

# CLOUDERA

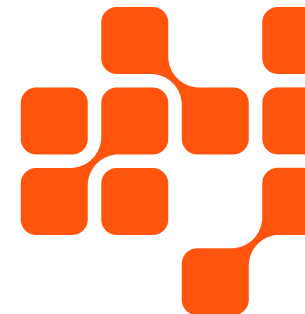
EBOOK

## Automated Data Lineage for Data Security



# Table of Contents

<b>Automated Data Lineage for Data Security</b>	<b>3</b>
<b>From Manual to Automated: The Power of Data Lineage</b>	<b>4</b>
<b>Intruder Alert!</b>	<b>5</b>
<b>Where's Your Weak Link?</b>	<b>6</b>
<b>Go With the Flow</b>	<b>7</b>
<b>Play By the Rules</b>	<b>8</b>
<b>Data With Integrity</b>	<b>9</b>
<b>Automated Data Lineage = an Unmatched Data Security Partner</b>	<b>10</b>



# Automated Data Lineage for Data Security

You spend your time involved with data, of course, but just for a moment, imagine that you're the head of a small physical security team for a distribution center with 1000 storerooms, 300 corridors and 100 employees on staff. You receive a tip-off that during work hours yesterday, one of your employees went into a storeroom that should have been off-access to them and tampered with some of your products. It's your job to discover and remediate the issue.

Of course, like a good security team, you have security cameras stationed in every corridor and every storeroom. So all you need to do is have your team of three people review all of yesterday's footage and identify the employee who was in the wrong place, at the wrong time, doing the wrong thing.

An easy job, right?

Well, yes... if you can deliver the answer next month, and your team won't have any other responsibilities until then.

Except that's never the reality. If you have a security breach, you need to identify and address the issue as soon as possible to stop damage from spreading. And your ongoing security work can't stop in the meantime.

Data security is no different.

If you're responsible for a modern data environment with multiple data systems and landscapes, hundreds of data pipelines and thousands of (or more) data assets, you just can't drive manual.





# From Manual to Automated: The Power of Data Lineage

Put yourself back in the swivel chair of the distribution center head of security. Instead of reviewing hours upon hours of video footage, you press a button and a program rapidly scans and processes the footage for you.

Within minutes, you're looking at a diagram of your complex's layout with a visual, color-coded overlay of all employee movements over yesterday's workday: a visualization of every single person's path from the time they came into the complex until the time they left. See something unusual? Just click on the path and get more details, like what that employee did in the room they were in at the time.

A security leader's dream, no?

While the technology isn't quite there to accomplish this kind of simple, elegant and effective visualization in the world of physical security, it is a reality in the world of data security.

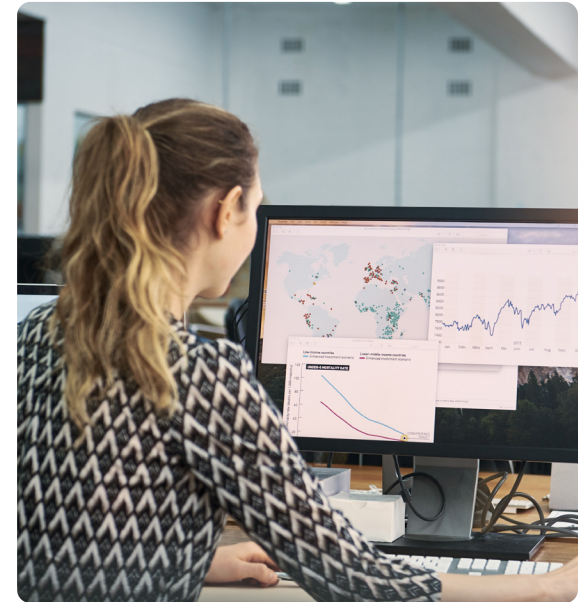
Data lineage is the mapping of data's movements from when it enters your data environment until it leaves - and everything that happens to it along the way.

- What data pipelines did it traverse?
- What transformations did it undergo?
- What systems did it impact?

Data lineage is a necessary component of any data security measure, implementation or investigation:

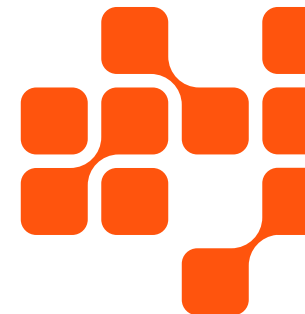
- Why are these numbers wrong - and was it an accident or tampering?
- Who has access to these data assets - and should they have that access?
- Where is the personal or sensitive information in our systems - and is it sufficiently protected?
- Has there been an unexpected change in our data pipelines - and which downstream assets and systems were affected?

To answer any of those questions requires detailed knowledge of the involved data's journey. Your team could trace the journey out manually... but that's such a waste of your valuable human resources.



Why spend days connecting the dots when automated data lineage can give you the complete picture in minutes so you can take swift, informed action?

Let's take a dive into some specific areas in which data lineage directly supports data security efforts.



# Intruder Alert!

## Threat Detection and Incident Response

Data lineage plays a vital role in detecting and responding to security threats and incidents promptly.

Since automated data lineage provides a clear understanding of the origin, movement and transformation of data across various systems and processes, it enables organizations to:

- Establish baseline data behavior
- Detect deviations or anomalies that might indicate security breaches or unauthorized activities
- Trigger real-time alerts to security teams for prompt investigation of suspicious data movements or transformations

If suspicions of a security incident are confirmed, automated data lineage is critical for thorough incident investigation and forensics and prompt response, providing or supporting the ability to quickly:

- Track down the source of an attack
- Trace the path of the compromised data
- Determine the scope and impact of the breach

- Understand the affected systems or applications
- Identify other potential areas of vulnerability
- Determine appropriate remedial measures
- Prioritize response efforts to contain the incident and restore security
- Reconstruct the sequence of events
- Gather evidence to support forensic analysis and legal proceedings
- Draw conclusions for future prevention

Data lineage enhances an organization's threat detection capabilities, supports incident response and aids in investigations. By leveraging data lineage, organizations can implement robust security measures, safeguard their sensitive data from unauthorized access and effectively deal with incidents when they occur.

# Where's Your Weak Link?

## Vulnerability and Risk Assessment

Data lineage facilitates vulnerability and risk assessments by providing insights into the potential risks associated with data flows and transformations.

Since automated data lineage creates intuitive visualizations of how data moves through an organization's systems, applications and processes, it enables organizations to:

- Easily understand the data's origins, transformations and destinations
- Identify weak points or vulnerabilities in their data infrastructure, data transfer mechanisms and data processing systems

With a clear and comprehensive view of the risks in the data environment, data lineage empowers organizations to:

- Proactively manage risk
- Prioritize security measures
- Patch vulnerabilities
- Mitigate potential threats

Data lineage provides organizations with the visibility, control and understanding they need to plan for the effective protection of their data systems and assets. By leveraging data lineage, organizations can easily identify potential security gaps and build targeted security plans and priorities.



# Go With the Flow

## Data Flow Control

Data lineage is crucial in empowering organizations to establish control over the flow of their data.

Since automated data lineage traces and presents easily understandable schematics of the end-to-end flow of data throughout an organization's data landscape, it enables the organization to:

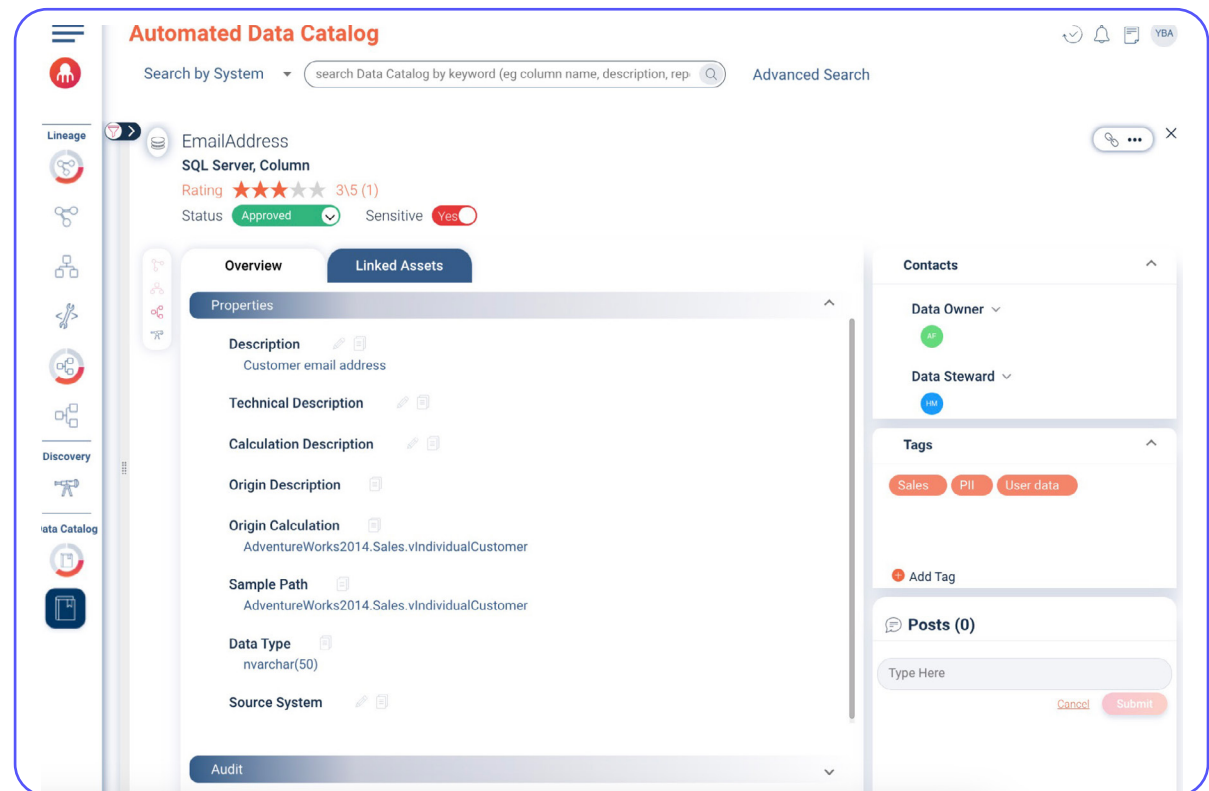
- Identify the systems, applications and individuals that have access to sensitive data
- Pinpoint critical points in the data flow where security measures need to be implemented

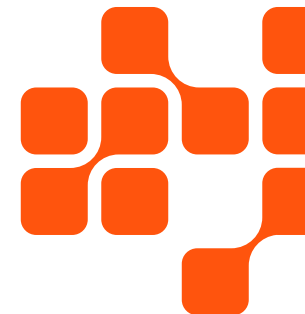
By providing a clear view of current data flow and access, data lineage allows organizations to efficiently:

- Set up appropriate access controls
- Define user privileges
- Enforce strong authentication measures
- Identify where encryption or masking is necessary
- Establish data loss protection mechanisms

By leveraging powerful tool for data lineage, organizations to smoothly and intelligently implement robust access control and authorization mechanisms, thereby preventing unauthorized access, insider threats, data breaches and data leakage.

**Leveraging Data Catalogs, you can assess who has access to sensitive assets.**





# Play By the Rules

## Compliance and Regulatory Requirements

Data lineage helps organize compliance with the stringent data security and privacy regulations achieve and demonstrate required by many industries and jurisdictions.

Since automated data lineage provides a detailed record of data movement and transformations, it facilitates compliance with regulatory requirements through enabling the organization to:

- Track how data is handled, stored and processed
- Understand data segmentation between systems and business units
- Easily locate the personal, sensitive or private data to which the regulations apply

When it comes to regulatory requirements, it is not enough for organizations to comply; they must be able to prove compliance through reporting, audits and regulatory assessments. Preparing such proof manually is a tremendous time and resource drain, and so automated data lineage is invaluable when it comes to:

- Showing how data is collected, stored, processed and shared
- Providing a clear trail of data handling practices
- Generating reports, visualizations and documentation required for audits
- Demonstrating adherence to legal and regulatory frameworks
- Reducing the risk of non-compliance penalties and reputational damage

The transparency provided by automated data lineage - and the speed at which it provides that transparency - facilitates compliance with data security regulations and frameworks, saving organizations immeasurable time and resources.



# Data With Integrity

## Data Quality Assurance

Data lineage contributes to data quality assurance efforts, ensuring the integrity of data within the organization's data landscape.

Since automated data lineage tracks and records the data's journey from origin to destination, including everything that happened to it along the way, it enables organizations to easily:

- Establish trusted data sources
- Identify data anomalies, discrepancies or inconsistencies
- Signal unauthorized modifications that could compromise data integrity
- Prevent data tampering
- Head off data poisoning attacks
- Verify the accuracy, consistency and reliability of data

Data lineage can be leveraged to ensure that data remains untampered and reliable throughout its lifecycle, thus enhancing an organization's overall data security.

Data lineage helps organizations comply with the stringent data security and privacy regulations they achieve and demonstrate required by many industries and jurisdictions.

Since automated data lineage provides a detailed record of data movement and transformations, it facilitates compliance with regulatory requirements through enabling the organization to:

- Track how data is handled, stored and processed
- Understand data segmentation between systems and business units
- Easily locate the personal, sensitive or private data to which the regulations apply

When it comes to regulatory requirements, it is not enough for organizations to comply; they must be able to prove compliance through reporting, audits and regulatory assessments. Preparing such proof manually is a tremendous time and resource drain, and so automated data lineage is invaluable when it comes to:

- Showing how data is collected, stored, processed and shared
- Providing a clear trail of data handling practices
- Generating reports, visualizations and documentation required for audits
- Demonstrating adherence to legal and regulatory frameworks
- Reducing the risk of non-compliance penalties and reputational damage

The transparency provided by automated data lineage – and the speed at which it provides that transparency – facilitates compliance with data security regulations and frameworks, saving organizations immeasurable time and resources.

# Automated Data Lineage = an Unmatched Data Security Partner

Data security in a modern data environment requires modern, automated tools that don't leave you waiting around for answers.

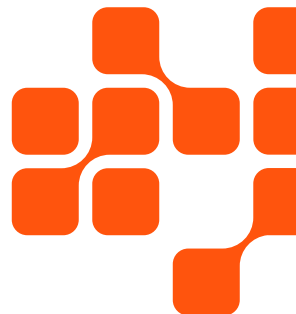
Automated data lineage gives you the full picture of where, when and how your data is moving. It enables you to spot vulnerabilities, catch anomalies and track issues to their source. It empowers you to take swift, informed action.

Automated data lineage does all the manual tracking work for you, freeing up your skilled human resources to do the truly strategic parts of data security.

If your goal is to efficiently maintain a secure data environment, automated data lineage is your partner in (stopping) crime.

Want a peek at how automated data lineage could boost your data security efforts?

[Schedule your demo now.](#)



# About Cloudera

Cloudera is the only true hybrid platform for data, analytics, and AI. With 100x more data under management than other cloud-only vendors, Cloudera empowers global enterprises to transform data of all types, on any public or private cloud, into valuable, trusted insights. Our open data lakehouse delivers scalable and secure data management with portable cloud-native analytics, enabling customers to bring GenAI models to their data while maintaining privacy and ensuring responsible, reliable AI deployments. The world's largest brands in financial services, insurance, media, manufacturing, and government rely on Cloudera to be able to use their data to solve the impossible—today and in the future.

To learn more, visit [Cloudera.com](https://cloudera.com) and follow us on [LinkedIn](#) and [X](#). Cloudera and associated marks are trademarks or registered trademarks of Cloudera, Inc. All other company and product names may be trademarks of their respective owners.



Cloudera, Inc. | 5470 Great America Pkwy, Santa Clara, CA 95054 USA | [cloudera.com](https://cloudera.com)

© 2024 Cloudera, Inc. All rights reserved. Cloudera and the Cloudera logo are trademarks or registered trademarks of Cloudera Inc. in the USA and other countries. All other trademarks are the property of their respective companies. Information is subject to change without notice. 0000-001 August 11, 2025