

data.all - User Guide

v1.4.0

None

None

Table of contents

1. Introduction	3
1.1 What is data.all?	3
1.2 Why did we built data.all?	3
1.3 How can data.all help data teams?	3
2. Main components	4
2.1 Administrate	4
2.2 Discover	4
2.3 Play	5
3. Administrate	6
3.1 Tenant and Organizations	6
3.2 Environments and Teams	9
4. Discover	21
4.1 Datasets	21
4.2 Tables and Folders	28
4.3 Centralized Catalog and glossaries	35
4.4 Shares	39
5. Play	51
5.1 Worksheets	51
5.2 Notebooks	53
5.3 ML Studio	57
5.4 Pipelines	59
5.5 Dashboards	69
6. Security	73
6.1 Data and metadata on data.all	73
7. Platform Monitoring	74
7.1 Observability	74
7.2 Platform usage	74
8. Labs	79
8.1 Hands-on Lab: Data Access Management with data.all teams	79

1. Introduction

This section defines what is data.all, what is the challenge that it is trying to overcome and the value that it can bring to your teams.

1.1 What is data.all?

A modern data workspace that makes collaboration among diverse users (like business, analysts and engineers) easier, increasing efficiency and agility in data projects ✨

1.2 Why did we built data.all?

Data teams can be diverse: analysts, scientists, engineers, business users. Diverse people, with diverse tools and skillsets — diverse "DNAs". All leading to chaos and resulting in titanic efforts spent in **Collaboration Overhead**.

Using data.all, any line of business within an organization can create their own isolated data lake, produce, consume and share data within and across business units, worldwide. By simplifying data discovery, data access management while letting more builders use AWS vast portfolio of data and analytics services, data.all helps more data teams discover relevant data and let them use the power of the AWS cloud to create data driven applications faster.

1.3 How can data.all help data teams?

Teams can easily DISCOVER AND UNDERSTAND data 🌐

data.all makes all your datasets easily discoverable! No more Slack messages saying "Where's that dataset?" or long email threads for approvals. With data.all, you can simply browse the data catalog.

Key Capabilities: [Discovery and Search](#), [Data Preview & Worksheets](#) and [Notebooks](#)

Teams can easily SHARE AND COLLABORATE with data 💬

Data practitioners spend 30-50% of their time finding and understanding data. data.all cuts that time by 95%. Your data team will be shipping 2-3 times more projects in no time.

Key Capabilities: [Data Profiling & Data Sharing](#) and [Subscriptions](#)

Teams don't have to worry about SECURING their data 🛡️

Don't lose sleep trying to figure out if your sensitive data is secure. Build ecosystems of trust, make your team happy, and let data.all manage governance and security behind the scenes.

Key Capabilities: [Granular Access Control](#)

2. Main components

This section introduces the main components of data.all which are divided in 3 groups. This is an overview, for more details please refer to their specific sections.

- **Administristrate:** used by team and data lake administrators to organise and manage teams and users inside data.all
- **Discover:** used by all users to contribute with data, search for data and share data.
- **Play:** once data is in data.all, all users can use these tools to work with data.

2.1 Administrate

2.1.1 Organizations

[Organizations](#) are high level constructs where business units can collaborate across different AWS accounts at once. An organization includes environments (see below). Organizations are abstractions, they **don't** contain AWS resources, consequently there is no CloudFormation stack associated with them.

Organizations usually correspond to whole organizations, organization divisions or a separated geographical region within an organization.

2.1.2 Environments

An [environment](#) is a workplace where a team can bring, process, analyze data and build data driven applications. This workspace is mapped to an AWS account in one region. It is possible to have more than one environment in the same AWS Account, however we recommend to stick to one environment - one account.

An environment usually corresponds to a business unit or a department. Inside an environment we add teams and assign them different levels of permissions.

2.1.3 Teams

A [team](#) corresponds to an IdP group that has been onboarded to data.all. A special case for the administration of data.all is the **Tenant**, an IdP group with high level application (tenant) permissions. As with IdP groups, users can belong to multiple teams.

Teams corresponds to real teams.

but really, what are teams?

Data in data.all is isolated at team level, meaning that all members of a team can access all team's datasets. Thus, a team is any group of users that can access the team's datasets. We can have bigger teams with generic data and project-based teams owning data that requires more restrictive access to only members of the project.

2.2 Discover

2.2.1 Datasets

A [dataset](#) is a representation of multiple AWS resources that helps users store data. When data owners create a dataset on data.all the following resources are created:

- Amazon S3 Bucket to store the data on AWS.
- AWS KMS key to encrypt the data on AWS.
- AWS IAM role that gives access to the data on Amazon S3.
- AWS Glue database that is the representation of the structured data on AWS.

Inside the dataset we can store structured data as tables or unstructured data in folders.

2.2.2 Catalog

data.all centralized [Catalog](#) is an inventory of datasets, tables, folders and dashboards. It contains metadata for each of the mentioned data assets and thanks to its search capabilities, users can filter based on type of data, type of asset, tags, region and on glossary terms.

We use the Catalog to search and discover data

2.2.3 Glossaries

A [Glossary](#) is a list of terms, organized in a way to help users understand the context of their datasets. For example, terms like "cost", "revenue", etc, can be used to group and search all financial datasets.

Glossaries are used to add meaning to data assets metadata facilitating and enhancing Catalog searching

2.2.4 Shares

A [Share](#) is an access request to a data asset. Users search and discover data in the catalog and for those data assets that belong to other teams, users can create a Share on behalf of a team (remember, data access: at team level!!!). Then, the owners of the asset can accept or reject the share.

We use Shares to collaborate and share data with other teams.

2.3 Play

2.3.1 Worksheets

Worksheets are AWS Athena sessions that allow us to query our datasets as if we were in the AWS Athena Query editor console.

2.3.2 Notebooks

Data practitioners can experiment machine learning algorithms spinning up Jupyter notebook with access to all your datasets. data.all leverages [Amazon SageMaker instance](#) to access Jupyter notebooks.

2.3.3 ML Studio

With ML Studio Notebooks we can add users to our SageMaker domain and open Amazon SageMaker Studio

2.3.4 Pipelines

In order to distribute data processing, data.all introduces data.all pipelines where: - data.all takes care of CI/CD infrastructure - data.all offers flexible pipeline blueprints to deploy AWS resources and a Step Function

2.3.5 Dashboards

In the Dashboard window we can start Quicksight sessions, create visual analysis and dashboards.

3. Administrate

3.1 Tenant and Organizations

data.all manages teams' permissions at four levels:

1. Tenant team
2. Organization
3. Environment (next section)
4. Teams (next section)

3.1.1 Tenant

data.all has a super user's team which is a group from your IdP that has the right to manage high level application (tenant) permissions for all IdP groups integrated with data.all.

This super user's team maps to a group from your IdP that's by default named "**DAAdministrators**", any user member of this group will be able to:

- create organizations
- manage tenant permissions on onboarded teams (IdP groups) as shown below.

Manage tenant permissions

As a user part of "**DAAdministrators**" on your IdP you can access the settings menu from the profile icon.

For example, Maria Garcia is not part of "**DAAdministrators**", therefore she sees nothing



On the other hand, Tenant user is part of this group and can navigate to **Admin settings**



In *Admin Settings*, the Tenant user can manage tenant permissions. In the following picture, the user is NOT granting the *DataScienceTeam* that John belongs to permissions to create an organization.

The screenshot shows the 'Team DataScienceTeam' configuration page. At the top, it says 'A Team is a group from your identity provider that has access to data.all. Administrators can manage permissions for each team.' Below this, the 'Team' dropdown is set to 'DataScienceTeam'. The 'Tenant Permissions' section contains a list of actions:

- Manage datasets (enabled)
- Manage Redshift clusters (enabled)
- Manage dashboards (enabled)
- Manage notebooks (enabled)
- Manage pipelines (enabled)
- Manage worksheets (enabled)
- Manage glossaries (enabled)
- Manage environments (enabled)
- Manage organizations (disabled)
- Manage pipelines (enabled)

At the bottom right is a large orange 'Save' button.

If the tenant revokes the permission of a team to manage an object, that team won't be able to perform any action on that particular object. For the given example, assuming that John only belongs to the *DataScienceTeam*, he is not able to create organizations:

The screenshot shows the 'Create a new organization' page. The 'Organization Name' field is populated with 'example'. The 'Organize' section shows 'Team' set to 'DataScienceTeam'. An error message at the top right reads: 'An error occurred (UnauthorizedOperation) when calling MANAGE_ORGANIZATIONS operation: User: john Doe@amazon.com is not authorized to perform: MANAGE_ORGANIZATIONS on dataall.'

3.1.2 Organizations

Organizations are high level constructs where business units can collaborate across many different AWS accounts at once. An organization includes environments and teams (see next section). Organizations are abstractions, they **don't** contain AWS resources, consequently there is no CloudFormation stack associated with them.

Organizations usually correspond to whole organizations, organization divisions or a separated geographical region within an organization.

Create an organization

 **Organization permissions**

Any user can create an organization as long as he or she belongs to a group with tenant permission "Manage Organizations" (see previous chapter, "Manage tenant permissions").

To create an organization, on the left pane select **Organization**, click **Create** and complete the following form.

Create a new organization

Admin > Organizations > Create  Cancel

Details	Organize
Organization Name	Team
Short description	Tags
Create Organization	

Field	Description	Required	Editable	Example
Organization name	Name of the organization	Yes	Yes	AnyCompany EMEA
Short description	Short description about the organization	No	Yes	AnyCompany EMEA region
Team	Name of the team managing the organization	Yes	No	EMEAAdmin
Tags	List of tags	No	Yes	fin,rnd,mark,sales

The next step to onboard your IdP groups is to link an environment and add teams, check [Link an environment](#) and [Add a team to an environment](#)

Edit and update an organization

On the organisation window we can check the organization metadata, as well as the environments and teams that belong to this organisation (we will come back to this in [Environments and teams](#)).

To edit the metadata of the organisation, click in **Edit** and update the information. Name, description and tags are editable, however the organisation team cannot be updated.

Delete an organization

 **Warning**

Make sure that you delete the organisation environments before deleting the organisation. Otherwise, orphan environments might run into conflicts.

To archive an organization, click on the **Archive** button next to the Edit button. A window with the previous warning will appear. If you want to go ahead and delete the organization, type *permanently archive* in the box and submit.

3.2 Environments and Teams

An environment is a **workplace** where a team can bring, process, analyze data and build data driven applications. Environments comprise AWS resources, thus when we create an environment, we deploy a CDK/CloudFormation stack to an AWS account and region. In other words, **an environment is mapped to an AWS account in one region, where users store data and work with data.**

210

 One AWS account, One environment

To ensure correct data access and AWS resources isolation, onboard one environment in each AWS account. Despite being possible, **we strongly discourage users to use the same AWS account for multiple environments.**

3.2.1 Bootstrap your AWS account

data.all does not create AWS accounts. You need to provide an AWS account and complete the following bootstrapping steps on that AWS account in each region you want to use.

1. Create AWS IAM role

data.all assumes a IAM role named **PivotRole** to be able to call AWS SDK APIs on your account. You can download the AWS CloudFormation stack from *data.all* environment creation form. (Navigate to an organization and click on link an environment to see this form)

2. Setup AWS CDK

`data.all` uses AWS CDK to deploy and manage resources on your AWS account. AWS CDK requires some resources to exist on the AWS account, and provides a command called `bootstrap` to deploy these specific resources.

Moreover, we need to trust data.all infrastructure account. data.all codebase and CI/CD resources are in the data.all **tooling account**, while all the resources used by the platform are located in a **infrastructure account**. From this last one we will deploy environments and other resources inside each of our business accounts (the ones to be boostraped).

To bootstrap the AWS account using AWS CDK, you need :

1. to have AWS credentials configured in `~/.aws/credentials` or as environment variables.
 2. to install cdk : `npm install -g aws-cdk`
 3. to run the following command :

```
cdk bootstrap --trust DATA.ALL_AWS_ACCOUNT_NUMBER -c @aws-cdk/core:newStyleStackSynthesis=true --cloudformation-execution-policies arn:aws:iam::aws:policy/AdministratorAccess aws://YOUR ENVIRONMENT AWS ACCOUNT NUMBER/ENVIRONMENT REGION
```

 Which account should I put in the command?

Let's check with an example: the **tooling account** is 111111111111 and data.all was deployed to the **infrastructure account** = 222222222222. Now we want to onboard a **business account** = 333333333333 in region eu-west-1. Then the cdk bootstrap command will look like: `cdk bootstrap --region eu-west-1 --account 333333333333`

3. Enable AWS Lake Formation

data.all relies on AWS Lake Formation to manage access to your structured data. If AWS Lake Formation has never been activated on your AWS account, you need to create a service-linked role, using the following command:

```
aws iam create-service-linked-role --aws-service-name lakeformation.amazonaws.com
```

**Service link creation error**

If you receive: An error occurred (InvalidInput) when calling the CreateServiceLinkedRole operation: Service role name AWSServiceRoleForLakeFormationDataAccess has been taken in this account, please try a different suffix. **You can skip this step, as this indicates the Lake formation service-linked role exists.**

4. Amazon Quicksight

This is an optional step. To link environments with **Dashboards enabled**, you will also need a running Amazon QuickSight subscription on the bootstrapped account. If you have not subscribed to Quicksight before, go to your AWS account and choose the Enterprise option as show below:

The screenshot shows a user interface for signing up for QuickSight. On the left, there is a blue icon depicting a computer monitor displaying a network graph. Next to it, the text reads: "Your AWS Account is not signed up for QuickSight. Would you like to sign up now?". Below this, the text "AWS Account" is followed by a redacted email address. At the bottom, a blue button labeled "Sign up for QuickSight" is visible. A note at the bottom states: "To access QuickSight with a different account, [log in again](#).

Create your QuickSight account		
Edition	<input checked="" type="radio"/> Enterprise	<input type="radio"/> Enterprise + Q Learn more
Team trial for 30 days (4 authors)*	FREE	FREE
Author per month (yearly)**	\$18	\$28
Author per month (monthly)**	\$24	\$34
Readers (pay-per-Session)	\$0.30 / session (max \$5)****	\$0.30 / session (max \$10)****
Additional SPICE per month	\$0.38 per GB	\$0.38 per GB
QuickSight Q regional fee	N/A	\$250 / mo / region
Personalized Q authoring workshop	N/A	Starting from \$199
Natural language query with QuickSight Q	N/A	INCLUDED
Single Sign On with SAML or OpenID Connect	✓	✓
Connect to spreadsheets, databases & business apps	✓	✓
Access data in Private VPCs	✓	✓
Row-level security for dashboards	✓	✓
Secure data encryption at rest	✓	✓
Connect to your Active Directory	✓	✓
Use Active Directory groups***	✓	✓
Send email reports	✓	✓
Embed QuickSight	✓	✓
Capacity-based pricing	✓	✓
Supported regions	Learn more	Learn more

After you've successfully subscribed to QuickSight, we need to trust *data.all* domain on QuickSight to enable Dashboard Embedding on *data.all* UI. To do that go to:

1. Manage QuickSight
2. Domains and Embedding
3. Put *data.all* domain and check include subdomains
4. Save

Account name: dataall-dataservice
Edition: Enterprise

Manage users
Manage groups
Your subscriptions
SPICE capacity
Account settings
Security & permissions
Manage VPC connections
Mobile settings
[Domains and Embedding](#)
Account customization
Single sign-on (SSO)

Manage domains for embedded dashboards
Embedded dashboards only work if they originate from domains you explicitly allow.

Domain: https://example.cloudfront.net/
 Include subdomains ①

Add

Domain	Include subdomains
No domains have been allow listed for embedding	

Username: Admin/dlpzx-lsengard
Community
Send feedback
English >
Ireland >
Tutorial videos
Help
Sign out

3.2.2 NEW Link an environment

Necessary permissions

Environment permissions

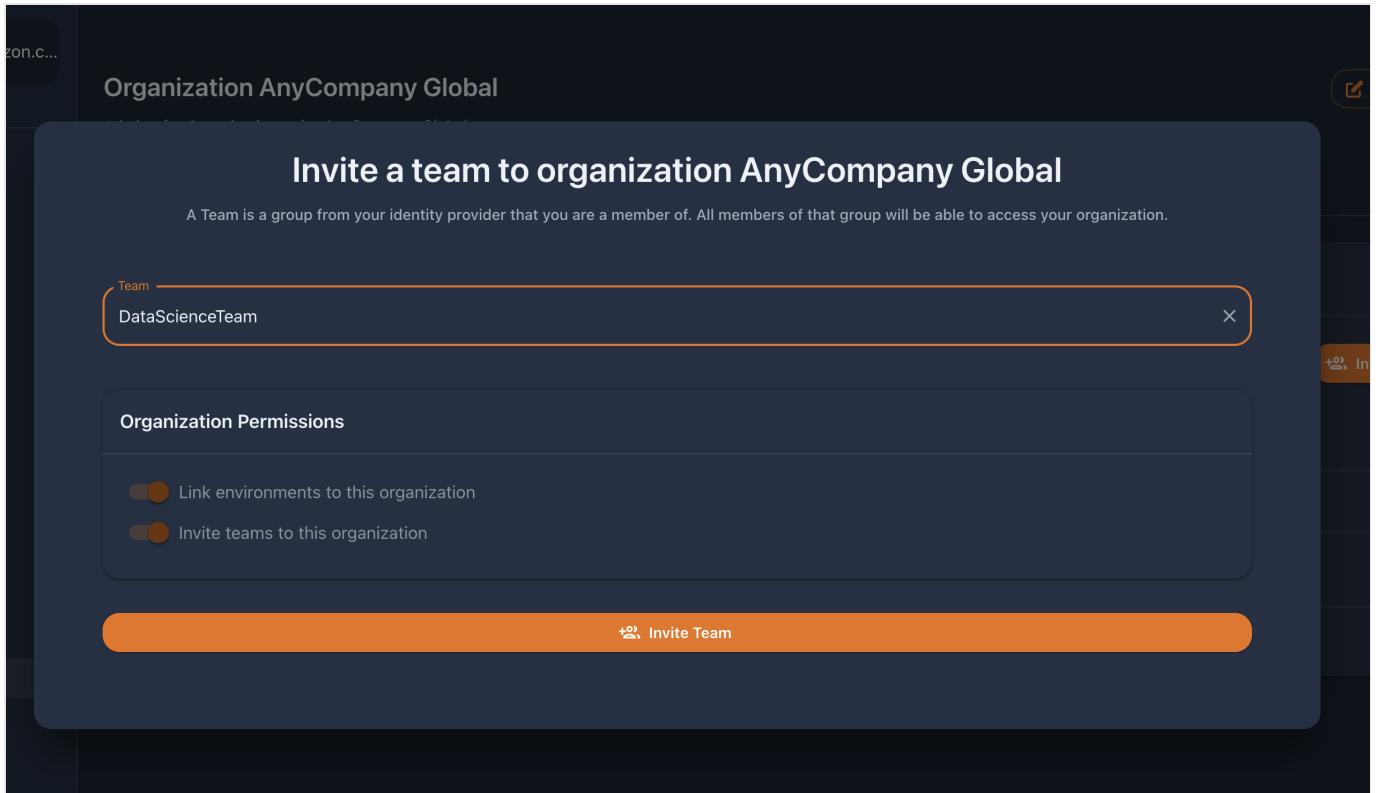
Only organization Administrator teams can link environments to the Organization. The Organization creator team is the by default Organization Administrator team, but users of this group can now invite other teams and grant them permission to manage organization teams, and link environment to the organization.

Managing organization teams can be done through the UI or APIs. From the UI, navigate to your organizations and click on the **Teams** tab.

OVERVIEW ENVIRONMENTS TEAMS

Name	Permissions	Actions
DAAAdministrators <small>ADMINS</small>	<small>ALL</small>	
DataAnalysisTeam		
MarketingTeam		
ResearchTeam		

Invite button opens a dialog that gives the organization creators the possibility to invite one of the IdP groups they belong to, which will appear in a dropdown when we click on **Teams**. They can also invite an IdP group that they don't belong to, as long as they type the exact group name (**case sensitive**):



You can check the Organization administrators teams in the Organization's **Teams** tabs and remove a team if necessary on the icon in the Actions column.

Name	Permissions	Actions
DHAdministrators ADMIN	ALL	
DataAnalyticsTeam	≡	⋮
DataScienceTeam	≡	⋮

Link environment

Once the AWS account/region is bootstraped and we have permission to link an environment to an organization, let's go! Navigate to your organization, click on the **Link Environment** button, and fill the environment creation form:

Field	Description	Required	Editable	Example
Environment name	Name of the environment	Yes	Yes	Finance
Short description	Short description about the environment	No	Yes	Finance department teams
Account number	AWS bootstraped account maped to the environment	Yes	No	111111111111
Region	AWS region	Yes	No	Europe (Ireland)
IAM Role name	Alternative name of the environment IAM role	No	No	anotherRoleName
Resources prefix	Prefix for all AWS resources created in this environment. Only (^[a-z-]*\$)	Yes	Yes	fin
Team	Name of the group initially assigned to this environment	Yes	No	FinancesAdmin
Tags	Tags that can later be used in the Catalog	Yes	Yes	finance, test
VPC Identifier	VPC provided to host the environment resources instead than the default one created by <i>data.all</i>	No	No	vpc-.....
Public subnets	Public subnets provided to host the environment resources instead than the default created by <i>data.all</i>	No	No	subnet-....
Private subnets	Private subnets provided to host the environment resources instead than the default created by <i>data.all</i>	No	No	subnet-....

Features Management

An environment is defined as a workspace and in this workspace we can flexibly activate or deactivate different features, adapting the workspace to the teams' needs. If you want to use Dashboards, you need to complete the optional fourth step explained in the previous chapter "Bootstrap your AWS account".

This is not set in stone!

Don't worry if you change your mind, features are editable. You can always update the environment to enable or disable a feature.

Click on Save, the new Environment should be displayed in the Environments section of the left side pane.

3.2.3 Manage your Environment

Go to the environment you want to check. You can find your environment in the Environments list clicking on the left side pane or by navigating to the environment organization. There are several tabs just below the environment name:

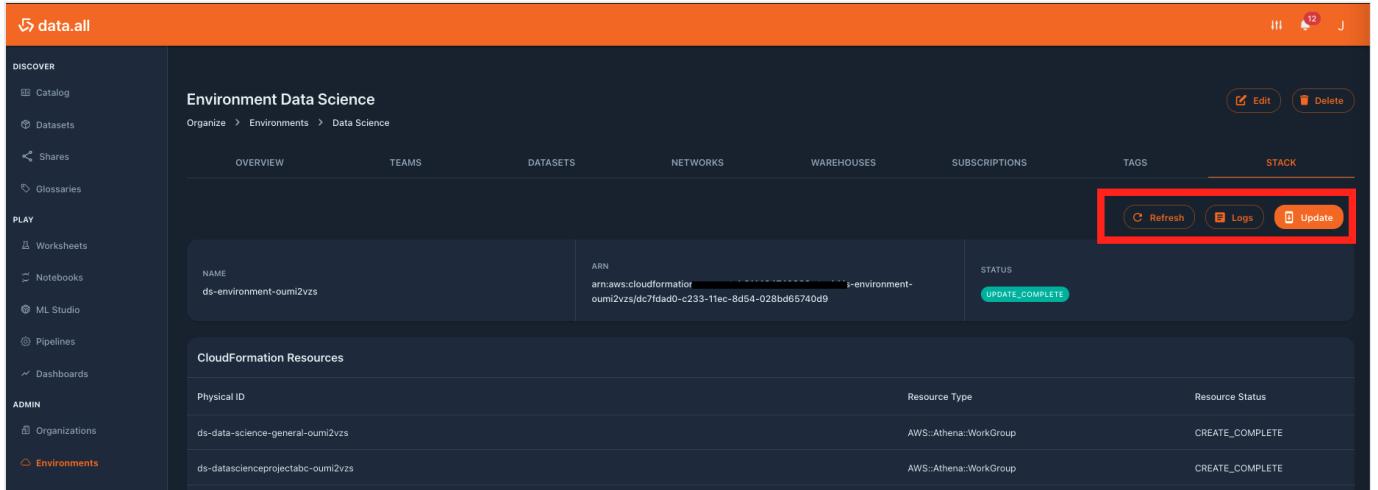
- Overview: summary of environment information and AWS console and credential access.
- Teams: list of all teams onboarded to this environment.
- Datasets: list of all datasets owned and shared with for this environment
- Networks: VPCs created and owned by the environment
- Warehouses: Redshift clusters imported or created in this environment
- Subscriptions: SNS topic subscriptions enabled or disabled in the environment
- Tags: editable key-value tags
- Stack: CloudFormation stack details and logs

Environment access

If **none** of the teams you belong to (IdP groups) has been onboarded to the environment, you won't be able to see the environment in the environments menu or in the organization environments list. **Check the "Manage teams" section**

Check CloudFormation stack

After linking an environment we can check the deployment of AWS resources in CloudFormation, click on the environment and then on the **Stack** tab. Right after linking an environment you should find something like the below picture.



Physical ID	Resource Type	Resource Status
ds-data-science-general-oumi2vzs	AWS::Athena::WorkGroup	CREATE_COMPLETE
ds-datasciencuprojectabc-oumi2vzs	AWS::Athena::WorkGroup	CREATE_COMPLETE

After some minutes its status should go from "PENDING" to "CREATE_COMPLETE" and we will be able to look up the AWS resources created as part of the environment CloudFormation stack. Moreover, we can manually trigger the update in case of change sets of the CloudFormation stack with the **Update** button.

Pro Tip

If something in the creation or update of an environment fails, we can directly check the logs by clicking the logs button. No need to navigate to the AWS console to find your logs!

After being processed (not in `PENDING`), the status of the CloudFormation stack is directly read from [CloudFormation](#).

Edit and update an environment

Find your environment in the Environments list or by navigating to the corresponding organization. Once in your selected environment, click on **Edit** in the top-right corner of the window and make all the changes you want.

Finally, click on **Save** at the bottom-right side of the page to update the environment.



Automatically updates the CloudFormation stack

Clicking on Save will update the environment metadata as well as the CloudFormation stack on the AWS account

Delete an environment

In the chosen environment, next to the Edit button, click on the **Delete** button.



Orphan *data.all* resources

A message like this one: "*Remove all environment related objects before proceeding with the deletion!*" appears in the delete display. Don't ignore it! Before deleting an environment, clean it up: delete its datasets and other resources.

Note that we can keep the environment CloudFormation stack. What is this for? This is useful in case you want to keep using the environment resources (IAM roles, etc) created by *data.all* but outside of *data.all*

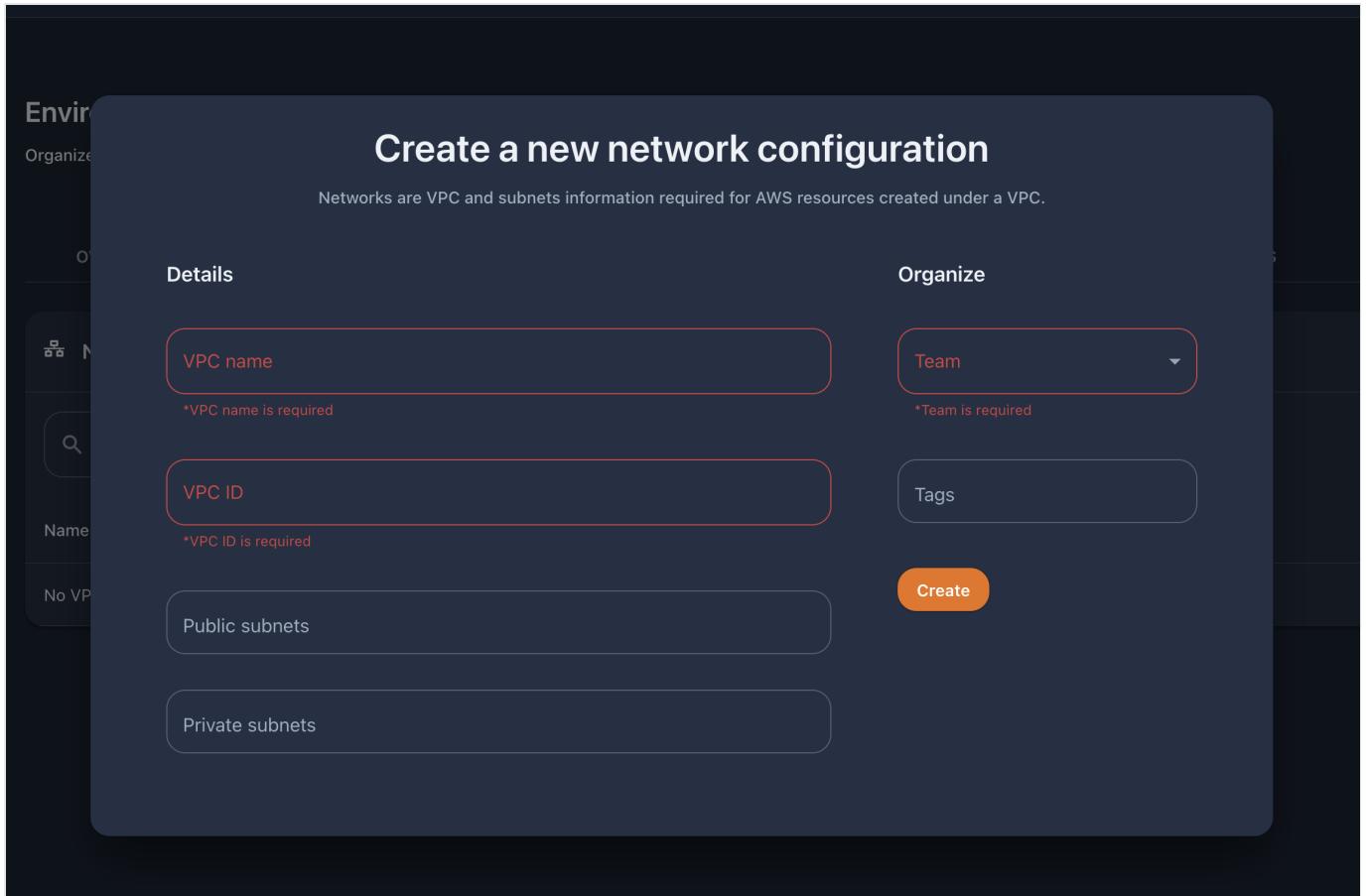
Create networks

Networks are VPCs created from *data.all* and belonging to an environment and team. To create a network, click in the **Networks** tab in the environment window, then click on **Add** and finally fill the following form.



I need an example!

What is the advantage of using networks from *data.all*?[MISSING INFO]



Create Key-value tags

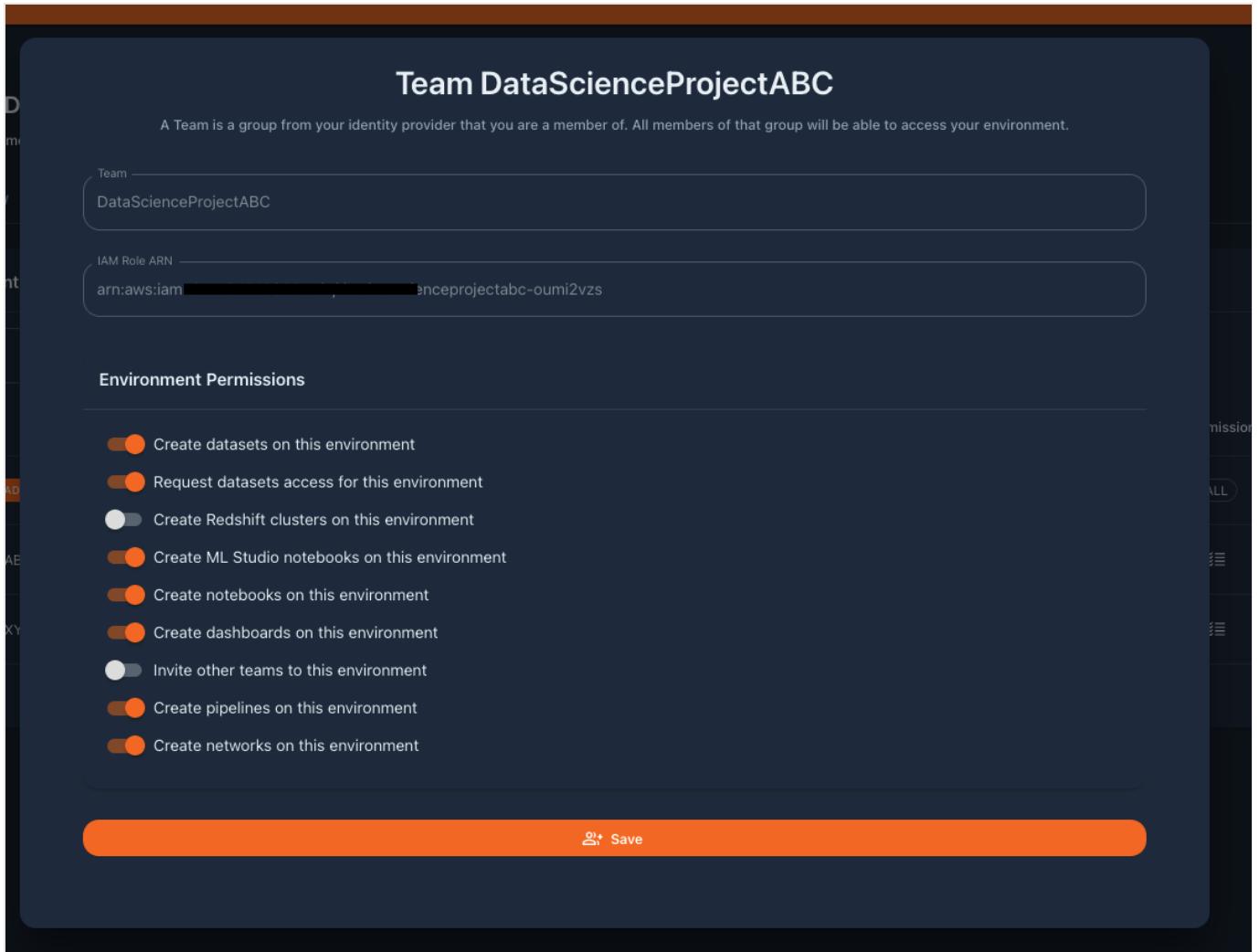
In the **Tags** tab of the environment window, we can create key-value tags. These tags are not *data.all* tags that are used to tag datasets and find them in the catalog. In this case we are creating AWS tags as part of the environment CloudFormation stack. There are multiple tagging strategies as explained in the [documentation](#).

3.2.4 Manage Teams

Environment creators have all permissions on the environment, and can invite other teams to the onboarded environment. To add an IdP group to an environment, navigate to the **Teams** tab of the environment and click on the **Invite** button.

The screenshot shows the "Environment Data Science" page with the "TEAMS" tab selected. At the top, there are buttons for "Edit" and "Delete". Below the tabs are buttons for "OVERVIEW", "DATASETS", "NETWORKS", "WAREHOUSES", "SUBSCRIPTIONS", "TAGS", and "STACK". A search bar and an "Invite" button are located at the bottom left. The main area displays a table with columns for "Name", "IAM Role", "Athena WorkGroup", "Permissions", and "Actions".

A display will allow you to customize the AWS permissions that the onboarded group will have, adapting to different types of users (data scientists, data engineers, data analysts, management). The customizable permissions can be enabled or disabled as appears in the following picture.



When the invitation is saved, the environment CloudFormation stack gets automatically updated and creates a new IAM role for the new team. The IAM role policies mapped to the permissions granted to the invited team (e.g., a team invited without "Create Redshift clusters" permission will not have redshift permissions on the associated IAM role). To remove a group, in the *Actions* column select the minus icon.

Automated permission assignment

Groups retrieved from the IdP are automatically granted all application high level permissions by default to accelerate the onboarding process.

Users will only be able to see the environments where a team that they belong to has been onboarded (either as creator of the environment or invited to the environment). In the following picture, John belongs to the *DataScienceTeam* that owns the *Data Science* environment, but on top of that he can access the *Data Analysis* environment because her team has been invited by Maria.

Pro tip!

You know whether you are `OWNER` or `INVITED` in an environment by checking **your Role** in that environment. This information appears in the picture in each environment box in the field "Role".

The screenshot shows the data.all interface with the 'Environments' tab selected. On the left, there's a sidebar with sections like 'Discover', 'Play', and 'Admin'. The main area displays two environments: 'Data Science' and 'Data Analysis'. Each environment card includes a thumbnail, owner information, a description field (which is empty), and a status section. For 'Data Science', the owner is 'johndoe@amazon.com' and the status is 'UPDATE_COMPLETE'. For 'Data Analysis', the owner is 'mariagarcia@amazon.com' and the status is also 'UPDATE_COMPLETE'. There are 'Learn More' buttons at the bottom of each card.

Difference between invited and owner

A team that has been invited to an environment has slight limitations, because well, it is not their environment! Invited teams cannot access the **Stack** tab of the environment because they should not be handling the resources of the environment. Same applies for **Tags** and **Subscriptions**. Other limitations come from the permissions that have been assigned to the team.

AWS access

data.all makes it easier to manage access to your AWS accounts. How? remember when we assigned granular AWS permissions to invited groups and this created an IAM role? (If not, check the "Manage teams" section). From the **Teams** tab of the environment we can assume our team's IAM role to get access to the AWS Console or copy the credentials to the clipboard. Both options are under the "Actions" column in the Teams table.

3.2.5 Manage Consumption Roles

Data.all creates or imports one IAM role per Cognito/IdP group that we invite to the environment. With these IAM roles data producers and consumers can ingest and consume data, but sometimes we want to consume data from an application such as SageMaker pipelines, Glue Jobs or any other downstream application. To increase the flexibility in the data consumption patterns, *data.all* introduces Consumption Roles.

Any IAM role that exists in the Environment AWS Account can be added to *data.all*. In the **Teams** tab click on *Add Consumption Role*

The screenshot shows the Data.all interface with the 'TEAMS' tab selected. The top navigation bar includes 'OVERVIEW', 'TEAMS' (selected), 'DATASETS', 'NETWORKS', 'SUBSCRIPTIONS', 'TAGS', and 'STACK'. Below the tabs, there are two main sections:

- Environment Teams**: Shows a table with columns: Name, IAM Role, Athena WorkGroup, Permissions, and Actions. One row is visible: SBResearch (ADMIN) with IAM Role arn:aws:iam::, Athena WorkGroup dataall-research-dev-2-y1q4tz5y, Permissions ALL, and Actions with AWS Lambda and CloudWatch Metrics icons.
- Environment Consumption IAM roles**: Shows a table with columns: Name, IAM Role, Role Owner, and Action. A search bar and an 'Add Consumption Role' button are at the top. The message 'No Consumption IAM Role added' is displayed below the table.

A window like the following will appear for you to introduce the arn of the IAM role and the Team that owns the consumption role. Only members of this team and tenants of data.all can remove the consumption role.

This is a modal dialog titled 'Add a consumption IAM role to environment Research-DEV-2'. It contains the following fields:

- Consumption Role Name: SageMaker DS role
- IAM Role ARN: arn:aws:iam::role/RoleSagemakerStudioUsers
- Owners: SBMarketingProjectABC

At the bottom is a large orange button labeled 'Add Consumption Role'.

Existing roles only

Data.all checks whether that IAM role exists in the AWS account of the environment before adding it as a consumption role.

Data Access

- By default, a new consumption role does NOT have access to any data in data.all.
- The team that owns the consumption role needs to open a share request for the consumption role as shown in the picture below.

4. Discover

4.1 Datasets

4.1.1 Datasets

In data.all, a Dataset is a representation of multiple AWS resources that helps users store data and establish the basis to make this data discoverable and shareable with other teams.

When data owners create a dataset the following resources are deployed on the selected environment and its linked AWS account:

1. Amazon S3 Bucket to store the data on AWS.
2. AWS KMS key to encrypt the data on AWS.
3. AWS IAM role that gives access to the data on Amazon S3.
4. AWS Glue database that is the representation of the structured data on AWS.

AWS Champion

data.all does all the infrastructure heavy lifting for data owners using **AWS CDK** and **AWS CloudFormation** service while following AWS deployment and security best practices.

Tables and Folders

Inside a dataset we can store structured data in tables and unstructured data in folders.

- Tables are the representation of **AWS Glue Catalog** tables that are created on the dataset's Glue database on AWS.
- Folders are the representation of an **Amazon S3 prefix** where data owners can organize their data. For example, when data is loaded, it can go to a folder named "raw" then after it's processed the data moves to a folder called "silver" and so on.

Dataset ownership

Dataset ownership refers to the ability to access, modify or remove data from a dataset, but also to the responsibility of assigning these privileges to others.

- **Owners:** When you create a dataset and associate it with a team, the dataset business ownership belongs to the associated team.
- **Stewards:** You can delegate the stewardship of a dataset to a team of stewards. You can type a name of an IdP group or choose one of the teams of your environment to be the dataset stewards.

Note

Dataset owners team is a required, non-editable field, while stewards are optional and can be added post the dataset has been created. If no other stewards team is designated, the dataset owner team will be the only responsible in managing access to the dataset.

Dataset access

In this case we are referring to the ability to access, modify or remove data from a dataset. Who can access the dataset content? users belonging to...

- the dataset owner team
- a dataset steward team
- teams with a share request approved to dataset content

Note

Dataset metadata is available for all users in the centralized data catalog.

4.1.2 NEW Create a dataset

On left pane choose **Datasets**, then click on the **Create** button. Fill the dataset form.

The screenshot shows the 'Create a new dataset' page in the data.all interface. The left sidebar has sections for DISCOVER (Catalog, Datasets), PLAY (Shares, Glossaries, Worksheets, Notebooks, ML Studio, Pipelines, Dashboards), and ADMIN (Organizations, Environments). The main area has tabs for 'Discover' and 'Datasets'. The 'Datasets' tab is selected. The 'Create a new dataset' form contains the following fields:

- Details**: Dataset name (text input), Short description (text area, 200 characters left).
- Deployment**: Environment (dropdown), Region (dropdown), Organization (dropdown).
- Classification**: Confidentiality (dropdown), Topics (dropdown), Tags (dropdown).
- Governance**: Owners (dropdown), Stewards (dropdown).
- Summary**: A table with columns: Field, Description, Required, Editable, Example.

A large orange 'Create Dataset' button is located at the bottom right of the form.

Field	Description	Required	Editable	Example
Dataset name	Name of the dataset	Yes	Yes	AnyDataset
Short description	Short description about the dataset	No	Yes	For AnyProject predictive model
Environment	Environment (mapped to an AWS account)	Yes	No	DataScience
Region (auto-filled)	AWS region of the environment	Yes	No	Europe (Ireland)
Organization (auto-filled)	Organization of the environment	Yes	No	AnyCompany EMEA
Owners	Team that owns the dataset	Yes	No	DataScienceTeam
Stewards	Team that can manage share requests on behalf of owners	No	Yes	FinanceBITeam, FinanceMgmtTeam
Confidentiality	Level of confidentiality: Unclassified, Official or Secret	Yes	Yes	Secret
Topics	Topics that can later be used in the Catalog	Yes, at least 1	Yes	Finance
Tags	Tags that can later be used in the Catalog	Yes, at least 1	Yes	deleteme, ds

4.1.3 Import a dataset

If you already have data stored on Amazon S3 buckets, data.all got you covered with the import feature. In addition to the fields of a newly created dataset you have to specify the S3 bucket and optionally a Glue database:

Field	Description	Required	Editable	Example
Amazon S3 bucket name	Name of the S3 bucket you want to import	Yes	No	importedBucket
AWS Glue database name	Name of the Glue database that you want to import	No	No	anyDatabase

The screenshot shows the 'Import a new dataset' interface. On the left is a sidebar with 'Discover' (Catalog, Datasets, Shares, Glossaries), 'PLAY' (Worksheets, Notebooks, ML Studio, Pipelines, Dashboards), and 'ADMIN' (Organizations, Environments). A 'User Guide' button is also present. The main area has tabs for 'Details', 'Classification', 'Deployment', and 'Governance'. Under 'Details', there are fields for 'Dataset name' and 'Short description'. Under 'Classification', there are dropdowns for 'Confidentiality', 'Topics', and 'Tags'. Under 'Deployment', there are dropdowns for 'Environment', 'Region', and 'Organization'. Under 'Governance', there are dropdowns for 'Team' and 'Stewards'. At the bottom right is an orange 'Import Dataset' button.

4.1.4 Navigate dataset tabs

When we belong to the dataset owner team

After creating or importing a dataset it will appear in the datasets list (click on Datasets on the left side pane). In this window, it will only be visible for those users belonging to the dataset owner team. If we select one of our datasets we will see the following dataset window:

The screenshot shows the 'Dataset January' window. At the top, there are buttons for 'Chat', 'AWS Credentials', 'S3 Bucket', 'Edit', and 'Delete'. Below is a navigation bar with tabs: 'DATA' (selected), 'OVERVIEW' (highlighted in orange), 'SHARES', 'UPLOAD', 'TAGS', and 'STACK'. The main content area is currently empty.

When we DON'T belong to the dataset owner team

How do we access a dataset if we don't have access to it? IN THE CATALOG! on the left pane click on Catalog, find the dataset you are interested in, click on it and if you don't have access to it, you should see only some of the tabs in comparison with the previous pic, something like:

The screenshot shows the 'Dataset Insights_1' window. At the top, there are buttons for 'Chat', 'AWS Credentials', 'S3 Bucket', 'Edit', and 'Delete'. Below is a navigation bar with tabs: 'DATA' (selected) and 'OVERVIEW' (highlighted in orange). The main content area is currently empty.

4.1.5 Edit and update a dataset

Data owners can edit the dataset by clicking on the **edit** button, editing the editable fields and saving the changes.

4.1.6 Delete a dataset

To delete a dataset, in the selected dataset window click on the **delete** button in the top-right corner. As with environments, it is possible to keep the AWS CloudFormation stack to keep working with the data and resources created but outside of data.all.

4.1.7 Check dataset info and access AWS

The **Overview** tab of the dataset window contains dataset metadata, including governance and creation details. Moreover, AWS information related to the resources created by the dataset CloudFormation stack can be consulted here: AWS Account, Dataset S3 bucket, Glue database, IAM role and KMS Alias.

You can also assume this IAM role to access the S3 bucket in the AWS console by clicking on the **S3 bucket** button. Alternatively, click on **AWS Credentials** to obtain programmatic access to the S3 bucket.

The screenshot shows the AWS Studio interface for a dataset named 'January'. The top navigation bar includes tabs for DATA, OVERVIEW (which is selected), SHARES, UPLOAD, TAGS, and STACK. The left sidebar has sections for Details, Governance & Classification, Topics, Tags, and Glossary terms. The main content area is divided into several sections: Details (URI: od779vcv, Name: January, Description: No description provided); Governance & Classification (Owners: DataScienceTeam, Stewards: DataScienceTeam, Classification: UNCLASSIFIED, Topics: Operations, Tags: prod); AWS Information (Account: [REDACTED], S3 bucket: arn:aws:s3:::ds-january-od779vcv, Glue database: arn:aws:glue:eu-west-1:[REDACTED]:database:ds_january_od779vcv, IAM role: arn:aws:iam:[REDACTED]:role/ds-january-od779vcv, KMS alias: arn:aws:kms:eu-west-1:[REDACTED]:alias:ds-january-od779vcv); and a bottom section with a status indicator 'UPDATE_COMPLETE'.

4.1.8 Fill the dataset with data

Tables

Quickly upload a file for data exploration

Users may want to experiment with a small set of data (e.g. a csv file). To create tables from a file, we first upload the file, then run the crawler to infer its schema, and finally, we read the schema by synchronizing the table. Upload & Crawl & Sync

1. Upload data: Go to the **Upload** tab of the dataset and browse or drop your sample file. It will be uploaded to the dataset S3 bucket in the prefix specified. By default, a Glue crawler will be triggered by the upload of a file, however this feature can be disabled as appears in the picture.

S3 Upload

Prefix: s3://ds-january-od779vcv/

Infer Schema: Enabling this will automatically start a crawler to infer your file schema.

Select file: Drop file [browse](#) through your machine

bestsellers_with_categories_2022_03_27.csv
64.01 KB

Remove All [Upload](#)

- Crawl data: the file has been uploaded but the table and its schema have not been registered in the dataset Glue Catalog database. If you have disabled the crawler in the upload, click on the **Start Crawler** button in the Data tab. If you just want to crawl one prefix, you can specify it in the Start Crawler feature.

Dataset January

Contribute > Datasets > January

DATA OVERVIEW SHARES UPLOAD TAGS STACK

Tables

Synchronize Start Crawler

Name	Database	Location	Actions
raw	ds_january_od779vcv	s3://ds-january-od779vcv/raw/	View Edit
videogames_sales	ds_january_od779vcv	s3://ds-january-od779vcv/videogames_sales/	View Edit
supermarket_sales	ds_january_od779vcv	s3://ds-january-od779vcv/supermarket_sales/	View Edit

Folders

+ Create

Name	S3 Location	Description	Actions
january-sales-pdfs	s3://ds-january-od779vcv/pdfs	PDF prints of sales reports	View Edit

- Synchronize tables: Once crawled and registered in the Glue database, you can synchronize tables from your dataset's AWS Glue database by using **Synchronize tables** feature in the Data tab. In any case, data.all will synchronize automatically the tables for you at a frequency of **15 minutes**.

You can preview your small set of data right away from data.all, check [Tables](#).

Ingest data

If you need to ingest larger quantities of data, manage bigger files, or simply you cannot work with local files that can be uploaded; this is your section!

There are multiple ways of filling our datasets with data and actually, the steps don't differ much from the upload-crawl-sync example.

- Crawl & Sync option: we can drop the data from the source to our dataset S3 bucket. Then, we will crawl and synchronize data as we did in the previous steps 2 and 3.
- Register & Sync option: we drop the data from the source to our dataset S3 bucket. However, if we want to have more control over our tables and its schema, instead of starting the crawler we can **register the tables** in the Glue Catalog and then click on Synchronize as we did in step 3.

How do we register Glue tables? There are numerous ways:

- manually from the [AWS Glue console](#) in your environment account
- Using [AWS Glue API](#), `CreateTable`.
- In a Glue Job leveraging Glue [PySpark DynamicFrame](#) class
- With [boto3](#)
- Or with [AWS Data Wrangler](#), Pandas on AWS.
- Also, you can deploy Glue resources using [CloudFormation](#)
- Or directly, [migrating from Hive Metastore](#).
- there are more for sure :)

Use data.all pipelines to register Glue tables

data.all pipelines can be used to transform your data including the creation of Glue tables using [AWS Glue pyspark extension](#). Visit the [pipelines](#) section for more details.

Folders

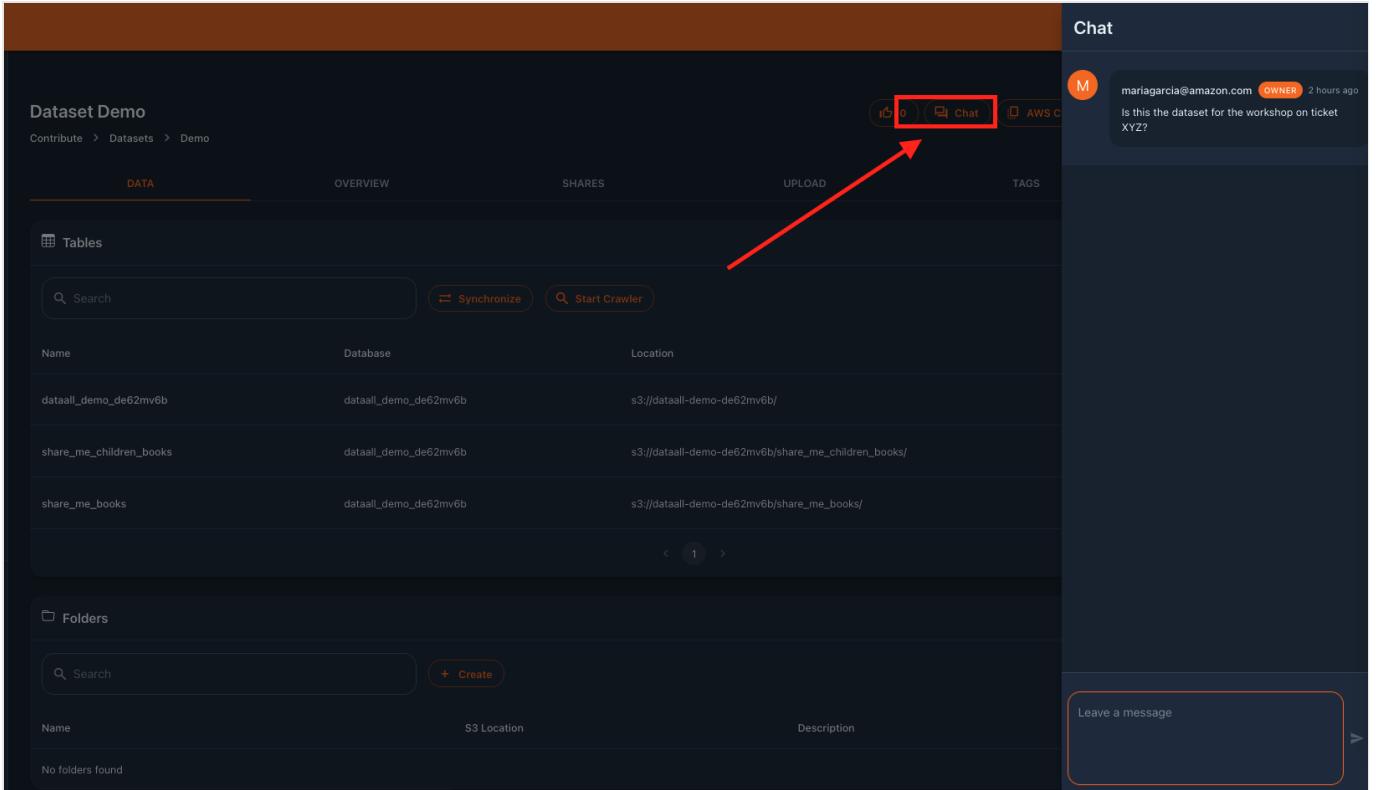
As previously defined, folders are prefixes inside our dataset S3 bucket. To create a folder, go to the **Data** tab and on the folders section, click on Create. The following form will appear. We will dive deeper in how to use folders in the [folders section](#).

The screenshot shows the AWS Glue Data Studio interface. The top navigation bar has tabs for DATA, OVERVIEW, SHARES, UPLOAD, and TAGS. The DATA tab is selected. On the left, there's a sidebar with sections for Tables and Folders. Under Tables, there's a search bar and a list of datasets: ds_january_od779vcv, raw, videogames_sales, supermarket_sales. Under Folders, there's a search bar and a list of folders: january-sales-pdfs. The main area is titled 'Create a new folder' and includes fields for 'Folder name' (with placeholder 'ds_january_od779vcv'), 'Amazon S3 prefix' (with placeholder 's3://ds-january-od779vcv/'), and 'Short description' (with placeholder 'PDF prints of sales reports'). A note above the fields says 'Creates an Amazon S3 prefix under the dataset bucket'. At the bottom of the dialog is an orange 'Create folder' button. Below the dialog, the list of existing datasets is shown again with columns for Name, S3 Location, Description, and Actions.

Name	S3 Location	Description	Actions
january-sales-pdfs	s3://ds-january-od779vcv/pdfs	PDF prints of sales reports	□ →

4.1.9 Leave a message in Chat

In the **Chats** button users can interact and leave their comments and questions on the Dataset Chat.



The screenshot shows the AWS DataAll interface for a dataset named "Dataset Demo". On the left, there's a sidebar with "Tables" and "Folders" sections. The main area displays a list of tables with columns for Name, Database, and Location. A red arrow points from the "Chat" button in the top navigation bar to a message card on the right. The message card shows a message from "mariagarcia@amazon.com" (OWNER) posted 2 hours ago, asking if this is the dataset for the workshop on ticket XYZ? Below the message card is a "Leave a message" input field.

Name	Database	Location
dataall_demo_de62mv6b	dataall_demo_de62mv6b	s3://dataall-demo-de62mv6b/
share_me_children_books	dataall_demo_de62mv6b	s3://dataall-demo-de62mv6b/share_me_children_books/
share_me_books	dataall_demo_de62mv6b	s3://dataall-demo-de62mv6b/share_me_books/

4.1.10 Create key-value tags

Same as in environments. In the **Tags** tab of the dataset window, we can create key-value tags. These tags are not data.all tags that are used to tag the dataset and find it in the catalog. In this case we are creating AWS tags as part of the dataset CloudFormation stack. There are multiple tagging strategies as explained in the [documentation](#).

4.2 Tables and Folders

4.2.1 Tables

In this section we will go through the different tabs in the Table window. We can reach this view:

1. by selecting a table from the data Catalog
2. or in the dataset view, in the **Tables** tab clicking on the arrow in the *Actions* column for the chosen table.

Name	Database	Location	Actions
raw	ds_january_od779vcv	s3://ds-january-od779vcv/raw/	Open →
videogames_sales	ds_january_od779vcv	s3://ds-january-od779vcv/videogames_sales/	Open →
supermarket_sales	ds_january_od779vcv	s3://ds-january-od779vcv/supermarket_sales/	Open →

Check table metadata

Also in the table window, go to the **Overview** tab where you will find the following information:

- URI: unique table identifier
- Name: name of the registered table in the Glue Catalog
- Tags
- Glossary terms
- Description
- Organization, Environment, Region, Team: inherited from the dataset
- Created: creation time of the table
- Status: INSYNC

Description, Tags and Glossary terms are not inherited!

If a dataset is tagged with Tags and Glossary terms, the child tables do not inherit these tags and terms. In the Overview tab, by clicking on **Edit** is where you can add them. Same applies for the description. Adding tags and terms to your tables will make them more discoverable in the Catalog.

Add or edit table metadata

Edit your table metadata by clicking on the **Edit** button.

Preview data

Data preview gives you the ability to preview a subset of the data available on data.all. Preview feature is available for data you own or data that's shared with you.

Just select a table and in the **Preview** tab you will find the results of an SQL select subset of the table.

Table supermarket_sales															 Chat	 Edit	 Delete
PREVIEW				OVERVIEW					COLUMNS					METRICS			
invoice id	branch	city	customer t...	gender	product line	unit price	quantity	tax %	total	date	time	payment	cogs	gross margi...	gross income	rating	
750-67-8...	A	Yangon	Member	Female	Health an...	74.69	7	26.1415	548.9715	1/5/2019	13:08	Ewallet	522.83	4.7619047...	26.1415	9.1	
226-31-3...	C	Naypyitaw	Normal	Female	Electronic ...	15.28	5	3.82	80.22	3/8/2019	10:29	Cash	76.4	4.7619047...	3.82	9.6	
631-41-31...	A	Yangon	Normal	Male	Home and...	46.33	7	16.2155	340.5255	3/3/2019	13:23	Credit card	324.31	4.7619047...	16.2155	7.4	
123-19-11...	A	Yangon	Member	Male	Health an...	58.22	8	23.288	489.048	1/27/2019	20:33	Ewallet	465.76	4.7619047...	23.288	8.4	
373-73-7...	A	Yangon	Normal	Male	Sports an...	86.31	7	30.2085	634.3785	2/8/2019	10:37	Ewallet	604.17	4.7619047...	30.2085	5.3	
699-14-3...	C	Naypyitaw	Normal	Male	Electronic ...	85.39	7	29.8865	627.6165	3/25/2019	18:30	Ewallet	597.73	4.7619047...	29.8865	4.1	
355-53-5...	A	Yangon	Member	Female	Electronic ...	68.84	6	20.652	433.692	2/25/2019	14:36	Ewallet	413.04	4.7619047...	20.652	5.8	
315-22-5...	C	Naypyitaw	Normal	Female	Home and...	73.56	10	36.78	772.38	2/24/2019	11:38	Ewallet	735.6	4.7619047...	36.78	8.0	
665-32-9...	A	Yangon	Member	Female	Health an...	36.26	2	3.626	76.146	1/10/2019	17:15	Credit card	72.52	4.7619047...	3.626	7.2	
692-92-5...	B	Mandalay	Member	Female	Food and ...	54.84	3	8.226	172.746	2/20/2019	13:27	Credit card	164.52	4.7619047...	8.226	5.9	
351-62-0...	B	Mandalay	Member	Female	Fashion ac...	14.48	4	2.896	60.816	2/6/2019	18:07	Ewallet	57.92	4.7619047...	2.896	4.5	

Rows per page: 100 ▾ 1–50 of 50 < >

Leave a message in Chat

As with datasets, in the **Chats** button users can interact and leave their comments and questions on the Table Chat.

Add column description

Metadata makes more sense when columns description fields are not empty. With data.all you can add columns description and avoid the pain of figuring out fields purpose.

Select one table and in the **Columns** tab, directly type the description in the Description column as shown in the picture.

Table supermarket_sales															 Chat	 Edit	 Delete										
PREVIEW				OVERVIEW					COLUMNS					METRICS													
Name	Type	Description																									
invoice id	string	No description provided																									
branch	string	No description provided																									
city	string	No description provided																									
customer type	string	premium, standard																									
gender	string	No description provided																									

 Synchronize

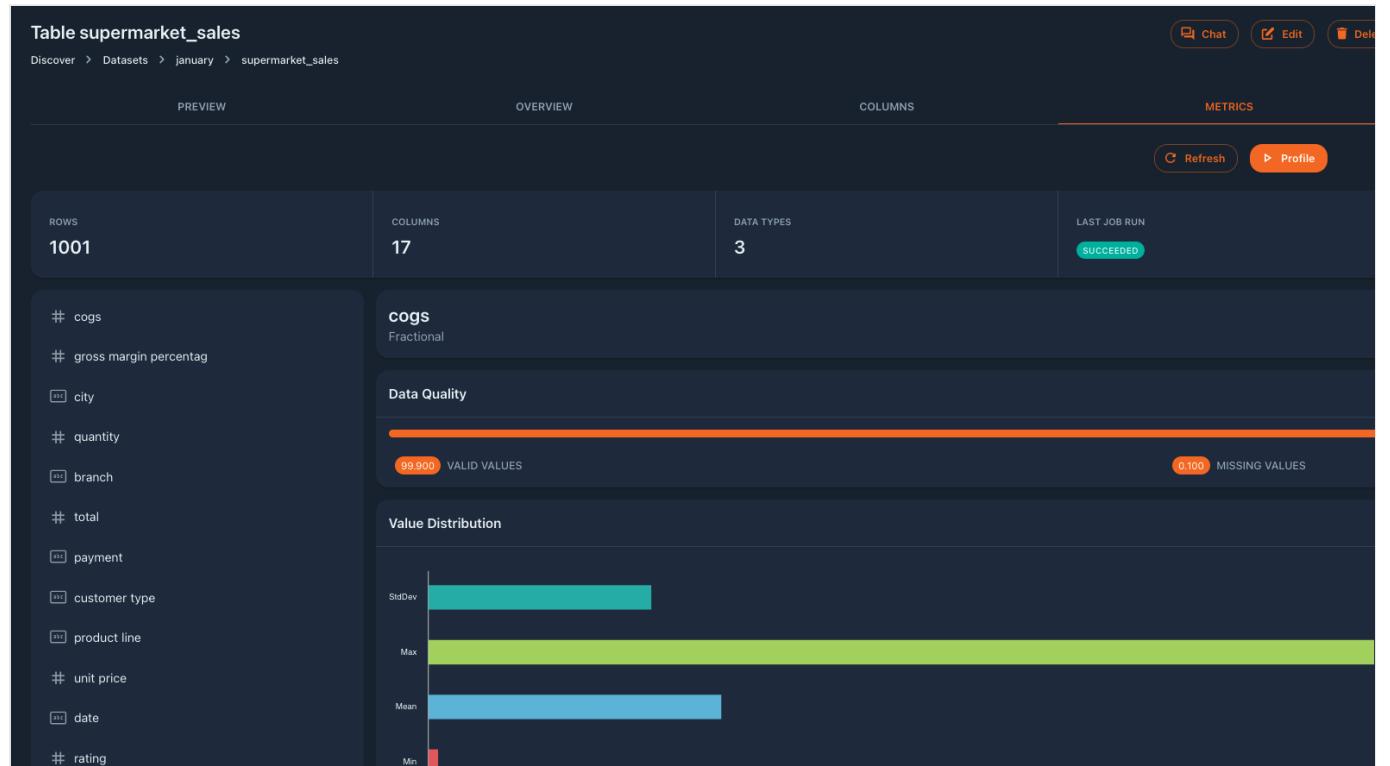
Profile data

Data profiling refers to the process of examining, analyzing, and reviewing the data available in the source by collecting statistical information about the data set's quality and hygiene. This process is called also data archaeology, data assessment, data discovery, or data quality analysis. Data profiling helps in determining the accuracy, completeness, structure, and quality of your data.

Data profiling in data.all involves:

- Collecting descriptive statistics like minimum, maximum, mean, median, and standard deviation.
- Collecting data types, along with the minimum and maximum length.
- Determining the percentages of distinct or missing data.
- Identifying frequency distributions and significant values.

By selecting the **Metrics** tab of your data table you can run a profiling job (click in the **Profile** button) , or view the latest generated data profiling metrics:



Delete a table

Deleting a table means deleting it from the data.all Catalog, but it will be still available on the AWS Glue Catalog. Moreover, when data owners delete a table, they are **not** deleting its data from the dataset S3 bucket. Teams with shared access to the dataset cannot delete tables or folders, even if they are shared.

It is possible to delete a table from the dataset **Tables** tab with the trash can icon next to each of the tables in the *Actions* column.

The screenshot shows the 'Tables' section of the AWS Glue Data Catalog interface. It lists five tables under the 'ds_january_od779vcv' dataset. The columns are 'Name', 'Database', 'Location', and 'Actions'. The 'Actions' column contains icons for 'Edit', 'Delete', and more. A red arrow points from the top right towards the 'Delete' icon for the first table.

Name	Database	Location	Actions
ds_january_od779vcv	ds_january_od779vcv	s3://ds-january-od779vcv/	
books_sales	ds_january_od779vcv	s3://ds-january-od779vcv/books_sales/	
raw	ds_january_od779vcv	s3://ds-january-od779vcv/raw/	
videogames_sales	ds_january_od779vcv	s3://ds-january-od779vcv/videogames_sales/	
supermarket_sales	ds_january_od779vcv	s3://ds-january-od779vcv/supermarket_sales/	

Another option is to go to the specific table (on the above picture click on the arrow icon next to the trash can icon). Click on the **Delete** button in the top right corner and confirm the deletion.

The screenshot shows the 'ds_january_od779vcv' table preview page. It includes sections for 'PREVIEW', 'OVERVIEW', 'COLUMNS', and 'METRICS'. A modal dialog titled 'Delete ds_january_od779vcv ?' is open, containing a warning message: '⚠️ Table will be deleted from data.all catalog, but will still be available on AWS Glue catalog.' and a red 'Delete' button. A red arrow points from the top right towards the 'Delete' button.

invoice id	branch	city	customer t...	gender	product line	unit price	quantity	tax %	total	date	time	payment	cogs	gross margin	gross income	rating	partition_
750-67-8...	A	Yangon	Member	Female	Health an...	74.69	7	26.1415	548.9715	1/5/2019	13:08	Ewallet	522.83	4.7619047...	26.1415	9.1	supermar...
226-31-3...	C	Naypyitaw	Normal	Female						9	10:29	Cash	76.4	4.7619047...	3.82	9.6	supermar...
631-41-3...	A	Yangon	Normal	Male						9	13:23	Credit card	324.31	4.7619047...	16.2155	7.4	supermar...
123-19-11...	A	Yangon	Member	Male						19	20:33	Ewallet	465.76	4.7619047...	23.288	8.4	supermar...
373-73-7...	A	Yangon	Normal	Male						9	10:37	Ewallet	604.17	4.7619047...	30.2085	5.3	supermar...
699-14-3...	C	Naypyitaw	Normal	Male						19	18:30	Ewallet	597.73	4.7619047...	29.8865	4.1	supermar...
355-53-5...	A	Yangon	Member	Female	Electronic ...	68.84	6	20.652	433.692	2/25/2019	14:36	Ewallet	413.04	4.7619047...	20.652	5.8	supermar...
315-22-5...	C	Naypyitaw	Normal	Female	Home and...	73.56	10	36.78	772.38	2/24/2019	11:38	Ewallet	735.6	4.7619047...	36.78	8.0	supermar...
665-32-9...	A	Yangon	Member	Female	Health an...	36.26	2	3.626	76.146	1/10/2019	17:15	Credit card	72.52	4.7619047...	3.626	7.2	supermar...
692-92-5...	B	Mandalay	Member	Female	Food and ...	54.84	3	8.226	172.746	2/20/2019	13:27	Credit card	164.52	4.7619047...	8.226	5.9	supermar...
351-62-0...	B	Mandalay	Member	Female	Fashion ac...	14.48	4	2.896	60.816	2/6/2019	18:07	Ewallet	57.92	4.7619047...	2.896	4.5	supermar...



An error occurred (ResourceShared) when calling `DELETE_DATASET_TABLE` operation: Revoke all table shares before deletion

To protect data consumers, if the table is shared you cannot delete it. The share requests to the table need to be revoked before deleting the table. Check the [Shares](#) section to learn how to grant and revoke access.

4.2.2 Folders

To open the Folder window you can either find your chosen folder in the Catalog or navigate to the dataset and then in the **Folders** tab click on the arrow in the **Actions** column of your folder:

Name	S3 Location	Description
january-sales-pdfs	s3://ds-january-od779vcv/pdfs	PDF prints of sales reports

Check folder and S3 metadata

The **Overview** tab of the folder contains folder metadata: - URI: unique folder identifier - Name: name of the folder, it is made out of the dataset name concatenated with the S3 prefix - Tags - Glossary terms - Description - Organization, Environment, Region, Team: inherited from the dataset - Created: creation time of the table

Details	
URI	3qwncc3rh
Name	january_sales_pdfs
Tags	demo
Glossary terms	-
Description	PDF prints of sales reports

S3 Properties	
S3 URI	s3://ds-january-od779vcv/pdfs/
S3 ARN	arn:aws:s3:::ds-january-od779vcv/pdfs/
Region	eu-west-1

CREATED BY	
john doe	john.doe@amazon.com
Organization	
Environment	
Region	eu-west-1
Team	DataScienceTeam
Created	5 days ago

Add or edit table metadata

Edit your folder metadata by clicking on the **Edit** button.

Description, Tags and Glossary terms are not inherited

Careful, those 3 fields are not synced with their dataset metadata. Just click on the **Edit** button of the folder to complete any missing information. This is especially useful to improve Catalog search of your folders.

Check the content of your folder

To check what kind of files does our prefix content, we can access the AWS S3 console on the **S3 Bucket** button of the Folder **Overview** tab.

Objects (5)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 Inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	airlines/	Folder	-	-	-
<input type="checkbox"/>	pdfs/	Folder	-	-	-
<input type="checkbox"/>	raw_airlines/	Folder	-	-	-
<input type="checkbox"/>	raw/	Folder	-	-	-
<input type="checkbox"/>	website/	Folder	-	-	-

Leave a message in Chat

Exactly the same as with tables. Allow your teams to discuss directly on the Folder Chat.

Delete a folder

Deleting folders is analogous to deleting tables. Deletion means deletion from the data.all Catalog and the content of the S3 prefix remains in the dataset S3 bucket. Only dataset owners can delete dataset folders.

The steps to delete a folder are exactly the same as with tables. You can either go to the dataset and in the **Folders** tab click on the can trash icon on the *Actions* column of the selected folder; or you can navigate to the Folder and click on the **Delete** button.



An error occurred (ResourceShared) when calling DELETE_DATASET_FOLDER operation: Revoke all folder shares before deletion

To protect data consumers, if the table is shared you cannot delete it. The share requests to the table need to be revoked before deleting the table. Check the [Shares](#) section to learn how to grant and revoke access.

The screenshot shows the AWS Data Studio interface for managing datasets. The top navigation bar indicates the path: Discover > Datasets > airlines > airlines_pdfs. Below this, there are two tabs: OVERVIEW (selected) and FEED.

OVERVIEW Tab Content:

- Details:**
 - URI: mMdoFv
 - Name: airlines_pdfs
 - Tags: -
 - Glossary terms: Canada
 - Description: registration forms
- S3 Properties:**
 - S3 URI: s3://ds-airlines-pfr3jm/pdfs/
 - S3 ARN: arn:aws:s3:::ds-airlines-pfr3jm/pdfs/
 - Region: eu-west-1
 - Account: [REDACTED]
- Actions:** S3 Bucket

Header Actions: Edit, Delete (highlighted with a red box and arrow)

Right Panel: A sidebar with user information and dataset metadata.

CREATED BY	john doe@amazon.com
Organization	
Environment	
Region	eu-west-1
Team	DataScienceTeam
Created	2 days ago

4.3 Centralized Catalog and glossaries

4.3.1 Catalog

In the Catalog we have a record with metadata for each dataset, table, folder and dashboard in data.all. Users come to this centralized Catalog to search and find data owned by other teams. Once users find a data asset they are interested in, they will create a [Share](#) request.

How do users find the data that they need?

Data needs to be discoverable, for this reason data.all Catalog offers a variety of filters that use business context to improve your search:

- **Type of data:** `dataset`, `table`, `folder` and/or `dashboard`
- **Tags:** tags of the data asset.
- **Topics:** filter by general topics created by the user.
- **Region:** AWS region where the data asset is located.
- **Classification:** `unclassified`, `official` and/or `secret`
- **Glossary:** filter datasets by the glossary terms created by users. This helps in two ways: It lets you narrow down results quickly using granular glossary terms like "sales", "profit", etc. Traditionally, a data glossary is just used to organize data. However, data.all uses it to power its search. This further encourages users to enrich and maintain the glossary regularly.

The screenshot shows the data.all interface with the 'Catalog' tab selected. The left sidebar includes sections for 'Discover' (Catalog, Datasets, Shares, Glossaries), 'PLAY' (Worksheets, Notebooks, ML Studio, Pipelines, Dashboards), and 'ADMIN' (Organizations, Environments). A prominent orange 'User Guide' button is at the bottom of the sidebar. The main area is titled 'Catalog' and shows a search bar with 'Search' placeholder text. Below the search bar, a message says 'No filters applied'. A dropdown menu provides filtering options: Type (Type, Tags, Topics, Region, Classification, Glossary). The results section displays eight dataset cards. Each card includes a thumbnail icon, the dataset name, the creator's email and creation date, a brief description, and metadata fields (Team, Environment, Region). Some cards also show a lock icon and a like count (e.g., 0 or 1).

Name	Creator	Created	Description	Team	Environment	Region
pdःfs	marlagarcia@amazon.com	3 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1
january_sales_pdfs	johndoe@amazon.com	5 days ago	PDF prints of sales reports	DataScienceTeam	data-science-general	euwest1
Insights_2	marlagarcia@amazon.com	5 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1
cannes_dates_3	marlagarcia@amazon.com	5 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1
raw	marlagarcia@amazon.com	3 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1
Demo	marlagarcia@amazon.com	2 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1
ds_johny_2_yngfwzpb	johndoe@amazon.com	5 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1
supermarket_sales	johndoe@amazon.com	5 days ago	No description provided	DataAnalysisTeam	data-analysis-general	euwest1

4.3.2 Glossaries

A Glossary is a list of terms, organized in a way to help users understand the context of their datasets. For example, terms like "cost", "revenue", etc, can be used to group and search all financial datasets.

The use of familiar terminology helps in quickly understanding the data and its background. It is a crucial element of data governance as it helps in bringing the business understanding closer to an organization's data initiatives.

On data.all, glossary terms can be attached to any dataset and can be leveraged to power quick and ease data discovery in the Catalog.

Spotlight

Glossaries are built hierarchically. They are made of categories and terms. This structure allows for glossaries from multiple domains to co-exist.

Term:

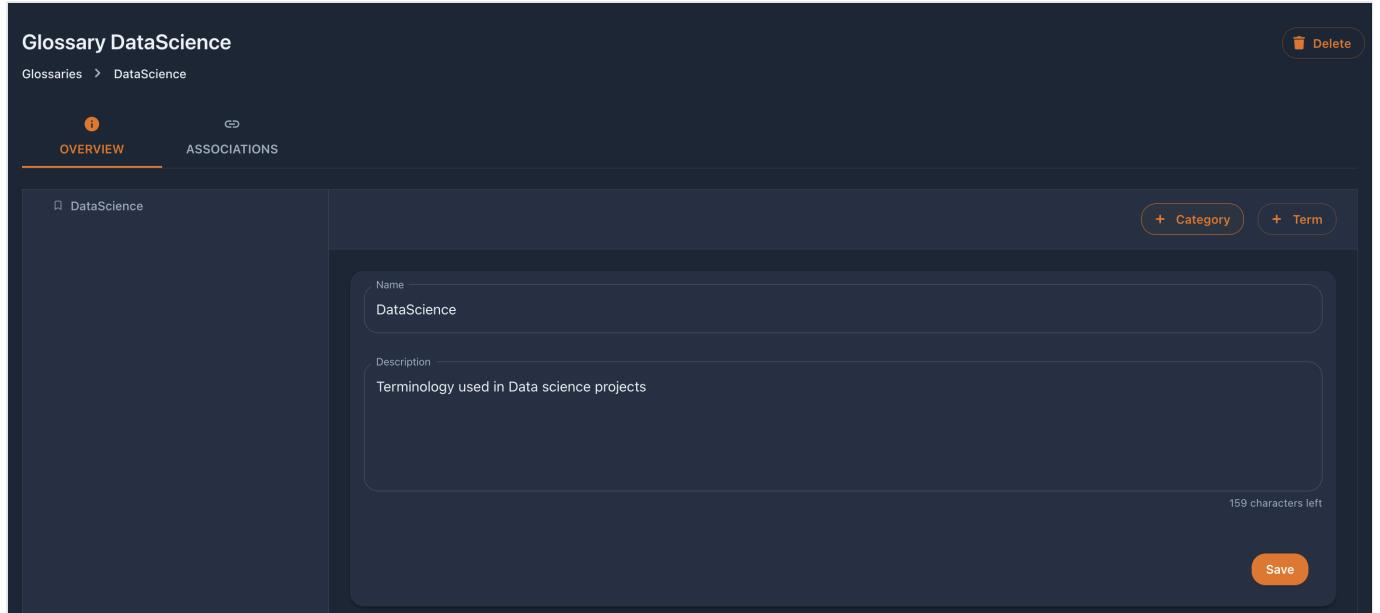
- A term is the lowest unit which is unique inside each glossary.
- It describes the content of the data assets in the most useful and precise way.
- It can exist independently, without belonging to any particular category or sub-category.

Category:

A category is used to group the terms of a similar context together. it is just a way of organizing terms.

Create a new glossary

1. Go to **Glossaries** menu on the left side pane.
2. Click on **Create**.
3. Fill the form and add a new glossary.



The screenshot shows a dark-themed user interface for creating a new glossary. At the top, a navigation bar displays 'Glossaries > DataScience'. On the right, there are 'Delete' and 'Save' buttons. Below the navigation, there are two tabs: 'OVERVIEW' (which is selected) and 'ASSOCIATIONS'. The main area contains a form for the 'DataScience' glossary. The 'Name' field is filled with 'DataScience'. The 'Description' field contains the text 'Terminology used in Data science projects'. There are two buttons at the bottom right: '+ Category' and '+ Term'. A character count of '159 characters left' is shown below the description field. The overall layout is clean and modern, with a focus on the input fields and navigation elements.

Add a category inside a Glossary

1. Click on the button "Add category" to add a new category.
2. Add a name and description to your category for better understanding.

Glossary DataScience

Glossaries > DataScience

OVERVIEW ASSOCIATIONS

- DataScience
- **Supervised Learning**
- ▽ Classification
- ▽ Regression

Name
Supervised Learning

Description
Terms included in supervised learning models. Supervised learning is when the model is getting trained on a labelled dataset. A labelled dataset is one that has both input and output parameters.

6 characters left

Save Delete

Add terms to a category

1. Click on the button "Add term" to add a new term to the category.
2. Give it an appropriate name and description.

Glossary DataScience

Glossaries > DataScience

OVERVIEW ASSOCIATIONS

- DataScience
- **Supervised Learning**
- ▽ **Classification**

Name
Classification

Description
Classification models are a subset of supervised machine learning . A classification models make predictions on DISCRETE data. A classification model reads some input and generates an output that classifies the input into some category.

-36 characters left

Save Delete

Remember!

The term will be used to recognize and filter the datasets. Hence, keep it short and precise.

Link your data with appropriate glossary terms

You can associate a glossary term to a dataset or a table. Go to a dataset click on "edit" and update the glossary terms field as shown below

Edit dataset January

Dataset name: January

Short description:

200 characters left

Classification

Confidentiality: Unclassified

Topics: Operations

Glossary Terms: Classification

Tags: prod

Deployment

Environment: Data Science

Region: eu-west-1

Organization: AnyCompany_EMEA

Governance

Team: DataScienceTeam

Stewards: DataScienceTeam

Save

Approve and Check all data related to a glossary

To see a list of all datasets and tables that have been linked with terms of a specific glossary, go to Glossaries and select the glossary. In the **Associations** tab it is possible to check the related data assets (target name), their types (e.g. dataset) and the specific term that they have used.

Important: Glossary owners need to approve the association. If it is not approved it won't be used as filter in the catalog.

Glossary DataScience

OVERVIEW ASSOCIATIONS

Term Associations

Term	Target Type	Target Name	Approval
Classification	Dataset	January	Approved

Delete

4.4 Shares

Teams can browse data.all catalog and request access for data assets. data.all shares data between teams securely within and environment and across environments without any data movement.

Datasets can contain tables and folders. Tables are Glue Tables registered in Glue Catalog. data.all uses (and automates) [Lake Formation sharing feature](#) to create access permissions to tables, meaning that no data is copied between AWS accounts.

Under-the-hood, folders are prefixes inside the dataset S3 bucket. To create sharing of folders in data.all, we create an S3 access point per requester group to handle its access to specific prefixes in the dataset.

Concepts

- Share request or Share Object: one for each dataset and requester team.
- Share Item refers to the individual tables and folders that are added to the Share request.

Sharing workflow

Requesters create a share request and add items to it. Both requesters and approvers can work on this `DRAFT` of the request and add and delete items to the request Draft. Items that are added go to the `PENDINGAPPROVAL` status.

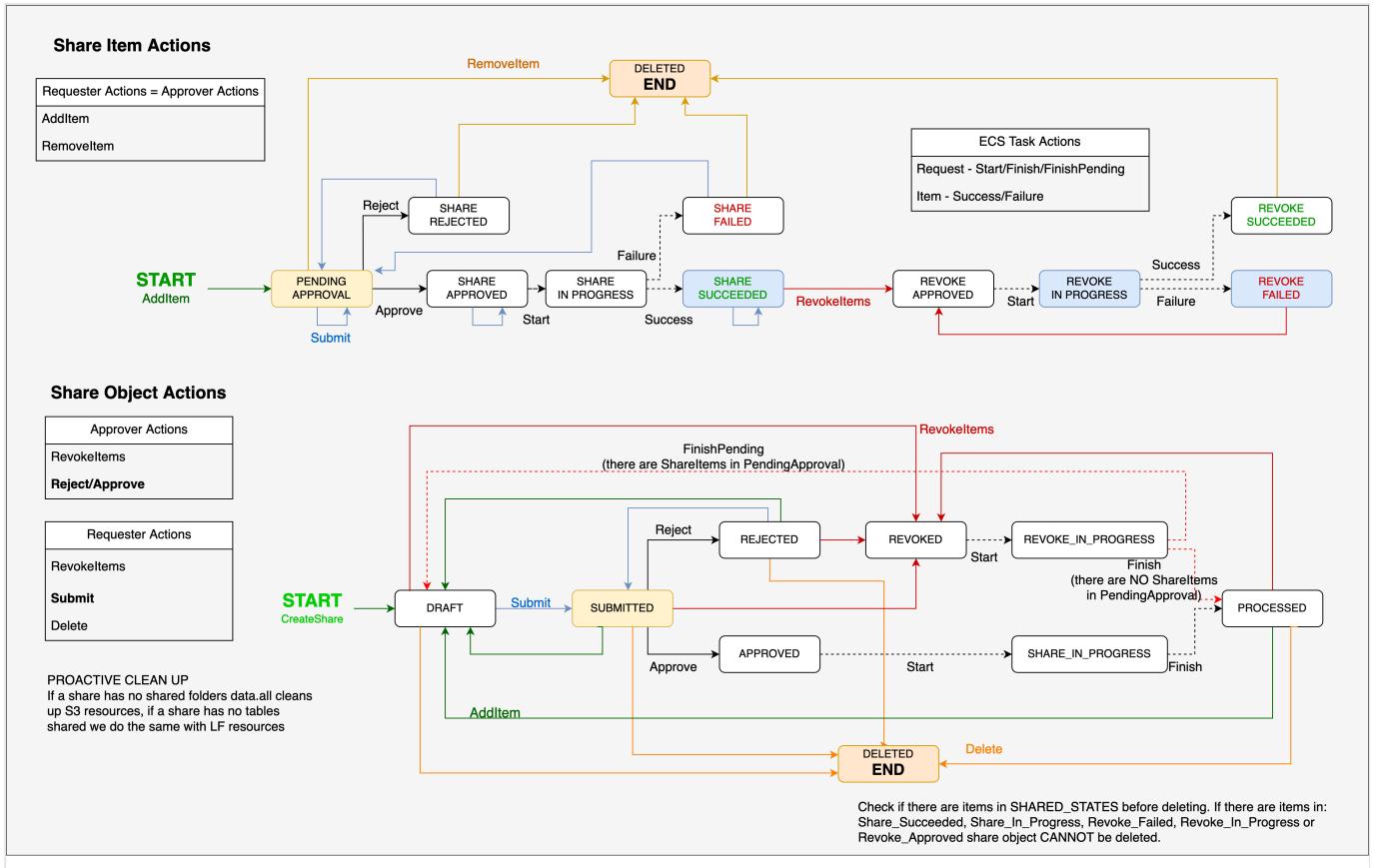
Once the draft is ready, requesters **submit** the request, which moves to the `SUBMITTED` status. Then, approvers **approve** or **reject** the request which will go to `APPROVED` or `REJECTED` status and its items to `SHARE_APPROVED` or `SHARE_REJECTED` correspondingly.

When the sharing task starts in the backend, both items and the share object move to `SHARE_IN_PROGRESS`. Once all items have been processed, the Share object is `PROCESSED` and each of the items is in either `SHARE_SUCCEEDED` or `SHARE_FAILED`. New items can be added to the share requests, the request will go back to `DRAFT` to be re-processed.

Both approvers and requesters can revoke access to shared items. They open the revoke items window and select which items should be revoked from the share request. The items move to `REVOKE_APPROVED` while the share is in `REVOKED` status.

While the revoking task is executing, the items and the request remain in `REVOKE_IN_PROGRESS` until the revoke is complete and items go to `REVOKE_FAILED` or `SUCCEEDED`. If there are share items in `PENDINGAPPROVAL` in the share request, it will go back to `DRAFT`. Otherwise, it will go to `PROCESSED`.

Requesters can delete the share request with the **delete** button. However, the request cannot contain any shared items. Users must revoke all shared items before deletion.



Create a share request (requester)

On left pane choose **Catalog** then **Search** for the table you want to access. Click on the lock icon of the selected data asset.

The screenshot shows the AWS Data Catalog interface with the following details:

- Left Sidebar:** DISCOVER, Catalog, Datasets, Shares, Glossaries, PLAY, Worksheets, Notebooks, ML Studio, Pipelines, Dashboards, ADMIN, Organizations, Environments.
- Top Bar:** data.all, New Dataset, Search bar.
- Current View:** Catalog > Catalog.
- Search Results:**
 - 5 results found in 446ms.
 - DatasetSB1-preupdate: by johndoe@amazon.com | created 13 days ago. Status: Energy, none. Lock icon highlighted with a red box and arrow.
 - folder1: by johndoe@amazon.com | created 13 days ago. Status: research-a, euwest1. Lock icon.
 - raw: by johndoe@amazon.com | created 13 days ago. Status: research-a, euwest1. Lock icon.
 - DatasetSB2: by johndoe@amazon.com | created an hour ago. Status: research-a, euwest1. Lock icon.
- Bottom Buttons:** User Guide, Like 0.

The following window will open. Choose your target environment and team and optionally add a *Request purpose*.

Request Access

Data access is requested for the whole requester Team or for the selected Consumption role. The request will be submitted to the data owners, track its progress in the Shares menu on the left.

Dataset name

DatasetSB2

Environment

Marketing-B

Requesters Team

SB2Marketing

Consumption Role (optional)

No additional consumption roles owned by this Team in this Environment.

Request purpose

200 characters left

► Send Request

johndoe@amazon.com | created 13

If instead of to a team, you want to request access for a Consumption role, add it to the request as in the picture below.

Request Access

Data access is requested for the whole requester Team or for the selected Consumption role. The request will be submitted to the data owners, track its progress in the Shares menu on the left.

Dataset name — DatasetSB2

Environment — Marketing-B ▾

Requesters Team — SB2Marketing ▾

Consumption Role (optional) — SagemakerPipelineMarketing23 [arn:aws:iam:...]

Request purpose
200 characters left

► Send Request

Finally, click on **Send Request**. This will create a share request or object for the corresponding dataset and if you have requested a table or folder it will add those items to the request. The share needs to be submitted for the request to be sent to the approvers.

4.4.1 Check your sent/received share requests

Anyone can go to the Shares menu on the left side pane and look up the share requests that they have received and that they have sent. Click on **Learn More** in the request that you are interested in to start working on your request.

Share Requests

Shares > Share Requests

RECEIVED SENT

mariagarcia@amazon.com
DRAFT | For datasetb2 | 2023-01-24 12:34:55.669354

Read access to Dataset: datasetb2 for Principal: SB2Marketing [arn:aws:iam:-----:role/dataall-marketing-b-t4zxupvl] from Environment: Marketing-B

Currently shared items: 0
Revoked items: 0
Failed items: 0
Pending items: 0

Learn More

< 1 >

4.4.2 Add/delete items

When you create a share request for a dataset, you still need to add the items (tables or folders) that you want to get access to. Initially the share request should be empty of items and in **DRAFT** state, it should look like the following picture.

Share object for datasetb2

Shares > Shares > datasetb2

Requested Dataset Details

Dataset: datasetb2
Dataset Owners: SB2Research
Dataset Environment: Research-A
Your role for this request: Requesters

REQUEST CREATED BY
mariagarcia@amazon.com

Principal: SB2Marketing [arn:aws:iam:-----:role/dataall-marketing-b-t4zxupvl]
Requester Team: SB2Marketing
Requester Environment: Marketing-B
Creation time: 2023-01-24 10:45:05.306815
Status: DRAFT

Shared Items

+ Add Item - Revoke Items

Type	Name	Status	Action
No items added.			

As appears in the picture, by clicking on **Add Item**, the following window will pop up to let you choose a specific table or folder in the dataset.

Add new item to share object datasetsb2

After adding an item, share object will be in draft status. Don't forget to submit your request !

Type	Name	Action
Folder	iot_files	+
Folder	images	+
Table	supermarkets	+
Table	books	+
Table	videogames	+

< 1 >

Note that the request is in `DRAFT` status and that the items that we add are in `PENDINGAPPROVAL`. They are not shared until the request is submitted and processed.

Share object for datasetsb2

Shares > Shares > [datasetsb2](#)

Refresh
 Submit
 Delete

Requested Dataset Details

Dataset **datasetsb2**

Dataset Owners **SB2Research**

Dataset Environment **Research-A**

Your role for this request **Requesters**

REQUEST CREATED BY
mariagarcia@amazon.com

Principal **SB2Marketing [arn:aws:iam::█████████████████████:role/dataall-marketing-b-t4zxu...**

Requester Team **SB2Marketing**

Requester Environment **Marketing-B**

Creation time **2023-01-24 12:34:55.669354**

Status **DRAFT**

Shared Items

Type	Name	Status	Action
Folder	iot_files	PENDINGAPPROVAL	
Table	supermarkets	PENDINGAPPROVAL	
Table	books	PENDINGAPPROVAL	

- 44/82 -

To remove an item from the request click on the **Delete** button with the trash icon next to it. We can only delete items that have not been shared. Items that are shared must be revoked, which is explained below.

4.4.3 Submit a share request (requester)

Once the draft is ready, the requesters need to click on the **submit** button. The request should be now in the `SUBMITTED` state. Approvers can see the request in their received share requests, alongside the current shared items, revoked items, failed items and pending items.

The screenshot shows the 'Share Requests' page with a 'RECEIVED' tab selected. A single share request is listed:

- Requester:** mariagarcia@amazon.com
- Status:** SUBMITTED | For datasetsb2 | 2023-01-24 12:34:55.669354
- Description:** Read access to Dataset: datasetsb2 for Principal: SB2Marketing [arn:aws:iam::█████████████████████:role/dataall-marketing-b-t4zxupvl] from Environment: Marketing-B
- Metrics (top right):**
 - Currently shared items: 0
 - Revoked items: 0
 - Failed items: 0
 - Pending items: 3
- Actions:** Learn More

4.4.4 Approve/Reject a share request (approver)

As an approver, click on **Learn more** in the `SUBMITTED` request and in the share view you can check the tables and folders added in the request. This is the view that approvers see, it now contains buttons to approve or reject the request.

The screenshot shows the 'Share object for datasetsb2' page. The top navigation is 'Shares > Shares > datasetsb2'. The main area has two sections:

- Requested Dataset Details:**
 - Dataset: datasetsb2
 - Dataset Owners: SB2Research
 - Dataset Environment: Research-A
 - Your role for this request: Approvers
- Request Created By:** mariagarcia@amazon.com

Principal	SB2Marketing [arn:aws:iam::█████████████████████:role/dataall-marketing-b-t4zxupvl]
Requester Team	SB2Marketing
Requester Environment	Marketing-B
Creation time	2023-01-24 10:45:05.306815
Status	SUBMITTED

Shared Items:

Type	Name	Status	Action
Folder	iot_files	PENDINGAPPROVAL	<input type="button" value="Delete"/>
Table	supermarkets	PENDINGAPPROVAL	<input type="button" value="Delete"/>

Actions: + Add Item, - Revoke Items

If the approvers **approve** the request, it moves to the `APPROVED` status. Share items IN `PENDINGAPPROVAL` will go to `SHARE_APPROVED`.

Research-A

Your role for this request
Approvers

Creation time: 2023-01-24 12:34:55.669354

Status: APPROVED

Shared Items

+ Add Item | - Revoke Items

Type	Name	Status	Action
Folder	iot_files	SHARE_APPROVED	<button>Delete</button>
Table	supermarkets	SHARE_APPROVED	<button>Delete</button>
Table	books	SHARE_APPROVED	<button>Delete</button>

< 1 >

Data.all backend starts a sharing task, during which, items and the request are in `SHARE_IN_PROGRESS` state.

Shares > Shares > datasetsb2

Share object for datasetsb2

Requested Dataset Details

Dataset: datasetsb2

Dataset Owners: SB2Research

Dataset Environment: Research-A

Your role for this request
Approvers

REQUEST CREATED BY: mariagarcia@amazon.com

Principal: SB2Marketing [arn:aws:iam:...:role/dataall-marketing-b-t4zxu...]

Requester Team: SB2Marketing

Requester Environment: Marketing-B

Creation time: 2023-01-24 12:34:55.669354

Status: SHARE_IN_PROGRESS

Shared Items

+ Add Item | - Revoke Items

Type	Name	Status	Action
Folder	iot_files	SHARE_SUCCEEDED	Revoke access to this item before deleting
Table	supermarkets	SHARE_IN_PROGRESS	<button>Delete</button>
Table	books	SHARE_APPROVED	<button>Delete</button>

When the task is completed, the items go to `SHARE_SUCCEEDED` or `SHARE_FAILED` and the request is `PROCESSED`.

The screenshot shows the 'Share object for datasetsb2' interface. At the top, there are 'Refresh' and 'Delete' buttons. Below that, the 'Requested Dataset Details' section shows the dataset name 'datasetsb2', its owner 'SB2Research', environment 'Research-A', and approvers. The 'Shared Items' section lists three items: 'iot_files', 'supermarkets', and 'books', all with a status of 'SHARE_SUCCEEDED'. There are 'Add Item' and 'Revoke Items' buttons at the bottom of this section.

Type	Name	Status	Action
Folder	iot_files	SHARE_SUCCEEDED	Revoke access to this item before deleting
Table	supermarkets	SHARE_SUCCEEDED	Revoke access to this item before deleting
Table	books	SHARE_SUCCEEDED	Revoke access to this item before deleting

If a dataset is shared, requesters should see the dataset on their screens. Their role with regards to the dataset is `SHARED`.

Datasets

Contribute > Datasets



Search



DatasetSB2

by johndoe@amazon.com

No description provided

Role

SHARED

Team

SB2Research

Tables

3

Folders

2

Status

CREATE_COMPLETE

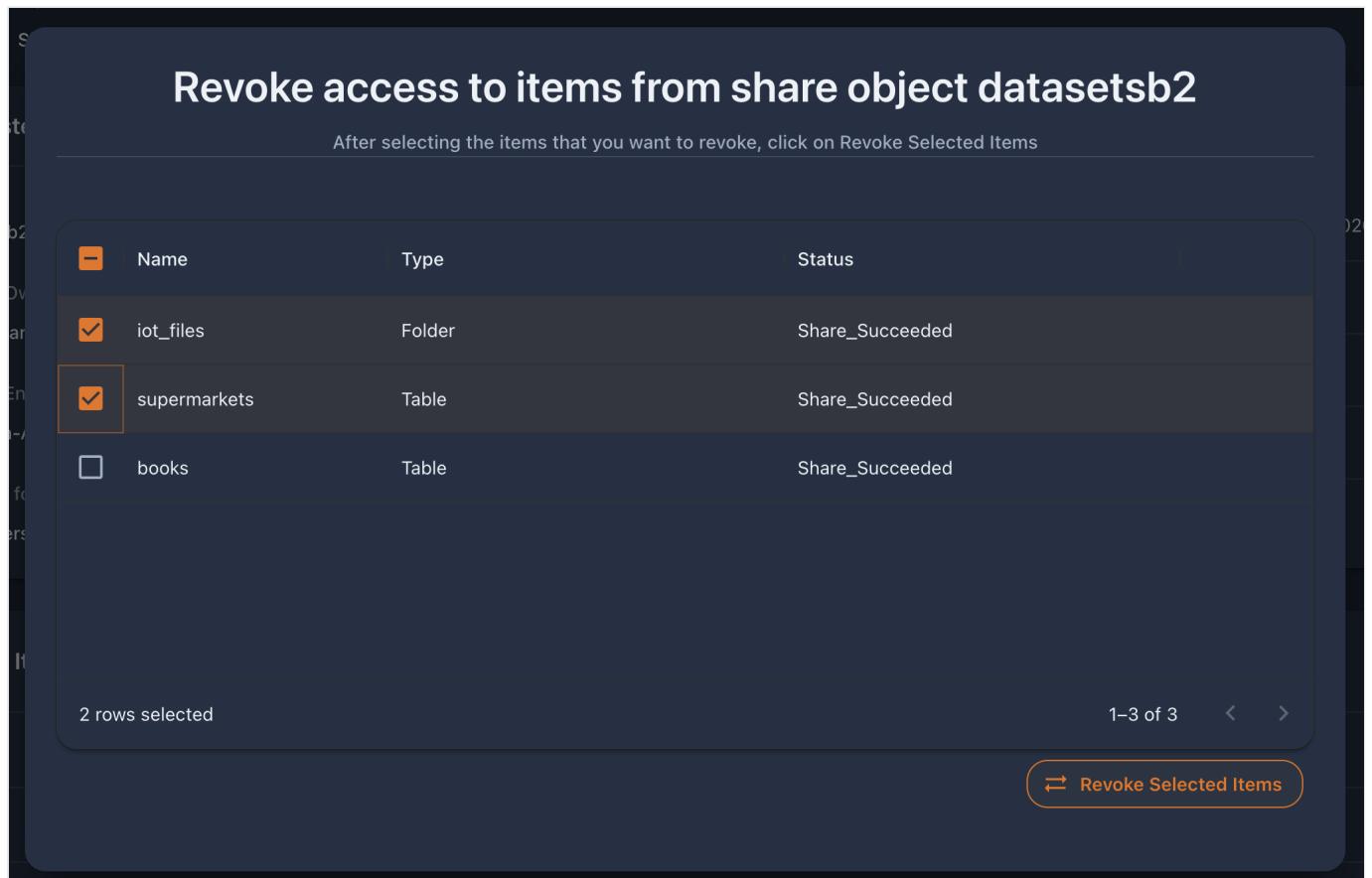
Learn More

0

4.4.5 Revoke Items

Both approvers and requesters can click on the button **Revoke items** to remove the share grant from chosen items.

It will open a window where multiple items can be selected for revoke. Once the button "revoke selected items" is pressed the consequent revoke task will be triggered.



Proactive clean-up

In every revoke task, data.all checks if there are no more shared folders or tables in a share request. In such case, data.all automatically cleans up any unnecessary S3 access point or Lake Formation permission.

4.4.6 Delete share request

To delete a share request, it needs to be empty from shared items. For example, the following request has some items in `SHARE_SUCCEEDED` state, therefore we receive an error. Once we have revoked access to all items we can delete the request.

The screenshot shows the 'Share object for datasetsb2' page in the AWS Data Catalog. At the top, there is an error message: 'An error occurred (UnauthorizedOperation) when calling Delete operation: This transition is not possible, Share_Succeeded cannot go to [Deleted]. If there is a sharing or revoking in progress wait until it is complete and try again.' Below the message, there are two main sections: 'Requested Dataset Details' and 'Shared Items'.

Requested Dataset Details:

- Dataset: datasetsb2
- Dataset Owners: SB2Research
- Dataset Environment: Research-A
- Your role for this request: Requesters

Shared Items:

Type	Name	Status	Action
Folder	iot_files	REVOKE_SUCCEEDED	<button>Delete</button>
Table	supermarkets	REVOKE_SUCCEEDED	<button>Delete</button>
Table	books	SHARE_SUCCEEDED	<button>Delete</button>

A red arrow points from the 'SHARE_SUCCEEDED' status of the 'books' table row to the 'Status' section in the 'Requested Dataset Details' panel, highlighting the status of the shared item.

4.4.7 Consume shared data

Data.all tables are Glue tables shared using AWS Lake Formation, therefore any service that reads Glue tables and integrates with Lake Formation is able to consume the data. Permissions are granted to the team role or the consumption role that has been specified in the request.

For the case of folders, the underlying sharing mechanism used is S3 Access Points. You can read data inside a prefix using the IAM role of the requester (same as with tables) and executing get calls to the S3 access point.

For example:

```
aws s3 ls arn:aws:s3:<SOURCE_REGION>:<SOURCE_AWSACCOUNTID>:accesspoint/<DATASETURI>--<REQUESTER-TEAM>/folder2/
```

5. Play

5.1 Worksheets

data.all offers a rich editor to write SQL queries and explore data. It is Athena on the backend that runs our queries on environments where our teams have been onboarded.

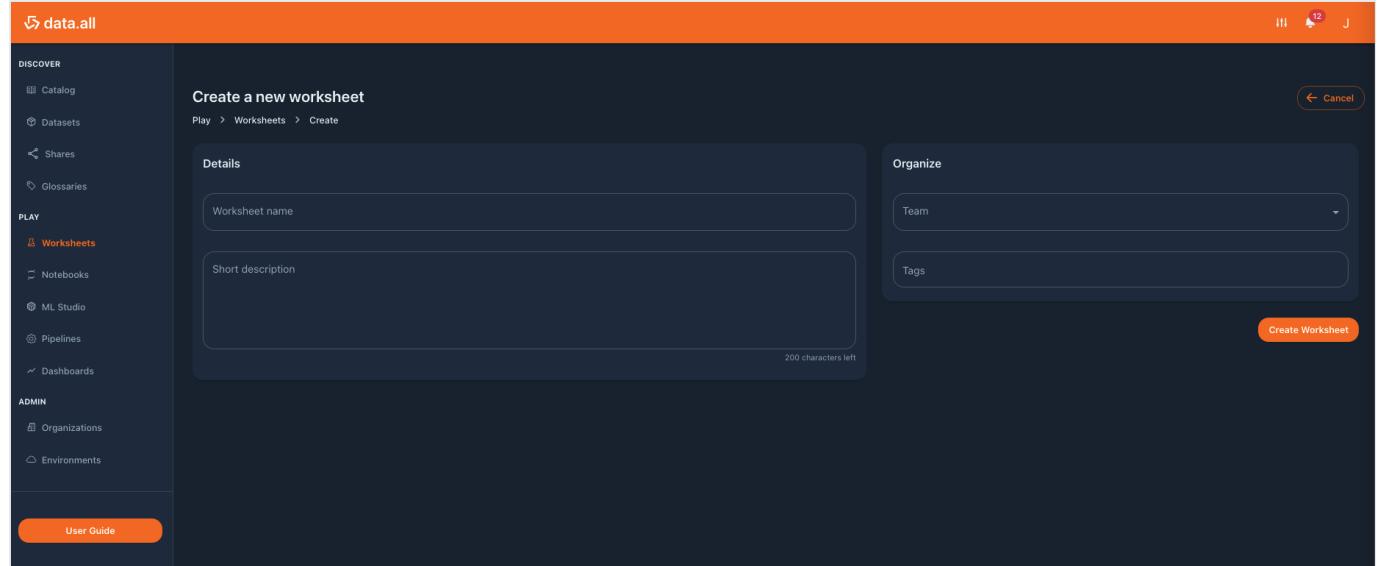
5.1.1 NEW Create a Worksheet

On the left pane under **Play** click on **Worksheets** to go to the Worksheet menu. Here you will find all Worksheets owned by your teams.

Shared queries = Seamless Collaboration

Check, learn from and collaborate with other members of your team to improve your analyses and get insights from your data, directly from data.all worksheets. No need to send queries by email, no need to create views :)

To create a new worksheet click on the **Create** button in the top right corner and fill the Worksheet form:



Field	Description	Required	Editable	Example
Worksheet name	Name of the worksheet	Yes	Yes	PalmDor
Short description	Short description about the worksheet	No	Yes	Query used to retrieve Palm D'or winners
Team	Team that owns the worksheet	Yes	No	DataScienceTeam
Tags	Tags	No	Yes	adhoc

No AWS resources

When we are creating a worksheet we are NOT deploying AWS resources. We don't provision clusters, we are not creating tables or views. We simply store the query in data.all database and we run it serverlessly on AWS Athena.

5.1.2 Edit worksheet metadata

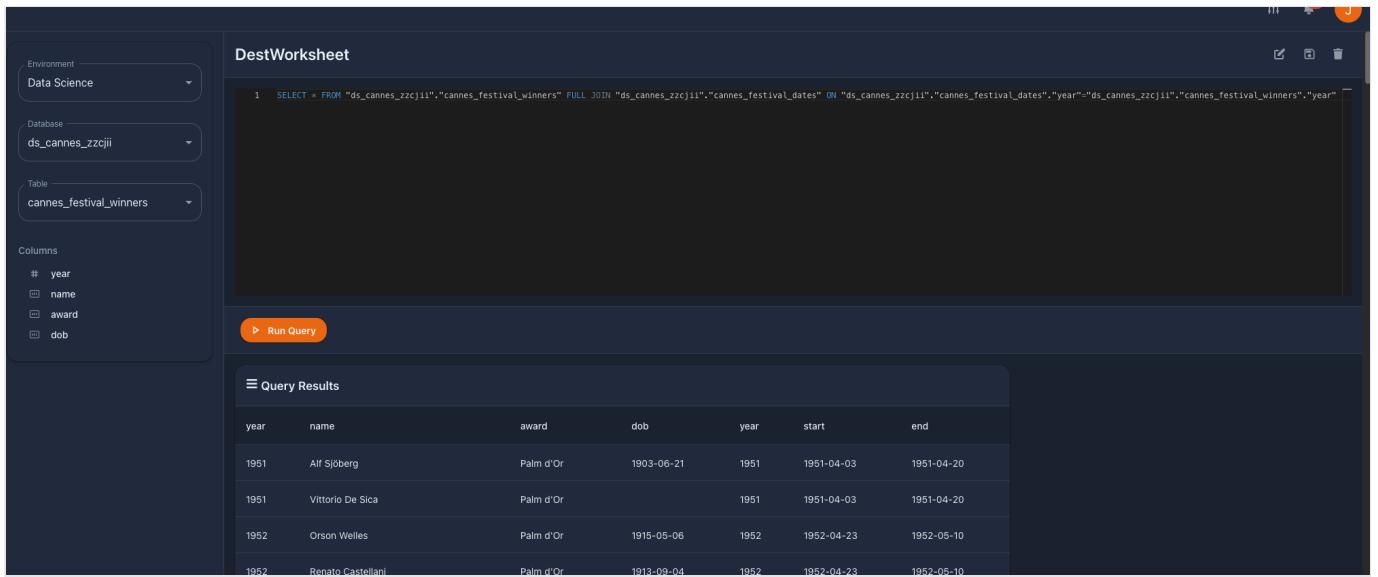
Select a worksheet and click on the pencil icon to edit the metadata of the worksheet. This includes worksheet name, description and tags. The ownership of the worksheet, its team, is not editable.

5.1.3 Delete a worksheet

Next to the edit button, there are 2 other buttons. To delete a worksheet click on the trash icon one. Worksheets are not AWS resources, they are a data.all construct whose information is stored in the data.all database. Thus, when we delete a worksheet we are not deleting AWS resources or CloudFormation stacks.

5.1.4 Write and save your queries

Select your worksheet and choose any of the environments, datasets and tables of your team to list column information. In the query editor write your SQL statements and click on **Run Query** to get your query results. Error messages coming from Athena will pop-up automatically.



The screenshot shows the DestWorksheet interface. On the left, there are dropdown menus for Environment (Data Science), Database (ds_cannes_zzcjii), and Table (cannes_festival_winners). Below these are sections for Columns (year, name, award, dob) and a list of columns (year, name, award, dob). The main area is titled "DestWorksheet" and contains a query editor with the following SQL code:

```
1 SELECT * FROM "ds_cannes_zzcjii"."cannes_festival_winners" FULL JOIN "ds_cannes_zzcjii"."cannes_festival_dates" ON "ds_cannes_zzcjii"."cannes_festival_dates"."year"="ds_cannes_zzcjii"."cannes_festival_winners"."year"
```

Below the query editor is a "Run Query" button. To the right of the main area is a "Query Results" section with a table:

year	name	award	dob	year	start	end
1951	Alf Sjöberg	Palm d'Or	1903-06-21	1951	1951-04-03	1951-04-20
1951	Vittorio De Sica	Palm d'Or		1951	1951-04-03	1951-04-20
1952	Orson Welles	Palm d'Or	1915-05-06	1952	1952-04-23	1952-05-10
1952	Renato Castellani	Palm d'Or	1913-09-04	1952	1952-04-23	1952-05-10

If you want to save the current query for later or for other users, click on the **save** icon (between the edit and the delete buttons).

More than just SELECT

Worksheets can be used for data exploration, for quick ad-hoc queries and for more complicated queries that require joins. As far as you have access to the joined datasets you can combine information from multiple tables or datasets. Check the [docs](#) for more information on AWS Athena SQL syntax.

5.2 Notebooks

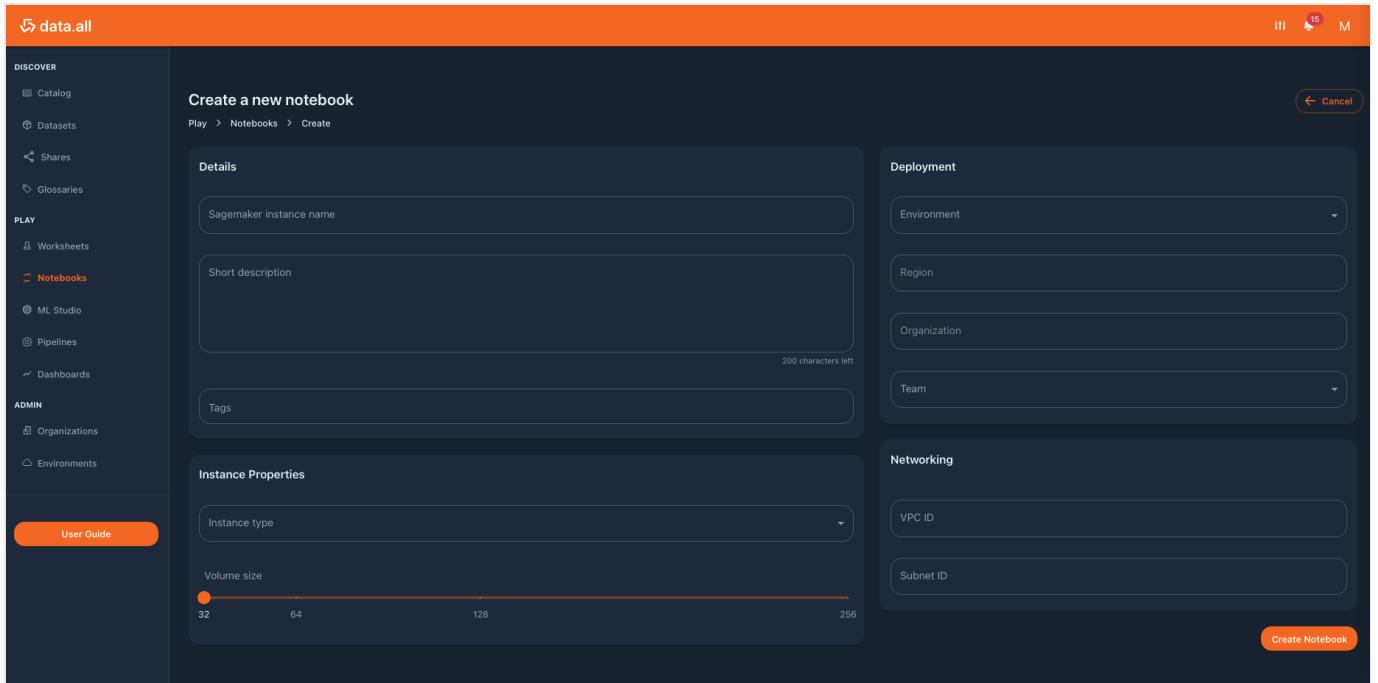
Data practitioners can experiment machine learning algorithms spinning up Jupyter notebook with access to all your datasets. `data.all` leverages [Amazon SageMaker instance](#) to access Jupyter notebooks.

5.2.1 Create a Notebook

Pre-requisites

To use Notebooks you need to introduce your own VPC ID or create a Sagemaker Studio domain inside a VPC (read the [docs](#)). Provisioning the notebook instances inside a VPC enables the notebook to access VPC-only resources such as EFS file systems.

To create a Notebook, go to Notebooks on the left pane and click on the **Create** button. Then fill in the following form:



Field	Description	Required	Editable	Example
Sagemaker instance name	Name of the notebook	Yes	No	Cannes Project
Short description	Short description about the notebook	No	No	Notebook for Cannes exploration
Tags	Tags	No	No	deleteme
Environment	Environment (and mapped AWS account)	Yes	No	Data Science
Region (auto-filled)	AWS region	Yes	No	Europe (Ireland)
Organization (auto-filled)	Organization of the environment	Yes	No	AnyCompany EMEA
Team	Team that owns the notebook	Yes	No	DataScienceTeam
VPC Identifier	VPC provided to host the notebook	No	No	vpc-.....
Public subnets	Public subnets provided to host the notebook	No	No	subnet-....
Instance type	[ml.t3.medium, ml.t3.large, ml.m5.xlarge]	Yes	No	ml.t3.medium
Volume size	[32, 64, 128, 256]	Yes	No	32

If successfully created we can check its metadata in the **Overview** tab. Unlike other data.all resources, Notebooks are non-editable.

Notebook Cannes exploration

Play > Notebooks > Cannes exploration

OVERVIEW TAGS AWS STACK

Details

URI: iloNk8

Name: Cannes exploration

Tags: -

Description: Some

Instance Properties

Instance type: ml.t3.medium

Volume size: 32 Go

VPC: [REDACTED]

Subnet: [REDACTED]

Instance Profile: arn:aws:iam:-----role/ds-data-science-rykpp5

CloudFormation Stack

- CREATED BY: john doe@amazon.com
- Organization: AnyCompany Global
- Environment: Data Science
- Team: DataScienceTeam
- Created: an hour ago
- Status: INSERVICE

5.2.2 Check CloudFormation stack

In the **Stack** tab of the Notebook, is where we check the AWS resources provisioned by data.all as well as its status. As part of the Notebook CloudFormation stack deployed using CDK, data.all will deploy:

1. AWS EC2 Security Group
2. AWS SageMaker Notebook Instance
3. AWS KMS Key and Alias

5.2.3 Delete a Notebook

To delete a Notebook, simply select it and click on the **Delete** button in the top right corner. It is possible to keep the CloudFormation stack associated with the Notebook by selecting this option in the confirmation delete window that appears after clicking on delete.

5.2.4 Open JupyterLab

Click on the **Open JupyterLab** button of the Notebook window to start writing code on Jupyter Notebooks.

jupyter

Open JupyterLab Quit Logout

Files Running Clusters SageMaker Examples Conda

Select items to perform actions on them.

Upload New

Name	Last Modified	File size
Untitled.ipynb	Running seconds ago	72 B

5.2.5 Stop/Start instance

As we briefly commented, `data.all` uses AWS SageMaker instances to access Jupyter notebooks. Be frugal and stop your instances when you are not developing. To do that, close the Jupyter window and click on **Stop Instance** in the Notebook buttons. It takes a couple of minutes, just refresh and check the Notebook Status in the overview tab. It should end up in `STOPPED`.



Save money, stop your instances

This feature allows users to easily manage their instances directly from `data.all` UI.

Same when you are coming back to work on your Notebook, click on **Start instance** to start the SageMaker instance. In this case the Status of the notebook should first be `PENDING` and once the instance is ready, `INSERVICE`.

5.2.6 Create Key-value tags

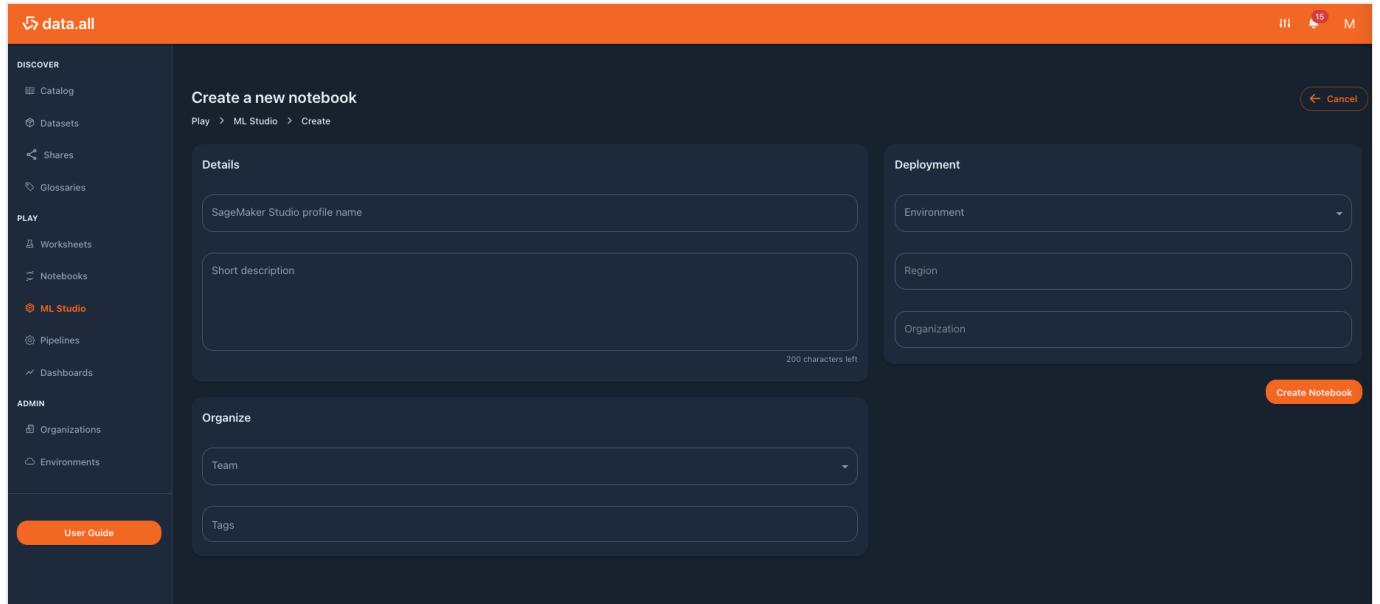
In the **Tags** tab of the notebook window, we can create key-value tags. These tags are not `data.all` tags that are used to tag datasets and find them in the catalog. In this case we are creating AWS tags as part of the notebook CloudFormation stack. There are multiple tagging strategies as explained in the [documentation](#).

5.3 ML Studio

With ML Studio Notebooks we can add users to our SageMaker domain and open Amazon SageMaker Studio

5.3.1 Create a ML Notebook

To create a new Notebook, go to ML Studio on the left side pane and click on Create. Then fill in the creation form with its corresponding information.



Field	Description	Required	Editable	Example
Sagemaker Studio profile name	Name of the user to add to SageMaker domain	Yes	No	johndoe
Short description	Short description about the notebook	No	No	Notebook for Cannes exploration
Tags	Tags	No	No	deleteme
Environment	Environment (and mapped AWS account)	Yes	No	Data Science
Region (auto-filled)	AWS region	Yes	No	Europe (Ireland)
Organization (auto-filled)	Organization of the environment	Yes	No	AnyCompany EMEA
Team	Team that owns the notebook	Yes	No	DataScienceTeam

5.3.2 Check CloudFormation stack

In the **Stack** tab of the ML Studio Notebook, is where we check the AWS resources provisioned by data.all as well as its status. As part of the CloudFormation stack deployed using CDK, data.all will deploy some CDK metadata and a SageMaker User Profile.

5.3.3 Delete a Notebook

To delete a Notebook, simply select it and click on the **Delete** button in the top right corner. It is possible to keep the CloudFormation stack associated with the Notebook by selecting this option in the confirmation delete window that appears after clicking on delete.

Notebook johndoe

Play > ML Studio > johndoe

OVERVIEW STACK

Details

URI: THKIJm

Name: johndoe

Tags: -

Description: some

CREATED BY: johndoe@amazon.com

Organization: AnyCompany Global

Environment: Data Science

Team: DataScienceTeam

Created: an hour ago

Status: NOTFOUND

5.3.4 Open Amazon SageMaker Studio

Click on the **Open ML Studio** button of the ML Studio notebook window to open Amazon SageMaker Studio.

Amazon SageMaker Studio

File Edit View Run Kernel Git Tabs Settings Help

Get started

Explore solutions, models, algorithms, and tutorials

SageMaker JumpStart

Solution: Detect malicious users and transactions →

Solution: Demand forecasting →

Go to SageMaker JumpStart →

Build models automatically

SageMaker Autopilot

Video: Get started with Autopilot →

Blog: Getting started with Autopilot →

New autopilot experiment →

Instantly prepare data for ML

SageMaker Data Wrangler

Blog: Getting started with Data Wrangler →

Blog: Predicting credit risk →

Start now →

ML tasks and components

New compilation job

Create a new compilation job. View compilation jobs

+ New feature group

Create a new feature group in the feature store to logically group and manage features. View feature store

+ New data flow

Prepare and visualize your data with SageMaker Data Wrangler. View data flows

+ New project

Organize ML components and automate MLOps with built-in or custom project templates. View projects

Notebooks and compute resources

Select a SageMaker Image: Data Science

Select a start-up script: No Script

+ Notebook Python 3

+ Console Python 3

+ Image Terminal Image Terminal

Utilities and files

Show Contextual Help

+ System terminal

+ Text File

+ Markdown File

5.4 Pipelines

Different business units might have their own data lake and ingest and process the data with very different tools: Scikit Learn, Spark, SparkML, AWS SageMaker, AmazonAthena... The diversity of tools and use-cases result in a wide variety of CICD standards which discourages development collaboration.

In order to distribute data ingestion and processing, data.all introduces data.all pipelines:

- data.all takes care of CICD infrastructure
- data.all integrates with [AWS DDK](#), a tool to help you build data workflows in AWS
- data.all allows you to define development environments directly from the UI and deploys data pipelines to those AWS accounts



Focus on value-added code

data.all takes care of the CICD and multi-environment configuration and DDK provides reusable assets and data constructs that accelerate the deployment of AWS data workflows, so you can focus on writing the actual transformation code and generating value from your data!

5.4.1 Multi-environment Pipelines

In some cases, enterprises decide to separate CICD resources from data application resources, which at the same time, need to be deployed to multiple accounts. Data.all allows users to easily define their CICD environment and other infrastructure environments in a flexible, robust way.

Let's see it with an example. In your enterprise, the Research team has 3 AWS accounts: Research-CICD, Research-DEV and Research-PROD. They want to ingest data with a data pipeline that is written in Infrastructure as Code (IaC) in the Research-CICD account. The actual data pipeline is deployed in 2 data accounts. First, in Research-DEV for development and testing and once it is ready it is deployed to Research-PROD.

Pre-requisites

As a pre-requisite, Research-DEV and Research-PROD accounts need to be bootstrapped trusting the CICD account (`-a` parameter) and setting the stage of the AWS account, the environment id, with the `e` parameter. Assuming `111111111111` = CICD account the commands are as follows:

- In Research-CICD (111111111111): `ddk bootstrap -e cicd`
- In Research-DEV (222222222222): `ddk bootstrap -e dev -a 111111111111`
- In Research-PROD (333333333333): `ddk bootstrap -e prod -a 111111111111`

In data.all we need to link the AWS accounts to the platform by creating 3 data.all Environments: Research-CICD Environment, Research-DEV Environment and Research-PROD Environment.

Creating a pipeline

data.all pipelines are created from the UI, under Pipelines. We need to fill the creation form with the following information:

- Name, Description and tags
- CICD Environment: AWS account and region where the CICD resources will be deployed.
- Team, this is the Admin team of the pipeline. It belongs to the specified CICD Environment where the pipeline is defined as IaC
- CICD strategy: This is the development strategy that determines the type of CICD Pipeline that is created by data.all. Currently the following 4 types are supported depending on your use case:
 - CDK Pipelines - Trunk-based**: A CICD pipeline based on [CDK Pipelines library](#). It defines a DDK Core construct which deploys Continuous Integration and Delivery for your app. Specifically, it provisions a stack containing a self-mutating CDK code pipeline to deploy one or more copies of your CDK applications using CloudFormation with a minimal amount of effort on your part.
 - CodePipeline - Trunk-based**: A CICD pipeline similar to CDK Pipelines and with a trunk-based approach but is not self-mutating.
 - CodePipeline - Gitflow**: A Gitflow branching strategy where each branch of the source repository has a corresponding CICD Pipeline that deploys resources for that branches environment.
 - GitHub Template**: This is a Bring-Your-Own-Template approach where users can specify they git clone path and deploy their own pipelines and IaC rather than using one of the previous 3 strategies.

Finally, we need to add **Development environments**. These are the AWS accounts and regions where the infrastructure defined in the CICD pipeline is deployed.



Environment ID = data.all environment stage

When creating the pipeline and adding development environments, you define the stage of the environment. The ddk bootstrap `-e` parameter needs to match the one that you define in the data.all UI. In our example, we bootstrapped with the parameters "dev" and "prod" and then we defined the stages as "dev" and "prod" correspondingly.

Create a new pipeline

Play > Pipelines > Create

Cancel

Details <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> Pipeline name <input type="text" value="My Pipeline"/> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> Short description <div style="height: 100px; border: 1px solid #ccc; width: 100%;"></div> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> Tags <div style="height: 40px; border: 1px solid #ccc; width: 100%;"></div> </div>	CICD <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> CICD Environment <input type="text" value="Research-CI_CD"/> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> Team <input type="text" value="Engineers"/> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> Region <input type="text" value="eu-west-2"/> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> Organization <input type="text" value="new"/> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> CICD strategy <div style="border: 1px solid #ccc; padding: 2px; display: inline-block;">CDK Pipelines - Trunk-based</div> </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> CDK Pipelines - Trunk-based </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> CodePipeline - Trunk-based </div> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 10px;"> CodePipeline - Gitflow </div>												
Development environments <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 10%;">Order</th> <th style="width: 20%;">Development Stage</th> <th style="width: 20%;">Environment</th> <th style="width: 10%;">Team</th> </tr> </thead> <tbody> <tr> <td>1</td> <td style="text-align: center;"><input type="button" value="dev"/></td> <td style="text-align: center;"><input type="text" value="Research-DEV"/></td> <td style="text-align: center;"><input type="text" value="Scientists"/></td> </tr> <tr> <td>2</td> <td style="text-align: center;"><input type="button" value="prod"/></td> <td style="text-align: center;"><input type="text" value="Research-PROD"/></td> <td style="text-align: center;"><input type="text" value="Scientists"/></td> </tr> </tbody> </table> <p>Add environment</p>		Order	Development Stage	Environment	Team	1	<input type="button" value="dev"/>	<input type="text" value="Research-DEV"/>	<input type="text" value="Scientists"/>	2	<input type="button" value="prod"/>	<input type="text" value="Research-PROD"/>	<input type="text" value="Scientists"/>
Order	Development Stage	Environment	Team										
1	<input type="button" value="dev"/>	<input type="text" value="Research-DEV"/>	<input type="text" value="Scientists"/>										
2	<input type="button" value="prod"/>	<input type="text" value="Research-PROD"/>	<input type="text" value="Scientists"/>										

Create Pipeline

CDK Pipelines Overview

CODECOMMIT REPOSITORY

When a pipeline is created, an AWS CodeCommit repository with the code of an AWS DDK application is created in the CICD environment AWS account. It contains an set up for a multi-account deployment, as explained in its [documentation](#).

In the deployed repository, data.all pushes a `ddk.json` file with the details of the selected development environments:

```
{
  "environments": {
    "cicd": {
      "account": "111111111111",
      "region": "eu-west-1"
    },
    "dev": {
      "account": "222222222222",
      "region": "eu-west-1",
      "resources": {
        "ddk-bucket": {"versioned": false, "removal_policy": "destroy"}
      }
    },
    "prod": {
      "account": "333333333333",
      "region": "eu-west-1",
      "resources": {
        "ddk-bucket": {"versioned": true, "removal_policy": "retain"}
      }
    }
  }
}
```

In addition, the `app.py` file is also written accordingly to the development environments selected in data.all UI. It will look similar to the following:

```
#!/usr/bin/env python3

import aws_cdk as cdk
from aws_ddk_core.cicd import CICDPipelineStack
from ddk_app.ddk_app_stack import DDKApplicationStack
from aws_ddk_core.config import Config

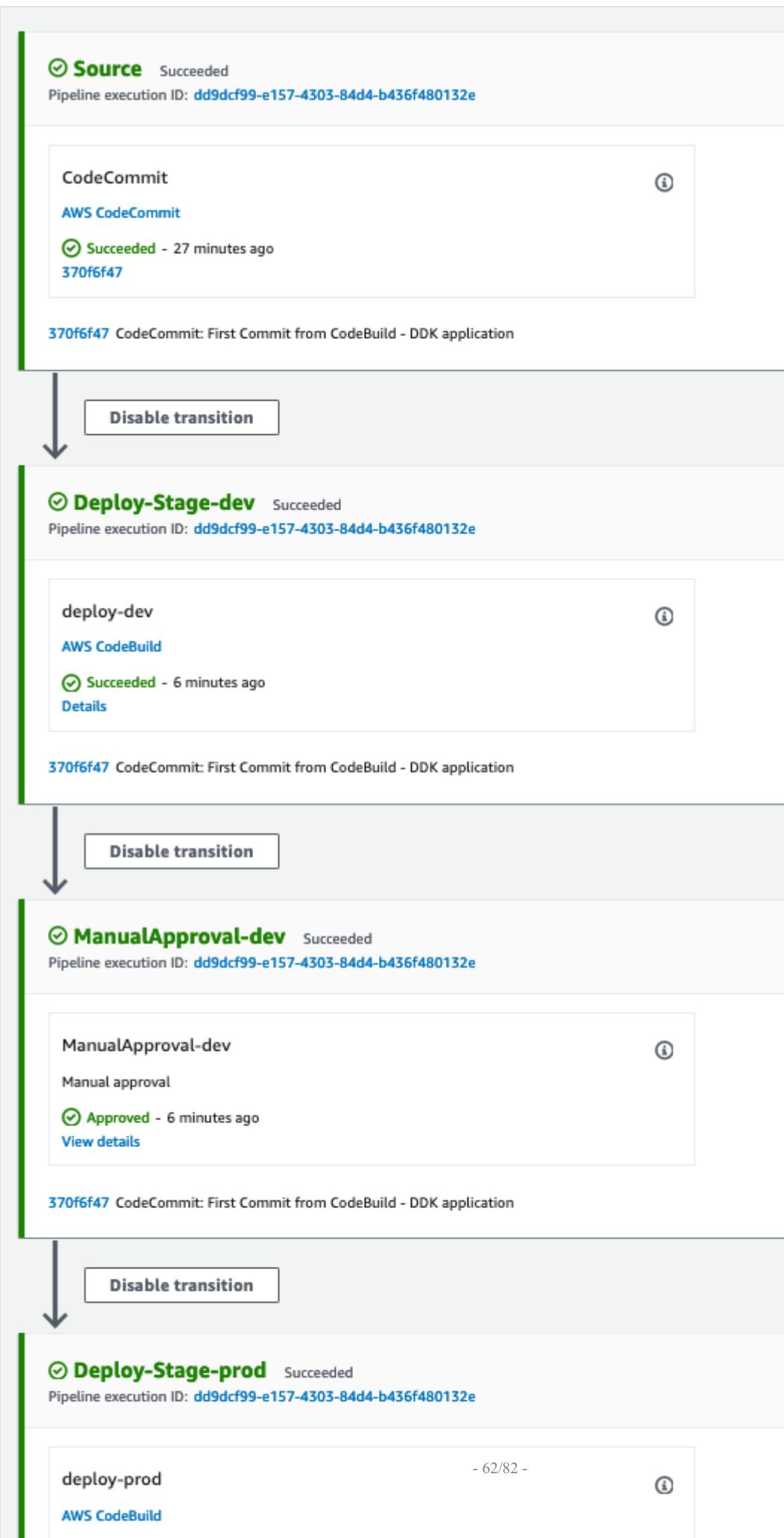
app = cdk.App()

class ApplicationStage(cdk.Stage):
    def __init__(self, scope, environment_id: str, **kwargs,):
        super().__init__(scope, f"dataall-{environment_id.title()}", **kwargs)
        DDKApplicationStack(self, "DataPipeline-PIPELINENAME-PIPELINEURI", environment_id)

config = Config()
(
    CICDPipelineStack(
        app,
        id="dataall-pipeline-PIPELINENAME-PIPELINEURI",
        environment_id="cicd",
        pipeline_name="PIPELINENAME",
    )
    .add_source_action(repository_name="dataall-PIPELINENAME-PIPELINEURI")
    .add_synth_action()
    .build().add_stage("dev", ApplicationStage(app, "dev", env=config.get_env("dev"))).add_stage("prod", ApplicationStage(app, "prod", env=config.get_env("prod")))
    .synth()
)
app.synth()
```

CICD DEPLOYMENT

data.all backend performs the first deployment of the CICD stack defined in the CodeCommit repository. The result is a CloudFormation template deploying a CICD pipeline having the aforementioned CodeCommit repository as source. This CodePipeline pipeline is based on the [CDK Pipelines library](#).



CodePipeline pipelines - Trunk-based or GitFlow Overview

For cases in which we do not want to use [CDK Pipelines library](#) we can instead use CodePipeline CICD Strategy which leverages the `aws-codepipeline` construct library.

CODECOMMIT REPOSITORY AND CICD DEPLOYMENT

When a pipeline is created, a CloudFormation stack is deployed in the CICD environment AWS account. It contains:

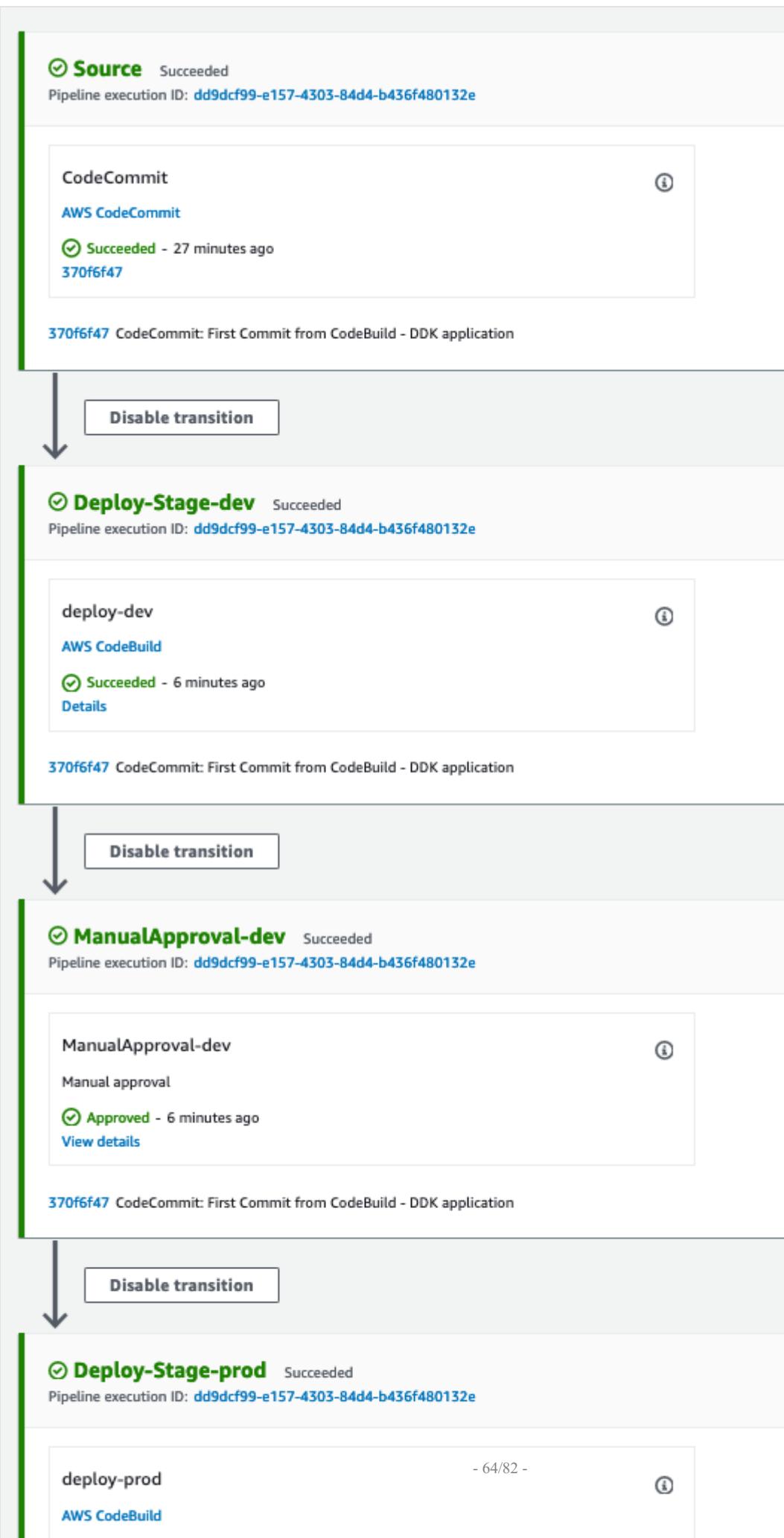
- an AWS CodeCommit repository with the code of an AWS DDK application (by running `ddk init`) with some modifications to allow cross-account deployments.
- CICD CodePipeline(s) pipeline that deploy(s) the application

The repository structure will look similar to:

Name
__pycache__
ddk_app
utils
app.py
cdk.json
ddk.json
deploy_buildspec.yaml
init_branches_deploy_buildspec.yaml
init_deploy_buildspec.yaml
README.md
requirements-dev.txt
requirements.txt
setup.py
source.bat
test.sh

The added `Multiaccount` configuration class allows us to define the deployment environment based on the `ddk.json`. Go ahead and customize this configuration further, for example you can set additional `env_vars`.

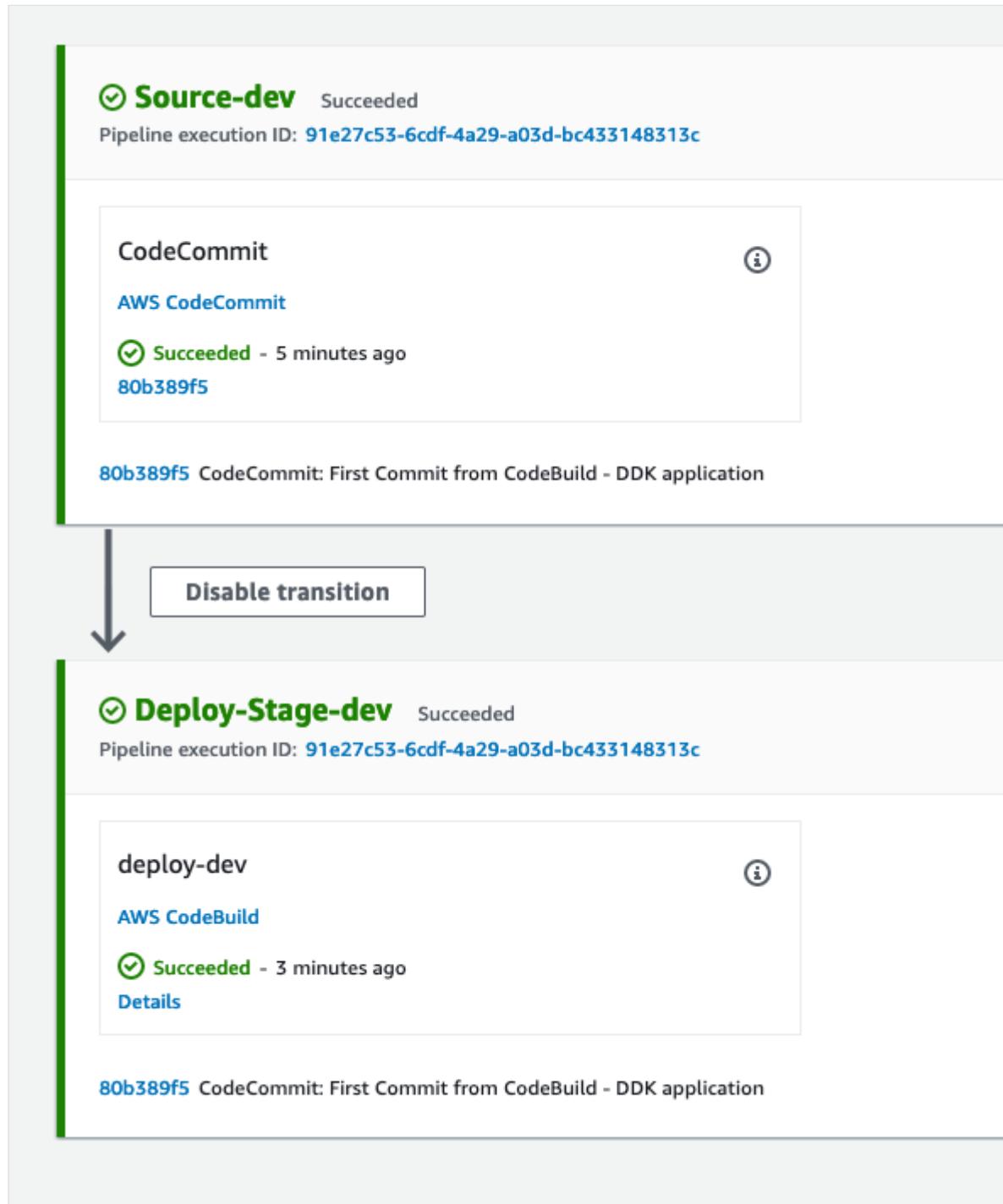
Trunk-based pipelines append one stage after the other and read from the main branch of our repository:



Gitflow strategy uses multiple CodePipeline pipelines for each of the stages. For example if you selected `dev` and `prod`:

Pipelines		Info										
Name	Most recent execution	Latest source revisions	Last executed									
dataall [REDACTED] gitflow [REDACTED] dev	⌚ In progress	CodeCommit - 80b389f5: First Commit from CodeBuild - DDK application	Just now									
dataall [REDACTED] gitflow [REDACTED] prod	⌚ In progress	CodeCommit - 80b389f5: First Commit from CodeBuild - DDK application	Just now									

The `dev` pipeline reads from the `dev` branch of the repository:



Github Template Pipelines Overview

This pipeline strategy takes a pre-defined IaC CDK Application that exists in a github repository and deploys the pipeline to be managed by data all. An AWS CodeCommit repository with the code of the github repository is created in the CICD environment AWS account.

NOTE: You may have to specify a access token in the HTTPS Clone Path of the Github Repository if the repository is private

data.all performs the initial deployment of this pipeline by running `cdk deploy` for the code now existing in AWS CodeCommit in the CICD environment. Adding development environments here is on the responsibility of the pipeline creator to align with the deployment environments specified in the cloned repository.

Create a new pipeline

Pipeline name: My Pipeline

Short description:

Tags:

CICD

- CICD Environment: consumerEnv
- Team: UserB
- Region: us-east-1
- Organization: TestOrg
- CICD strategy: GitHub Template
- GitHub Template Clone Path: HTTPS_GITHUB_CLONE_PATH

Development environments

Order	Development Stage	Environment	Team
1	dev	consumerEnv	UserB

Add environment

Create Pipeline

5.4.2 Editing a Data All Pipeline

For users who would like to promote their pipeline deployments to new environments managed by data all, you can do so by first bootstrapping the new environment(s) to be deployed to (as mentioned in the [Pre-requisites](#)) and then adding and/or editing the development environments.

Based on pipeline use case, editing a data all pipeline's development environments will:

- CDK Pipelines:** On update, the `ddk.json` and `app.py` will be edited to update the new development environment information. The self-mutating, CICD Pipeline will trigger and deploy to the new environments based on the source CodeCommit repository changes.
- CodePipelines - Trunk-based:** On update the `ddk.json` will be edited. A new `cdk deploy` will run to update the CICD CloudFormation Stack for the AWS CodePipeline to add the new stages required for the additional environment deployment(s) (as well as manual approval steps between stages in the code pipeline). You will see these updates to the CICD stack in CloudFormation of the CICD environment.
- CodePipelines - Gitflow:** On update the `ddk.json` will be edited. A new `cdk deploy` will run to update the CICD CloudFormation Stack to add the new AWS CodePipelines required for the additional environment deployment(s). You will see these updates to the CICD stack in CloudFormation of the CICD environment.
- Github Template Pipelines:** Editing development environments **will NOT** re-deploy the application or update the CodeCommit repository. Editing of template pipeline's development environment(s) is the responsibility of the pipeline creator for proper data all pipeline management.

5.4.3 Which development strategy should I choose?

CDK pipelines - Trunk-based

- The `CDK-pipelines` construct handles cross-account deployments seamlessly and robustly. It synthesizes CDK stacks as CloudFormation stacks and performs the deployment cross-account. Which means that we don't manually assume IAM roles in the target accounts, all is handled by CDK :)
- It also allows developers to modify the CICD stack as it is self-mutating. It is easy to customize having several typical CodePipeline stages out-of-the-box. For example, developers can add monitoring, tests, manual approvals directly in the repository with single-line changes.

CodePipeline pipelines - Trunk-based or GitFlow

- The `aws-codepipelines` construct uses AWS CodePipelines directly. We are able to define any type of CICD architecture, such as in this case Trunk-based and GitFlow.
- Developers working on the pipeline cannot modify the CICD pipeline
- Cross-account deployments require specific definition of the environment in the code.

Github Template Pipelines

1. The aforementioned pipeline strategies do not align with your desired pipeline architecture
2. You already have pipelines IaC written in AWS CDK and ready to be deployed rather than creating pipeline(s) and developing from scratch

Summary

CDK pipelines are recommended for flexibility and for a robust cross-account application deployment, whereas CodePipeline pipelines are recommended if you need to provide an immutable pipeline architecture or if you want to implement a GitFlow strategy.

5.4.4 Cloning the repository

Pre-requisites:

1. Install git: `sudo yum install git`
2. Install pip: `sudo yum -y install python-pip`
3. Install git-remote-codecommit: `sudo pip install git-remote-codecommit`

Clone the repo:

1. Get the AWS Credentials from the AWS Credentials button in the Pipeline overview tab.
2. Clone the repository with the command in the overview tab.

CICD

- CICD Environment: Research-CI_CD
- AWS Account: [REDACTED]
- Team: Research
- Repository name: dataall-multi-5-7th7hbmz
- Development Strategy: trunk
- Git clone: `git clone codecommit:eu-west-2://dataall-multi-5-7th7hbmz`
- AWS Credentials** button

Development environments

ID	Name	Type	Region	Last Update
1	dev	Research-DEV	Research	[REDACTED]
2	prod	Research-PROD	Research	[REDACTED]

STACK

CREATED BY: john Doe (john.doe@amazon.com)

- Team: Research
- Created: 2 hours ago
- Status: PENDING

Details

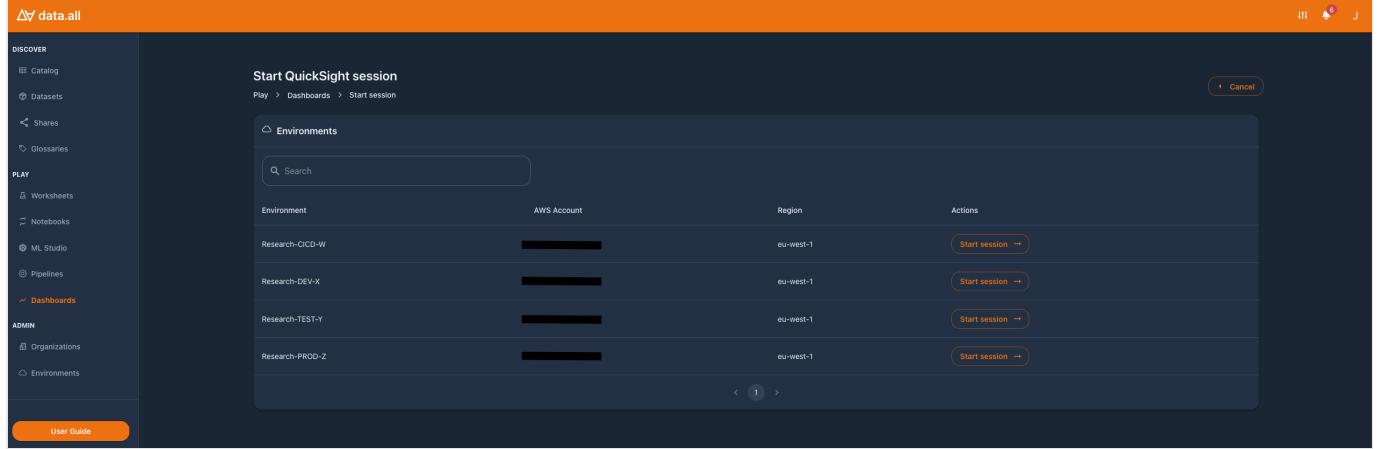
- URI: 7th7hbmz
- Name: Multi-5
- Tags: [REDACTED]
- Description: No description provided

5.5 Dashboards

Data.all connects with Amazon Quicksight to allow users to quickly visualize and analyse their data.

5.5.1 Start Quicksight session

Go to the Dashboards menu of data.all and click on the orange "Quicksight" button in the top right corner. It will redirect you to the following page, in which you can start a Quicksight session in one of your environment accounts. If *Dashboards* are not enabled in the environment, an error message will appear on the screen.



Environment	AWS Account	Region	Actions
Research-CI_CD-W	[REDACTED]	eu-west-1	<button>Start session →</button>
Research-DEV-X	[REDACTED]	eu-west-1	<button>Start session →</button>
Research-TEST-Y	[REDACTED]	eu-west-1	<button>Start session →</button>
Research-PROD-Z	[REDACTED]	eu-west-1	<button>Start session →</button>

5.5.2 Import a dashboard

Our user has been working on a Dashboard in Quicksight and wants to register it and make it available in data.all. The first step is to copy the Dashboard ID to your clipboard. You can find this ID in the Quicksight URL.



Now, go back to data.all and in the Dashboards menu, click on the *Import* button in the top-right corner. Fill in the following form and paste the dashboard ID correspondingly.

Import a QuickSight dashboard

Play > Dashboards > Import Cancel

Details

- Dashboard name
- QuickSight dashboard identifier
- Short description

Deployment

- Environment
- Region
- Organization
- Team

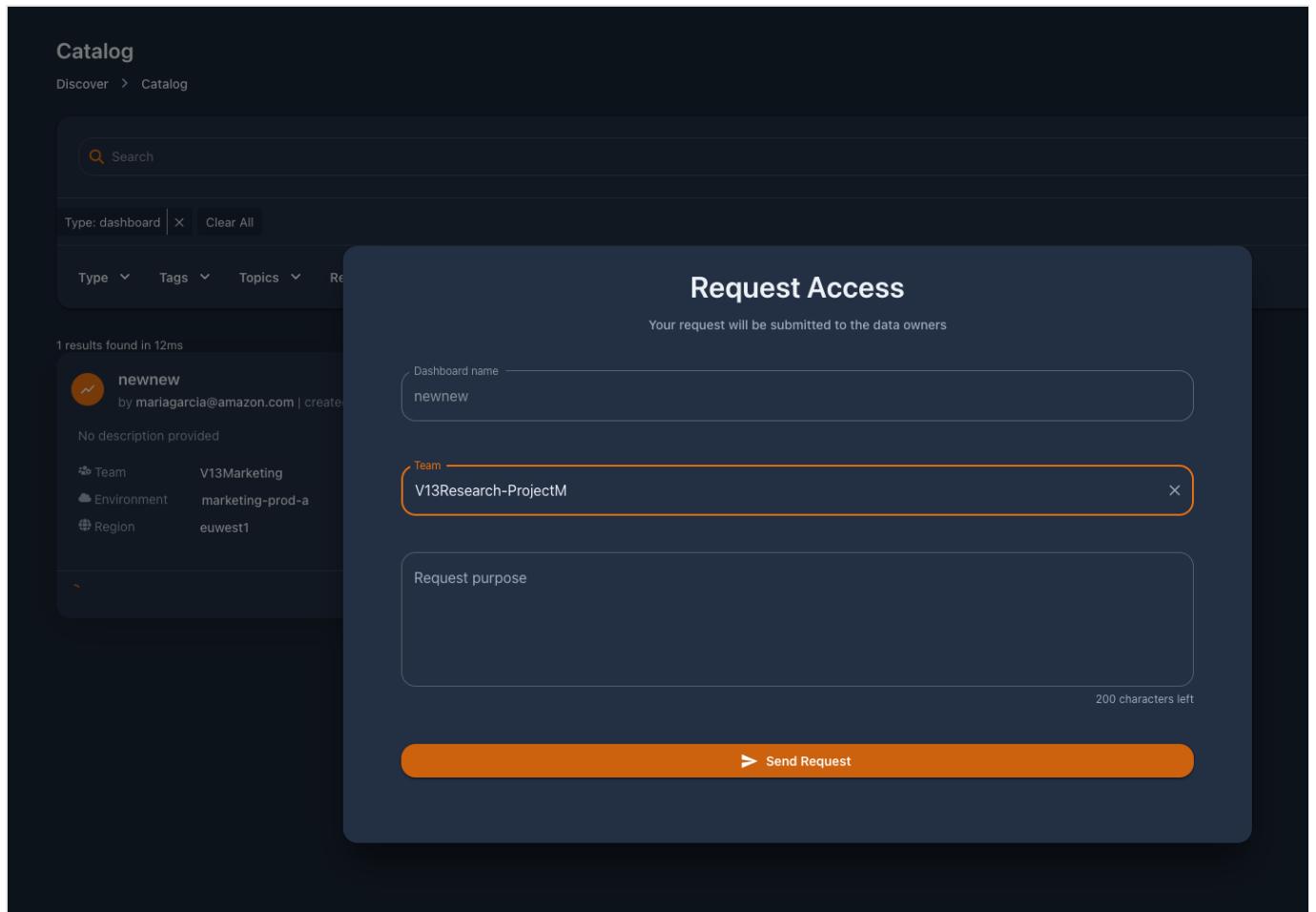
Organize

- Tags
- Glossary Terms

Import Dashboard

5.5.3 Share a dashboard

Once a dashboard is imported, it is catalogued in the central Catalog. Users can go to the Catalog, filter by dashboard and request access to a dashboard as shown in the picture below.



Once the request is open, the Dashboard owner team can accept or reject the request in the Dashboard *Shares* tab. If the request is approved, the requesters' team can visualize the Dashboard directly from data.all UI.

⚠ Session capacity pricing

To be able to embed the Dashboard using anonymous sessions, Session Capacity pricing needs to be enabled in the source Quicksight account. A Quicksight administrator in the AWS account needs to go to the **Manage Quicksight** menu and enable capacity pricing.

QuickSight

- Manage users
- Manage groups
- Manage assets
- Your subscriptions**
- SPICE capacity
- Account settings
- Security & permissions
- Manage VPC connections
- Mobile settings
- Domains and Embedding
- Account customization
- Single sign-on (SSO)
- KMS keys

Paginated Reports allows you to build highly formatted multi-page reports.

Create, schedule, share highly formatted multi-page reports and schedule data exports at scale.

[Get Paginated Reports](#)

Ask natural language questions on your data with Q

Q enables your business users to simply type in their ad-hoc business questions in plain English and get instant answers in QuickSight as a visualization.

[Get Q add-on](#)

Manage Subscriptions

Readers

With different plans to choose from, you can find one that works best for you. [Learn more](#)

Per reader pricing



Good for

- Predictable pricing for BI use cases
- Accounts with many repeat readers

Up to \$5/reader/month, \$0.30/session. Add lots of readers, pay only when they're active.

This is your current plan

Session capacity pricing



Good for

- Scalable pricing for embedded analytics and apps
- Large scale BI deployments with infrequent per reader usage
- Giving users access without the need for provisioning
- Embedding in public websites, internal applications and more

Get started for as little as \$250/month

[Get monthly subscription](#)

Lower per session pricing and a free [AWS Data Lab consult](#)

[Get annual subscription](#)

*1 session = 30 mins. from login

Authors

You've got 0 annual author subscriptions

[Add more authors](#)

6. Security

Details on data.all security-first approach to building a modern data workspace

6.1 Data and metadata on data.all

6.1.1 Data virtualization

data.all is a fully virtualized solution that does not involve moving data from existing storage layers.

Any queries run on the data.all are pushed to existing processing layers (e.g. directly to your database, warehouse, or a processing layer such as Athena or Presto on top of S3).

6.1.2 Data and metadata storage

Data and metadata collected and created by data.all are stored in applications and databases within the customer's VPC (virtual private cloud). This includes information for data previews and queries, data quality, metadata, and user data.

6.1.3 Data previews and queries

data.all gives users the ability to see sample data previews for a datasets and results for any queries run on data.all.

In both cases, the request is pushed upstream to the original data source, and a 100-row sample of the result is provided to data.all users.

6.1.4 Data quality profile

Users can generate data quality metrics with the click of a button on data.all. Once generated, these metrics are stored in PostgreSQL on the customer's VPC.

6.1.5 Metadata

Dataset metadata, including metadata generated on data.all, is stored across Elasticsearch and Aurora PostgreSQL.

Elasticsearch is used to optimize search on the product, and Aurora PostgreSQL acts as a persistence backend.

6.1.6 User data

Data on users, roles, and IdP groups is stored in a PostgreSQL database.

Any user data transmitted over the internet is SSL-encrypted over HTTPS.

6.1.7 Authentication

The data.all authentication process is based SAML 2.0-based login. data.all can also integrate into organizations' existing SAML 2.0-based SSO authentication systems.

7. Platform Monitoring

As an administrator of data.all I want to know the status of data.all. In this section we will focus on the following aspects of monitoring:

- Platform observability
- Platform usage

7.1 Observability

It refers to the infrastructure of data.all, the frontend and backend.

7.1.1 AWS CloudWatch

As part of the deployment, data.all deploys observability AWS resources with CDK and ultimately in CloudFormation. These include AWS CloudWatch Alarms on the infrastructure: on Aurora DB, on the OpenSearch cluster, on API errors... Operation teams can subscribe to a topic on Amazon SNS to receive near real time alarms notifications when issues are occurring on the infrastructure.

7.1.2 AWS CloudWatch RUM

Additionally, if we enabled CloudWatch RUM in the config.json file when we deployed data.all we will be able to collect and view client-side data about your web application performance from actual user sessions in near real time.

7.2 Platform usage

I want to know how my teams are using the platform. Inside this category we answer questions such as "*how many environments or datasets are in data.all?*".

7.2.1 RDS Queries

The first option is to query the RDS metadata database that contains all the information regarding environments, datasets and other data.all objects. You need access to the data.all infrastructure account, in which you will: 1) Navigate to RDS Console 2) Connect with secrets manager ARN 3) Get this ARN from AWS Systems Manager Parameter Store (search for "aurora") 4) Run SQL statements to extract insights about the usage of the platform

7.2.2 Quicksight enabled monitoring

When we deployed data.all, we can configure optional monitoring of Quicksight, this is the `enable_quicksight_monitoring` parameter. If enabled, we allow AWS Quicksight to establish a VPC connection with our RDS metadata database in that account. We modify the security group of our Aurora RDS database to communicate with Quicksight, then we can use AWS Quicksight to create rich dynamic analyses and dashboards based on the information on RDS. Once the deployment is complete you need to follow the next steps:

1) Pre-requisite: Quicksight Enterprise Edition We need to subscribe to Quicksight and allow data.all domain to embed dashboards, follow the instructions in the step 4 of the [Linking environment section](#).

2) Create Quicksight VPC connection

Follow the steps in the [documentation](#) and make sure that you are in the same region as the infrastructure of data.all. For example, in this case Ireland region.

The screenshot shows the AWS QuickSight user profile page. At the top, it displays the user's name, "Admin/dlpzx-Isengard", followed by a dropdown arrow. Below this, there is a summary section with the following details:

- Username:** Admin/dlpzx-Isengard
- Account name:** dlpzx-demo-dataall

Below this summary is a "Manage QuickSight" button. The main menu on the left includes the following options:

- Community
- Send feedback
- English > (with a globe icon)
- Ireland > (with a location pin icon)
- Tutorial videos
- Help
- Sign out

On the far left, there is a vertical sidebar with a checkmark icon and a timestamp "6:54". At the bottom of the main content area, it says "Showing 1 - 1 of 1 users."

To complete the set-up you will need the following information:

- VPC_ID of the RDS Aurora database, which is the same as the data.all created one. If you have more than one VPC in the account, you can always check this value in AWS SSM Parameters or in the Aurora database as appears in the picture:

RDS > Databases > dataall-dev-db

dataall-dev-db

Summary

DB cluster ID dataall-dev-db	CPU <div style="width: 14.19%;">14.19%</div>	Info Available	Current capacity 4 capacity units
Role Serverless	Current activity	Engine Aurora PostgreSQL	Region & AZ eu-west-1

Connectivity & security

Endpoint & port

Endpoint
dataall-dev-db.clust...west-1.rds.amazonaws.com

Networking

VPC
dataall-second-cicd-stack/dat...stage/backend-stack/Vpc/VP... (vpc-...)

Security

VPC security groups
dataall-dev-aurora-sg (sg-...pf)
Active

- Security group created for Quicksight: In the VPC console, under security groups, look for a group called <resource-prefix>-<envname>-quicksight-monitoring-sg. For example using the default resource prefix, in an environment called prod, look for dataall-prod-quicksight-monitoring-sg.

3) Create Aurora data source We have automated this step for you! As a tenant user, a user that belongs to DAADministrators group, sign in to data.all. In the UI navigate to the **Admin Settings** window by clicking in the top-right corner. You will appear in a window with 2 tabs: Teams and Monitoring. In the Monitoring tab, introduce the VPC connection name that you created in step 2 and click on the *Save* button. Then, click on the *Create Quicksight data source* button. Right now, a connection between the RDS database and Quicksight has been established.

Settings

Administration > Settings

MONITORING

Prerequisites

1. Enable Quicksight Enterprise Edition in AWS Account = 733017067868. Check the user guide for more details.
2. Create a VPC Connection between Quicksight and RDS VPC. Check the user guide for more details.

Create the RDS data source in Quicksight

3. Introduce or Update the VPC Connection ID value in the following box:
XXXXXX Save
4. Click on the button to automatically create the data source connecting our RDS Aurora database with Quicksight
Create Quicksight data source +

Get insights in Quicksight

5. Go to Quicksight to build your Analysis and publish a Dashboard. Check the user guide for more details.
Start Quicksight session →
6. Introduce or update your Dashboard ID
XXXXXXXXXXXX Save

4) Customize your analyses and share your dashboards Go to Quicksight to start building your analysis by clicking on the *Start Quicksight session* button. First, you need to create a dataset. Use the **dataall-metadata-db** data source, this is our connection with RDS.

Datasets

Upload a file (.csv, .tsv, .clif, .xlsx, .json)

Salesforce Connect to Salesforce

S3 Analytics

S3

Athena

RDS

Redshift Auto-discovered

Redshift Manual connect

MySQL

PostgreSQL

ORACLE

SQL Server

Aurora

MariaDB

Presto

Spark

Teradata Provided by Teradata

Snowflake

Amazon OpenSearch Ser... Successor to Amazon Elasticsearch Ser...

Exasol

GitHub

Twitter

Jira

ServiceNow

FROM EXISTING DATA SOURCES

dataall-metadata-db Updated a day ago

Use this dataset in an analysis (check the docs [customization of analyses](#)) and publish it as a dashboard (docs in [publish dashboards](#))

Not only RDS

With Quicksight you can go one step further and communicate with other AWS services and data sources. Explore the documentation for cost analyses in AWS with Quicksight or AWS CloudWatch Logs collection and visualization with Quicksight.

5) Bring your dashboard back to data.all Once your dashboard is ready, copy its ID (you can find it in the URL as appears in the below picture)

Not only RDS

With Quicksight you can go one step further and communicate with other AWS services and data sources. Explore the documentation for cost analyses in AWS with Quicksight or AWS CloudWatch Logs collection and visualization with Quicksight.

https://eu-west-2.quicksight.aws.amazon.com/sn/dashboards/**dash**

Back in the data.all Monitoring tab, introduce this dashboard ID. Now, other tenants can see your dashboard directly from data.all UI!

Settings

Administration > Settings

TEAMS MONITORING

Prerequisites

1. Enable Quicksight Enterprise Edition in AWS Account = 733017067868. Check the user guide for more details.
2. Create a VPC Connection between Quicksight and RDS VPC. Check the user guide for more details.

Create the RDS data source in Quicksight

3. Introduce or Update the VPC Connection ID value in the following box:
XXXXXX Save
4. Click on the button to automatically create the data source connecting our RDS Aurora database with Quicksight
Create Quicksight data source +

Get insights in Quicksight

5. Go to Quicksight to build your Analysis and publish a Dashboard. Check the user guide for more details.
Start Quicksight session →
6. Introduce or Update your Dashboard ID
XXXXXXXXXXXXXX Save

8. Labs

8.1 Hands-on Lab: Data Access Management with data.all teams

This document is a step-by-step guide illustrating some functionalities of the data.all "Teams" feature. This guide is far from exhaustive and mainly focuses on how users can share data across environment and teams. After completing it, you are free to continue exploring data.all and the functionalities it provides.

8.1.1 ⓘ Scope of this guide

To follow this guide, you will need:

- An AWS account (#111111111111) where data.all is deployed. Your version of data.all must support the "Teams" feature
- An AWS account that will be used as a data.all environment for the data platform team (#222222222222)
- An AWS account that will be used as a data.all environment for the data science team (#333333333333)

The scenario you will implement in this guide is the following. The data platform team owns a dataset. Data scientists are interested by the content of this dataset for their analysis. There are however two different types of data scientists, that are members of two different teams: data science team A and data science team B.

A data scientist from team A will request access to the data platform dataset for its team. A data platform user will then accept the request, thus granting team A read-only access to the dataset. Team B does not have access to the dataset from the data platform team.

Then a user in team A creates a dataset in the data science environment. We will check that users in team B does not have access to this data.

You will go through the following steps to implement this scenario:

1. Create users and groups in Cognito
2. Create the Organisation and the Environment for the data platform team
3. Create the Dataset for the data platform team and upload some data
4. Create the Organisation and the Environment for the data science team
5. Invite team A and team B to the data science environment
6. Share data platform data with team A in the data science environment
7. Create a Dataset managed by team A in the data science account

Here is an illustration of the scenario:

1. Create users and groups in Cognito

First, you need to create users and groups from the Cognito console. This happens in the account where the infrastructure of data.all is deployed (#111111111111). You will later use these users to connect to data.all. Go to the AWS console and create five groups and four users as follow:

Cognito group

- **DAAdministrators:** group for data.all administrators
- **DataPlatformAdmin:** group for data platform admin team
- **DataScienceAdmin:** group for data science admin team
- **TeamA:** first category of data scientists
- **TeamB:** second category of data scientists

Cognito users

- **data.alladmin:** create this user and add it to both in the DAAdministrators and DataPlatformAdmin groups. This user will be able to manage permissions of all teams in data.all (tanks to the DAAdministrators group membership) and will be able to create resources for the data platform team
- **ds-admin:** create this user and add it to the DataScienceAdmin group. This user will create resources for the data science team
- **ds-a:** create this user and add it to TeamA
- **ds-b:** create this user and add it to TeamB



When creating users, you will need to **provide both the name of the user and its email.**

After creating users and assigning them to groups in Cognito, you end-up with the following situation:

2. Create the Organisation and the Environment for the data platform team

We will start by creating the resources for the data platform team. Log into data.all with the user in the **DataPlatformAdmin** group. Create an Organisation for the Data Platform team. Make sure that the DataPlatformAdmin team manages this organization.

Now link a new environment to this organisation. You can do this by clicking on **Environment** and then **Link Environment**

When onboarding a new AWS account as an environment in data.all, you usually need to make some operations in the account first. The UI lists these operations for you: bootstrapping the AWS account and creating the data.allPivotRole notably. You will have to go through these operations if it is the first time you use this AWS account to create an environment in data.all. Then, create the environment by providing a name, the account ID (#222222222222) and the Team managing it (**DataPlatformAdmin**).

Wait until the stack is deployed successfully. You can check the status of the stack in the **stack** tab of the environment. Once the status is **create_complete**, create a new dataset in this environment. You can do it from the **Contribute** window

Deploy this dataset in the environment you have just created. Also make sure that the **DataPlatformAdmin** team owns this dataset (Governance section):

Wait until the dataset is created successfully. You can check the status in the **stack** tab of the dataset. Once the status is **create_complete**, you can start uploading some data from the **upload** tab:

From this window, you are able to upload files in your dataset. When uploading files, you can ask for a crawler running automatically in your dataset, thus populating a glue database. To make sure the crawler will work, please upload a csv file of your choice. Insert any name you want in the **prefix** section. This will be the name of your Glue table.

Click on the **upload** button. This puts your file in the S3 bucket related to your data.all dataset. It also launches the Glue crawler populating the Glue database. Leave some time for the crawler to run and click on the **Tables** tab. If the crawler ran successfully, clicking on the **synchronize** button will display your table. At this point, feel free to explore your table from the data.all user interface.

We have completed all the tasks on the data platform side. This included the creation of the organisation, the environment, the dataset and the upload of a csv file. This is an illustration of where we are in the process:

Note: As you may have already noted down at the beginning of this guide, the data platform user is also part of the **DAAdministrators** group. Being part of this group enables this user to manage the permissions of all the other teams in data.all. To do so, click **Setting**. This provides the list of teams for which you can manage the permissions.

Click on the icon next to the team's name to manage its permissions

This opens a new window from where you can manage all permissions of the team.

3. Create the Organisation and the Environment for the data science team

You will now create data.all resources for the data science team. Log into data.all with the user in the **DataScienceAdmin** group. Create an Organisation for the data science team. Make sure that the **DataScienceAdmin** team manages this organization.

Now link a new environment to this organisation. Provide a name for this environment, the AWS account ID (#333333333333), and the team owning it (**DataScienceAdmin**).

You now have an organisation and an environment managed by the **DataScienceAdmin** team. The next step is to invite team A and team B to this data science environment. This will enable data scientists from team A and team B to access the environment.

4. Invite Team A and Team B to the data science environment

With the user in DataScienceAdmin team, select the data science environment and click on the **Teams** tab. You can invite other teams in your environment with the **invite** button.

This opens a new window asking you to indicate the name of the team you want to invite. You can also manage the permissions this team will have in your environment. Use this **invite** button to invite **TeamA** and **TeamB** in your data science environment.

Users from team A and team B now have access to your environment

5. Share data platform data with Team A in the data science environment

Log into data.all with user in **TeamA**. This user does not own any data, but wants to access data of the data platform team. Go on data.all **Data Catalog** tab. This shows all the datasets and tables you can request access to. There are different tools you can use in order to find the data you are looking for (you can find more information about these in the data.all documentation):

- Directly typing the name of the dataset or the table in the search bar
- Use tags or topics associated to the datasets
- Use data.all Glossary

In this case, data scientist in team A wants to access **mydpdata** uploaded by the data platform team. Use the search bar to find the data. When you see the table you want, click on **Request access**.

This button opens a new window where you can configure your request. When you share a dataset or a table in data.all, the share occurs at an environment and team level. You therefore need to indicate for which environment and for which team you make the request. In this case, the user in team A wants to access data in the data science environment. Fill the request accordingly.

When you click on **Send Request**, this does not directly send the request to the data platform team. It rather creates a Draft that you can still edit in the **Collaborate** tab, under **Sent**. Click on the **Submit** button to send the request

Now re-open a new data.all window connected as the user in the **DataPlatformAdmin** team. This team owns the dataset and is therefore responsible of accepting access requests. It is possible to delegate this right to other teams using **Data Stewards** but we did not set this up in this guide. Under **Collaborate** and **Received**, you can find all the access requests received by the data platform team. Locate the request you just made with the user in Team A. If you want to know more about this request (who is making it, for which table in the dataset,...), click on **Learn More**. If you agree to grant access, click on **Approve**.

This action triggers an ECS task that updates the permissions of the table in Lake Formation. Users in Team A are now able to access the data platform data. Let us verify it.

Re-open data.all connected as the user in **TeamA**. You can first visit the **Contribute** tab where you will see the dataset that has been shared with team A.

Quick reminder: The data platform team agreed to share the **mydpdata** table with TeamA in the data science environment called **DSENV**.

As a conclusion, the table **mydpdata** is accessible from the environment **DSENV**, through a role only team A can assume. Team A users can assume this role directly from the data.all user interface. Select the data science environment and go under the **Teams** tab. You will then see all the teams that have access to the environment. Find TeamA line and click on the AWS logo.

This opens a new window in the AWS console. The AWS account is the one you associated to the data science environment earlier in this guide (#333333333333). Also note that you are assuming a role specific to your team. Use the search bar to get to the Athena console. In the Athena Query editor, you will be able to see under **Database** the dataset shared by the data platform team. The name of this database is a concatenation of "dh" (for data.all), the name of the dataset (dpdataset) and random characters to ensure unicity. Under **Tables**, you can now see **mydpdata** which you can query using with Athena.

Explanation: When the data platform team uploaded the csv file under the **mydpdata** prefix, the crawler created a new Glue table called **mydpdata** in the AWS account associated to the data platform environment (#333333333333). When the data platform team accepted to share **mydpdata** with team A in the data science environment, it triggered an ECS task that updated the Lake Formation (AWS service managing data access) settings in both the data platform and data science

environments. It updated the settings in a way that allows the IAM role of team A in the data science environment to read the **mydpdata** table stored in the data platform environment. In short, only the role of team A in the data science environment is able to read **mydpdata** table (in addition to data platform team of course). This is a **read-only** access, and the data is not moved from the data platform environment to the data science environment.

You can repeat the same thing to check that team B does not have access to the data. Log into data.all with a user in **TeamB** and select the data science environment. Under the **Teams** tab, click on the AWS logo to connect to the AWS console assuming the role of **TeamB**. Go to the Athena Query Editor and you will see that you won't be able to see data shared with team A.

At this stage of the guide, you should better understand how data sharing cross account works. The graph below illustrates where we are in the original scenario.

6. Create a Dataset managed by team A in the data science account

In the previous section of this guide, you went through an example of how you can share data across environment and teams. In this section, we will focus on the creation of datasets in a single account. Team A will create a dataset in the data science environment. We will make sure that other teams invited to the data science environment (teamB) are not able to access the dataset of team A.

Open data.all with a user in **TeamA**. In the **Contribute** section, create a new dataset in the data science environment. Make sure that **TeamA** owns this dataset.

When the dataset is fully created, upload a csv file from the **upload** tab of the dataset. Upload this file under a prefix named **datateama** to create a new Glue table with the same name. After uploading the file, wait a few minutes to let the crawler do its job. In the **Tables** section, click on **Synchronize** to display your new table.

Now that the data is uploaded, team A is able to access the data as it is registered as the owner of the dataset. However, team B is not able to read the data even if it has access to the environment. If you log into data.all with the user in team B, you won't be able to see the **TeamADataset** in the **Contribute** section. In addition, you will find below two screenshots of the Athena console. In the first screenshot, we assume the role of **TeamA** in the data science environment (process already explained in the previous section). In the second screenshot, we assume the role of **TeamB** in the data science environment. When assuming the role of team A, we can see the team A dataset in the **database** section, and also the **datateama** table. We can then query the data with Athena. However, when assuming the role of team B in the data science environment, we are not able to see any dataset. This is because in this guide, we have not created or shared any dataset with team B. Team B is thus unable to query the data of team A.

This last section illustrated how you can use teams to manage data access in a single environment. You have reached the end of the guide that illustrated some capabilities that data.all brings. Now that you got the basis, feel free to explore all the other things you can do with your data.

Cleanup

When you are done with this guide, you delete your data.all resources (dataset, environment, organization). This also automatically deletes the Cloudformation stacks created in your AWS accounts.