cb6f16beb
HW state: OK
Node state: MANAGED
Management IP: 10.200.100.173/64
S3 auth signing type: Default
Created: 2025-06-03 23:59:31 UTC
```



```
Modified:                2025-06-05 05:29:37 UTC
System services:
    Object type:          Data Mover
    ID:                   553b23ef-0709-567b-80f0-2be3fe68a852
    Internal ID:          536870923
    Admin state:          Up
    Oper state:           Down
    Oper state reason:    Failed port tests:
[10.200.100.173:9095(Connection refused), 10.200.100.173:9096(Connection refused)]
    Port:                 9095
```

The above error can also be caused by the pd-di service not installed, correctly configured, or started. You might also see:

```
Oper state reason:    Failed port tests: [10.200.101.3:9095(No route to host),
10.200.101.3:9096(No route to host)]
```

When you fix firewall and pd-di service issues, it may take a few seconds for Anvil to recognize the change.

## Appendix 2 - Use Tier 0 for Checkpointing

### Checkpointing Workflow Example

1. **Checkpoint Initiation:** The application triggers a checkpoint, at which point each node creates a file on the Hammerspace global shared storage system. Based on Service Level Objectives set up in Hammerspace, the storage node selected to back each file is the NVMe that is local to each node, thus making it possible to bypass the NFS protocol and networking stack as described above.
2. **Local Write Completion:** Once the write is complete, the GPU resumes computation without delay.
3. **Asynchronous Replication:** Hammerspace detects the new checkpoint file and begins replicating it to designated storage tiers based on policies.
4. **Data Availability:** The checkpoint is now safely stored in multiple locations, ensuring it can be used for recovery if needed.

### Tier 0 Checkpoint Analysis

This analysis quantifies the benefits of using Tier 0 storage. Specifically, leveraging local NVMe storage within compute nodes—for checkpointing, compared to traditional methods that rely on external networked storage systems, even those connected via high-speed 800Gb Ethernet (800GbE).

We will use the NVIDIA A100 GPU for our calculations, consider a 1,000-node cluster with 8,000 GPUs, 100 Petabytes of storage capacity, and 1 TB/sec of aggregate throughput.

### Estimating Checkpoint Data Size

System Configuration:

- **Compute Node:** NVIDIA DGX A100 or HGX system
- **GPUs per Node:** 8 NVIDIA A100 GPUs
- **GPU Memory per GPU:** 80 GB (also available in 40 GB variants)
- **Total GPU Memory per Node:** 8 GPUs × 80 GB = 640 GB
- **CPU Memory:** Assume 256 GB per node (can vary)
- **Total Memory to Checkpoint:** GPU Memory + Relevant CPU Memory

#### Checkpoint Data Size Estimate

600 GB



## Estimating Checkpoint Time When Writing to External Storage

Considerations when Writing Checkpoints to External Storage:

- **Network Overheads:** Protocol overheads, congestion, and latency reduce effective bandwidth.
- **Shared Infrastructure:** Multiple nodes checkpointing simultaneously can saturate the network and storage array.
- **Effective Bandwidth per Node:** Often significantly less than theoretical maximum. Let's conservatively estimate 1 GB/s per node.

### Time to write a 600 GB checkpoint to external storage

600 GB/1 GB per second = **600 seconds**

## Estimating Checkpoint Time When Writing to Tier 0 Local NVMe Storage

Local NVMe Storage Specifications:

- **NVMe Devices per Node:** 8 NVMe drives
- **NVMe Interface:** PCIe Gen5
- **Bandwidth per NVMe Device:** Approximately 14 GB/s (real-world write performance)
- **Total Aggregate Bandwidth:** 112 GB/s per node
- **Effective Bandwidth Assuming 90% Efficiency:** 100.8 GB/s per node

### Time to write a 600 GB checkpoint to Tier 0 storage

600 GB/100.8 GB per second = **5.95 seconds (6)**



## Glossary

Term/Component	Description
<b>Anvil Nodes</b>	A Hammerspace metadata server. Responsible for managing all the metadata for the Global File System, and hosts the API and Management GUI services. Manages data policies, objectives, and failover coordination.
<b>Cloud Mover (CM)</b>	A Hammerspace data service responsible for uploading and downloading data from cloud and object storage. This service can run on a Hammerspace DSX node or a supported Linux server.
<b>DI (Data Instantiator)</b>	Responsible for moving data from one NFS (Network File System) location to another NFS location, also known as an "NFS to NFS data mover". This service can run on a Hammerspace DSX node or a supported Linux server.
<b>Data Services Nodes (DSX)</b>	The Hammerspace data services node. Includes functional capabilities like Portal, Mover, Cloud Mover, and Store.
<b>HSTK</b>	Hammerspace Toolkit is a BASH command line tool that provides access to Hammerspace metadata and other functionality for clients.
<b>Instance</b>	A managed copy of a file placed on a Hammerspace volume. A file can have one or many instances for performance or protection reasons, but this is typically invisible to applications and users.
<b>Linux Storage Server (LSS)</b>	Servers running Linux that are specifically configured and optimized to serve as storage nodes within a larger system.
<b>Mover</b>	A stateless DSX service that moves files non-disruptively between file storage volumes.
<b>NFS (Network File System)</b>	A distributed file system protocol that allows users on a client computer to access files over a network as if they were stored locally.
<b>RDMA (Remote Direct Memory Access)</b>	Architecture where data can be transferred directly between computer memory and storage devices (or between two storage devices) without involving the CPU or operating system of the sending or receiving systems.





<b>SMB (Server Message Block)</b>	A network communication protocol for providing shared access to files, printers, and serial ports between nodes on a network.
<b>Tier 0 Nodes</b>	GPU compute nodes (e.g., NVIDIA DGX) with local NVMe used as high-performance storage.