

Une expérience des conteneurs sur HPC. Mise en œuvre de la chaîne de traitement iota2 et de ses applications connexes

AUDA Yves¹, AZEMA Philippe¹, BIGOT Jérôme¹, CABANAS Laurent⁴, ESCOBAR Joan⁴, GAZEN Didier⁴, GONDET Etienne⁵, GEGOUT Pascal¹, GRIPPA Manuela¹, HAGOLLE Olivier², HIERNAUX Pierre¹, INGLADA Jordi², KERGOAT Laurent¹, PENA LUQUE Santiago², Arthur VINCENT³

1. GET - Géosciences Environnement Toulouse, France
 2. CESBIO - Centre d'Etude Spatiale de la Biosphère, Toulouse, France
 3. CS-SI - Toulouse, France



Résumé

L'analyse d'image satellite requiert des puissances de calcul et des espaces de taille supérieure à ceux disponibles sur les ordinateurs personnels. Par exemple, un suivi annuel de l'occupation du sol sur une région de 200 x 200 km à l'aide d'images Sentinel-2 demande de traiter 5 TO de données. Seul un cluster équipé des dernières technologies dispose des espaces de stockage et des puissances de calcul adaptés à ces volumes de données.

La chaîne de traitement mise en œuvre est constituée d'un ensemble de logiciels relevant des Systèmes d'Information Géographique et de l'analyse des données (iota2, Orfeo Toolbox, GRASS, R du CRAN...) dont l'installation nécessite des environnements spécifiques (bibliothèques gdal, r-proj4, openCV...). La multiplicité des versions des logiciels et bibliothèques, leurs incompatibilités et leurs fréquentes évolutions s'avèrent un casse-tête pour les administrateurs système des clusters. Une solution est de recourir à la conteneurisation.

Nous présentons une solution qui est mise en œuvre sur poste de travail Ubuntu, puis déployée sur le cluster du laboratoire Géosciences Environnement de Toulouse géré par le laboratoire d'Aérodynamique (Dell R740 Dual-Socket Xeon Gold 6154 18 cœurs Skylake SP@3.0Ghz per processor for a total of 36 cores and 192GB of memory). Le système de conteneur utilisé est singularity ; L'ordonnanceur est slurm.

Les avantages et inconvénients de ces solutions techniques sont discutés en terme de performance, et fonctionnalité.

Objectifs

La chaîne iota2 est développée par le CESBIO/CNES. Elle consiste à interpoler selon une grille temporelle et spatiale régulière des images satellite (Sentinel-2, Landsat 7, Landsat 8 et Sentinel-1) puis à analyser cette grille par machine learning. Cette chaîne s'appuie pour une grande partie sur des modules de l'Orfeo ToolBox.

Cette chaîne constitue un cas d'école pour tester la performance d'un environnement de calcul en terme de performance et de portabilité des codes.

Architecture du cluster

Les nœuds du laboratoire Géosciences Environnement Toulouse sont mutualisés par leur intégration au cluster géré par le Laboratoire d'Aérodynamique. Le nœud utilisé est composé d'un Dell R740 Dual-Socket Xeon Gold 6154 18 cœurs Skylake SP@3.0Ghz par processeur pour un total de 36 cœurs et 192 GB de mémoire.

Le système d'exploitation est OpenSuSE, L'ordonnanceur de tâches est slurm.

Architecture logicielle

La chaîne iota2 est disponible sous la forme d'un paquet conda. Les capacités des moyens de calcul sont optimisés par le multithreading de la plupart des traitements.

Nous avons intégré cet environnement conda dans un conteneur singularity pour réaliser la mise au point de notre

Les données

Dans cet exemple, les données Sentinel-2 pour la saison 2018 (28 dates) qui couvrent deux tuiles (T31PDR et T31PDQ) occupent 500 GB d'espace disque. Les tuiles sont produites par le pôle thématique surfaces continentales Theia. Les variables analysées comprennent les bandes spectrales natives de Sentinel-2 (bleu, vert, 3, rouge, NIR, SWIR1, SWIR2) et des néo-canaux (NDVI, indice de brillance, STD).

363 relevés terrain fournissent une description de l'occupation du sol en 19 classes relatives aux cultures, à la couleur du sol... En post-analyse ces 19 classes sont regroupées en 6 classes sur des critères fonctionnels décrivant l'usage du sol.

Application

Une première expérience menée sur d'autres données montre sans surprise que le cluster est 24 fois plus rapide que l'ordinateur de bureau (2 coeurs, i7 2,6Ghz).

Pour les données de l'exemple, le temps de calcul de notre configuration de la chaîne iota2 (fig. 1) sur le cluster pour 11 runs est de 7h30.

L'interprétation des résultats utilise les sorties de la chaîne iota2 rapatriées sur l'ordinateur personnel. La matrice de confusion (fig. 2) montre que 70% des pixels sont correctement classés, valeur qui atteint 90% après regroupement en 6 classes. L'image classée (fig. 3) souligne les structures spatiales (culture, jachère, plan d'eau, village, arbre, pâturage libre).

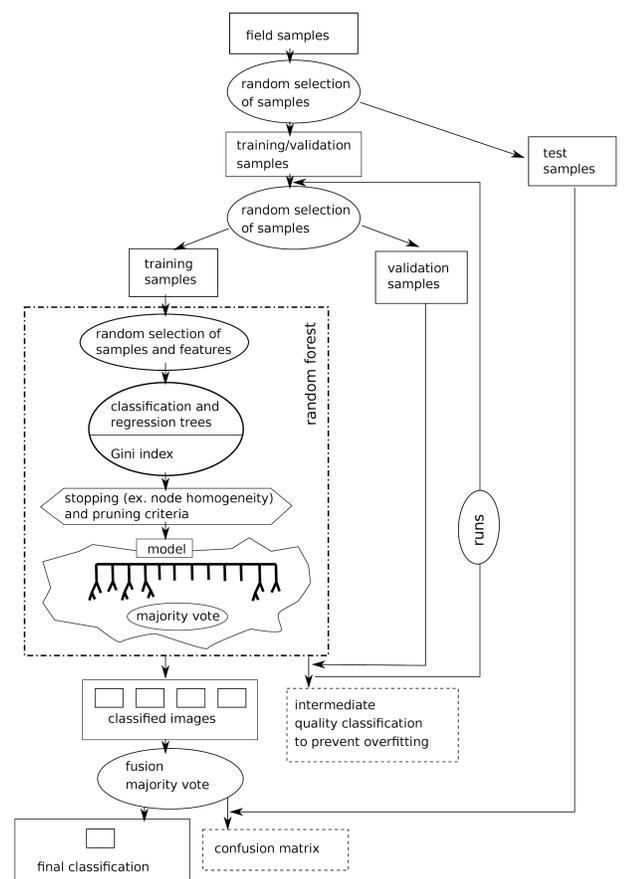


Figure 1. La configuration de la chaîne iota2 utilisée pour analyser une série temporelle d'images satellite par Random Forest.

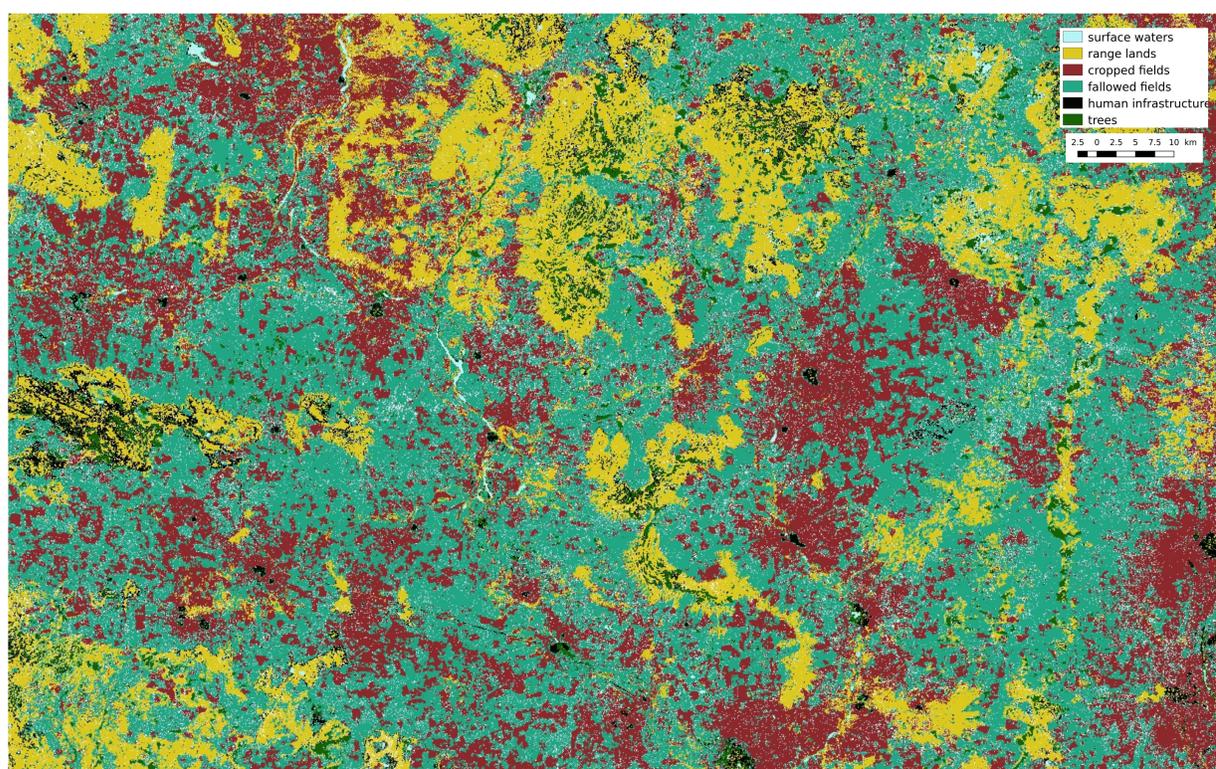


Figure 3. Image classée par la chaîne iota2 après regroupement en 6 classes.

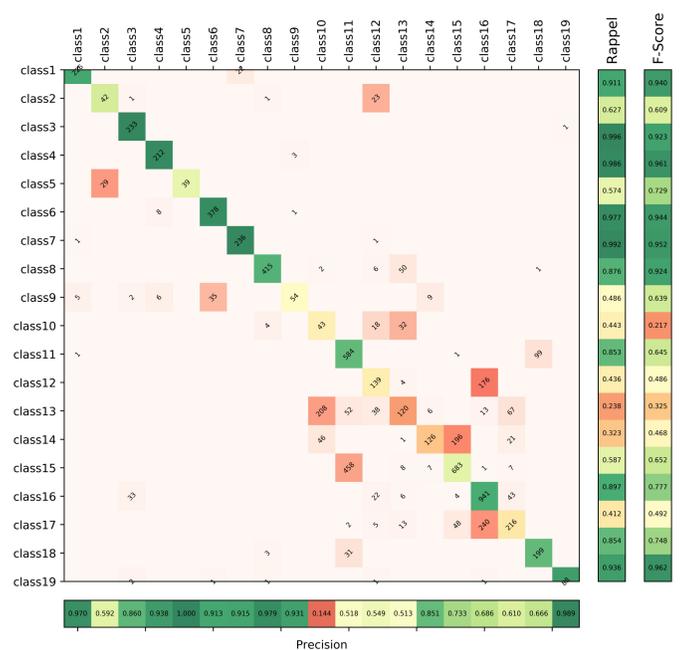


Figure 2. Matrice de confusion calculée pour les 19 classes.

Intérêts de cette solution

- La mise au point des traitements est réalisée dans un environnement convivial sur un ordinateur personnel.
- L'utilisation des containers assure la portabilité par simple copie sur tout cluster possédant singularity. Cette solution est particulièrement intéressante quand les chaînes de traitement comprennent de nombreuses bibliothèques. Il est même parfois nécessaire de créer plusieurs containers dans le cas d'incompatibilité de versions de bibliothèque !
- La solution d'environnement virtuel conda facilite énormément l'installation des applications. Il est possible de choisir le paquet adapté à une architecture particulière. Tensorflow est disponible pour les architectures AVX, AVX2, AVX512 optimisé par les compilateurs Intel, GPU optimisé par Cuda.