

Introduction à iRODS

Formateurs : Emmanuel Medernach (IPHC) et Jérôme Pansanel (IPHC)
avec l'appui du groupe iRODS de France Grilles

iRODS

Un *middleware* pour la gestion des données

- Open Source (licence BSD)
- Supportant plusieurs milliers d'utilisateurs et de groupes
- Permettant d'accéder, de gérer et de partager des données stockées sur différents types de stockage
- Facilitant l'accès à des ressources hétérogènes (Unix, S3, DDN, HPSS, ...), à travers un seul espace de nom (zone)
- Permettant de contrôler finement les données grâce à un moteur de règles et un ensemble de micro-services (réplication, vérification des types, ...)

Une solution puissante qui rend possible

- La virtualisation de l'accès au stockage
- La gestion de plusieurs péta-octets de données
- Le transfert parallèle pour les données volumineuses
- La recherche des données (méta-données)
- L'automatisation des processus grâce aux règles et aux micro-services
- La sécurisation des données grâce à la réplication et la gestion des accès

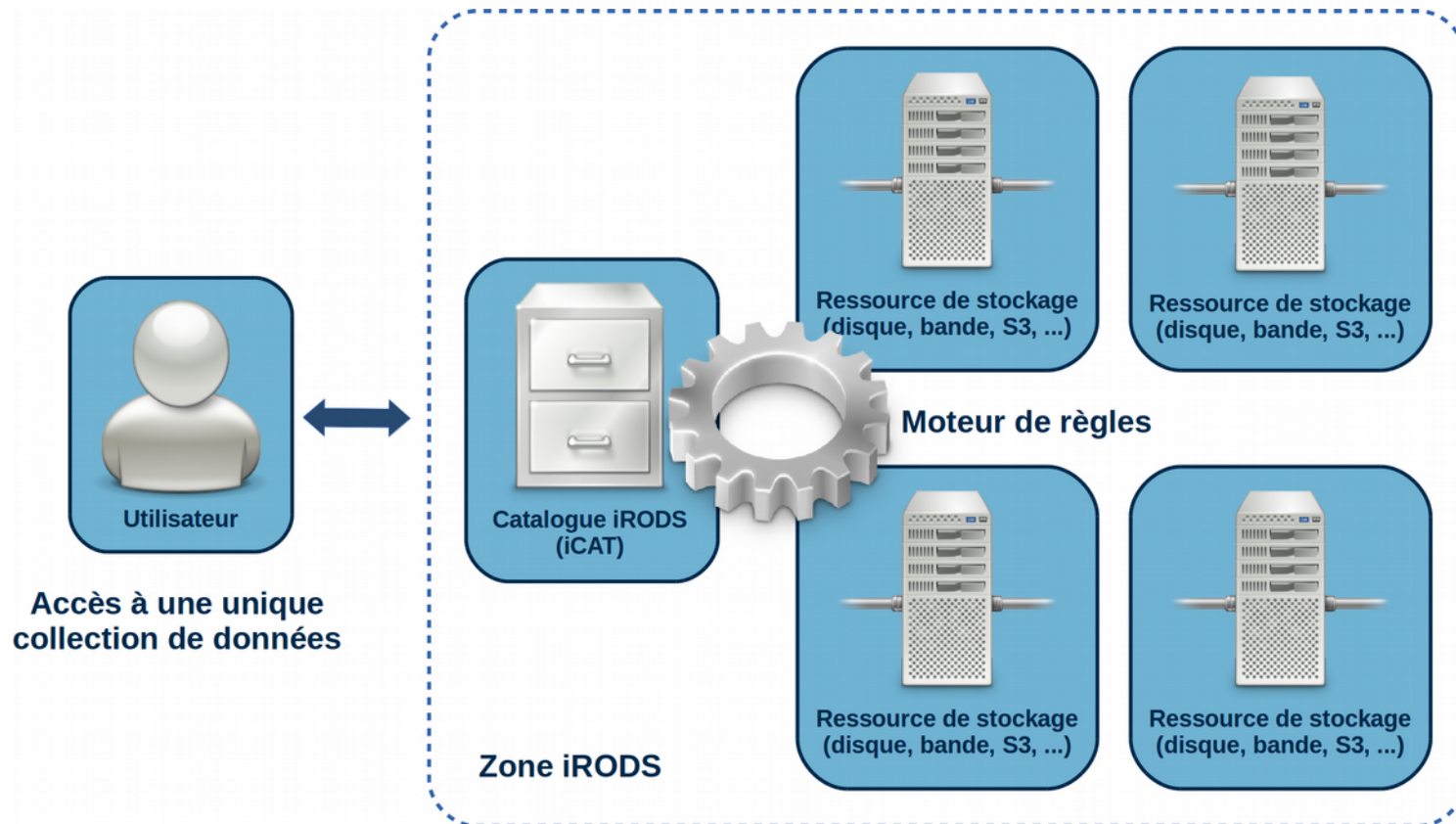
Introduction

Un peu d'histoire :

- 1995 – Démarrage du projet SRB (*Storage Resource Broker*), un système de stockage géographiquement distribué (*data grid*). Ce projet est piloté par DICE et UCSD.
- 2006 – Début le développement d'iRODS, successeur de SRB. Il intègre les concepts de SRB, mais avec une licence libre et un outil de gestion de règles.
- 2013 – Création du consortium iRODS avec les acteurs historiques et RENCI.
- Le consortium intègre de nombreux nouveaux membres (Bayer, DDN, EMC, IBM, Novartis, Seagate, ...)

The logo for iRODS, featuring the lowercase letter 'i' in a teal color and the uppercase letters 'RODS' in a dark grey color, all in a bold, sans-serif font.

Fonctionnement



iRODS

Une zone iRODS

- Comporte un catalogue de méta-données (iCAT)
- Une base de données stockant les méta-données
- Un ensemble de ressource de stockage
- La version des éléments dans une zone doit être homogène

Le serveur iCAT stocke les informations

- Sur la zone
- Les données et les méta-données associées
- Le système de fichier virtuel
- La configuration des ressources
- La base des utilisateurs

iCAT supporte actuellement les bases PostgreSQL, MySQL et Oracle.

iRODS

Objet de données (Data Object) et méta-données

- Le terme objet de données (ou élément de données) est une représentation logique d'une donnée stockée sur une ou plusieurs ressources de stockage
- Ces objets sont stockés dans des Collections(~ répertoire)
- La hiérarchie est séparée par des '/'
- Par exemple, la collection racine de la zone **tmpZone** est /tmpZone
- Une hiérarchie de base est la suivante :

```
/tmpZone/  
|-- home  
`-- trash
```
- Le chemin complet inclus le nom de la zone
- iRODS peut stocker des informations complémentaires (méta-données) à propos des objets de données (clé*, valeur*, unité)
- Il est possible de retrouver des données en effectuant des recherches par méta-données (CLI, Web GUI)

Les règles

Le moteur de règle

- Automatisation des processus
- Langage spécifique (*C-like*)
- Programmation procédurale avec fonctions, commentaires, opérations sur les types, ...
- Possibilité d'utiliser des expressions régulières
- Possibilité d'utiliser le *Language Integrated General Query* (LIGQ) pour rechercher directement des informations dans le catalogue
- Utilisation de micro-services
- Depuis iRODS v4.2, possibilité d'utiliser le moteur Python

```
HelloWorld {  
    writeLine("stdout", "Hello, world!");  
}
```

Cas d'étude

Points liés aux données

- Où sont stockées les données ?
- Faut-il les centraliser ou les distribuer (avantage / inconvénient) ?
- Faut-il stocker plusieurs copies des données ?
- Quelle quantité de données doit être stockée ? Quelle évolution à prévoir ?
- Quelle est le format des fichiers (ouvert / propriétaire) ?
- Est-ce que les données sont confidentielles ?

Points liés au réseau

- Quel est le débit des réseaux ? Est-il stable dans le temps ?
- Quelle confiance pouvons-nous avoir ?

Points liés aux ressources

- Est-ce que le projet de déploiement intègre un budget pour la partie logicielle ?
- Quelle est la position des décideurs par rapport à l'Open Source ?
- Quelles sont les ressources disponibles pour l'acquisition, la gestion et la maintenance du service ?

Cas d'étude

Points liés à l'organisation

- Quelle est l'organisation / gouvernance du projet
- Quelles sont les libertés décisionnelles et opérationnelles des participants du projet ?
- Qui est responsable ? Qui pilote ?

Points liés aux utilisateurs

- Quels types d'utilisateurs accèdent aux données ? Quels sont les privilèges à leur accorder ?
- Combien d'utilisateurs utiliseront le service ? Est-ce que ce nombre va croître, et dans quelle proportion ?
- Où sont situés les utilisateurs ? Quel support devons-nous leur apporter ?

Besoin en matériel

iRODS est une solution légère

- Les paquets iRODS occupent moins de 100 Mo (mais les dépendances ont besoin de plus d'espace)
- Serveur iCAT avec minimum 2 Go de RAM
- La taille de la base de données est dimensionnante pour le choix du serveur
- Dans un environnement de production, il est conseillé d'utiliser une base de données (sous forme cluster) en dehors du service
- Pour une ressource de 40 To, un bi-pro AMD avec 32 Go de RAM permet de saturer deux liens 1 Gb/s

Alternatives

Outils avec des fonctionnalités proches

- OneData – <https://onedata.org>
- Dspace – <http://dspace.org>
- DIRAC File Catalog – <http://diracgrid.org/>