

Comment faire fonctionner le cluster

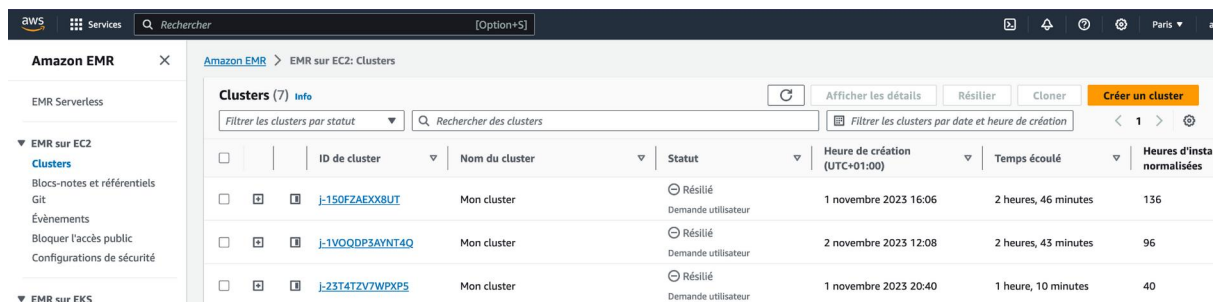
Créer/reproduire le cluster EMR

Pour créer un cluster avec EMR :

- aller sur la page EMR en tapant « EMR » dans la barre de recherche
- choisir une option parmi celles proposées ; pour notre projet, nous avons opté pour un EMR s'exécutant sur Amazon EC2
- cliquer sur « Créer un cluster »

S'ouvre alors la page de configuration de l'EMR.

Pour reproduire un EMR précédemment créé puis résilié, aller dans la rubrique **EMR sur EC2 > Clusters** (cf. image ci-dessous), cocher le cluster que l'on veut reproduire, puis cliquer sur cloner. S'ouvre alors la page de configuration de l'EMR.



Pour un EMR identique à celui que j'ai créé, choisir les réglages suivants :

Attention : lorsque l'on tente de cloner un cluster qui avait été créé avec ces réglages, certains se remettent automatiquement sur les valeurs par défaut et doivent être redéfinis. J'ai coloré les réglages en question en rouge.

- version d'EMR : emr-6.14.0
> La version la plus récente au moment où j'ai créé mon EMR. Lorsque de nouvelles versions seront disponibles, il est possible que ce réglage fasse partie de ceux qui reviennent automatiquement sur une autre valeur au moment de cloner le cluster.
- choisir l'offre « Custom » et sélectionner les packages Spark, Hadoop, JupyterHub et TensorFlow
- comme option du système d'exploitation, choisir « Version Amazon Linux »
- configuration du cluster : sélectionner « Groupes d'instances »
et « m5.xlarge » pour le type d'instance EC2
- dimensionnement et mise en service du cluster : sélectionner « Définir manuellement la taille du cluster »
taille de l'instance(s) : entrer « 1 » pour l'unité principale, et « 2 » pour la tâche 1.

- Pour l'amorçage, charger sur S3, dans le bucket du projet, un fichier .sh qui contienne ces lignes :

```
sudo python3 -m pip install -U setuptools
sudo python3 -m pip install -U pip
sudo python3 -m pip install wheel
sudo python3 -m pip install pillow
sudo python3 -m pip install pandas==1.2.5
sudo python3 -m pip install pyarrow
sudo python3 -m pip install boto3
sudo python3 -m pip install s3fs
sudo python3 -m pip install fsspec
sudo python3 -m pip install jupyterlab
sudo python3 -m pip install jupyter
sudo python3 -m pip install notebook
sudo python3 -m pip install pyspark
sudo python3 -m pip install tensorflow
```

Puis, dans la rubrique Actions d'amorçage de la page de configuration du cluster, cliquer sur « Ajouter », donner un nom au fichier d'amorçage et entrer son emplacement sur S3 (s3://nom-du-bucket/nom-du-fichier si le fichier se trouve au niveau racine, dans aucun dossier).

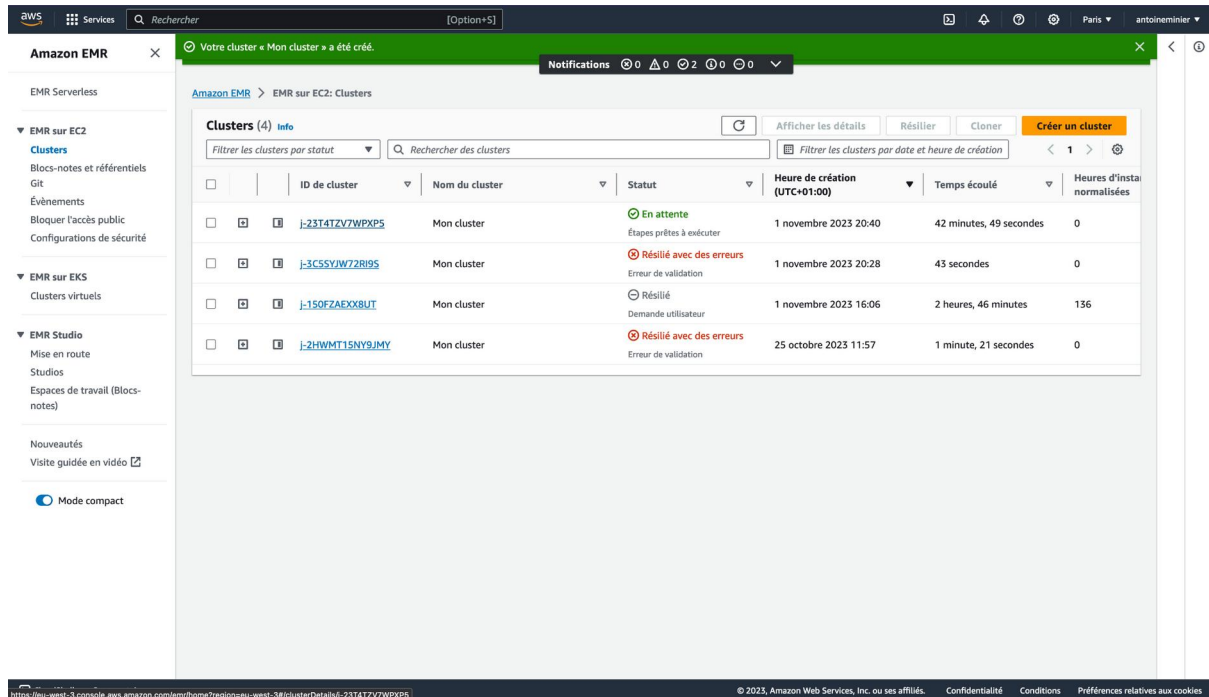
- Paramètres logiciels : sélectionner « Entrer la configuration » et entrer le code suivant dans l'espace dédié :

```
[
{
  "Classification": "jupyter-s3-conf",
  "Properties": {
    "s3.persistence.bucket": "aminier-p8",
    "s3.persistence.enabled": "true"
  }
}
]
```

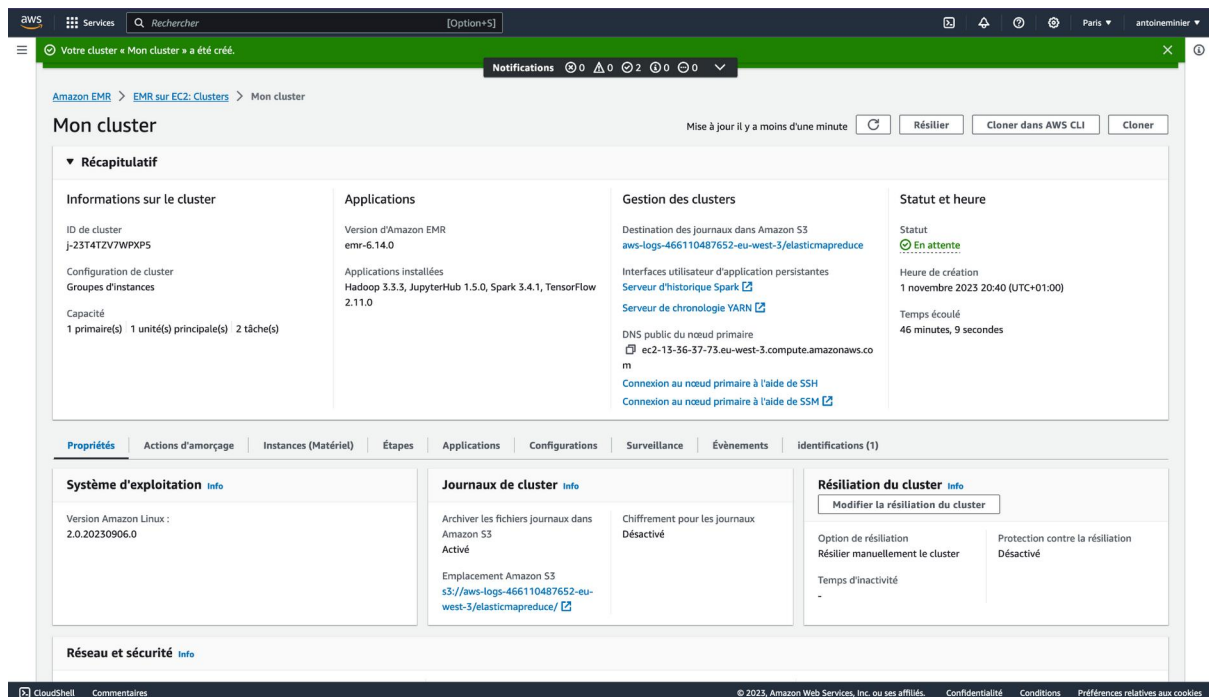
- Configuration de sécurité et paire de clés EC2 : dans « Paire de clés Amazon EC2 pour SSH sur le cluster » cliquer sur parcourir et sélectionner la paire de clés que l'on aura préalablement créée (pour cela, suivre les indications de la rubrique dédiée du cours Réalisez des calculs distribués sur des données massives, chapitre [Déployez un cluster de calculs distribués](#))
- Rôle Identity and Access Management (IAM) :
 Dans « Fonction du service Amazon EMR », sélectionner « Créez un fonction du service » et non pas « Choisir une fonction du service existant ».
 Puis dans « Profil d'instance EC2 pour Amazon EMR », sélectionner « Choisir un profil d'instance » et non pas « Choisir un profil d'instance existant ». Juste en-dessous, dans « Accès au compartiment S3 », sélectionner « Tous les compartiments S3 de ce compte avec accès en lecture et en écriture ».

Utiliser le cluster EMR une fois qu'il est activé

Sur la page des clusters EMR (cf. l'image ci-dessous), cliquer sur l'ID du cluster que l'on veut utiliser.



On arrive sur la page du cluster (cf. l'image ci-dessous).



Par défaut, on se trouve dans l'onglet « Propriétés » ; cliquer sur l'onglet « Applications ». Descendre à la rubrique « Interfaces utilisateur d'application sur le nœud primaire », puis cliquer sur “Activer une connexion SSH” en haut à droite de la rubrique.

Connexion au nœud primaire à l'aide de SSH

Propriétés

Actions d'amorçage

Instances (Matériel)

Étapes

Applications

Configurations

Surveillance

Événements

Identifications (1)

Interfaces utilisateur d'application

Les applications installées sur votre cluster Amazon EMR publient des interfaces utilisateur en tant que sites web. Vous pouvez les utiliser pour surveiller l'activité du cluster.

Interfaces utilisateur d'application sur le cluster

Les interfaces utilisateur sur le cluster sont disponibles uniquement pendant l'exécution de votre cluster. Utilisez les liens suivants pour démarrer. Pour accéder à toutes les interfaces utilisateur d'application, configurez le tunneling SSH.

Interfaces utilisateur d'application persistantes

Les interfaces utilisateur persistantes ne nécessitent pas de tunneling SSH. Elles sont hébergées hors du cluster et sont disponibles pendant 30 jours après la fin d'une application.

Interfaces utilisateur d'application en direct

Ces interfaces utilisateur d'application sur cluster sont disponibles sans tunneling SSH.

Interfaces utilisateur d'application

Interface utilisateur du serveur d'historique Spark

Interfaces utilisateur d'application sur le nœud primaire

Celles-ci nécessitent l'activation du tunneling SSH.

Activer une connexion SSH

Application	URL de l'interface utilisateur
Gestionnaire de ressources	http://ec2-13-36-37-73.eu-west-3.compute.amazonaws.com:8088/
JupyterHub	https://ec2-13-36-37-73.eu-west-3.compute.amazonaws.com:9443/
Nom du nœud HDFS	http://ec2-13-36-37-73.eu-west-3.compute.amazonaws.com:9870/
Serveur d'historique Spark	http://ec2-13-36-37-73.eu-west-3.compute.amazonaws.com:18080/

Interface utilisateur d'application sur les nœuds principaux et nœuds de tâches

Application	URL de l'interface utilisateur
Gestionnaire de nœuds	http://ec2-000-000-000-000.compute-1.amazonaws.com:8042/
Nœud de données HDFS	http://ec2-000-000-000-000.compute-1.amazonaws.com:9864/

Copier la commande SSH surlignée en gris.

Activer une connexion SSH

Les applications EMR publient des interfaces utilisateur sous forme de sites Web hébergés sur le nœud primaire. Pour des raisons de sécurité, ces sites Web ne sont disponibles que sur le serveur Web local du nœud primaire.

Pour accéder aux interfaces Web, vous devez établir un tunnel SSH avec le nœud primaire à l'aide d'une redirection de port dynamique ou locale. Si vous utilisez la redirection de port dynamique, vous devez également configurer un serveur proxy pour afficher les interfaces Web. [En savoir plus](#)

Étape 1: Ouvrez un tunnel SSH vers le nœud primaire Amazon EMR.

Windows

Mac/Linux

1. Ouvrez une fenêtre de terminal. Sur Mac OS X, sélectionnez **Applications > Utilities > Terminal**. Sur les autres distributions Linux, le terminal se trouve généralement dans **Applications > Accessoires > Terminal**.

2. Pour établir un tunnel SSH avec le nœud primaire à l'aide de la redirection de port dynamique, entrez la commande suivante. Remplacez `~/aminier-p8-ec2-rsa.pem` par l'emplacement du fichier de clé privée (.pem) utilisé pour lancer le cluster.

```
ssh -i ~/aminier-p8-ec2-rsa.pem -ND 8157 hadoop@ec2-15-188-48-170.eu-west-3.compute.amazonaws.com
```

Remarque : le port 8157 utilisé dans la commande est un port local non utilisé sélectionné au hasard.

3. Saisissez yes (oui) pour ignorer l'avertissement de sécurité.

Étape 2: Configurez un outil de gestion de proxy.

Fermer

contenant la paire de clé pour la connexion SSH.

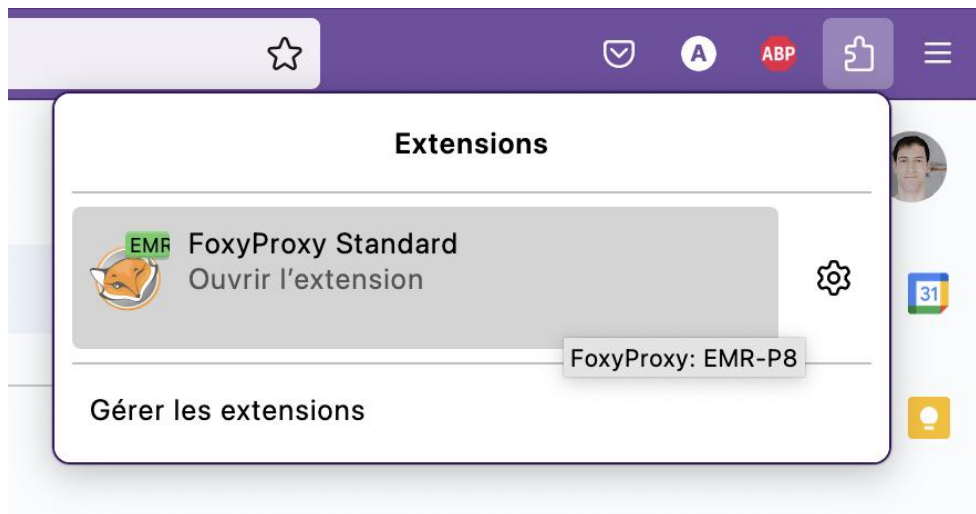
Puis coller la commande précédemment récupérée d'Amazon pour la connexion SSH.
Enlever de cette commande ~/ et le N dans -ND pour n'avoir plus, dans notre exemple, que :

```
ssh -i aminier-p8-ec2-rsa.pem -D 8157 hadoop@ec2-13-36-37-73.eu-west-3.compute.amazonaws.com
```

Changer aussi le port selon ce que l'on souhaite ; par exemple 5000 à la place de 8157.
Exécuter la commande, répondre **yes** à la question qui s'affiche en retour.

Si l'on utilise Firefox, installer l'extension FoxyProxy, et la régler selon les indications données dans le notebook donné en ressource pour le P8 (si l'on n'utilise pas Firefox, utiliser un autre proxy). Indiquer le même port que celui que l'on aura choisi d'utiliser pour la commande SSH.

Une fois installée et réglée, dans Firefox, dans les extensions, cliquer sur FoxyProxy.



Puis activer le proxy en cliquant sur le nom que l'on aura défini lors du réglage, en vert (dans l'exemple, « EMR-P8 »).



Revenir ensuite sur la page du cluster, et cliquer sur le lien de l'url qui correspond à JupyterHub.

Passer les avertissements de sécurité.

S'ouvre alors l'interface JupyterHub avec le contenu du dossier Jupyter de notre bucket S3.

Si des identifiants sont d'abord demandés, entrer les identifiants par défaut : `XXXX` en nom d'utilisateur, et `jupyter` en mot de passe.

Une fois le cluster résilié, désactiver le proxy.