# cloudera®

# MULTITENANCY AND THE ENTERPRISE DATA HUB:

Best practices and architectures for multiple services and user communities in Apache Hadoop
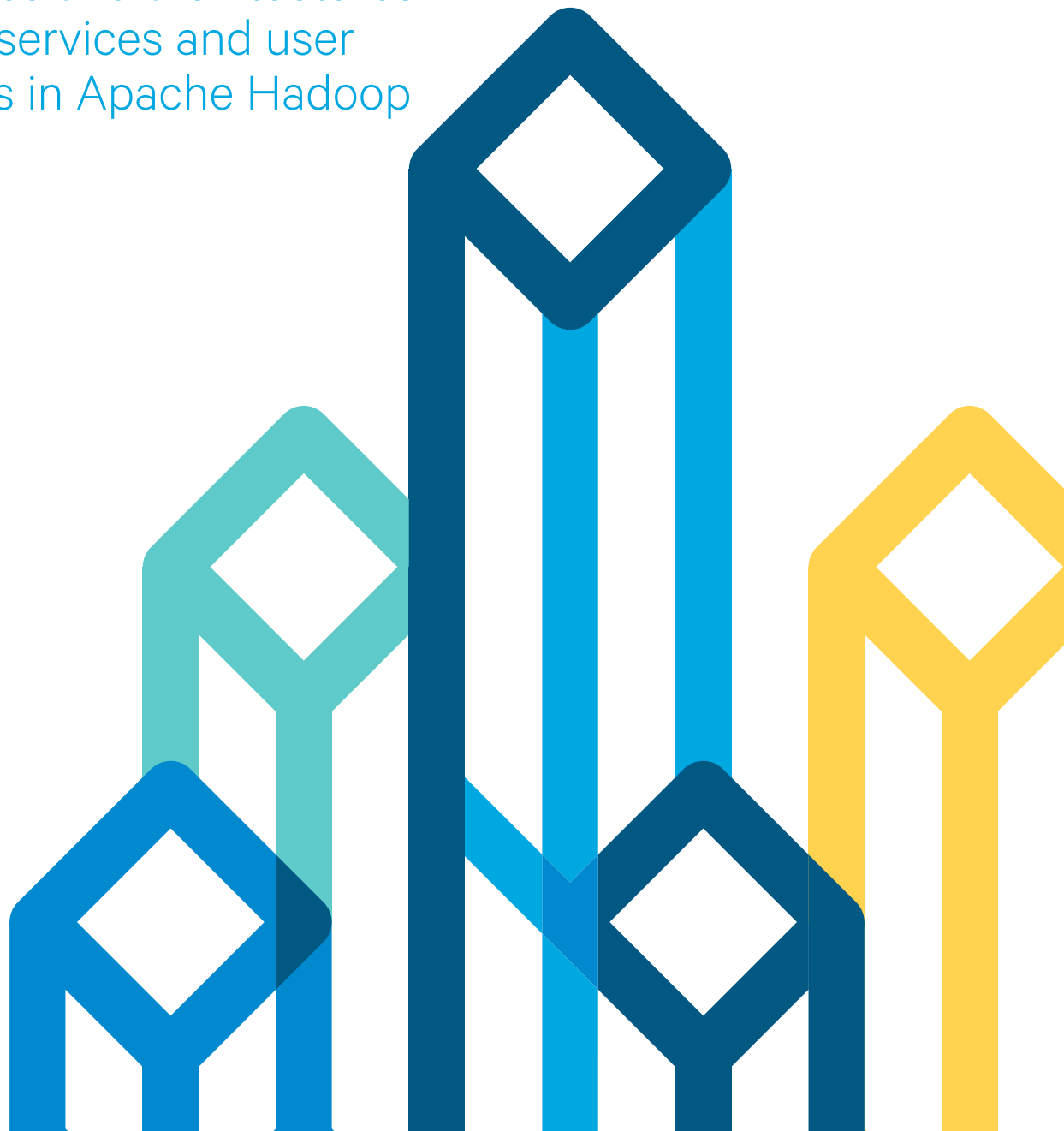
# cloudera®

**Table of Content**

# cloudera®

Multitenancy in an enterprise data hub lets a business share the collective resources of an Apache Hadoop cluster across its user communities without impacting business SLAs and capabilities or security and privacy needs.

## Introduction

Information-driven enterprises increasingly turn to organization-wide information repositories to provide the necessary breadth and depth of data and analysis to answer the business queries of today and tomorrow. A enterprise data hub (EDH), powered by Apache Hadoop, is this next step in the evolution of data management. With an enterprise data hub, organizations get a single central system to store and work with all data of all types and formats, with the flexibility to run a variety of enterprise workloads - including batch processing, analytic SQL, stream processing, enterprise search, and machine learning. These capabilities are built to integrate with existing systems and come with the security, governance, data protection, and management required by modern enterprise data management infrastructures. Yet IT administrators must consider multitenancy for their user communities as the enterprise data hub emerges as the necessary center in this data infrastructure and as enterprise users shift to the EDH as their primary engine for business.

Multitenancy generally refers to a set of features that enable multiple business users and processes to share a common set of resources, such as an Apache Hadoop cluster. via policy rather than physical separation, yet without negatively impacting service level agreements (SLA), violating security requirements, or even revealing the existence of each party. This paper discusses the strategies for developing an appropriate multitenant architecture as well as the features and capabilities Cloudera offers to enable multitenant administration for an enterprise data hub.

## Business Objectives for Multitenant Environments

Realizing an EDH offers enterprises a platform and building blocks for a number of broad solutions that address specific problems or improve efficiencies within the business. For example, an EDH is a cornerstone to self-service business intelligence (BI) applications due to the flexibility of its storage substrate to handle all types of data, unstructured to structured, in volumes typically unreachable with conventional systems while providing access to this trove of data for existing and familiar BI tools via ODBC connections which utilize the analytic SQL capabilities inherent to the EDH.

Such an application can have broad adoption within an organization, and thus the range and degree of data sources and volumes as well as the variations in SQL usage and application requires balancing several key factors when developing and operating a multitenant environment. Therefore, enterprise IT teams must consider several strategic business objectives in their architectural choices and cluster design for an EDH.

### Data Sharing

Businesses often turn to an enterprise data hub to build a single repository that manages data from multiple users and departments. By creating such a data store, an organization can promote effective sharing of relevant data sets without having to manage the duplication of data into private or standalone repositories. Enterprises can employ this approach as a way to mitigate "shadow IT" projects that stem from frustrations with lengthy and slow operational overhead and data processing that typically exist with siloed and separated data sets and applications.

MULTITENANCY AND THE
ENTERPRISE DATA HUB:
**Best practices and architectures for multiple services and user communities in Apache Hadoop**

WHITE PAPER

3

## Consolidated Operations

Enterprises employ a multitenant environments as a way to amortize administrative overhead across different groups by consolidating operations, which results in better and consistent experiences at a lower cost. The proliferation of data and sources that characterize big data exacerbate the complexity of managing distinct and separate systems for the range of computing capabilities required by business, from analytic SQL to search to machine learning.
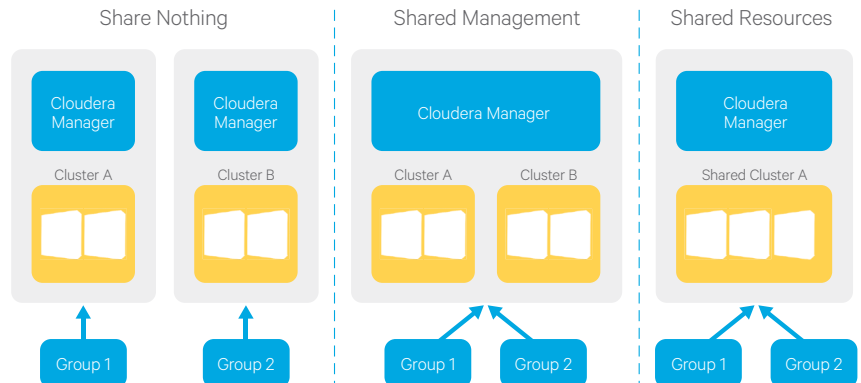
## Increased Performance and Utilization

Enterprises can multiplex of the resources of larger EDH installations across a broad range of users and processes. This approach can achieve better utilization of cluster resources and better performance for individual requests through dynamic resource allocation and optimized, just-in-time execution.

Using these objectives along with specific business use cases and the needs of the serviced applications, IT teams can consider different models and approaches of multitenancy that are more appropriate for their organization.

## Standard Isolation Models of an EDH

Before diving into the elements of a multitenant environment, organizations should first recognize the fundamental EDH architectures and decision framework for choosing the right approach and foundation for their particular goals and operations.

There are three key architectural models that articulate the various approaches to strategic cluster design.
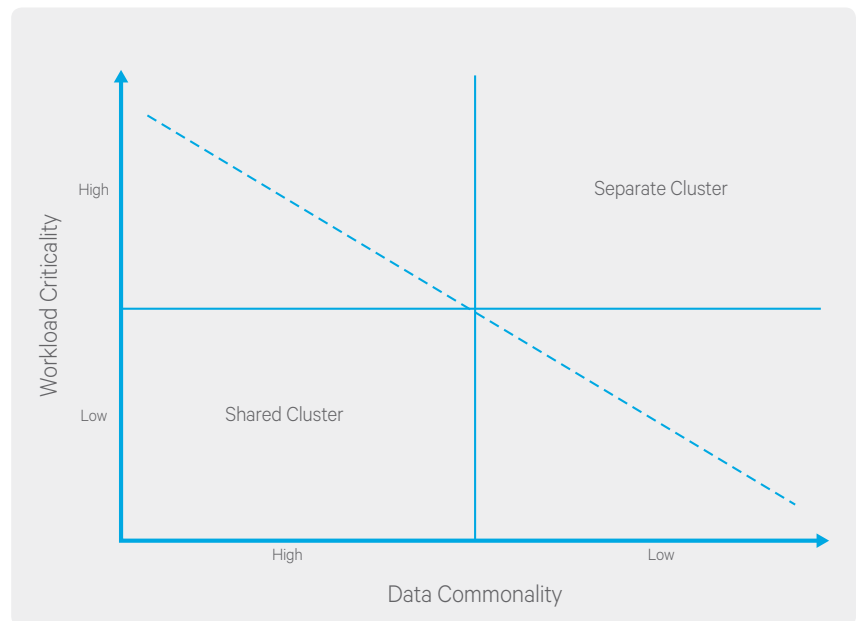


## Share Nothing

In a *share nothing* architecture, both management and data are completely separated and nothing is shared between clusters. This architecture does not provide the benefits of multitenancy - shared data sets and innovation, streamlined operations, and aggregate performance — but IT teams may find it appropriate based on specific operational realities and governance policies. For example, geographically disparate data centers or specific company contracts or policies, such as the legal concerns surrounding a business conglomerate, may force an enterprise IT team to employ this model. Another common example is security and data privacy mandates, like those found within the European Union, that prevent centralized administrative access for different groups and data sets and restrict the transfer of these data sets across geographical boundaries.

## Shared Management

A *shared management* model offers organizations the benefits of reduced operational overhead yet without actual cluster resource and data sharing between groups. This approach is a middle ground, granting some of the benefits of multitenancy while maintaining isolation at the cluster level. Generally speaking, this is the preferred choice for environments where full multitenancy is not appropriate. For example, enterprises commonly employ this model for purpose-built, high-priority clusters that cannot risk any performance issues or resource contention due to their sensitivity to latency, such as an ad serving system or a retail personalization, "next offer" engine. While a multitenant EDH always faces the risks of misconfigured applications or malicious users and a well-designed EDH has the tools and frameworks to address these situations, this model strongly mitigates these risks at the cost of data sharing and resource pooling.

## Shared Resources

The shared resource model embraces full multitenancy with all the accretive benefits, from consolidated management to shared data and resources, and represents the desired end state for many EDH operators. For example, a biotechnology firm can harness the entire corpus and insight of research, trial data, and individual perspectives from all of its research teams and other departments by employing a full multitenant EDH to house and analyze its information, greatly accelerating innovation and improving sharing through transparency and accessibility.



A high level framework for evaluating the appropriate architecture for an EDH environment.

## The Balance of Criticality and Commonality

As mentioned, enterprise IT teams often discover that multitenancy is not necessarily a good fit for their mission-critical workloads running uniquely tailored data sets. On the other hand, multitenancy can be extremely useful for less critical workloads that employ shared data sets by reducing unnecessary or burdensome data duplication and synchronization. For most situations where business units share data, but the specific workloads are critical to the organization, the overarching business priorities and SLA goals drive the choice of a multitenant or isolated architecture. For some, the risk of latency and resource contention weighs heavily on their performance goals, and thus would suggest a shared management model, while others consider data visibility as paramount, such as for fraud detection and insider threat analysis, and would lean towards a shared resources model.

MULTITENANCY AND THE
ENTERPRISE DATA HUB:
**Best practices and architectures
for multiple services and user
communities in Apache Hadoop**

WHITE PAPER

5

# cloudera®

Apache Sentry (incubating) is an independent security module that integrates with SQL engines like Apache Hive and Cloudera Impala and with Cloudera Search. Sentry delivers advanced authorization controls to enable multi-user applications and cross-functional processes for enterprise data sets and provides enterprise IT and business teams with:

> Role-based administration and delegation

> Data classification for commingled data sets

> Improved regulatory compliance

> Expanded user reach and data ROI

> Visual policy management through Hue

Learn more about Sentry at http://www.cloudera.com/content/cloudera/en/products-and-services/cdh/sentry.html

MULTITENANCY AND THE ENTERPRISE DATA HUB: Best practices and architectures for multiple services and user communities in Apache Hadoop

WHITE PAPER

6

## Elements of a Multitenant Cluster Architecture

When the business data and workload mix accommodates a multitenant environment, IT teams need to address three critical facets when enabling a Hadoop-powered EDH as the platform of the system: security and governance, resource isolation and management, and chargeback and showback capabilities.

### Security

Security for Hadoop is clearly critical in both single tenant and multitenant environments as it establishes the foundation for trusted data and usage among the various actors of the business environment. Without such trust, enterprises cannot rely on the resources and information when making business-critical decisions, which in turn undermines the benefits of operational consolidation and the decreased friction of shared data and insights. Cloudera's EDH provides a rich set of tools and frameworks for security. Key elements of this strategy and its systems include:

> **Authentication**, which proves users are who they say they are;

> **Authorization**, which determines what users are allowed to see and do; and

> **Auditing**, which determines who did what, and when.

> **Data Protection**, which encrypts data-at-rest and in-motion

Cloudera's EDH offers additional tools like network connectivity management as well as data masking. For further information on how IT teams can enable enterprise-grade security measures and policies for multitenant clusters, please refer to the Cloudera white paper, *Securing Your Enterprise Hadoop Ecosystem*[1].

In the context of multitenant administration, security requirements should also include the following capabilities.

**Security Management Delegation**

A central IT team tends to become the bottleneck in granting permissions to individuals and teams to specific data sets when handling large numbers of data sources with different access policies. Organizations can use Apache Sentry (incubating), the open source role-based access control (RBAC) system for Hadoop, to delegate permissions management for given data sets. Using this approach, local data administrators are responsible for assigning access for those data sets to the appropriate individuals and teams.

**Auditor Access**

For most large multitenant clusters, audit teams typically need access to data audit trails - for example, to monitor usage patterns for irregular behavior like spikes in request access to credit card details or Social Security data and other sensitive information - yet without having full access to the cluster and its resources and data; enterprise IT teams often to adhere to the best practice of "least privilege" and restrict operational access to the minimal data and activity set required. For these cases, Cloudera Navigator provides a data auditor role that partitions the management rights to the cluster so that administrators can grant the audit team access only to the informational data needed and thus mitigate the impact to operations and security. This approach also answers the common request of audit teams to simplify and focus the interaction model of their applications.

**Data Visibility Policies**

Data visibility, in particular for the cluster administrator, is another security requirement that is prominent in most multitenant environments, especially those under strict compliance policies or regulations. Typical security approaches encrypt data, both on-disk and in-use, such that only users with the correct access can view data and even administrators without proper access cannot view data stored on Hadoop. Cloudera Navigator provides data encryption and enterprise-grade key management with encrypt and key trustee, out-of-the-box.

> "Fair scheduling is a method of assigning resources to jobs such that all jobs get, on average, an equal share of resources over time... Unlike the default Hadoop scheduler, which forms a queue of jobs, [the Fair Scheduler] lets short jobs finish in reasonable time while not starving long jobs. It is also an easy way to share a cluster between multiple of users. Fair sharing can also work with job priorities - the priorities are used as weights to determine the fraction of total compute time that each job gets."

Learn more about the Fair Scheduler at http://hadoop.apache.org/docs/r1.1.2/fair_scheduler.html

MULTITENANCY AND THE ENTERPRISE DATA HUB:
Best practices and architectures for multiple services and user communities in Apache Hadoop
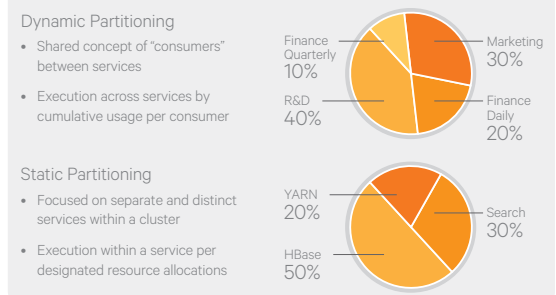
WHITE PAPER

7

## Resource Isolation & Management

IT administrators must manage another crucial aspect of running a multitenant environment: facilitating the fair and equitable usage of finite cluster resources across disparate users and workloads. Typically, this is a result of aggregating resources to drive improved performance and utilization - a key business driver for multitenancy - and multiple groups within the organization will finance the operations of this resource pool to meet this goal. As an outcome of many of these financing models, EDH administrators require systems to grant proportional access to this pool based on the proportion of payment. In addition, a successful multitenant environment employs these tools and frameworks to let users meet SLAs for critical workloads even in the presence of unpredictable usage stemming from multiple, simultaneous workloads and ill-constructed or misconfigured processes.

### Resource Management

The utilitarian batch processing engine of Hadoop, MapReduce, provides a scheduler framework that administrators can configure to ensure multiple, simultaneous user jobs share physical resources. More specifically, many production environments have

**Dynamic Partitioning**
- Shared concept of "consumers" between services
- Execution across services by cumulative usage per consumer

Finance Quarterly 10%
R&D 40%
Marketing 30%
Finance Daily 20%

**Static Partitioning**
- Focused on separate and distinct services within a cluster
- Execution within a service per designated resource allocations

YARN 20%
HBase 50%
Search 30%

successful implementations of the Fair Scheduler, and these environments provide maximal utilization while enforcing SLAs through assigned resource minimums.

### Dynamic Partitioning

With the advent of computing capabilities such as Cloudera Impala and Cloudera Search and a growing ecosystem of partner applications built to take advantage of CDH and the elements of Cloudera's EDH, cluster faculties and data are increasingly shared with systems that do not operate within the original management framework of MapReduce. Thus, a resource management solution must take the full range of these systems into account. To address this challenge, Cloudera's EDH and the underlying Hadoop platform, CDH, ship with the YARN resource management framework. In this context, YARN becomes a building block for computing engines to coordinate consumption and usage reservations to ensure resources are fairly allocated and used. This approach is sometimes referred to as *dynamic partitioning*. Currently, Impala, MapReduce, and other well-designed YARN applications participate in dynamic partitioning in CDH.

IT administrators should also consider, with respect to the scheduler capability, how best to regulate user access to specified allocations (i.e. 'pools') of resources. For example, IT teams may wish to balance allocations between the processing needs for their marketing team's near real-time campaign dashboards and their finance department's SLA-driven quarterly compliance and reporting jobs. These administrative needs also extend to multiple applications within a single group. For example, the finance team must balance their quarterly reporting efforts with the daily expense report summaries. To achieve these goals, Hadoop and YARN support Access Control Lists for the various resource schedulers, thus ensuring that a user (or application) or group of users may only access a given resource pool.

### Static Partitioning

While dynamic partitioning with YARN offers the IT administrator immense flexibility from a resource management perspective, IT teams do operate applications that are not built on the YARN framework or require hard boundaries for resource allocation in order to separate them fully from other services in the cluster. Typically, these applications are purpose-built and by design do not permit this degree of resource flexibility.

To satisfy such cases, Cloudera's EDH, through Cloudera Manager, supports a static partitioning model, which leverages a technology available on modern Linux operating systems called *container groups (cgroups)*. In this model, IT administrators specify policies within the host operating system to restrict a particular service or application to a given allocation of cluster resources. For instance, the IT administrator may choose to partition a cluster by limiting an Apache HBase service to a maximum of 50% of the cluster resources and allotting the remaining 50% to a YARN service and its associated dynamic partitioning in order to accommodate the business SLAs and workloads handled by each of these services.

### Quota Management

While resource management systems ensure appropriate access and minimum resource amounts for running applications, IT administrators must also carefully govern cluster resources in terms of disk usage in a multitenant environment. Like resource management, disk management is a balance of business objectives and requirements across a range of user communities.

The underlying storage substrate of a Hadoop-based EDH, the Hadoop Distributed File System (HDFS), supports two quotas mechanisms that administrators can tune to manage space usage by cluster tenants.

**Disk Space Quotas:** Administrators can set disk space limits on a per-directory basis. This quota prevents users from accidentally or maliciously consuming too much disk space within the cluster, which can impact the operations of HDFS, similara to other file systems.

**Name Quotas:** Name quotas are similar to disk quotas and administrators can employ them to limit the number of files or subdirectories within a particular directory. This quota helps IT administrators optimize the metadata subsystem (NameNode) within the Hadoop cluster.

Name quotas and disk space quotas are the primary tools for administrators to make sure that tenants have access only to an appropriate portion of disk and name resources, much like the resource allocations mentioned in the previous section, and cannot adversely affect the operations of other tenants through their misuse of the shared file system.

### Monitoring and Alerting

Resource and quota management controls are critical to smooth cluster operations, yet even with these tools and systems, administrators have to plan for unforeseen situations such as an errant job or process that overwhelms an allotted resource partition for a single group and requires investigation and possible response, such as a poorly constructed SQL query within a self-service BI application.

Cloudera Manager provides Hadoop administrators a rich set of reporting and alerting tools that can be used to identify dangerous situations like low disk space conditions; once identified, Cloudera Manager can generate and send alerts to a network operations center (NOC) dashboard or an on-call resource via pager for immediate response.

## Chargeback and Showback

Another common requirement for multitenant environments is the ability to meter the cluster usage of different tenants. As mentioned, one of the key business drivers of multitenancy is the aggregation of resources to improve utilization and performance, and the multiple participants will build internal budgets to finance this resource pool. In many organizations, IT departments use the metered information to drive showback or chargeback models and illustrate compliance.

Cloudera Manager offers IT teams historical and trending disk and CPU usage, and with this information —which can be exported in common formats, such as Microsoft Excel, to financial modeling applications—can provide a strong foundation for an internal chargeback model. These metering capabilities can also facilitate capacity planning and accurate budgeting for growth of the shared platform, thus ensuring that IT teams allocate sufficient resources in line with cluster demand.

## The Hidden Potential of Multitenancy

While this paper has examined the elements necessary to protect and isolate tenants in a shared enterprise data hub, enterprises business and IT teams should note that multitenancy in Hadoop provides rich opportunities to discover new value about data and usage by elements and people within the organization. These opportunities can lead to many significant and interesting insights, such as:

Awareness of expired or unused data, so IT teams can retire data sets that are no longer accessed by groups in the organization due to age and thus provide little to no value and, more commonly, remove data that has exceeded a contractual expiration date.

Insight into the key parties that use specific data sets across an organization and identify shared access patterns that may provide opportunities for further cross-group collaboration within an organization.

Better data quality from shared ownership of specific data sets by mitigating the need to replicate data sets and the effects of the inevitable divergence and resulting errors between replicas.

Improved accuracy of organizational metadata attributes such as "most frequent accessor" and "data owner," thus providing clear identity of the correct points of contact in the organization for provenance and usage.

These are just some of the potential benefits to having a single, well managed, enterprise data hub. Clearly, multitenancy represents not only mutual protection and consolidated operations, but shared insights and improved collaboration for data and business usage.

## Conclusion

Hadoop has quickly evolved from a single-workload, data processing platform into the foundation of a comprehensive information repository serving a wide variety of user communities and applications, data sets and business processes. This evolution in information management — the next generation of data management — is the enterprise data hub. IT teams and administrators demand that Hadoop, as a cornerstone to the enterprise data hub and the multiple audiences, data, and capabilities enabled therein, offer the controls necessary for a multitenant platform: security and data protection, scheduling and resource isolation, and quota management and usage tracking.

Cloudera is making it easier for enterprises to deploy, manage, and monitor multitenancy Hadoop environments and to realize the full potential of an enterprise data hub. Multitenancy in Hadoop has matured rapidly, with developments such as Apache Sentry for role-based authorization; Cloudera Manager for a single, comprehensive interface to control, configure, and monitor all aspects of a cluster, from resource management to chargeback details; and Cloudera Navigator for audit, data lineage, encryption and key management, and data discovery capabilities. With Cloudera, enterprises can more easily reap the benefits that come from a multitenant enterprise data hub.

# cloudera®

## About Cloudera

Cloudera is revolutionizing enterprise data management by offering the first unified Platform for Big Data, an enterprise data hub built on Apache® Hadoop™. Cloudera offers enterprises one place to store, process and analyze all their data, empowering them to extend the value of existing investments while enabling fundamental new ways to derive value from their data. Only Cloudera offers everything needed on a journey to an enterprise data hub, including software for business critical data challenges such as storage, access, management, analysis, security and search. As the leading educator of Hadoop professionals, Cloudera has trained over 40,000 individuals worldwide. Over 800 partners and a seasoned professional services team help deliver greater time to value. Finally, only Cloudera provides proactive and predictive support to run an enterprise data hub with confidence. Leading organizations in every industry plus top public sector organizations globally run Cloudera in production. www.cloudera.com.

---

**cloudera.com**

1-888-789-1488 or 1-650-362-0488
Cloudera, Inc. 1001 Page Mill Road, Palo Alto, CA 94304, USA