

Lake Formation: Data Security e Data Governance mediante LF-TBAC

9 Novembre 2021 - 10 min. read

AWS Glue

AWS Identity and Access Management (IAM)

AWS Lake Formation

Data Security and Governance

Database

Introduzione

Il Big Data è cresciuto rapidamente come modo per descrivere le informazioni ottenute da fonti eterogenee quando diventa incredibilmente difficile gestirle in termini di **varietà, veridicità, valore, volume e velocità**. Tuttavia, esso può essere considerato il "nuovo petrolio a causa del potenziale di generare valore aziendale".

Senza adeguate procedure di governance o di gestione della qualità, i data lake possono trasformarsi rapidamente in "paludi" di dati ingestibili. I data engineer sanno perfettamente che i dati di cui hanno bisogno vivono in questi data lake, ma senza una chiara strategia di governance dei dati, non saranno in grado di trovarli, fidarsi o utilizzarli in modo adeguato ai loro bisogni.

Una sfida molto comune è **mantenere** la governance, il **controllo degli accessi** sugli utenti che operano sul Data Lake e la protezione delle informazioni sensibili.

Le aziende devono centralizzare la governance, il controllo degli accessi e applicare una strategia supportata da servizi managed per controllare in modo granulare l'accesso degli utenti ai dati.

Affrontare queste problematiche richiede tipicamente due approcci: *manuale*, **più flessibile** ma **complesso**; *gestito* che richiede che la propria **soluzione si adatti a**

standard specifici, ma in cambio **elimina tutte le complessità di gestione** per gli sviluppatori.

Questo articolo guiderà il lettore attraverso la configurazione del proprio Data Lake mediante Lake Formation, mostrando tutte le sfide che devono essere affrontate durante il processo mostrando un occhio particolare alla sicurezza e alla governance attraverso l'approccio LF-TBAC.

Il controllo degli accessi basato su tag, in breve **TBAC**, è un modo sempre più diffuso per risolvere queste sfide, applicando vincoli basati su tag associati a risorse specifiche.

Quindi, senza ulteriori indugi, andiamo a dare un'occhiata!

Che cos'è l'accesso tramite TBAC

Il controllo degli accessi basato su tag consente agli amministratori di risorse abilitate per IAM di creare policy di accesso, basate su tag esistenti, associati a risorse idonee.

I provider cloud gestiscono i permessi sia degli utenti che delle applicazioni tramite delle policy, documenti con regole che fanno riferimento a risorse. Applicando i tag a tali risorse, è possibile definire condizioni di allow/deny semplici ed efficaci.

L'utilizzo dei tag di gestione degli accessi può ridurre il numero di criteri di accesso necessari all'interno di un account cloud, fornendo anche un modo semplificato per concedere l'accesso a un gruppo eterogeneo di risorse.

Perché S3 da solo non basta

S3, come la maggior parte dei servizi AWS, **sfrutta i *principal* di IAM per la gestione degli accessi**, il che significa che è possibile definire quali parti di un bucket (file e cartelle o prefissi) un singolo principal IAM può leggere o scrivere; tuttavia non è possibile limitare ulteriormente l'accesso IAM a parti specifiche di un oggetto, né a determinati segmenti di dati archiviati all'interno di oggetti.

Ad esempio, supponiamo che i dati della nostra applicazione siano archiviati come una raccolta di file parquet suddivisi per paese in cartelle diverse.

È possibile vincolare un utente ad accedere solo agli utenti appartenenti a un determinato paese. Tuttavia, non c'è modo di impedire loro di leggere le informazioni anagrafiche (es. nome utente e indirizzo) memorizzate come colonne nel file parquet.

L'unico modo per impedire agli utenti di accedere a informazioni sensibili sarebbe crittografare le colonne prima di scrivere i file su S3, il che può essere **lento**, **ingombrante** e aprire un "vaso di Pandora" per quanto riguarda **l'archiviazione delle chiavi**, la condivisione e, infine, la **disattivazione delle chiavi**.

Inoltre, **concedere l'accesso a entità esterne utilizzando i principal IAM è spesso di per sé un problema non banale**.

Fortunatamente, AWS **offre una soluzione tutto-in-uno al problema delle autorizzazioni di Data Lake su S3**: facciamo entrare in gioco AWS Lake Formation!

AWS Lake Formation è un servizio completamente gestito che semplifica la creazione, la protezione e la gestione dei data lake, automatizzando molti dei complessi passaggi manuali necessari per crearli.

Lake Formation fornisce anche **il proprio modello di autorizzazioni, che è ciò che vogliamo esplorare in dettaglio, che potenzia il classico modello di autorizzazioni AWS IAM**.

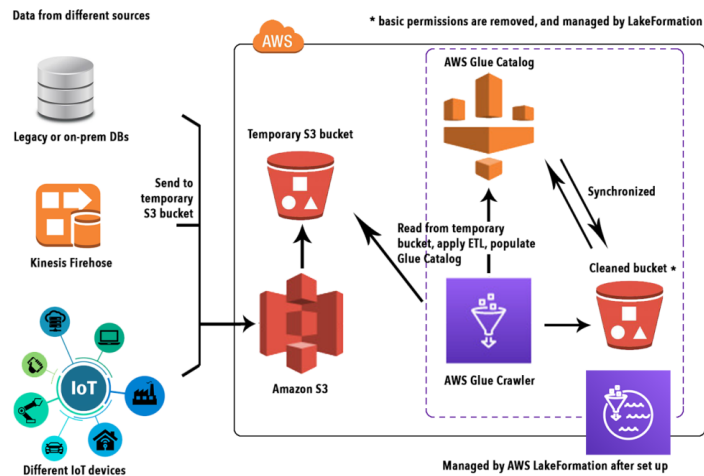
Questo modello di autorizzazioni definito in modo centralizzato, consente l'accesso granulare ai dati archiviati nei data lake attraverso un semplice meccanismo di concessione/revoca.

Quindi, sfruttando la potenza di Lake Formation, vorremmo dimostrare, con una soluzione semplice, come affrontare le suddette sfide di S3; proseguiamo!

Utilizziamo l'approccio TBAC in LakeFormation

Per accompagnare il lettore nella comprensione del motivo per cui AWS Lake Formation può essere una buona scelta nell'affrontare le complessità della gestione di un DataLake, abbiamo preparato un semplice tutorial su come migrare dati eterogenei.

Dai database locali legacy a S3, creando anche un catalogo Lake Formation per gestire la pulizia dei dati, le autorizzazioni e ulteriori operazioni.



Schema di esempio del nostro tutorial

Migrazione di dati on-prem via AWS Glue

Il primo passo per la creazione di un Data Lake è ovviamente quello di recuperare, trasformare e inserire i dati. In questo semplice esempio, abbiamo utilizzato un set di dati utente simulato da un database MySQL. La colla AWS è il modo naturale per connettersi all'origine dati eterogenea, dedurre il loro schema importare e trasformare i dati e infine scriverli su S3 [come spiegato in dettaglio qui](#).

Dopo che i dati sono stati caricati in un bucket S3 temporaneo, è necessario creare **un database in Lake Formation** per connettersi a un **Glue Crawler** ed eseguirlo sul prefisso S3 per popolare un Glue Catalog per i dati.

Vai alla **console AWS Lake Formation**, nella pagina *Database* nella **scheda Data catalog**, e inserisci un nome per il database e il tuo percorso su S3.

AWS Lake Formation > Databases > Create database

Create database

Database details
Create a database in the AWS Glue Data Catalog.

Database
Create a database in my account.

Resource link
Create a resource link to a shared database.

Name
articolo-lakeformation

Location - optional
Choose an Amazon S3 path for this database, which eliminates the need to grant data location permissions on catalog table paths that are this location's children
s3://lakeformation-article-blog

Description - optional
articolo-lakeformation
Descriptions can be up to 2048 characters long.

Default permissions for newly created tables
This setting maintains existing AWS Glue Data Catalog behavior. You can still set individual permissions, which will take effect when you revoke the Super permission from IAMAllowedPrincipals. See [Changing Default Settings for Your Data Lake](#).

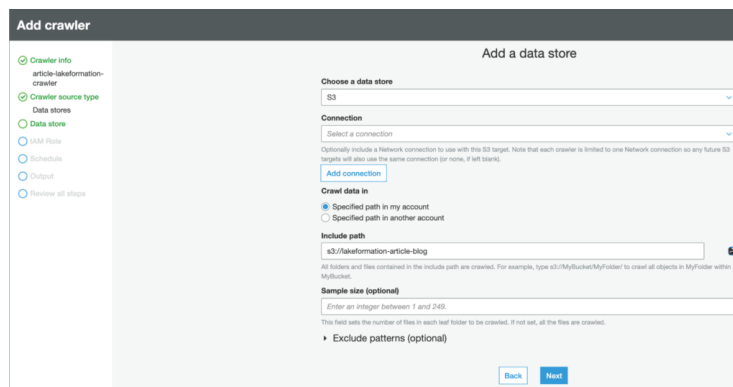
Use only IAM access control for new tables in this database

Creare un nuovo Database da Lake Formation

Nota: la creazione di un database direttamente da Lake Formation assicura che le autorizzazioni siano correttamente associate ad esso; avremmo potuto fare la stessa cosa anche da AWS Glue, ma avremmo avuto bisogno di uno sforzo aggiuntivo per modificare le autorizzazioni per i passaggi successivi.

Dopo che il database è stato creato, abbiamo bisogno di creare un Glue Catalog, che è un metastore contenente lo schema (schema-on-read) dei dati salvati su S3 (di solito come file parquet).

Avere uno schema Glue è **necessario per configurare i livelli di accesso da AWS Lake Formation per il proprio Data Lake su S3**. Per farlo, basta creare un crawler e collegarlo allo stesso percorso S3 del database e impostare **quel DB come output del crawler**.



The screenshot shows the 'Add crawler' wizard in the AWS Glue console, specifically the 'Add a data store' step. On the left, a navigation pane shows the progress: 'Crawler info' (completed), 'Crawler source type' (completed), 'Data stores' (current step), 'Data store' (pending), 'IAM Role' (pending), 'Schedule' (pending), 'Output' (pending), and 'Review all steps' (pending). The main area is titled 'Add a data store' and contains the following fields and options:

- Choose a data store:** A dropdown menu with 'S3' selected.
- Connection:** A dropdown menu with 'Select a connection' selected. Below it, a note states: 'Optionally include a Network connection to use with this S3 target. Note that each crawler is limited to one Network connection so any future S3 targets will also use the same connection (or none, if left blank)'. There is an 'Add connection' button.
- Crawl data in:** Two radio buttons: 'Specified path in my account' (selected) and 'Specified path in another account'.
- Include path:** A text input field containing 's3://lakeformation-article-blog'. Below it, a note states: 'All folders and files contained in the include path are crawled. For example, type s3://MyBucket/MyFolder to crawl all objects in MyFolder within MyBucket.'
- Sample size (optional):** A text input field with a placeholder 'Enter an integer between 1 and 243'. Below it, a note states: 'This field sets the number of files in each leaf folder to be crawled. If not set, all the files are crawled.'
- Exclude patterns (optional):** A section for defining patterns to exclude.

At the bottom right, there are 'Back' and 'Next' buttons.

Setup di un AWS Glue Crawler standard

Per utilizzare il crawler, è necessario un ruolo IAM, ma fortunatamente AWS ha un step per questo nella procedura guidata di creazione del crawler:



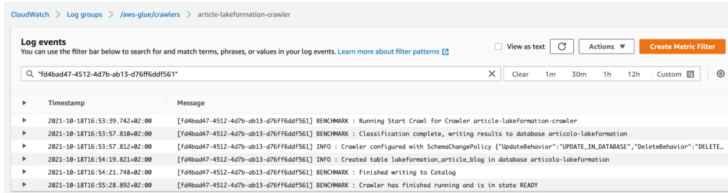
The screenshot shows the 'Choose an IAM role' step in the AWS Glue console. It features three radio buttons for selection:

- Update a policy in an IAM role
- Choose an existing IAM role
- Create an IAM role

Below the radio buttons, there is an 'IAM role' section with an information icon. It contains a text input field for the role name, which is 'AWSGlueServiceRole-lake-formation-article'. Below the input field, a list of permissions is shown: 's3://lakeformation-article-blog'. At the bottom, there are 'Back' and 'Next' buttons.

Come creare uno IAM role da utilizzare per il Crawler

Non appena il Crawler è creato, e i dati importati nel catalogo, siamo pronti per il prossimo step.

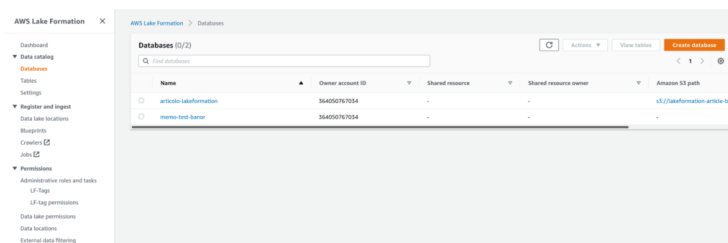


Cloudwatch Logs dimostra che il Crawler ha funzionato correttamente

AWS Lake Formation

Avendo a disposizione un catalogo di dati di Glue, è il momento di configurare Lake Formation per gestire finalmente le autorizzazioni di accesso degli utenti.

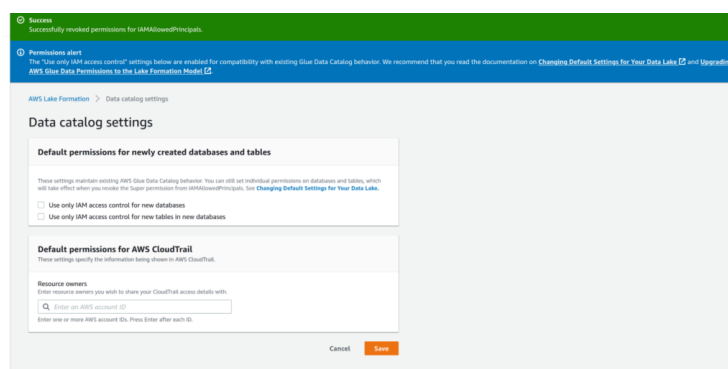
Per fare ciò, iniziamo andando nella **dashboard di Lake Formation e rimuovendo i consueti permessi di accesso su S3**.



Dashboard di Lake Formation

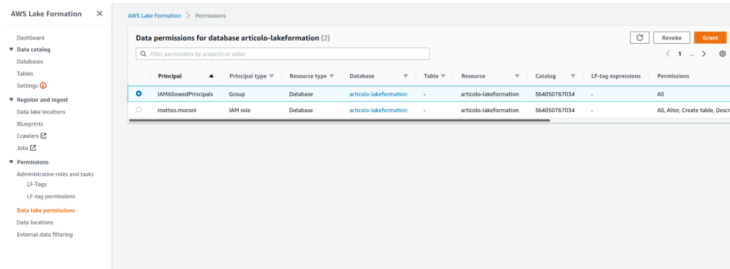
Quindi possiamo andare su *Data Catalog Settings* e deselezionare *Use only IAM access control for new databases* e *Use only IAM access control for new tables in new databases*.

Per impostazione predefinita, l'accesso alle risorse del catalogo dati e ai percorsi su Amazon S3 è controllato esclusivamente dalle policy di AWS Identity and Access Management (IAM), deselezionando i valori rendiamo effettive le **autorizzazioni** individuali di Lake Formation.



Settaggi di Lake Formation data catalog: dobbiamo disabilitare entrambi i flag Use only

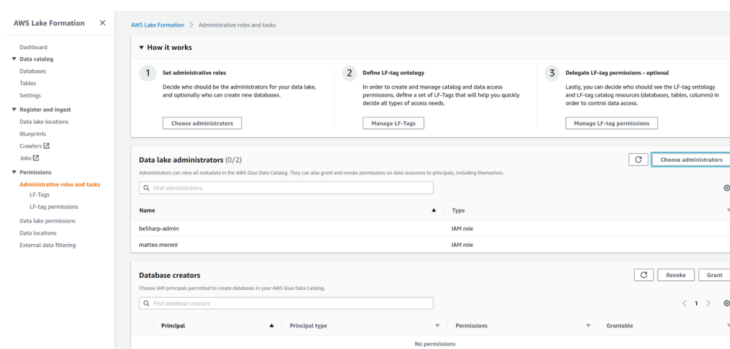
Una volta che le **responsabilità di accesso sono state delegate a Lake Formation**, possiamo rimuovere l'accesso per il gruppo IAM standard *IAMAllowedPrincipals*, nella scheda *Permissions* del data lake; selezioniamo dunque **permission of the IAM group** e facciamo click su *Revoke*.



Rimuoviamo le autorizzazioni standard IAMAllowedPrincipals

L'utente che crea il data lake sarà anche elencato in questa sezione con i privilegi di amministratore; se si desidera invece che l'utente mantenga l'accesso ai dati è possibile lasciare l'autorizzazione così com'è, altrimenti è possibile **revocare o limitare l'autorizzazione all'utente/ruolo**.

*Nota: se è necessario aggiungere un Data lake administrator principal, è possibile farlo accedendo ad Administrative roles and tasks e aggiungendo un **amministratore del Data Lake**.*



Aggiungiamo admin e db creator dalla console

Una volta completati tutti questi passaggi, è il momento di iniziare a definire i tag Lake Formation (da ora in poi **LF-Tags**), che verranno utilizzati per limitare l'accesso al data lake.

Dalla pagina *LF-Tag* sotto la scheda *Permissions* **creiamo un nuovo LF-Tag**, come chiave utilizziamo *level*, come valori invece *private*, *sensitive* e *public* separati da virgola proprio come in figura. Facciamo infine click su **Add LF-tag**.

Add LF-Tag [Learn More](#) ✕

LF-Tags have a key and one or more values that can be associated with data catalog resources. Tables automatically inherit from database LF-tags, and columns inherit from table LF-tags.
Example: Key = Confidentiality | Values = private, sensitive, public

Key
level
Key string must be less than 128 characters long, and cannot be changed once LF-tag is created.

Values
Type a single value and select [Enter] or specify multiple values separated by commas.
Add

private ✕ sensitive ✕ public ✕
Enter up to 15 values; each value must be less than 256 characters long.

Cancel **Add LF-tag**

Creazione di un LF-Tag

Ora, una volta creati, come possiamo utilizzare questi tag per imporre il controllo degli accessi? Per prima cosa andiamo nella sezione database e **selezioniamo il nostro database**, creato all'inizio del tutorial. Nelle azioni del database, possiamo selezionare il tag appena creato e il livello di autorizzazione.

Di solito, lasciamo aperto l'accesso al database e limitiamo le autorizzazioni in base alla tabella e ai campi, ma questa scelta di solito varia a seconda del database. Nel nostro esempio, assegniamo il livello public all'intero database di esempio.

Edit LF-Tags: articolo-lakeformation [Learn More](#) ✕

LF-Tags
After they are associated with catalog resources, LF-Tags allow you to create scalable permissions.

Assigned keys Values

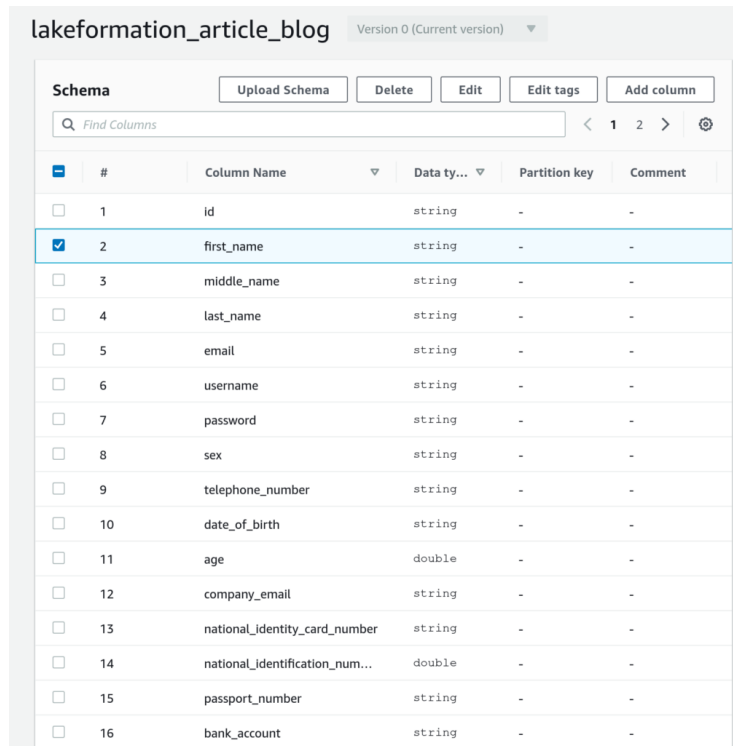
level ✕ public ▼ Remove

Assign new LF-Tag
You can add 49 more LF-tags.

Cancel **Save**

Modifichiamo gli LF-Tag per un intero database

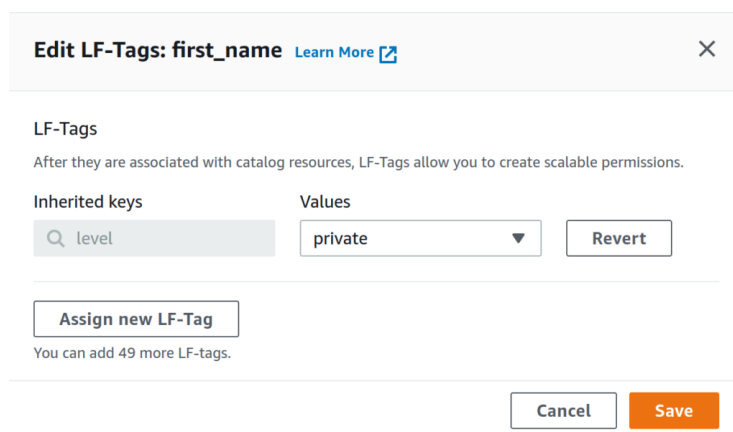
Ora, se vogliamo **restringere l'accesso alle colonne contenenti informazioni personali nella tabella utenti**, possiamo accedere alla tabella da modificare, selezionare una colonna specifica e modificare il suo LF-tag da **public** a **private** (come mostrato nelle figure successive).



The screenshot shows the 'lakeformation_article_blog' schema interface. At the top, there are buttons for 'Upload Schema', 'Delete', 'Edit', 'Edit tags', and 'Add column'. Below these is a search bar labeled 'Find Columns'. The main part of the interface is a table with columns: #, Column Name, Data ty..., Partition key, and Comment. The table lists 16 columns, with the second column, 'first_name', selected (indicated by a blue checkmark in the first column of the row).

#	Column Name	Data ty...	Partition key	Comment
1	id	string	-	-
2	first_name	string	-	-
3	middle_name	string	-	-
4	last_name	string	-	-
5	email	string	-	-
6	username	string	-	-
7	password	string	-	-
8	sex	string	-	-
9	telephone_number	string	-	-
10	date_of_birth	string	-	-
11	age	double	-	-
12	company_email	string	-	-
13	national_identity_card_number	string	-	-
14	national_identification_num...	double	-	-
15	passport_number	string	-	-
16	bank_account	string	-	-

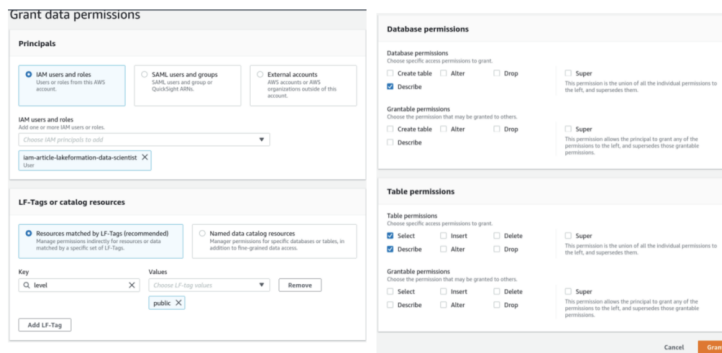
Schema del database di esempio nel quale selezioniamo una colonna



The screenshot shows a dialog box titled 'Edit LF-Tags: first_name'. It contains a search bar with 'level' entered. Below the search bar, there is a dropdown menu with 'private' selected. To the right of the dropdown is a 'Revert' button. At the bottom of the dialog, there is a 'Cancel' button and a 'Save' button.

Editiamo gli LF-Tag per i permessi sulle colonne

Ora dobbiamo solo definire quali IAM principals (ad esempio il nostro user di prova) avranno accesso ad uno specifico livello di LF-Tag. Per fare questo, andiamo su *Data lake permissions* e **diamo i permessi ad un utente, ruolo o gruppo IAM accesso alle risorse taggate con uno specifico LF-Tag.**

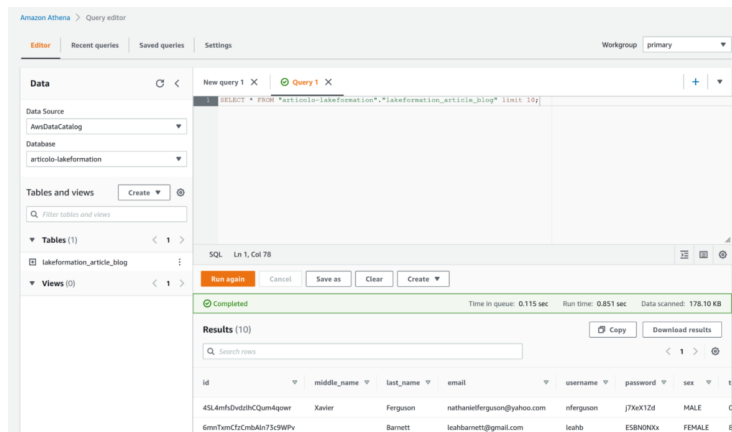


Diamo un'autorizzazione di read

Questo esempio mostra come dare ad un utente l'accesso a tutte le risorse etichettate con "level": "public".

L'utente di esempio potrà così vedere tutti i nostri database ad eccezione dei dati personali etichettati come privati. Un altro utente potrà avere accesso sia alle informazioni pubbliche che a quelle private, aggiungendo il livello *private* nella sezione LF-Tag o, in generale, modificando i tag delle colonne in base alle proprie esigenze.

Ora possiamo interrogare la tabella del database usando il nostro utente di prova che, in base al nostro set di autorizzazioni, non è in grado di vedere la colonna `first_name` (che è etichettata come privata).



Athena viene usato per fare query ai dati e dimostrare che la colonna `first_name` non è mostrata nella select perchè taggata come private

Come mostrato nella figura, siamo riusciti a negare al nostro utente di prova il diritto di vedere una colonna "riservata" a nostra scelta.

Vorremmo incoraggiare l'utente a sperimentare aggiungendo o rimuovendo anche le opzioni *describe* e *select* dai permessi LF-Tag nella sezione Data Lake, per vedere come sia possibile negare anche la visualizzazione di interi database o tabelle specifiche.

Nota: a partire dal 3 novembre 2021 per migliorare la sicurezza, AWS Lake Formation ha anche aggiunto il supporto per gli endpoint VPC gestiti tramite AWS PrivateLink per accedere a un data lake in una Virtual Private Cloud.

Feature in preview: sicurezza livello di riga

Lake Formation è ancora un servizio giovane, quindi c'è molto margine di miglioramento. AWS lavora costantemente per aumentare le funzionalità dei suoi servizi e Lake Formation non fa eccezione.

AWS Lake Formation consente già di impostare policy di accesso per nascondere i dati, ad esempio una colonna con informazioni riservate, agli utenti che non dispongono dell'autorizzazione per visualizzare tali dati.

La sicurezza a livello di riga si aggiungerà a ciò, consentendo di impostare criteri a livello di riga oltre ai criteri a livello di colonna.

Un esempio potrebbe essere l'impostazione di una policy che dia a un data scientist l'accesso solo ai dati dell'esperimento contrassegnati con un ID specifico.

Un altro aspetto interessante sarebbe condividere lo stesso Data Lake per set di dati diversi, riducendo così i costi e gli sforzi di gestione.

Concludendo

In questo articolo abbiamo visto come possiamo sfruttare la potenza dei servizi AWS per lo storage e l'analisi dei dati per affrontare la sfida imposta dai Big Data, in particolare come gestire l'accesso, le autorizzazioni e la governance.

Abbiamo dimostrato che i crawler di AWS Glue possono recuperare in modo efficace dati non strutturati da repository temporanei, che si tratti di database come RDS o on-premise, o storage di oggetti come S3, e ottenere uno schema per popolare un catalogo di Glue.

Abbiamo visto che partendo da S3 e da un metadata store è possibile creare un Lake Formation Catalog sopra S3, interamente gestito da AWS, per ridurre drasticamente lo

sforzio di gestione per impostare e amministrare un Data lake.

Abbiamo visto brevemente cos'è una metodologia Tag-Based Access Control (TBAC) e come può essere efficacemente utilizzata per gestire accessi e permessi.

Abbiamo dimostrato che AWS Lake Formation può applicare policy IAM e regole TBAC per concedere o limitare l'accesso a utenti e gruppi anche per riga/colonna. Abbiamo dimostrato che con Lake Formation e AWS Glue, possiamo oscurare i dati sensibili a committenti specifici.

Abbiamo descritto gli LF-Tags in dettaglio, con un semplice tutorial. Infine, abbiamo parlato della sicurezza a livello di riga.

Per concludere, possiamo dire che per le sfide relative ai Big Data e alle corrette soluzioni di storage, con un occhio di riguardo alle questioni di sicurezza e governance, le scelte possibili da fare sono sempre due: fai da te o optare per una soluzione gestita.

In questo articolo abbiamo scelto una soluzione gestita per mostrare tutti i vantaggi di un approccio più rigido al problema. Pur essendo una soluzione meno flessibile ad adattarsi, offre un servizio più aderente alle best practices di AWS e meno oneri nell'amministrazione e nella governance.

Come sempre, sentiti libero di commentare nella sezione sottostante e contattaci per qualsiasi dubbio, domanda o idea! Ci vediamo su **Proud2beCloud** tra un paio di settimane per un'altra storia emozionante!



Alessandro Gaggia

Head of software development di beSharp, Full-Stack developer, mi occupo di garantire lo stato dell'arte di tutta la nostra codebase. Scrivo codice in quasi ogni linguaggio, ma prediligo Typescript. Respiro Informatica, Game design, Cinema, Fumetti e buona cucina. Disegno per passione!



Matteo Moroni

DevOps e Solution Architect di beSharp, mi occupo di sviluppare soluzioni SaaS, Data Analysis, HPC e di progettare architetture non convenzionali a complessità divergente. Appassionato di informatica e fisica, da sempre lavoro nella prima e ho un PhD nella seconda. Parlare di tutto ciò che è tecnico e nerd mi rende felice!

Copyright © 2011-2021 by beSharp srl - P.IVA IT02415160189