

Utiliser le Datalake => utiliser le protocole S3



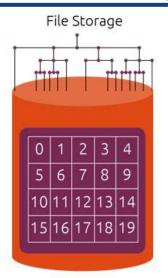


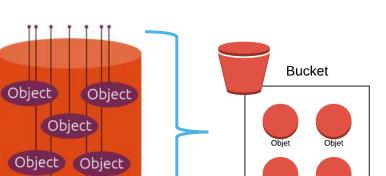
AWS S3 (Simple Storage Service): stockage objet

Object Storage

Object

Object

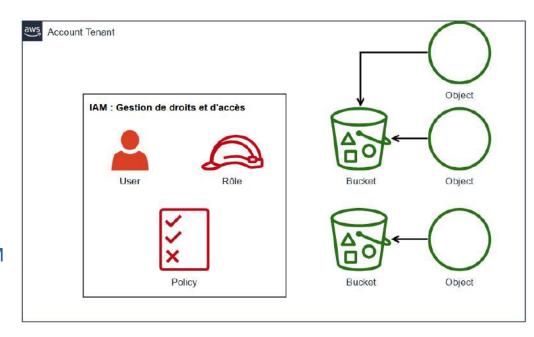






Concepts Importants

- Organisation des données :
 - Tenant / Projet (Account au sens AWS)
 - Buckets à l'intérieur des tenants
 - Des objets placés à plat au sein des buckets
- Les accès via entités IAM :
 - Policy : les politiques d'accès
 - User/Group : utilisateurs/groupes IAM authentifiés avec une paire de clé
 - Rôle: Entité détenant des droits d'accès à des ressources (bucket/objet) et pouvant être endossé temporairement



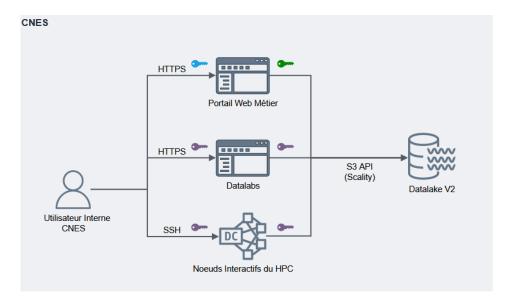


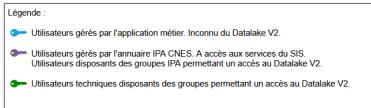
Credentials S3 == paire de clés

- aws_access_key_id
- aws_secret_access_key
- aws_session_token (si utilisation d'un rôle)

Différents modes d'authentification

- Via credentials permanents:
 - IAM User (pas lié au compte IPA)
 - À destination des projets / producteurs de données
- Via credentials temporaires (max 12h):
 - En endossant un rôle grâce à son compte IPA
 - · À destination de tout les utilisateurs









Les services AWS

Uniquement 3 services AWS sont disponibles sur le Datalake (gestion stockage S3)



Gestion de ressources
[bucket et object]



Gestion de entités IAM
[user, group, role et policy]
Chaque account possède son propre IAM



Gestion de credentials temporaires avec une durée de 12h

Le Datalake respecte les APIs standards AWS

la documentation AWS S3 est donc une excellente base pour commencer



Accès au Datalake

- Un endpoint : https://s3.datalake.cnes.fr

- Une région : us-east-1

S'authentifier au Datalake

- Paire de clés aws_access_key_id et aws_secret_access_key
- Token aws_session_token (si on utilise un rôle)

Configuration des credentials

- Variables d'environnement AWS_ACCESS_KEY_ID AWS_SECRET_ACCESS_KEY AWS_SESSION_TOKEN
- Profils en configurant les fichiers ~/.aws/credentials et ~/.aws/config
- Directement dans le client (NON recommandé)

Les credentials permettent d'accéder et d'avoir la visibilité aux ressources d'un seul Account.

Les accès sont limités aux autorisations fournis via les politiques (IAM policy).

Par défaut, aucune autorisation (implicit deny).



Le Datalake du CNES : Usage du S3

Quel est le périmètre d'un utilisateur ?

Dans son périmètre



- Connaitre les principes de la gestion des droits
- Pouvoir demander si son besoin le nécessite des clés permanentes (Read-only, Read/write)



- Connaitre les principes
- Utiliser les scripts présents dans /softs/tools/datalake/ ou librairie assumerole



- Choisir son client S3
- Paramétrer ses authentifications (paire de clé)
- Manipuler les objets
- Tier 3 : Restaurer les objets et superviser le tier de l'objet

Hors périmètre (équipe Datalake)

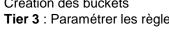
Création de IAM User, Rôle ou policy

Paramétrer les rôles

- en gestion d'accès (Policy)
- en gestion d'autorisation pour les comptes IPA

Création des buckets

Tier 3 : Paramétrer les règles ILM





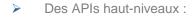




Clients S3 et particularités du Datalake

Toujours utiliser un client S3 reconnu ou librairies de traitements

Utiliser un client S3 qui permet de manipuler les données



- Simuler les commandes POSIX (ls, cp ...)
- Uploads/downloads partitionnés
- Des APIs bas niveaux : get-object, put-object, head-object, listobject...



- très performant en lecture mais comparativement peu performant en écriture.
- organisation des données : metadonnées personnalisable, tags ...
- formats de données plus ou moins adaptés (en terme de performance) pour le S3

SDKs



































Les opérations autorisées (GetObject, PutObject ...) dépendent des autorisations acquises par l'entité IAM (credentials) utilisée

