

# 4 ESSENTIAL STEPS FOR MANAGING SENSITIVE DATA



# SPEAKERS



**Balaji Ganesan**

CEO, Privacera



**Srikanth Venkat**

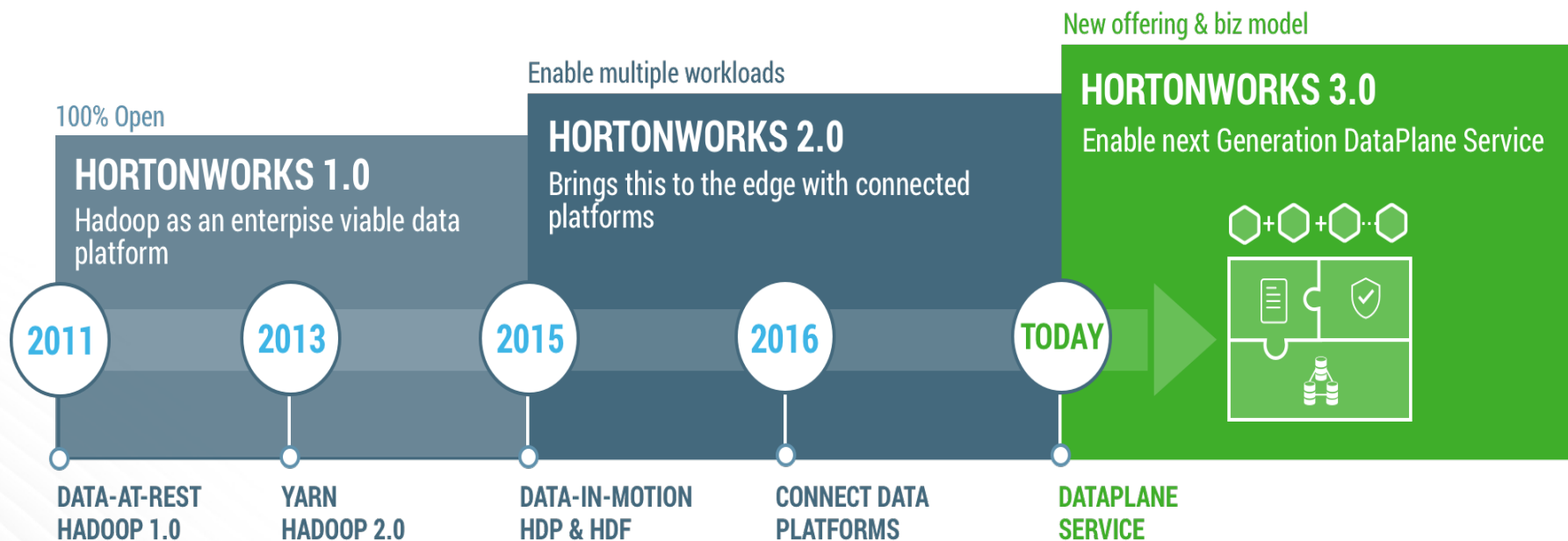
Senior Director, Product  
Management, Hortonworks

# AGENDA

- ▶ Hortonworks Introduction
  - ▶ Security & Governance with Hortonworks
  - ▶ Sensitive Data Management Challenges
  - ▶ Hortonworks DataPlane Service
  - ▶ Demo (Data Steward Studio)
- ▶ Privacera Introduction
  - ▶ 4 steps in managing sensitive data
  - ▶ Representative scenarios & solutions
  - ▶ Demo (Privacera)
- ▶ Wrap up

# About Hortonworks:

Enabling the Modern Data Architecture through consistent and continuous innovation



# Apache Ranger

## Authorization

- Centralized platform to define, administer and manage security policies consistently across Hadoop components
  - HDFS, Hive, HBase, YARN, Kafka, Solr, Storm, Knox, NiFi, Atlas
- Extensible Architecture
  - Custom policy conditions, user context enrichers
  - Easy to add new component types for authorization

## Ranger KMS

- Store and manage encryption keys
- Support HDFS Transparent Data Encryption
- Integration with HSM
  - Safenet LUNA

## Auditing

- Central audit location for all access requests
- Support multiple destination sources (HDFS, Solr, etc.)
- Real-time visual query interface

# Dynamic Row Filtering & Column Masking: Apache Ranger with Apache Hive

**User 1: Joe**  
Location : US  
Group: Analyst

Users from US Analyst group see data for US persons with CC and National ID (SSN) as masked values and MRN is nullified

EU HR Policy Admins can see unmasked but are restricted by row filtering policies to see data for EU persons only

**User 2: Ivanna**  
Location : EU  
Group: HR

Country	National ID	CC No	MR N	Name
US	xxxxx3233	4539 xxxx xxxx xxx	null	John Doe
US	xxxxx7465	5391 xxxx xxxx xxx	null	Jane Doe



**Ranger Policy Enforcement**  
Query Rewritten based on Dynamic Ranger Policies: Filter rows by region & apply relevant column masking

Country	National ID	Name	MRN
Germany	T22000129	Ernie Schwarz	876452830A

**Original Query:**  
SELECT country, nationalid,  
name, mrm FROM  
ww\_customers

**Original Query:**  
SELECT country, nationalid,  
ccnumber, mrm, name FROM  
ww\_customers

Country	National ID	CC No	DOB	MRN	Name	Policy ID
US	232323233	4539067047629850	9/12/1969	8233054331	John Doe	nj23j424
US	333287465	5391304868205600	8/13/1979	3736885376	Jane Doe	cadsd984
Germany	T22000129	4532786256545550	3/5/1963	876452830A	Ernie Schwarz	KK-2345909



# Apache Atlas : Open Metadata & Governance



## Vision:

*Metadata-driven foundational governance services for enterprise data ecosystem*

- Open frameworks and APIs
- Agile and secure collaboration around data and advanced analytics
- Reduce operational costs while extracting economic value of data

(lineage), owner,

provide enterprise

your data catalog out

metadata repositories

new data is created or

and classify the data

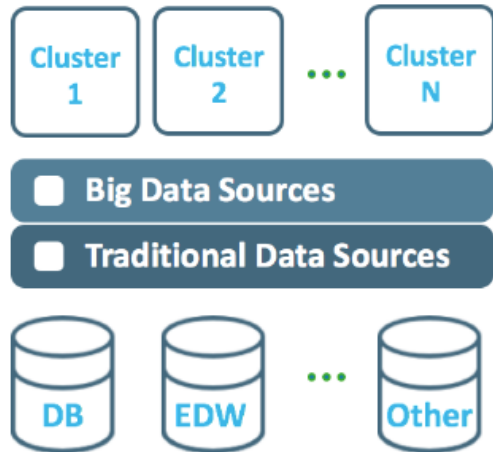
quickly and efficiently,

ge to help others

omatically

Predefined standards for glossaries, data schemas, rules and regulations

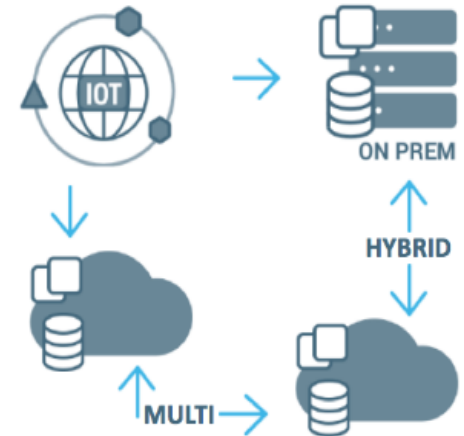
# Next Generation Data Problems



**Data Is Spread Across Multiple Clusters and Data Sources**



**Store & Analyze Data From ERP/CRM, Systems, IoT/ Mobile Devices, Social Media, Geo Location etc.**



**Some data is on-premise, rest in the cloud.**

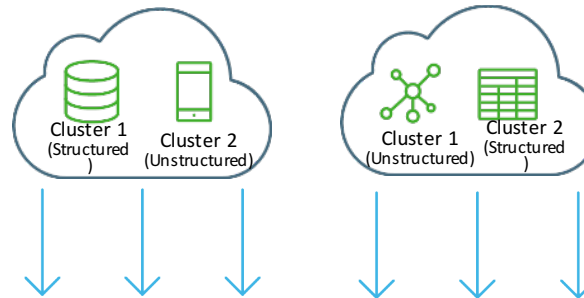
**Moving data from cloud to on-premise & vice versa**

**Moving data between different clouds**



# What If...

In the Cloud

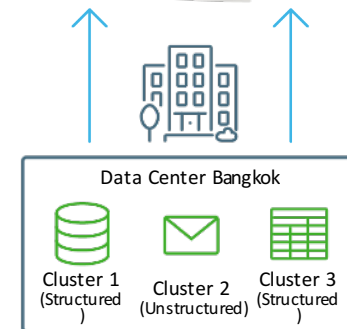
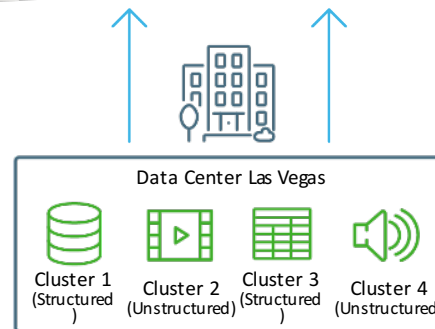
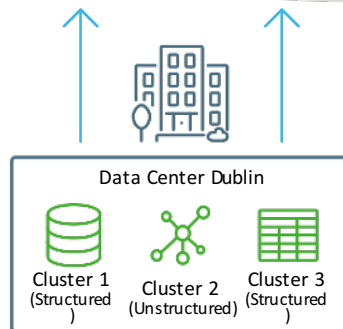


Aware of  
Data Sources

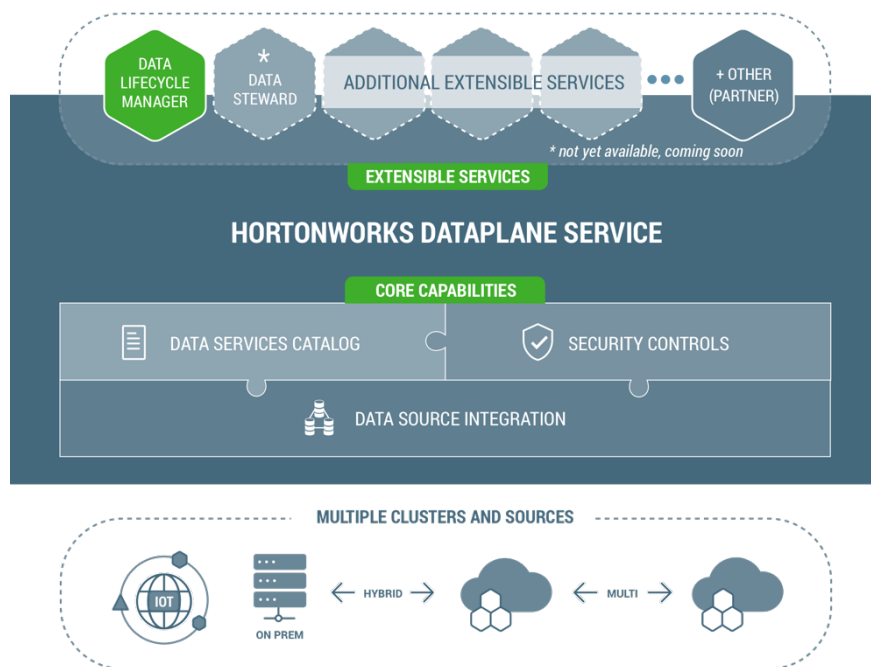
Enable  
New Services

Unified  
Security &  
Governance  
Model

On Premises



# What is Hortonworks DataPlane Service?

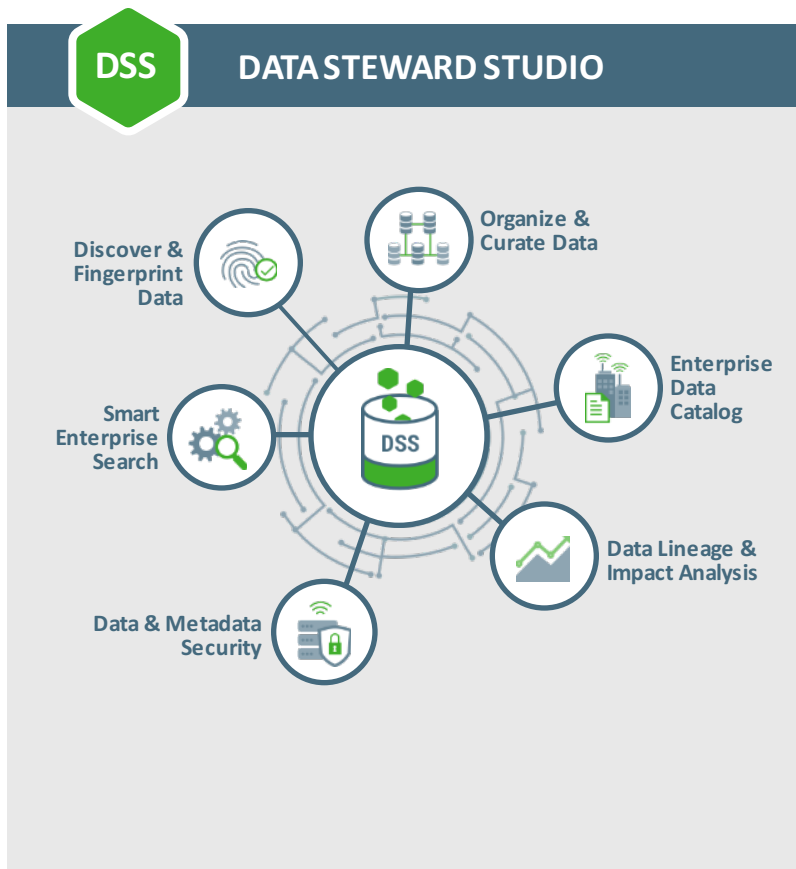


## Hortonworks DataPlane Service

a platform with extensible data management services for:

- ❑ Addressing compliance and regulatory requirements for enterprise
- ❑ Providing consistent security & governance across data landscape
- ❑ Enabling centralized management of data assets
- ❑ Responsible data sharing and collaboration

# Hortonworks DataPlane Service: Extensible Services



## Data Steward Studio (DSS)

Suite of capabilities that allows users to understand, secure, and govern data across enterprise data lakes

Ensure consistent security and governance for data assets **across tiers**

- Curate, discover and organize data assets based on business classifications, purpose, protections, relevance, etc.
- Govern proper usage and lineage of data assets to identify schema, classification and view lineage/data supply chain
- Understand and audit data asset security and use for anomaly detection, forensic audit/compliance & proper control mechanisms

**...all across multiple types and tiers of data**

*Technical Preview Available*



## Data Steward / Asset Details

Datalake  
cl1

Database  
hortoniabank

DETAILS

LINEAGE

POLICY

AUDIT

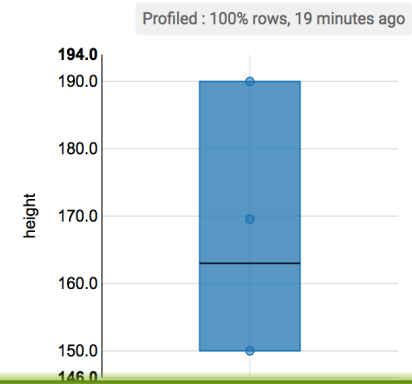
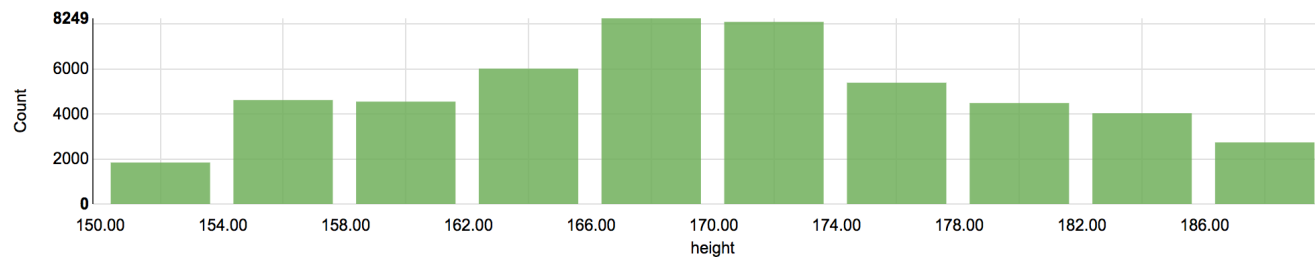
PROPERTIES

TAGS

SCHEMA

Next Profiler Schedule in : 41 minutes

	Name	Type	Unique Values *	Null Values	Max	Min	Mean	Comment
	streetaddress	string	49790	0				
	middleinitial	string	26	0				
	emailaddress	string	49970	0				
	age	int	67	0	85	19	52.17764	
	height	int	41	0	190	150	169.56316	



## Data Steward Studio (DSS)

**CONSUMABILITY:** Understand shape of Hive column data with statistical profiler, example: Profile shows box plot and histogram for distribution of column values

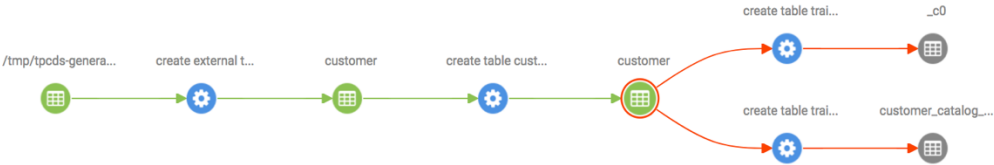


CUSTOMER

HIVE

Datalake prod_east	Database tpcds_bin_partitioned_orc_2	# of Rows 144000
-----------------------	---	---------------------

DETAILS   LINEAGE   POLICY   AUDIT



customer ✕

owner: centos

temporary: false

lastAccessTime: 29 Sep 2017 06:54:35 PM

aliases:

qualifiedName: tpcds\_text\_2.customer...

description:

viewExpandedTex

tableType: EXTERNAL\_TABLE

createTime: 29 Sep 2017 06:54:35 PM

name: customer

comment:

partitionKeys:

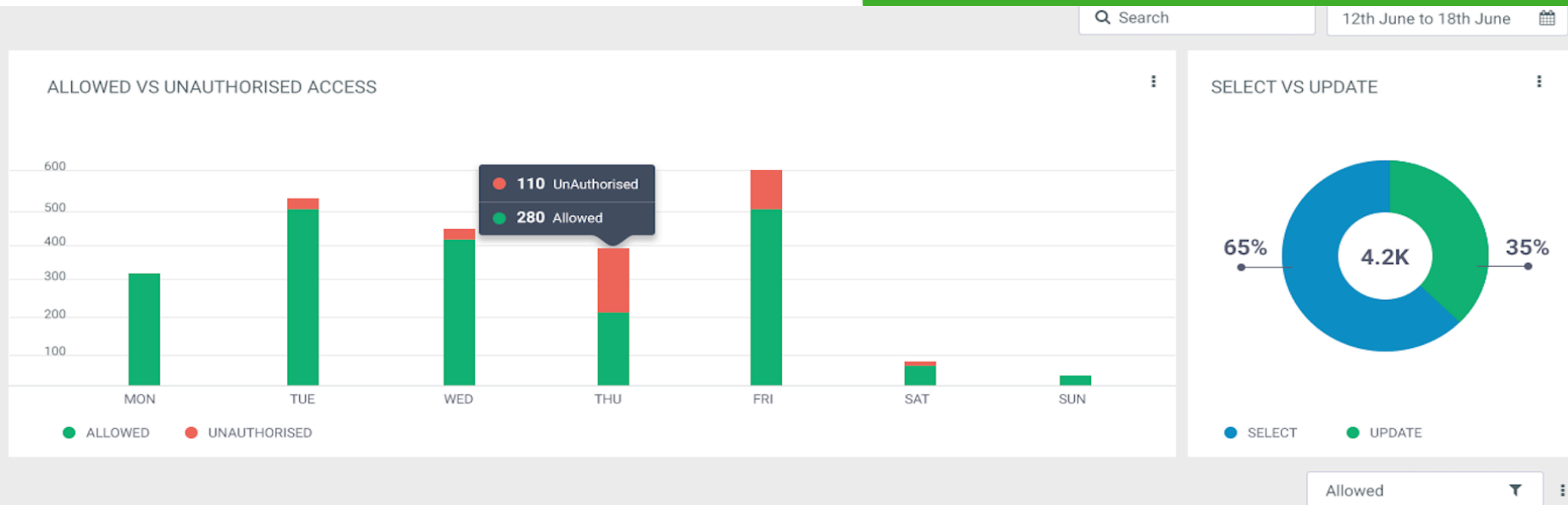
profileData: typeName: hive\_table\_pr...

db: tpcds\_text\_2

**CONSUMABILITY:** Data lineage shows complete chain of custody and downstream dependencies for an asset!



**CONSUMABILITY:** Audit Profiler shows both summarized views & patterns of access for a data asset.



Policy ID	Event Time	User	Access Type	Result	Access Enforcer	Client IP
19	07/28/2017 07:45:09 GMT	admin	UPDATE	ALLOWED	ranger-acl	172.27.25.135
19	07/28/2017 07:42:42 GMT	admin	UPDATE	ALLOWED	ranger-acl	172.27.25.135
19	07/28/2017 06:31:03 GMT	admin	CREATE	ALLOWED	ranger-acl	172.27.25.135



WEBINAR

---

# PRIVACERA INTRODUCTION

# ABOUT PRIVACERA

PLATFORM FOR DISCOVERING AND MANAGING SENSITIVE DATA

• • • • • • • • • •

PARTNERS



**cloudera**

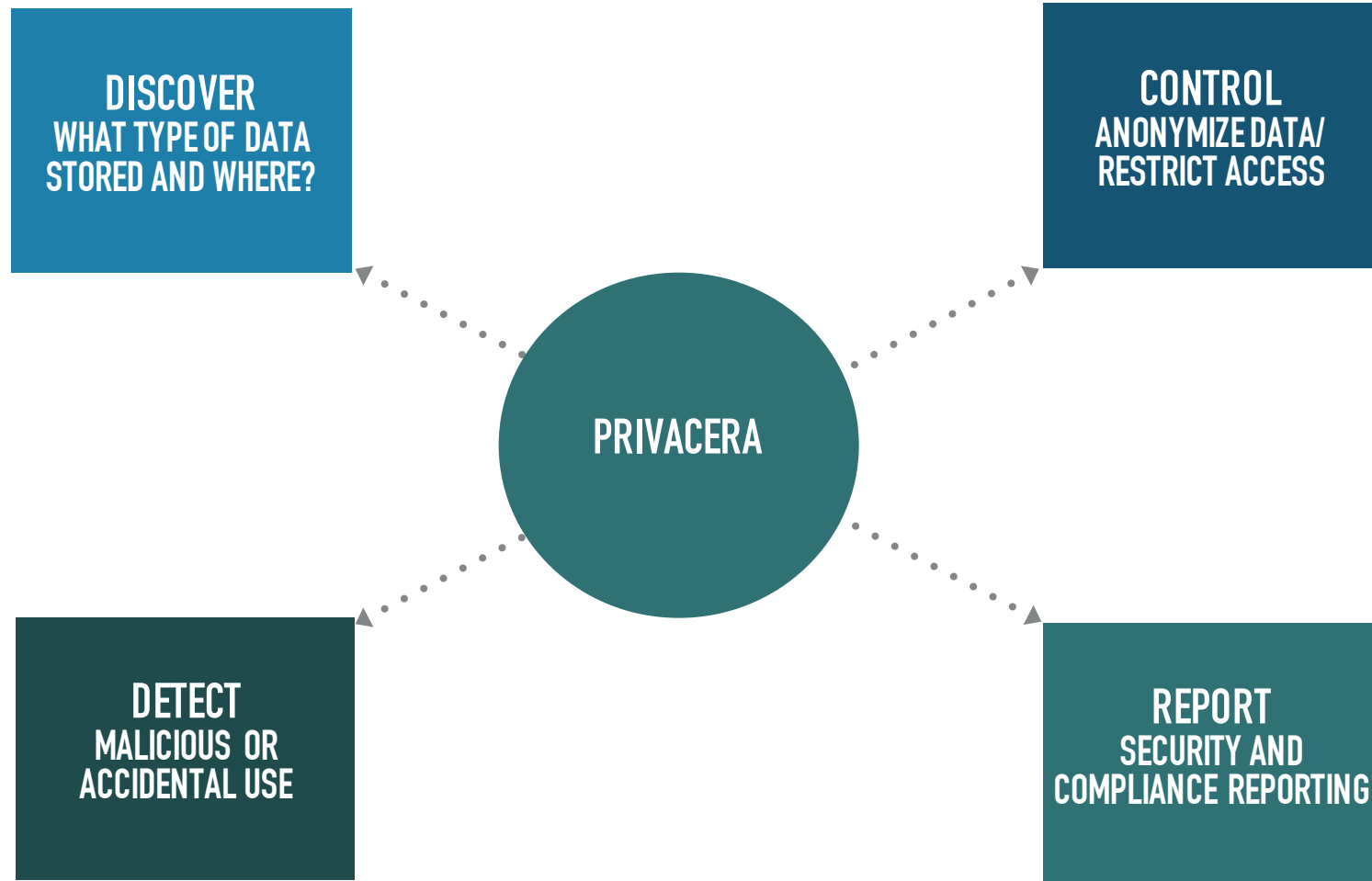
**splunk**>

• • • • • • • • • •

GLOBAL



# PLATFORM TO MANAGE SENSITIVE DATA



WEBINAR

---

# STEPS TO MANAGE SENSITIVE DATA

# 4 STEPS FOR MANAGING SENSITIVE DATA

**DATA  
DISCOVERY**



**ACCESS  
CONTROL**

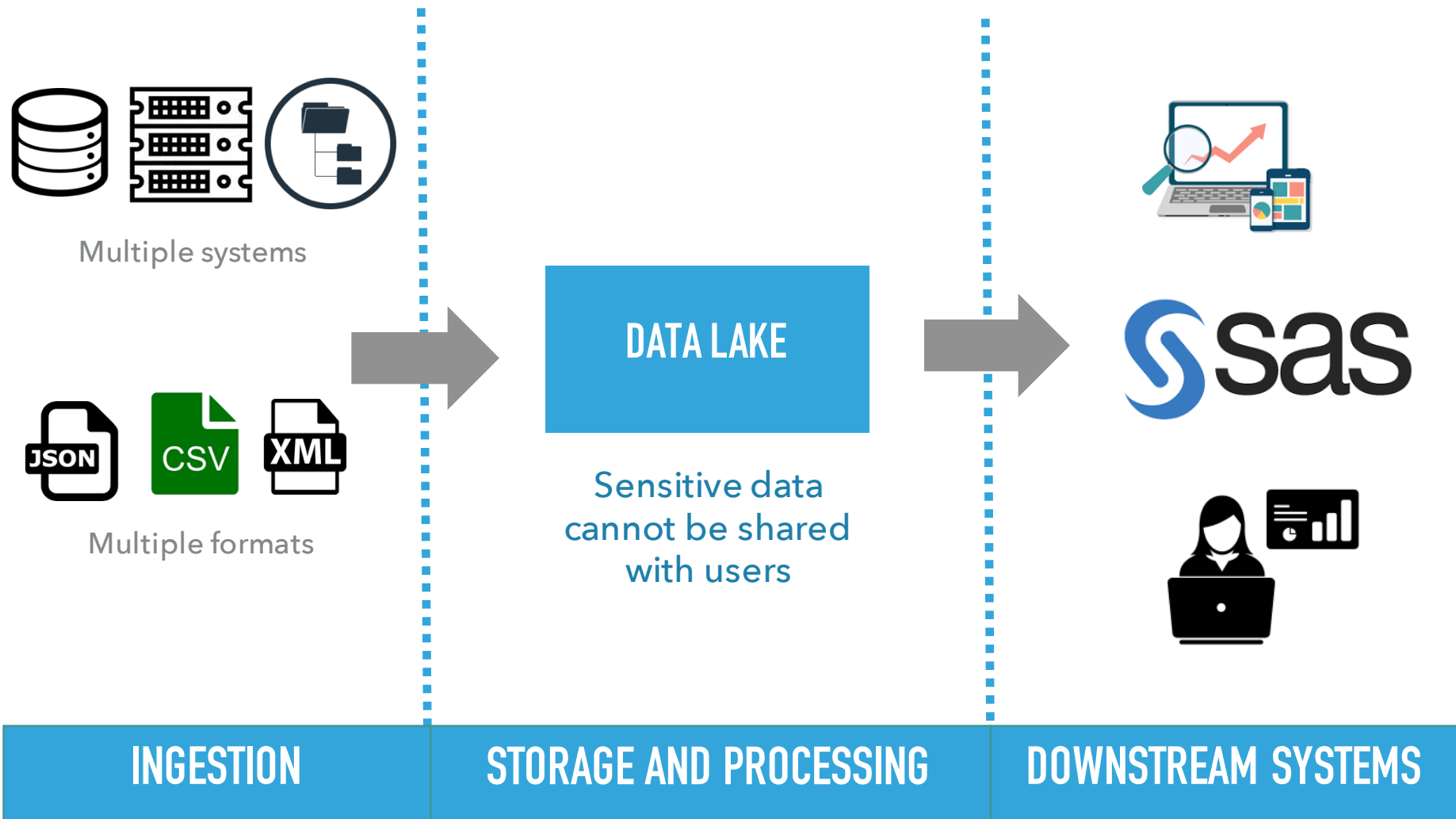


**ANONYMIZATION**



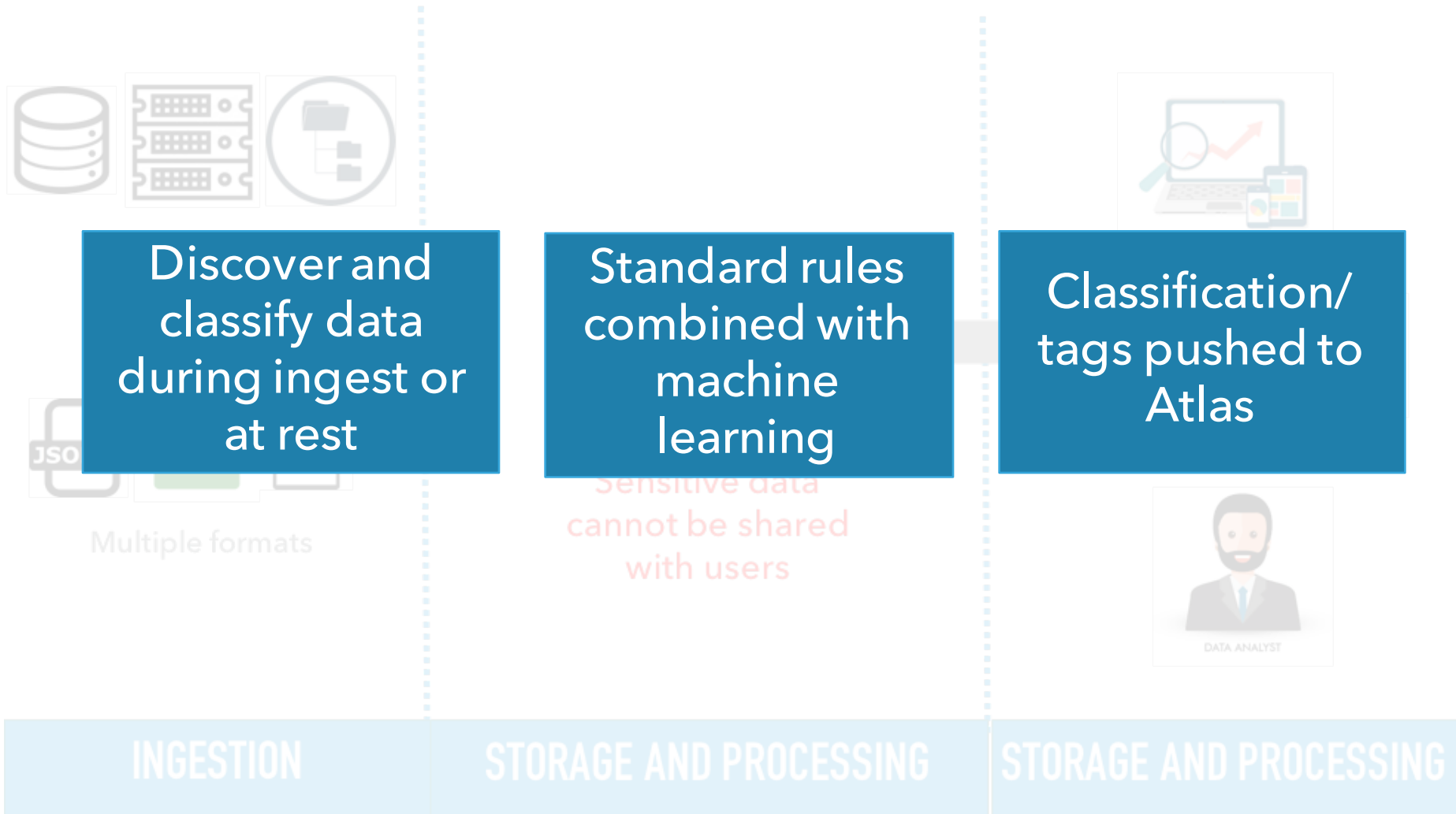
**MONITORING**

# REPRESENTATIVE SCENARIO – FINANCIAL SERVICES





# SOLUTION – PRIVACERA AUTOMATED DATA DISCOVERY



# REPRESENTATIVE SCENARIO – HEDGE FUND



Stock Info

**CONFIDENTIAL**

Proprietary  
Confidential data



**DATA LAKE**

Access to sensitive  
data is restricted



Data Scientist

**INGESTION**

**STORAGE AND PROCESSING**

**DOWNSTREAM SYSTEMS**

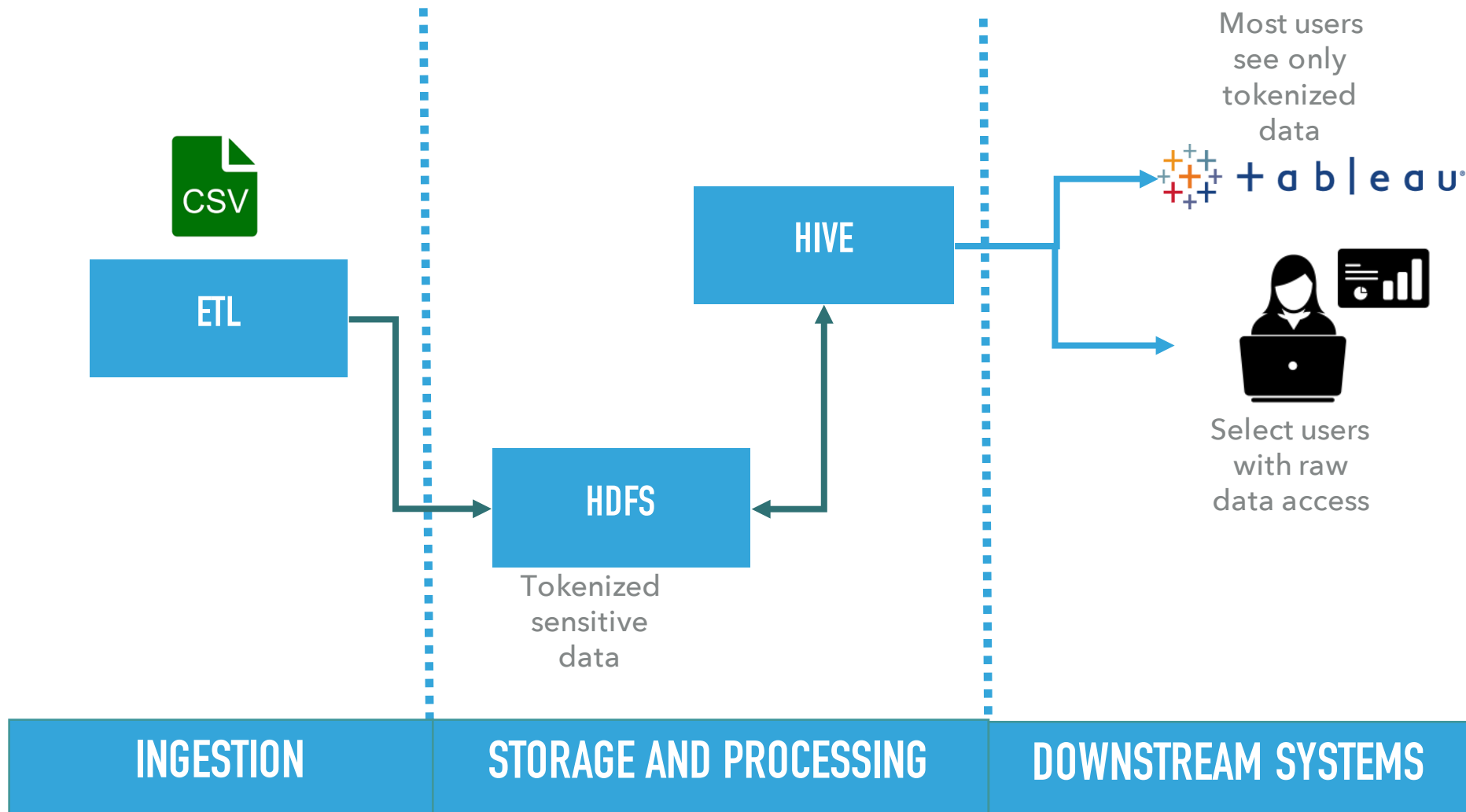
## SOLUTION – TAG BASED ACCESS CONTROL

Simplify policies  
by managing at  
tag level

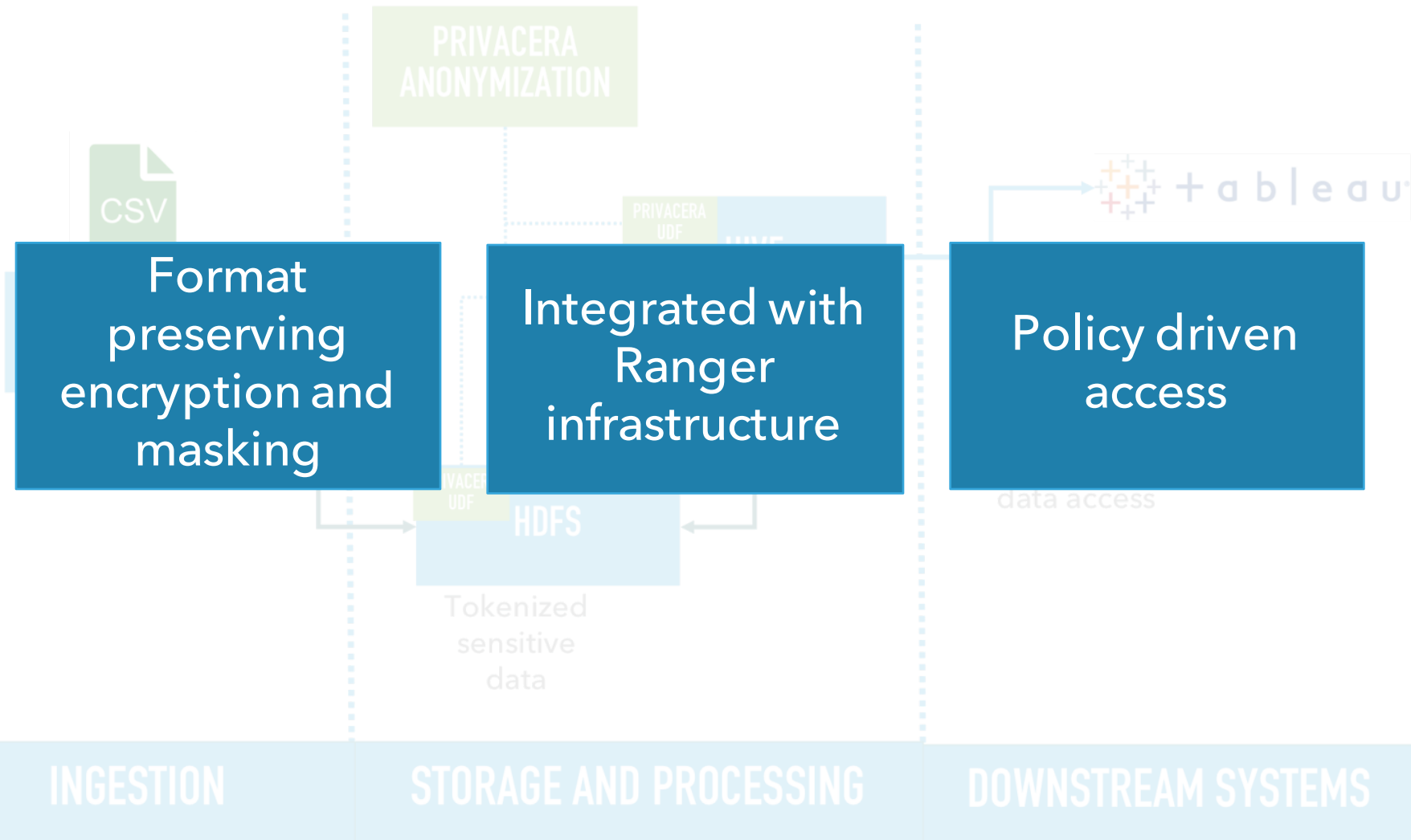
Tag attributes  
such as  
expiration date

Metadata  
updated by  
Privacera

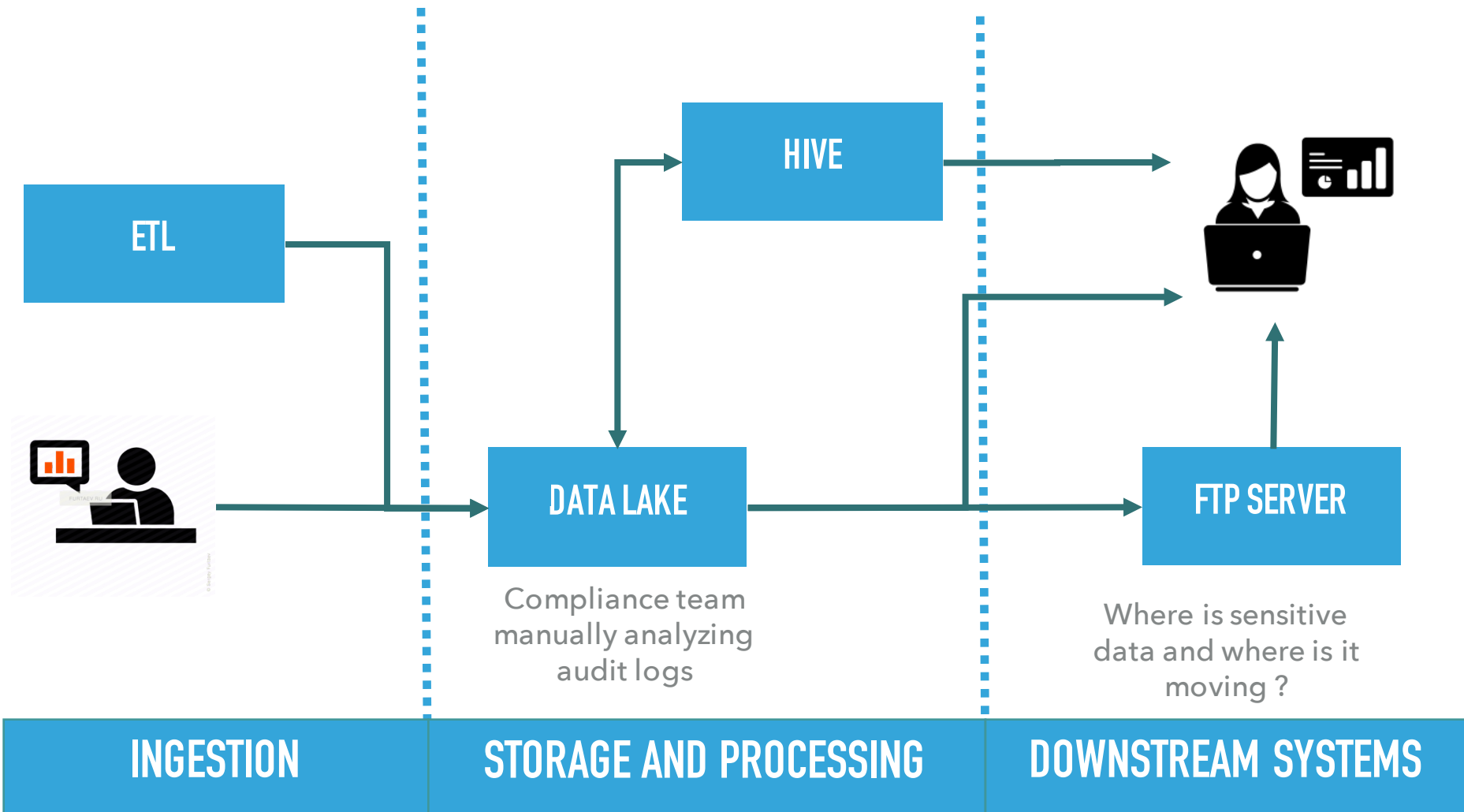
# REPRESENTATIVE SCENARIO – HEALTHCARE



# SOLUTION – PRIVACERA ANONYMIZATION

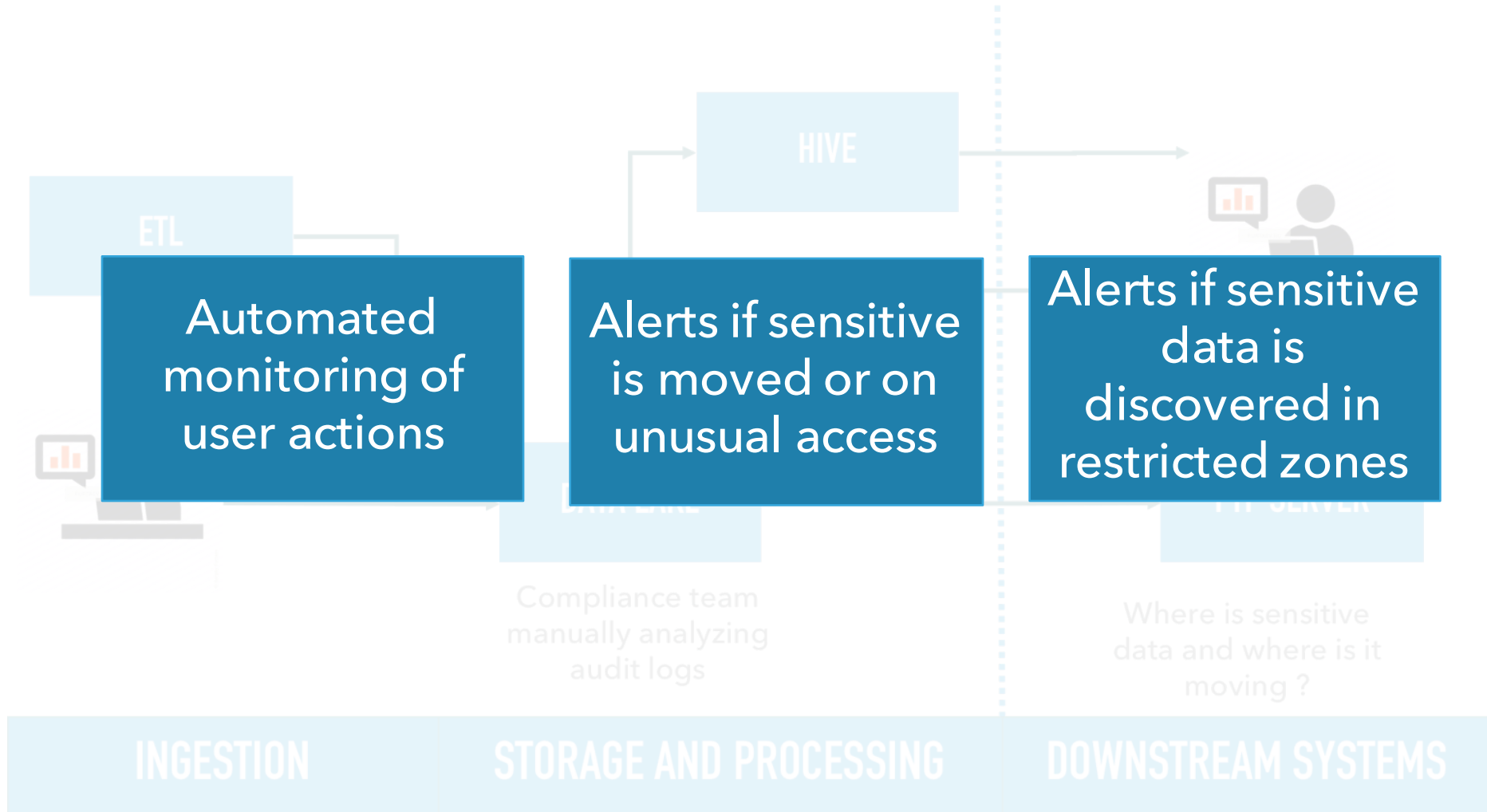


# REPRESENTATIVE SCENARIO – FINANCIAL SERVICES





# SOLUTION – PRIVACERA MONITORING



WEBINAR

---

DEMO

# SUMMARY

- ▶ Understand your data before expanding your data lake
- ▶ Invest in automated classification and centralized metadata
- ▶ Manage access to user by data classification
- ▶ Anonymize data to reduce exposure
- ▶ Monitor the use of data, “trust but verify”.
- ▶ Data plane provides next generation for tools for hybrid data infrastructure

[questions@privacera.com](mailto:questions@privacera.com)

---

**QUESTIONS ?**