



WHITEPAPER

# An Inside Look into the Talena Architecture

# INTRODUCTION

Big Data has matured beyond the confines of the research lab, and now serves as the foundation for high-value, business-oriented applications. Yet while most enterprises have happily embraced specialized Big Data technologies for gathering, organizing, and extracting meaning from this information, far too many continue to rely on haphazard and potentially risky methods for safeguarding these assets and dispensing them to software developers and testers.

Talena's always-on Big Data management platform closes this gap, delivering affordable, flexible, enterprise grade information protection and distribution to bolster these new mission-critical applications. This paper describes the challenges of fully utilizing yet defending Big Data resources, explains why an updated approach was needed, and supplies technical details about the Talena solution.

The intended audience includes:

- Database administrators
- Architects
- IT operations
- Software developers
- Quality assurance professionals
- DevOps teams

# BIG DATA PROTECTION AND DISTRIBUTION DEMANDS A BETTER STRATEGY

Transactional software solutions such as Customer Relationship Management (CRM) and Enterprise Resource Planning (ERP) used to be the primary sources of business information that was typically hosted in relational databases. This is no longer the case: a diverse and ever-expanding assortment of new technologies is spawning enormous amounts of raw data that are driving entirely new types of applications.

In an effort to keep pace with the volume, variety, and velocity of all this supplemental information – collectively labeled as Big Data – businesses are rolling out a diversified collection of applications built on scale-out file systems and databases such as:

- Hadoop-related ecosystem
  - » Hadoop File System (HDFS)
  - » Hive
  - » HBase
  - » Impala
  - » Tez
  - » Spark
- Massively Parallel Processing (MPP) data warehouses
  - » HP Vertica
  - » IBM Netezza
- NoSQL databases
  - » Cassandra
  - » Couchbase
  - » HBase
  - » MongoDB

Until fairly recently, most businesses primarily employed these Big Data platforms for research or proof-of-concept projects. This meant that secure information access and robust data protection were relatively minor considerations. However, everything changes once an organization begins to merge these assets into their core application portfolio: it's no longer acceptable to risk a system outage, data loss, or security breach.

Although the majority of organizations have longstanding tools and techniques for protecting conventional information as well as sharing it with software developers and testers, there are significant challenges where Big Data is concerned.

## Data protection issues

Conventional backup and recovery techniques - and supporting products - remain ideal for their original purposes, but are a mismatch for the Big Data technologies listed earlier.

Many of these new platforms supply built-in replication, which entails distributing multiple data copies onto distributed servers. While this approach reduces the possibility of outright data loss, it paradoxically serves to propagate user or application-driven data corruption: damage quickly spreads to all replicated copies. This means that it's also essential to implement a proper backup/recovery strategy for Big Data.

While the majority of these new platforms ship with their own dedicated backup and restore utilities, these are laden with substantial drawbacks:

- They're myopically focused on a single technology, and thus don't recognize the heterogeneous Big Data portfolio that prevails in many enterprises
- They're driven by command line interfaces (CLI), which demands manual interaction, making both backup and restore processes cumbersome and error-prone
- Automation – which is a fundamental characteristic of dependable data protection processes – demands extensive scripting
- They're restrictive in terms of their capabilities, thus minimizing their effectiveness as companies scale their data infrastructure.

Scripting is a particularly inefficient way to protect data: it siphons off valuable IT talent from mainline business responsibilities, and produces a continually growing inventory of brittle, maintenance-intensive assets.

## Data distribution challenges

To produce effective applications, software developers and testers need prompt access to meaningful amounts of representative production information from Big Data repositories. Faced with these requirements, enterprises have generally resorted to one of two approaches: create fabricated test data, or write scripts to transfer data from production to development and testing environments.

Organizations that rely on artificially constructed data confront the risk of delivering solutions that are divorced from reality, and are prone to software quality issues once they are placed into production and encounter authentic information.

For those enterprises that do permit information extraction from production Big Data systems in support of new software development and testing, the task of writing extraction scripts hampers and delays the entire new application creation process. In fact, according to a comprehensive study commissioned by Talena, 90% of organizations defer application rollouts waiting for data. What's more worrisome – and possibly exposes the business to legal consequences – is the potential for inadvertently divulging highly sensitive Personally Identifiable Information (PII) to unauthorized users.

Given these information protection and allocation deficiencies, it's clear that a fresh, comprehensive approach was necessary to support the Big Data platforms that are now core enterprise resources. Founded and led by seasoned technology industry veterans, Talena has delivered the first solution tailor-made for Big Data:

- It's designed to manage petabytes of Big Data economically
- It's built on a scale-out architecture and leverages commodity servers or virtual machines (VMs)
- It provides a number of enterprise-grade capabilities that naturally fit into your data center or cloud environment

Talena provides the first “always-on” big data management software solution to help companies protect valuable data assets and iterate rapidly on their business-critical applications through capabilities such as backup & recovery, test/dev management, and archiving.

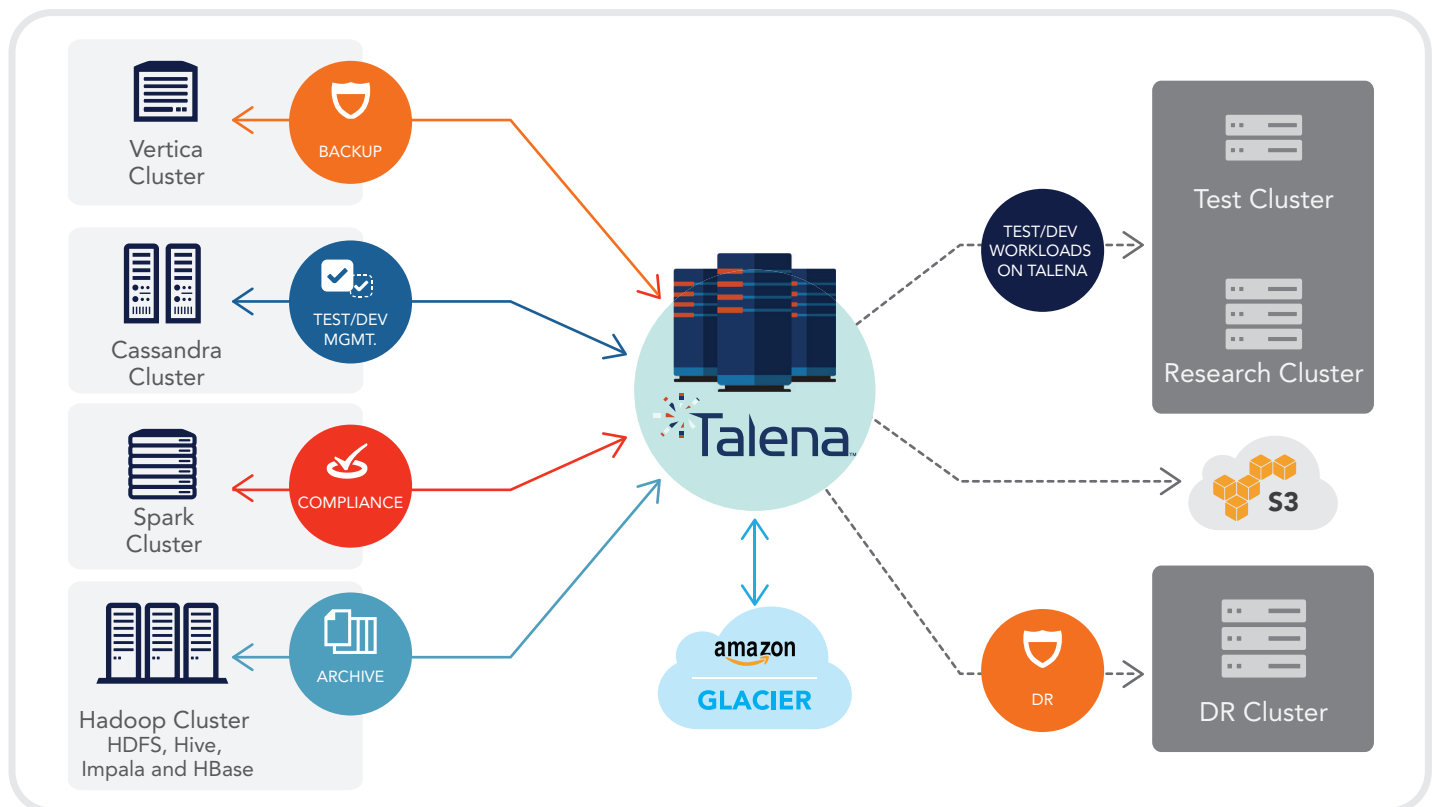
**The next section describes Talena's technical architecture and merits.**

# INTRODUCING THE TALENA ALWAYS-ON DATA MANAGEMENT PLATFORM

Recognizing the inadequacy of existing methods for protecting and distributing Big Data to software development teams, Talena set out to devise a highly focused solution by following three core philosophies:

1. Design for enterprise Big Data
2. Adapt to dynamic Big Data environments
3. Deliver a low total cost of ownership

The following three sections describe how the Talena architecture attained these objectives.



# 1. Design for enterprise Big Data

For many years, the process of backing up enterprise data has been comprised of a full weekly backup, followed by daily incremental backups. This time-tested approach falls flat in Big Data environments that are commonly measured in petabytes: it's simply too hardware-intensive, laborious, time-consuming, and error prone to conduct weekly full backups.

Furthermore, data restoration using this method meant applying the last full backup, and then restoring all the incremental backups to attain the user-specified restore point. This is a very slow procedure that's ill suited for Big Data.

## Incremental-only backups

Talena's incremental-only backup strategy is much more appropriate for Big Data technologies. It uses snapshots to ascertain what's changed since the previous incremental backup, and then only backs up relevant information. Incremental data is immediately materialized onto Talena servers, and a restore point is created. This fully materialized restore point – driven by the incremental backups – can be used to return data to the production clusters as-is, without requiring any additional work.

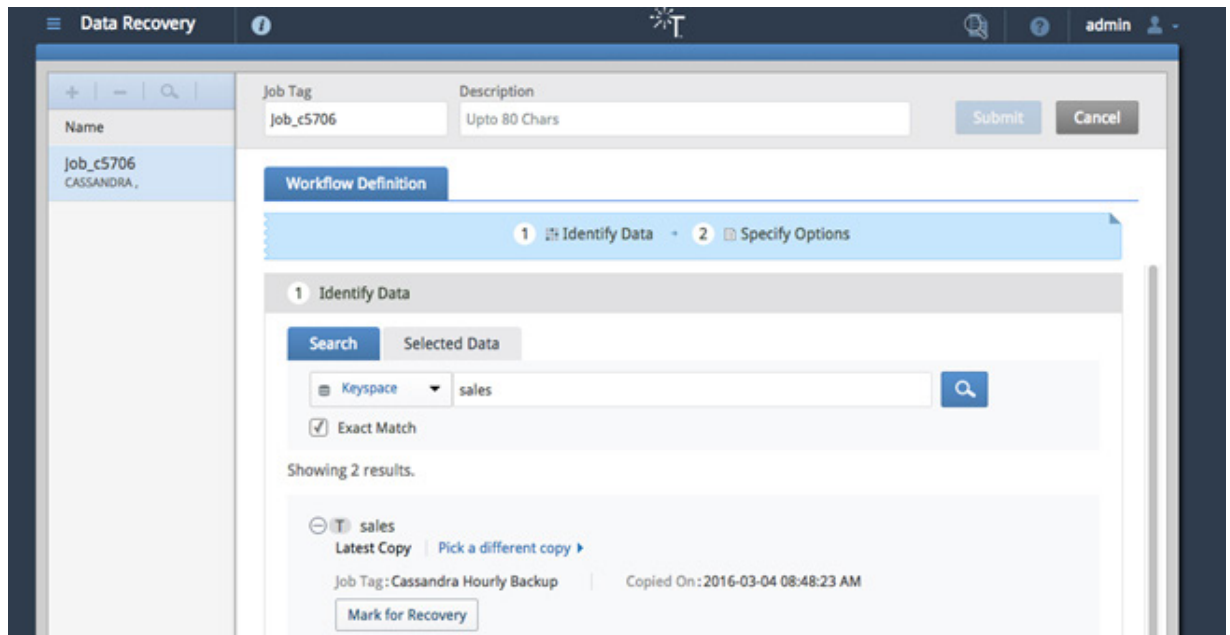
Talena's Big Data-aware backup and restore architecture thrives when applied to the specialized platforms listed earlier. Rather than creating a "one size fits all" product, Talena invested considerable time and resources to create an "application-aware" solution that not only backs up Big Data but also incorporates metadata such as file/directory attributes and table/database schemas. This essential information is retained along with the restore point, and upon a user's request is then restored back to the production cluster with the data itself.

Table 1 itemizes the objects that Talena assesses when performing a backup:

| TECHNOLOGY  | OBJECTS EVALUATED BY TALENA                                  |
|-------------|--|
| Hadoop/HDFS | Files and directories – including permissions and attributes |
| Hive        | Databases, tables, partitions, and metastore                 |
| Impala      | Databases, tables, partitions, and metastore                 |
| Cassandra   | Keyspaces and tables   |
| HP Vertica  | Databases, schemas, tables, and catalog                      |

## Flexible data restoration

Talena's recoveries are fast, granular, and capable of restoring data to any previous point in time no matter how large the backup data set may be. For example, although a Cassandra restore point may contain hundreds of tables and keyspaces, Talena's FastFind technology lets users search for a specific keyspace and table, as well as mutations for that specific restore point. Similarly, FastFind can restore a single Hive partition out of the hundreds of partitions that this type of table may have.



Since Talena uses an incremental-forever technique to instantly materialize backed-up data, recovery can take place very quickly: all that's necessary to perform a complete restore is the most recent checkpoint, since it already reflects all previous incremental backups. In addition, the flexible Talena data recovery capabilities extend to restoring data to a cluster with a different topology.

## Making data available to software developers and testers

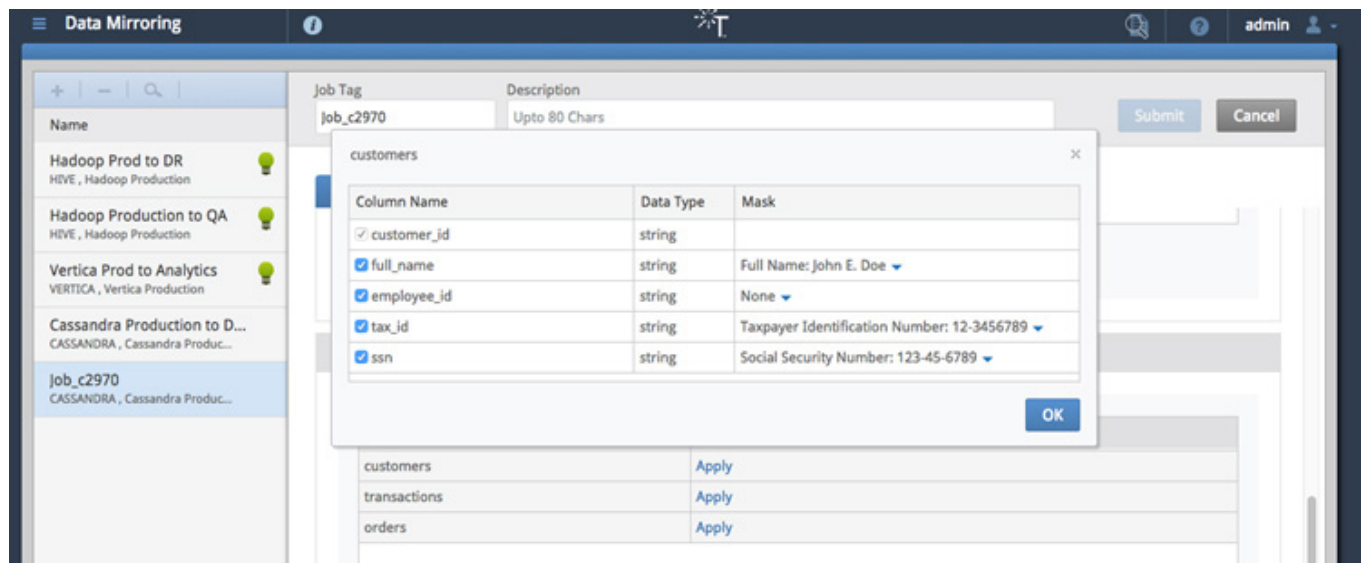
Protecting information is only part of Talena's mission: as described earlier, software developers and testers need prompt and regular access to the production data that's stored in specialized Big Data platforms. This is a critical requirement to ensure that the new breed of applications that they're building will match user expectations. However, those enterprises that grant carte blanche access to this information confront a very serious risk of divulging Personally Identifiable Information (PII).

To help eliminate this unpleasant possibility, Talena uses its deep platform-specific knowledge to offer sophisticated data masking that lets organizations freely dispense Big Data yet protect PII. Talena understands the data formats of each platform, so it's able to combine



the details from the information source's schema and raw files to appropriately mask the relevant data. Currently supported masks include:

- Full name
- Social security number
- Taxpayer identification number
- Employee identification number
- Credit card number



Talena's data masks are always consistent. For example, a given social security number will always cause the same mask to be generated, and this coherence will occur across all tables containing columns with social security numbers. This preserves the statistical properties of the original data, and makes accurate analysis possible. Since making a full copy of production data for development and testing isn't always necessary or practical, Talena lets users specify a meaningful sample size. It then uses this parameter to automatically drive the task of creating representative data in the target environment.

To further protect sensitive data, Talena's patent-pending algorithm is completely one-way: there is no technique - even for Talena - to reverse-engineer masked information. In addition, its data masking is stateless so there are no intermediate files or encryption technologies to maintain and protect.

Finally, regardless of the exact data distribution configuration, Talena offers self-service capabilities that let authorized users set up scheduled or ad-hoc production data extracts into development or testing sandboxes.

## Robust security

Talena implements industrial-strength access, security, and encryption capabilities, including:

- Kerberos and LDAP authentication
- User and role management
- Support for encrypted file systems
- SSL for data transmission between Talena and primary clusters

## 2. Adapt to dynamic Big Data environments

It's highly unusual to encounter a static, homogenous Big Data implementation. Instead, most organizations deploy a unique blend of tools and technologies that were meticulously selected to meet their particular needs. Furthermore, these configurations are constantly changing by incorporating additional Big Data platforms and adding, adjusting, or dropping data nodes. Talena was engineered to thrive in every type of landscape, with an intelligent set of features that appreciates how Big Data is actually utilized in the enterprise.

### Heterogeneous Big Data platform support

Each Big Data product offer highly targeted capabilities that enterprises earmark for achieving specific goals. For example, the use case for a Hadoop installation is very different than for a MongoDB instance. This means that it's likely that a business will acquire an assortment of these platforms to thoroughly address their Big Data needs.

Consequently, any product that's intended to strengthen an organization's collective Big Data implementation must address all of the deployed technologies. Talena attains this objective by blending a loosely-coupled core architecture, deep insight into each supported product, and platform-specific data movers that completely utilize each of the vendor-supplied APIs. This design principle also makes it straightforward for Talena to integrate new Big Data solutions based on customer requirements.

## Topology independent backup and restore

Backups performed using traditional solutions or Big Data vendor-supplied utilities are intolerant of topology alterations: in other words, any changes to the exact blend of servers and storage media that was present during the backup process can make a subsequent restore difficult - if not impossible - to carry out. Since the vast majority of Big Data environments are extremely fluid, this introduces a very tangible peril of permanent data loss because previous backups are likely to become useless.

Talena follows a much more Big Data-friendly strategy. First, it dynamically acquires the topology of the destination cluster at the time of the restore. It then uses this structure to reshard the backed-up data to match the topology of its destination. This is a much smarter, more powerful mechanism for restoring information: it has no requirement for an exact match between the backup and restore topologies, nor does it require Talena to preserve the production cluster's topology at the time of backup.

In addition, customers may elect to recover tables and keyspaces to the original source or to an alternate cluster. The destination cluster size is also independent and may be adjusted during the restore process.

## Performance and scalability

Since scale-out – which leverages the power of commodity servers to increase capacity – is a fundamental Big Data architectural philosophy, it's natural that the same must hold true for the technologies that are tasked with protecting and distributing this information.

Each software-defined, storage-agnostic Talena node is deployed on cost-effective commodity hardware, using directly attached disk drives to maintain backed-up information. This eliminates any reliance on more expensive Network Attached Storage (NAS) or Storage Area Network (SAN) components, although customers are welcome to make use of these devices if desired. It also frees Talena from mandating a specific hardware configuration.

All nodes in the Talena cluster establish connections to one or more nodes in the primary cluster, so that data can be transferred concurrently. Adding more nodes in Talena increases the data ingestion rate from the production cluster, which is a completely parallelized and highly scalable technique. Horizontal scalability is also a major attribute of Talena's catalog, which is a cornerstone of data restoration. The catalog is capable of tracking and versioning millions of objects, and offers full search capabilities.

### 3. Deliver a Low Total Cost of Ownership (TCO)

From the beginning, Talena elevated TCO to be a major influence on its architecture – a critical approach for the open source-driven Big Data environments that it would serve. This resulted in a series of design decisions that delivered a robust yet cost-effective solution to market.

#### Software-defined platform

Talena's paramount decision to avoid any hardware mandates – such as proprietary, expensive appliances – resulted in numerous money-saving benefits:

- It can be installed on bare metal servers or virtual machines
- It's available on popular enterprise Linux distributions:
  - » RHEL
  - » Ubuntu
  - » Oracle Linux
  - » Centos
- It utilizes inexpensive, commodity disk drives
- It's capable of running in the cloud, internal data center, or any combination
- It's able to be instantiated with a small number of servers, and automatically rebalances its clusters as more servers are added
- It provides administrators with a "single pane of glass" that lets them support multiple data sources and use cases, as well as manage all aspects of their Talena environment
- It integrates with Nagios agents that transmit alerts regarding disk space, backup failures, and other issues, along with daily email notifications about job status

Finally, many customers elect to deploy Talena using the same operating system as their production data nodes. This homogenizes their operating system-level software infrastructure, helping to reduce maintenance efforts and potential security issues.

#### Agentless architecture

Customarily, enterprise backup solutions have mandated installing software agents on all data nodes. These components are then tasked with the job of transferring information. This tactic simply won't work in ever-changing Big Data environments: administrators would quickly be overwhelmed with maintenance and monitoring responsibilities, and these far-flung agents would also introduce inherent security risks.

Instead of forcing customers to install agents throughout their production landscape, Talena's agentless architecture uses the already-optimized public interfaces supplied by the Big Data platform vendors as its cross-system communication pipeline.

### **Storage efficiency**

Given how large Big Data environments can get, it's no surprise that storage expenditures can quickly outpace expectations, especially when including the outlays that enable sufficient information backups. To help keep this overhead to a minimum, Talena employs a smoothly integrated pipeline of storage reduction techniques.

Talena's global, data-aware, variable-length deduplication engine is the initial component in this workflow. It begins the storage reduction process by identifying the data that is to be deduplicated. This information may be stored in one of many formats, such as compressed files (e.g. GZ, Snappy, LZO, and so on) or application-specific structures (e.g. RCFile, Parquet and ORC for Hive and Impala or SSTable for Cassandra). Once deduplicated and compressed, the output is saved onto the file system and erasure coded to increase its durability.

The upshot is that Talena offers significant storage savings yet is ready to quickly restore information to its original state.

# SUMMARY & NEXT STEPS

As Big Data continues its march towards becoming a mainline IT resource, protecting this critical information and making it available for the builders of new applications will require fresh technologies and tactics. The Talena Big Data Management Platform manifests a set of intelligent design decisions that help enterprises deal with the three most well agreed-upon traits of Big Data environments:



**Volume:** Talena's incremental-forever backups are capable of gracefully handling even the biggest data quantities



**Variety:** Talena supports the most popular Big Data platforms, NoSQL databases, and data warehouses



**Velocity:** As workloads grow, all that's required to achieve linear scalability is to add more Talena nodes

To set up a Talena demonstration visit [www.talena-inc.com](http://www.talena-inc.com)

## About Talena

Talena provides the first data availability management software expressly engineered for Big Data technologies, enabling always-on data across the full lifecycle of applications so they can meet internal and external service level agreements. Our data availability management software is uniquely designed to optimize the test/dev management, backup, recovery, and archive functions so that our customers can iterate rapidly on their applications, prevent data loss, and minimize compliance risk. Based in Silicon Valley, we are backed by Canaan Partners, Intel Capital, ONSET Ventures, and Wipro Ventures.

For more information, visit [www.talena-inc.com](http://www.talena-inc.com).



**Please contact us for more information at [info@talena-inc.com](mailto:info@talena-inc.com) or visit us at [www.talena-inc.com](http://www.talena-inc.com).**



Talena, Inc. | 2860 Zanker Road,, Suite 109, San Jose CA 95134 | +1.408.649.6338 | [talena-inc.com](http://talena-inc.com)

© 2016 Talena, Inc. All rights reserved. Talena and the Talena logo are trademarks of Talena in the US and in other countries. Information subject to change without notice. All other trademarks and service marks are property of their respective owners.