

 **MMUTA** +  **databricks**

E-BOOK

A Guide to Automated Data Security

In Databricks Using Immuta

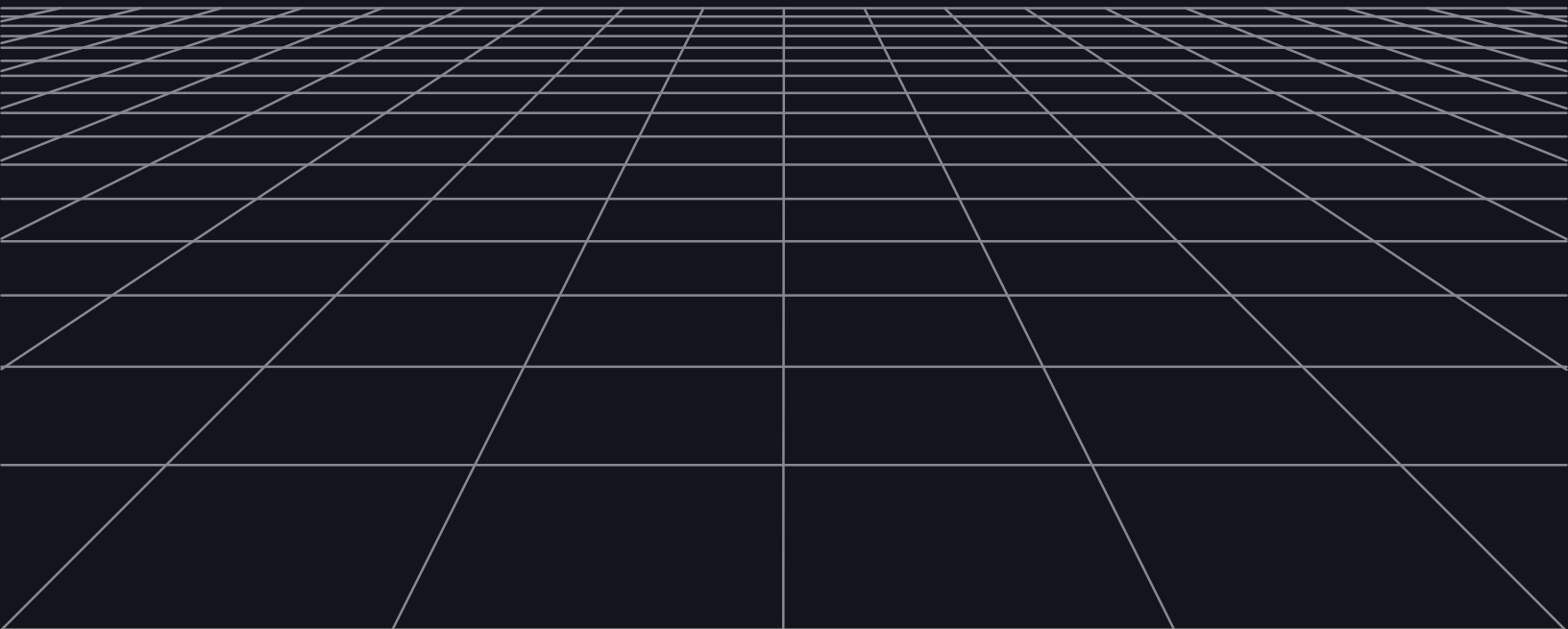


Table of Contents

Introduction	3
Automated Data Discovery & Classification	5
What is it?	5
What problems does it solve for data teams?	6
How is it implemented?	7
Dynamic Data Security & Access Controls	8
What are they?	8
What problem do they solve for data teams?	11
How are they implemented?	12
Continuous Data Security Monitoring	14
What is it?	14
What problem does it solve for data teams?	15
How is it implemented?	16
Conclusion	17

Introduction

Data-driven insights are no longer optional — they're necessary. Organizations that leverage data to make critical business decisions have a definitive competitive advantage. While consumers today appreciate and often pay more for customized products and services, delivering that experience often requires the use of sensitive personal data.

Sensitive data is the most valuable asset for analytics and data science. But if data engineers and architects can't proactively secure this resource while maximizing its utility for real-time access by analysts and data scientists, is it still as valuable?

Centralizing data and making it discoverable and actionable are key to data-driven innovation, personalized user experiences, and revenue growth. But the advent of new [data privacy regulations](#) like California's CCPA and Europe's GDPR, in addition to increasingly stringent internal data rules and security policies, have put a tax on data engineering teams, who are charged with translating these rules and regulations into executable policies so that data consumers can gain access to critical

data for innovation. This complicated and time-consuming process — which can lead to personal liability if a data leak or breach occurs — can delay or halt advanced analytics and data science projects, which require fast access to data.

To maintain the balance between [data security](#) and speed to data, Immuta integrates with Databricks to enable data teams to secure and scale lakehouse access by automating data discovery and classification, providing comprehensive secure protection, and continuous data monitoring. This white paper outlines those capabilities and considerations for when each should be used, so that data teams can unlock even more use cases with Databricks and Immuta.

Why might data teams need Immuta for Databricks?

- They're responsible for managing complex policies across many tables — as well as the ensuing role explosion.
- They need to empower stakeholders with business context of data use for self-service data access management.
- They're unable to easily prove compliant data use with corporate rules and regulations, or respond to investigations.
- They're unable to implement global policies and access controls for each platform in their cloud data ecosystem, so must manually do it for each individual platform.
- They're unable to dynamically restrict access based on time, geography, purpose, data sharing agreements, or other scenarios that may arise.
- They're responsible for consistent security and auditing across multiple data platform technologies, but don't have a centralized way to manage the process or ensure its uniformity, which causes confusion and frustration.
- They're expected to stay up-to-date on evolving regulations and to implement sufficient data security measures accordingly — otherwise they could be held personally liable for leaks or breaches.
- They spend their time managing case-by-case data access within an organization, as opposed to delivering new innovations.
- They need to be able to mask data while preserving its original format for non-production use.
- They're expected to enable compliant, secure collaboration on data sets without inadvertently granting unauthorized access or hindering analytics initiatives.
- Their policies are difficult to scale because they're unable to draw upon existing organizational glossaries or underlying metadata.

Automated Data Discovery & Classification

What is it?

Manual processes are simply not practical for an increasingly cloud-based and decentralized data ecosystem. Users rely on Databricks – in addition to their other cloud platforms – for fast, automated data storage and compute. Why should data teams expect any less from their data security solution?

With Immuta and Databricks, automated data security and privacy controls are enforced using dynamic attribute-based access controls to serve as additional buffers against unauthorized access, data leaks, and re-identification. But before these controls can even be enforced, data first needs to be discovered and classified appropriately. Only after data discovery is completed and sensitive data is correctly tagged can security policies be built to effectively protect it.

Sensitive Data Discovery & Tagging

Manually tagging sensitive data as it is uploaded to Databricks is a significant time commitment for data teams, in addition to introducing the risk of human error. As data becomes available faster than ever before, it's easy to see how this process can quickly become unmanageable for data engineers and architects, who could also potentially be held liable for any sensitive data that slips through the cracks and ends up in the wrong hands.

Immuta leverages Databricks Unity Catalog APIs to monitor for schema and user changes, and enriches user metadata by discovering and tagging sensitive data, inferring additional information like PII, tagging data with external information, and leveraging Unity lineage for tag propagation. This allows Databricks users to discover the sensitive data that is entered into their data ecosystem without requiring manual effort. New data sources are scanned and classified

using a combination of more than sixty prebuilt classifiers and any custom or domain-specific classifiers.

Data identified as being sensitive can be assigned tags that correspond to access control policies which are dynamically applied at query time. Immuta automatically tags sensitive fields, like PII, personal data, or PHI, in addition to enabling data engineers and architects to create tags as needed. Tags can also map to data privacy protection laws such as CCPA and GDPR, which in turn correspond to Immuta's templated starter policies that automate policy creation and ensure compliance with regulations. This means data engineers and architects can seamlessly ingest sensitive data and apply regulatory-compliant policies without fear of personal liability or wasted time.

What problem does it solve for data teams?

According to Immuta's report on Data Engineering & Operations for Analytics, the most challenging aspect of managing a data pipeline for data teams is data masking and security, followed by auditing and monitoring. The least challenging part? Extract and load.

Maintaining an ETL pipeline and GRANTS can be complicated and cumbersome for data teams when data discovery, classification, and governance aren't built into the process. This is because data engineers and architects must transform raw data into "clean" data before it can be utilized by analysts, data scientists, or any other data consumers. The report found that this task overwhelmingly falls on the shoulders of data engineers – more than half of survey respondents said transform is done by data engineers, either before or after load. So, in order to make data operational, data engineers must apply comprehensive security and privacy controls as part of transform. This model also sets data engineers and architects up for having to make copies of data and apply controls to them manually.

Without an integrated system that proactively discovers and classifies sensitive data, these time-intensive, manual processes may or may not result in sufficiently protected information.

Immuta's automated security and privacy controls make this task much faster and more secure for Databricks users. Incorporating automated data discovery and classification into the ETL process with Immuta's native Databricks integration eases the burden of combing through and identifying potentially sensitive data, and streamlines the process of updating pipelines and GRANTS as data users and policies change over time.

When data consumers run a query in Databricks, Immuta's dynamic controls are applied at run time based on how data has been discovered and tagged. This captures the most current data policies and permissions, simplifies the ETL process, and eliminates the need for manual processes and data copies. Dynamic, automated controls that reduce the risk of human error, missed sensitive data, and re-identification are a more efficient, secure approach for data teams.

How is it implemented?

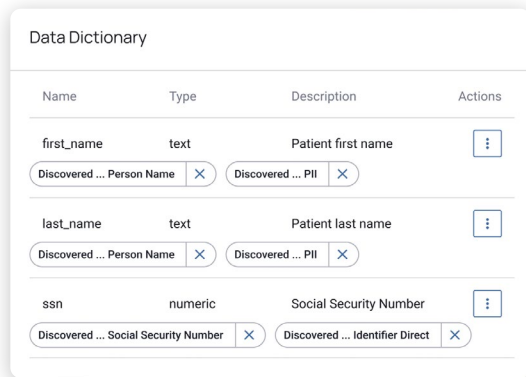
To understand the benefits of data discovery and classification with Databricks and Immuta in practice, use medical records as an example.

In this scenario, a data engineer receives a data set of medical records and is responsible for preparing it for use by a medical center's billing department and the city's public health department. The data set contains protected health information (PHI) including name, address, insurance information, and social security number, and the data engineer must implement policies and privacy controls that restrict data visibility by user attribute and purpose, while achieving HIPAA compliance.

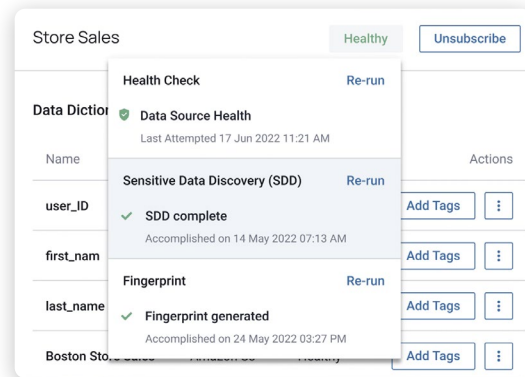
As the data set is uploaded to Databricks, Immuta's sensitive data discovery tags and classifies PHI so the data engineer can assign it to access control policies. The data engineer can then build policies in plain language that restrict data returned from a query based on user attributes. So, an accountant from the billing department may be able to view patients' full addresses and insurance information in order to process claims, but the data engineer can suppress diagnoses with a value of `redacted`.

Meanwhile, an analyst from the public health department may be able to view diagnoses, but the data engineer suppressed patients' names, social security numbers, and insurance information, and generalized addresses by only showing the first three numbers of a zip code – for example, `021XX`. This would allow the analyst to see diagnoses by general area, so they can make public health recommendations, without identifying any particular individuals with a positive diagnosis.

Using automated controls based on classified and tagged data, the data engineer in this scenario preserves the original data set within Databricks – without making copies – and ensures data consumers see only the information relevant to their functional need.



Name	Type	Description	Actions
first_name	text	Patient first name	[Info]
Discovered ... Person Name [X] Discovered ... PII [X]			
last_name	text	Patient last name	[Info]
Discovered ... Person Name [X] Discovered ... PII [X]			
ssn	numeric	Social Security Number	[Info]
Discovered ... Social Security Number [X] Discovered ... Identifier Direct [X]			



Store Sales Healthy Unsubscribe

Health Check Re-run

Data Dictionary

- ✓ Data Source Health
Last Attempted 17 Jun 2022 11:21 AM
- Sensitive Data Discovery (SDD) Re-run
 - ✓ SDD complete
Accomplished on 14 May 2022 07:13 AM
- Fingerprint Re-run
 - ✓ Fingerprint generated
Accomplished on 24 May 2022 03:27 PM

Columns: user_ID, first_name, last_name, Boston Sto... [Add Tags] [Info]

Dynamic Data Security & Access Controls

What are they?

As organizations continue to invest resources into migrating data to the cloud with powerful platforms like Databricks, they face an array of common challenges. They desire fast and easy data access for their varied – and likely growing – pool of data users, but need to ensure that this data is effectively and consistently secured against leak or breach.

Securing data in modern cloud data ecosystems can be an extremely complex process. With more data sources, users, and platforms introduced into the data stack, it becomes difficult to create policies that protect sensitive data wherever it lives and moves. On top of this, further decentralization can add increased risk of data breaches and leaks. This creates an ultimatum: companies can either halt the breakneck progress of their data use by aggressively locking it down, or leave it completely exposed to the dire effects of misuse.

Simply safeguarding access to raw data inputs is no longer sufficient to protect personal data. Data teams must also consider what information can be inferred about an individual from a model's behavior or API output, as well as any risks that may arise from publishing a data set. Dynamic data security and access controls protect data in a manner that maximizes privacy and utility, enabling both internal or external data use with a reduced possibility of attacks on the privacy of individuals. Databricks users can implement a range of dynamic security and privacy controls through Immuta:

Attribute-Based Access Control (ABAC)

Immuta's data access and security policies are built and applied based on user attributes. Attribute-based access control (ABAC) is an approach to data security that permits or restricts data access based on assigned user, object, action, and environmental attributes. In contrast to role-based access control (RBAC), which relies on the privileges specific to one role for data protection, ABAC has multiple dimensions on which to apply access controls. This makes attribute-based access control a highly dynamic model because policies, users, and objects can be provisioned independently, and policies make access control decisions when the data is requested.

Similar to how Databricks separates compute and storage, Immuta separates access control policy logic from the data platform. Within Databricks, Immuta extracts access control policy logic from both compute and storage, eliminating the need for data engineers and architects to spend time and energy recreating access control policies for each individual data tool. The flexibility to utilize the same access control policies – whether role-, attribute-, or purpose-based – consistently across distributed data storage and compute providers is a powerful capability, particularly in satisfying data use compliance, changing business rules, or working within a cloud data ecosystem that includes Databricks.

Data Masking

Data masking is the process of removing or obscuring identifiers, altering existing sensitive information in a data set to make a fake—but still convincing—version of it. This allows sensitive data to be stored and accessed, while maintaining the anonymity and safety of the information involved. This type of privacy-enhancing measure is crucial for organizations looking to secure their most sensitive data, and can facilitate measures like secure internal and external data sharing and protecting sensitive personal health information (PHI) and financial data.

Masking can be applied using a range of techniques, most commonly by way of generalization or suppression. Generalization assigns the same broad value for any given attribute — for example, replacing an attribute, `hair color`, with the value `any`. Suppression works the opposite way, by removing values entirely or replacing them with a constant — for instance, replacing an attribute, `hair color`, with the value `redacted`.

Differential Privacy

The sheer amount of personal information being collected and used in today's environment means that no single piece of data exists in a vacuum. Differential privacy aims to mathematically limit an outsider's ability to confidently use the output of an analysis to make inferences about its input. This allows individuals providing their personal data to credibly deny their participation in the input.

Differential privacy requires the data analysis mechanism to give the same answers with similar probabilities over any pair of databases that differ by a single row. In other words, Immuta injects noise into the data analysis in order to render inference attacks nearly impossible. This way, an individual may claim that the output of the mechanism came from a database that did not include their data.

Randomized Response

While differential privacy enables people to credibly deny their participation in a data input, randomized response — also known as local differential privacy — makes it possible for participating individuals to credibly deny the contents of their participation records. This approach allows data subjects to answer sensitive or potentially embarrassing questions confidentially.

Like differential privacy, randomized response employs randomization to enhance privacy; however, unlike differential privacy, randomization is applied prior to submission and formal constraints are applied to the randomized substitutions. This

means that any chosen substitution must be nearly — though not necessarily exactly — as likely to arise from any given input. As a result, all potential inputs look plausible to an attacker wishing to undo the randomized substitution.

Since the randomized response technique is applied prior to the data leaving a device, data subjects are assured protection from the moment of submission. This protection remains privatized — even in the case of subsequent breach.

k-Anonymization

k-Anonymization is the data equivalent of hiding in a crowd; the more people – or in this case, data points – that are present and generally similar, the harder it is to pick out the details that can identify individuals. This approach reduces re-identification risk by anonymizing indirect identifiers, thereby destroying the signal of data.

In k-anonymization, K represents instances of tuples in a data set. A data set is k-anonymous when attributes within it are generalized or suppressed until each row is identical to at least $k-1$ other rows. Therefore, the higher the value of k , the lower the

re-identification risk. Just as the larger the crowd is, the less likely you'll find exactly the person you're looking for, k-anonymization works particularly well with large data sets. However, lines of data may have to be redacted if there isn't enough data to anonymize indirect identifiers.

k-Anonymization can help transform, analyze, and share secure data at scale, making it an important privacy enhancing technique for Databricks users dealing with large sets of sensitive data.

What problem do they solve for data teams?

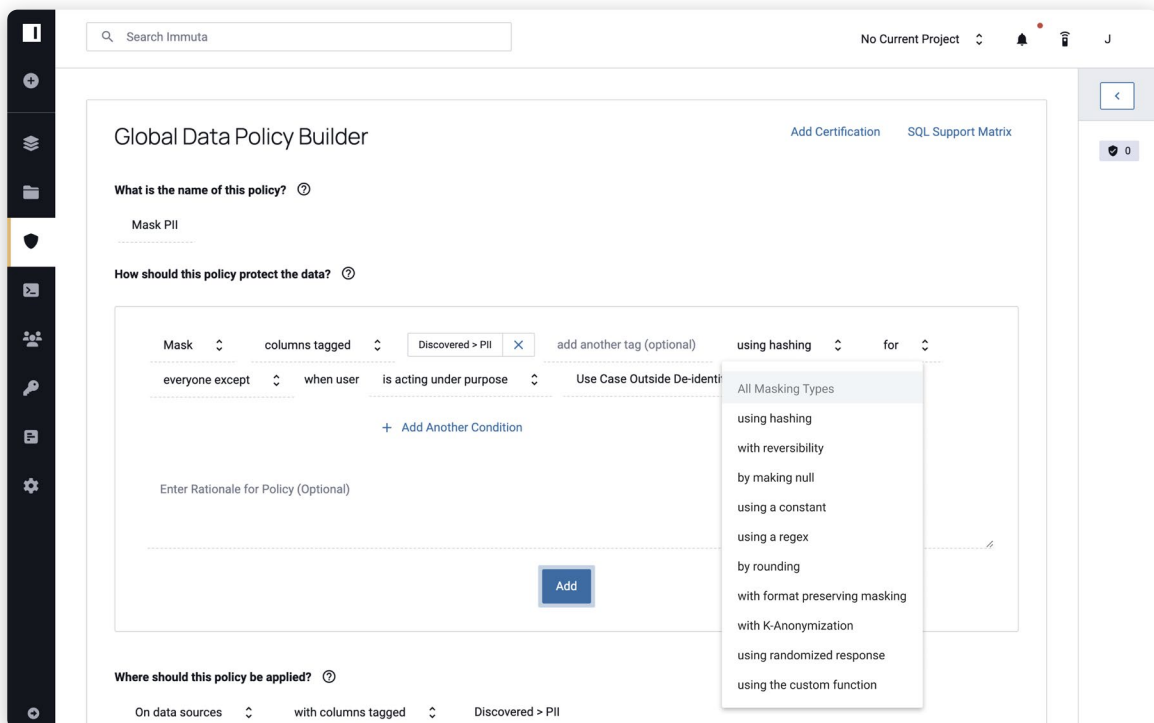
Data teams are often responsible for walking the line between data utility and security. The inherent problem is that creating and implementing controls that maximize both security and utility is highly complex and risky. Effective implementation that mathematically guarantees against re-identification often requires a PhD in applied mathematics.

Although the dynamic security and privacy controls shield some data fields, they are better able to preserve data's utility than some other methods.

In the k-anonymized health data example, suppressing [age] and [zip code] reduces the risk of re-identification by nulling the data, but a statistically relevant number of [age] and [zip code] combinations can still be analyzed to observe diagnosis trends by age and gender cohorts, for instance. For researchers at WorldQuant Predictive analyzing COVID-19 spread, this technique enabled data sharing from multiple data sources that helped streamline collaboration and quickly generate predictive models to form hypotheses without inadvertently exposing sensitive PHI or identifying information in the process.

Read More: [Explore WorldQuant Predictive's Immuta journey in this case study.](#)

A key benefit of dynamic security and privacy controls like attribute-based policies, k-anonymization, randomized response, and differential privacy is that they simultaneously reduce the risk of re-identification, maximize data utility and privacy, and enable data sharing between teams. Without these techniques, data teams are hindered in their ability to unlock collaboration and speed to data access without compromising data integrity, security, and compliance. Immuta's natively enforced advanced security capabilities mean Databricks users can access and share critical data — including sensitive data — quickly, efficiently, and securely.

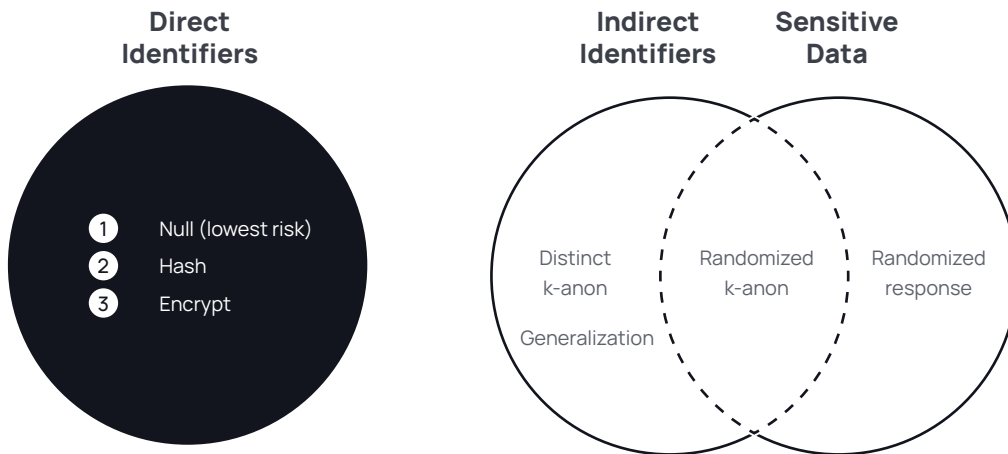


How are they implemented?

To understand how dynamic security and privacy controls apply to and help de-identify specific types of data, it helps to divide data into segments: that which contains direct identifiers, indirect identifiers and sensitive data.

It's important to first understand the two types of identifiers – in other words, the personal information that can be used to help identify an individual.

1. Direct identifiers are the pieces of personal information that are unique to an individual and can be used in isolation to identify that single person. For this reason, direct identifiers are highly sensitive and strictly regulated. Examples of direct identifiers include social security numbers, passport numbers, taxpayer identification numbers, full facial images, and medical record numbers.
2. Indirect identifiers, also known as quasi-identifiers, are pieces of personal information that are not unique to a single person. While indirect identifiers cannot alone be used to identify a specific individual, they are still considered sensitive because they can often be combined with other information to single somebody out. Examples of indirect identifiers include height, ethnicity, hair color, car make and model, and occupation.



As this figure demonstrates, these privacy enhancing technologies (PETs) can buffer protection for indirect identifiers and sensitive data, and vastly reduce the opportunity for inference or linkage attacks.

These advanced controls have grown in popularity across industries to enable secure data sharing. Data teams are well versed in the stringent requirements – and enforcements – associated with regulations like GDPR, HIPAA, and CCPA, which explicitly state that data must be protected such that the chances of an individual being re-identified remain as close to unfeasible as possible. History has shown that personal data is easy to re-identify if not adequately protected – Harvard professor Latanya Sweeney found that 87% of the population can be re-identified simply using birth date, zip code, and gender.

For Databricks users, leveraging personal data like credit card transactions and cell phone location data to drive timely analytics is a necessity, but one that comes at a high risk. By law, a Databricks table containing PHI, such as genders, zip codes, and credit card numbers, must have PETs applied before it can be shared with data scientists and analysts.

Once the table has been registered in Databricks with Immuta, a data engineer can create a new data

source and Databricks connection. Then, automated data discovery and tagging flags identifiers and sensitive data in the set. Using Immuta's plain language policy builder, the data engineer can enforce one or more dynamic security and privacy controls. For example, k-anonymity can be enforced using the masking or suppression method on columns tagged as `[gender]` and `[zip code]`.

With this policy selected, Immuta scans the Databricks table to calculate the statistics required for k-anonymization. A fingerprint service runs a query against Databricks to collect counts for each possible group of values in the data source and produces custom predicates for each column. The data engineer can either use this automated return's minimum group size, k, or can manually specify a value for k. To protect identity data, the predicates only contain a whitelist of values visible to users.

This secured data set is exposed as a table in Databricks so data analysts and scientists can access and query the table. Since the controls are enforced natively on read from Databricks, the underlying data remains unmodified and not copied, and policies are applied to the plan that Spark builds for a user's query from the Notebook.

Continuous Data Security Monitoring

What is it?

Data engineers and architects are often asked by security and compliance teams, “who accessed this data over the past two months?” On the surface, this seems like a simple question; in reality, it’s much more complicated.

Continuous data security monitoring and auditing are necessary parts of the data pipeline management process, but in a multi-cloud compute platform environment they are difficult to execute. Disparate compute layers mean security and privacy controls may be enforced in the compute layer, passed through to storage or even implemented in analytical applications. Various data consumer roles and their corresponding permissions add an extra layer of complexity.

Databricks users can bypass tedious manual monitoring and auditing processes with Immuta’s natively integrated capabilities. All audit logs and information, as well as cluster queries to the cloud provider, are done with system accounts so that data usage across cloud compute platforms – not just within Databricks or your cloud provider services – can be captured consistently. This includes audits

of not just data queries, but also of all policy actions being taken in Immuta, such as changing policies or subscribers to a table. Additionally, Databricks data teams can implement purpose restrictions that trigger consent workflows and simplify the process of monitoring and auditing data usage.

Immuta’s continuous data monitoring, unified audit logs, automated reports and purpose-based access controls provide granular snapshots of which data consumers accessed specific data sources, when, and for what purpose, as well as changes to data over time. As a result, legal and compliance stakeholders always have quick access to insights that prove compliant data usage across the organization.

What problem does it solve for data teams?

As personal and sensitive data has become more widely available and collected by organizations, data subjects have become increasingly aware of and concerned with how their information is being used.

As personal and sensitive data has become more widely available and collected by organizations, data subjects have become increasingly aware of and concerned with how their information is being used. Regulations have followed accordingly and now continuous data security monitoring and auditing is not an option – it's a requirement.

Yet, as multi-cloud compute platform adoption becomes the new normal, data teams' ability to monitor and audit data use consistently across platforms like Databricks and others is exponentially more complicated. It also broadens the risk landscape for sensitive data, as it moves between more tools and platforms for storage and analysis. In turn, there are a wider array of places that data breaches could potentially occur, creating a need for increased monitoring and breach detection capabilities. This has the potential to substantially limit acceptable use cases for data in cloud analytics. Furthermore, it makes the process of responding to security and compliance requests time-consuming, if not impossible.

Databricks users avoid these concerns and roadblocks with Immuta's integrated detection, monitoring, and auditing capabilities. Data teams can efficiently provide legal and compliance teams with detailed logs and reports at the data level, to provide full transparency about what data was accessed, by whom, when, and for what purpose. This automates the otherwise-laborious process of formalizing proof of compliance to adhere to federal, state, and industry standards and regulations, which in turn bolsters collaboration between data and compliance teams and mitigates concerns about legal enforcement and personal liability. Data security posture management (DSPM) is also enabled through an analysis engine with both sensitivity level and risk profile indicators to address and mitigate any potential threats.

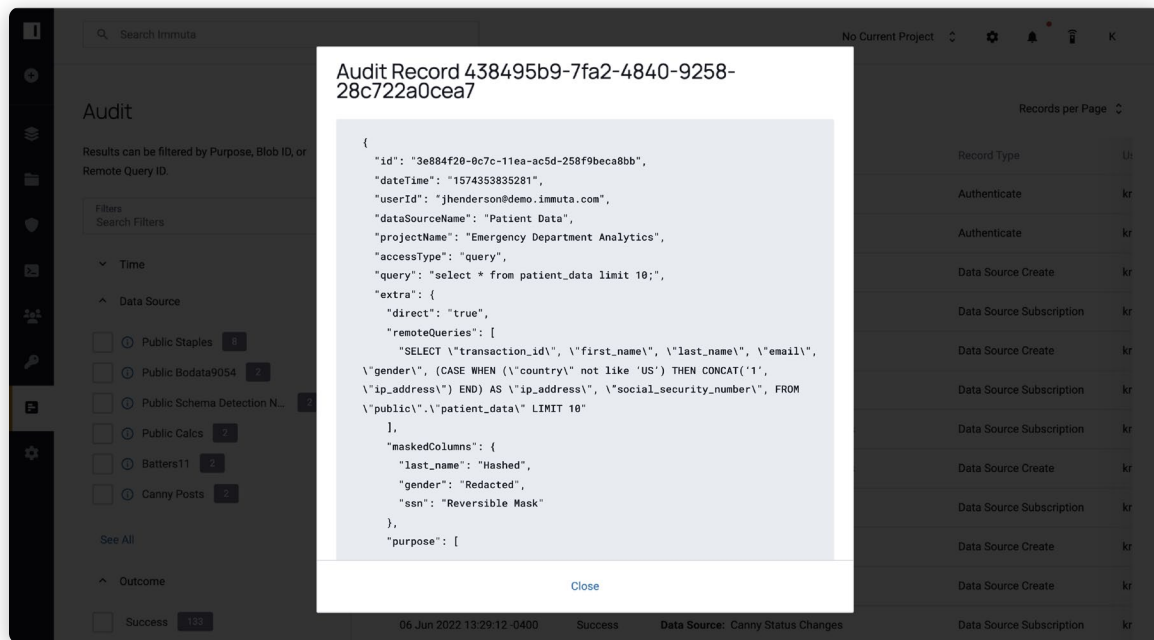
How is it implemented?

The GDPR requires data consumers to have an acceptable purpose for accessing data, and consequently, that data teams can provide evidence of that purpose.

This is a complicated process because, without the right tools, purpose is difficult to capture. Yet, GDPR and other regulations, including HIPAA, HITRUST, and SOC-2, are not optional and data teams can only avoid incurring noncompliance enforcement with comprehensive audits. Layered on top of regulations are employment and industry standards or contractual agreements, among others.

Immuta captures all details of where data is stored, who owns it, under what conditions it can be

accessed, when it was added, and how recently it was queried. Immuta's integration with Databricks enables full monitoring of all user activity, policy related activity and history, compliance and anomaly reports, and alerts and notifications. In addition to auditing, the plain language policy builder also means business stakeholders can assess the strength of policies in meeting compliance standards without having to rely on IT.



Conclusion

Centralized, secure, self-service data access is the best way to maximize data's impact. In today's rapidly evolving market, that can be the secret weapon to gaining a critical competitive edge. Whether analyzing data to inform public health policy decisions or to predict consumer behavior, data consumers need secure access to data – as fast as possible.

This isn't an option without the right combination of data science and automated data security capabilities. Immuta's integration with Databricks enables organizations to automatically secure sensitive data for analytical use in industries ranging from insurance to transportation. With automated data discovery and classification, dynamic security and access controls, and continuous data security monitoring and auditing, Databricks and Immuta together enable data teams to maximize data utility and security.

In fact, Databricks customers that take advantage of Immuta's core, native capabilities experience results such as 100x faster data access, a 93x reduction in data policies, and a 4x increase in data utilization. This means teams are able to accomplish more and unlock more data-driven outcomes in Databricks when Immuta is natively implementing dynamic data security and privacy measures.

To find out what you can accomplish when you combine the power of Databricks and Immuta, schedule a demo today. →

About Immuta

Immuta enables organizations to unlock value from their cloud data by protecting it and providing secure access. The Immuta Data Security Platform provides sensitive data discovery, security and access control, data activity monitoring, and has deep integrations with the leading cloud data platforms. Immuta is now trusted by Fortune 500 companies and government agencies around the world to secure their data. Founded in 2015, Immuta is headquartered in Boston, MA.

About Databricks

Databricks is the lakehouse company. More than 9,000 organizations worldwide – including Comcast, Condé Nast, and over 50% of the Fortune 500 – rely on the Databricks Lakehouse Platform to unify their data, analytics and AI. Databricks is headquartered in San Francisco, with offices around the globe. Founded by the original creators of Apache Spark™, Delta Lake and MLflow, Databricks is on a mission to help data teams solve the world's toughest problems. To learn more, follow Databricks on [Twitter](#), [LinkedIn](#) and [Facebook](#).

