

Aster Data *n*Cluster: Online Backup

Businesses have witnessed an exponential increase in data volumes in the last few years and continue to see them grow at a rapid pace. At the same time, today's organizations are adopting data-driven, decision-making strategies and analytics-intensive applications that place more stringent availability requirements on data warehouses. Backup technology used in traditional databases is unable to match the multi-terabyte and petabyte scale commonly seen in today's systems. Database backups, if done concurrently with queries, often don't complete in the dedicated "backup window" and cause significant performance slow-down. Similarly, recovery processes in traditional database systems have become unacceptably slow, given the large volumes of data that need to be recovered in case of a failure.

Not only do today's data warehousing and analytically-intensive applications need a scalable architecture for the database itself, they also need a backup and recovery architecture that can scale to meet data protection requirements.

Aster Data *n*Cluster delivers the first data-analytics server, a massively parallel (MPP) database with an integrated analytics engine. Aster Data *n*Cluster Online Backup is a part of "Always On" capabilities of *n*Cluster. It leverages *n*Cluster's "Always Parallel" architecture to provide powerful disk-based backup and recovery capabilities. Using massive parallelism, these capabilities can effectively handle the large data volumes and strict backup and recovery timelines that today's data-driven applications need.

The following aspects form the core of *n*Cluster Online Backup architecture:

- Massive parallelization for providing backup and restore speeds required at the petabyte scale
- Online operation so backup and restores don't require system shutdown
- Incremental backups so the time taken for every backup does not depend on the amount of historical data
- Compressed backups to reduce the cost of storage and improve network performance during data transfer

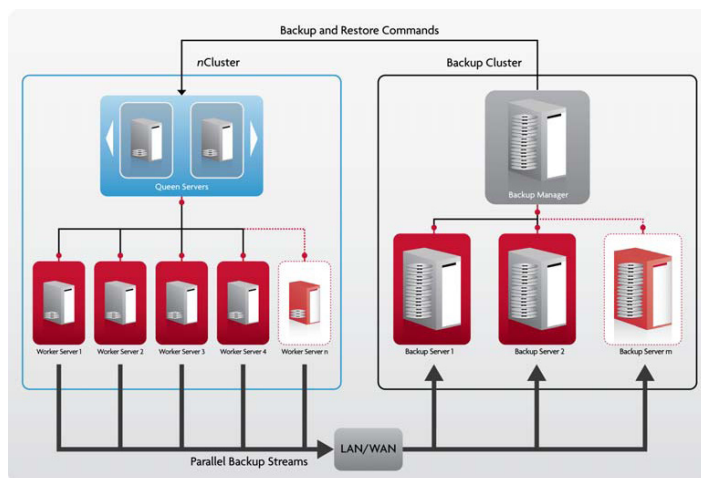


Figure 1: Aster Data *n*Cluster Online Backup

Quick Overview

Aster Data *n*Cluster leverages its "Always Parallel" architecture to deliver high-performance backup and restore capabilities required by terabytes- to petabytes-scale data management and data-driven applications. Using *n*Cluster Online Backup, data can be backed up without requiring a dedicated "backup window." *n*Cluster Online Backup includes powerful capabilities to support different backup policies of an IT department.

Highlights

- Massively parallel backup and recovery for high performance
- Online backup and recovery to eliminate downtime
- Backup storage on dedicated Backup servers or SAN/NAS
- Full and incremental backups
- Cluster-level and table-level backups
- Scheduled backups
- Compressed backups
- *n*Cluster Backup Terminal for managing backup and recovery tasks
- Low-cost, disk-based backups using high-density Backup servers

"With Aster Data's Online Backup we can take full and incremental back-ups while the system is running, data being loaded, queries being executed and the best part, it is also highly scalable for optimal backup performance."

Lenin Gali, Director of Business Intelligence
ShareThis



Backup Architecture

nCluster Online Backup is based on a unique architecture that leverages a cluster of dedicated Backup servers, called Backup Cluster. The Backup Cluster architecture can leverage massive parallelism of Backup servers to provide high-performance backup and recovery. For disaster recovery, Backup Cluster can reside in a location which is geographically separated from *nCluster*. Multiple *nCluster* databases can send data to a single Backup Cluster, providing an opportunity to consolidate backups and reduce costs. Data from *nCluster* can be backed up to two types of storage targets:

- **Backup Cluster** – The Backup servers themselves can store backups using their direct-attached storage, with disk mirroring for enhanced data protection. Storage-heavy servers can provide high-density storage at a low cost per terabyte, an important consideration for data volumes associated with analytically-intensive applications and data warehousing. Backup Cluster is built on incremental scalability principles similar to the *nCluster* architecture—more servers can be incrementally provisioned to add backup capacity. Thus, backup storage costs can be managed in a granular manner as data volumes grow. From Backup servers, backup files can be moved to tapes or virtual tape libraries (VTLs), if required by an organization's IT policies.
- **Network Storage** – *nCluster* Online Backup also provides the flexibility to store backup on network storage (SAN/NAS), if required by an organization's IT policies. If this option is chosen, Backup servers can run backup and restore processes in a massively parallel manner while using the network storage to store backup data.

Backup and Recovery Process

To manage backup and recovery tasks, one of the servers in Backup Cluster is designated as Backup Manager. Backup Manager runs *nCluster* Backup Terminal software (and also stores data, like other Backup servers), which provides a functionally rich interface for backup and restore tasks. To start a backup, Backup Manager first contacts the Queen server in *nCluster*, which, in turn, sends the backup request to the Worker servers. Next, Worker servers directly connect to the Backup servers and stream data in a massively parallel manner. Backup Manager and the Queen server do not get involved in the actual data transfer so they do not become performance bottlenecks. It should be noted that the number of Backup servers is not determined by the number of Worker servers, as the data transfer is done in a many-to-many fashion.

Restoration of backup data can also be initiated using *nCluster* Backup Terminal. During restoration, data is streamed in a massively parallel manner directly from the Backup servers to the Worker servers, providing faster time to recovery. For recovery, both online and offline modes are supported, depending on the parameters selected at the time of backup. Using table-level online recovery, data can be recovered while the system is available for business.

nCluster Online Backup also provides disaster recovery capabilities—full backups can be restored to a secondary site running *nCluster* software. Disaster recovery in *nCluster* is fast, as large volumes of data can be quickly transferred by leveraging massive parallelization and backup compression.

Types of Backup

nCluster provides the flexibility to manage backups according to an IT organization's administrative policies. Backups can be made at the level of the *nCluster* massively parallel database or at the level of individual tables. *nCluster* Online Backup supports both full and incremental backups of data. For table-level backup, full backup is currently supported.

A full backup includes all data, metadata, log files, etc. Incremental backup includes only the changes since last backup, minimizing data transferred over the network as well as backup storage requirements.

Scheduled Backups

nCluster Online Backup includes the ability to run scheduled backups according to pre-defined backup policies, minimizing manual intervention for routine backups. Administrators can schedule backups in any of the following three ways:

- Using the scheduler provided with *nCluster* backup utility
- Using an OS scheduler (e.g. CRON jobs)
- Using third-party schedulers

Backup Compression

nCluster Backup Compression provides powerful compression capabilities. Data can be compressed by Worker servers before it is sent to Backup servers. This minimizes backup storage requirements as well as data transfer over the network, speeding up both backup and recovery processes. Compression ratios in the range of 3X to 12X are typical, depending on the nature of data being backed up.

Backup and Recovery Performance

The unique architecture of *nCluster* Online Backup enables very high levels of performance. As data is streamed from a large number of Workers, tens, hundreds, or even thousands of processor cores can send data to Backup Cluster in parallel. On the receiving end, Backup servers can process data in parallel, maximizing throughput. Similarly, the recovery process also benefits from massive parallelism, providing high performance and shorter time to recovery. Compression of backups improves network performance during backup as well as recovery processes, providing further speed-up. Such high-performance architecture of *nCluster* Online Backup lets organizations protect data efficiently, even at the petabyte scale.

Backup Management

nCluster Online Backup provides a rich set of backup and restore functionality through *nCluster* Backup Terminal. Administrators can manage the following tasks using *nCluster* Backup Terminal:

- Adding and removing Backup servers from Backup Cluster for capacity and performance management
- Starting, pausing, resuming, and cancelling backup and restore jobs
- Monitoring the status of ongoing backup and restore jobs
- Scheduling backup jobs
- Deleting backups
- Viewing the history of backups and restores
- Monitoring total and used backup storage capacity (at Backup servers or SAN/NAS, as the case may be)

Reduce the Cost of Backup

For today's big data management and data-driven applications, the cost of managing backups can be high. *nCluster* Online Backup helps organizations reduce the cost of backup in the following ways:

- **Use of inexpensive commodity servers for backup storage** – *nCluster* Backup Cluster can use inexpensive commodity servers with direct attached storage to minimize the cost of backups. Organizations can select cost-effective server configurations according to their backup requirements. For example, storage-heavy commodity servers with 48 SATA disks (500 GB to 1 TB per disk) can provide low cost per terabyte for storing backups. In addition, Backup Cluster can be incrementally scaled-out to increase storage capacity at low costs, according to the rate of data growth. This helps organizations manage storage costs in small incremental steps.
- **Backup Compression** – As mentioned above, *nCluster* Backup Compression can provide significant reduction in backup storage requirements and also help improve backup and restore performance by minimizing network traffic.

About Aster Data

Aster Data is a proven leader in big data management and big data analysis for data-driven applications. Aster Data's *nCluster* is the first MPP data warehouse architecture that allows applications to be fully embedded within the database engine to enable ultra-fast, deep analysis of massive data sets. Aster Data's unique "applications-within" approach allows application logic to exist and execute with the data itself. Termed a "Data-Analytics Server", Aster Data's solution effectively utilizes Aster Data's patent-pending SQL-MapReduce together with parallelized data processing and applications to address the big data challenge. Companies using Aster Data include Coremetrics, MySpace, comScore, Akamai, Full Tilt Poker, and ShareThis. Aster Data is headquartered in San Carlos, California and is backed by Sequoia Capital, JAFCO Ventures, IVP, and Cambrian Ventures, as well as industry visionaries including David Cheriton, Ron Conway, and Rajeev Motwani. For more information please visit www.asterdata.com, or call 1.888.Aster.Data.