



# Le Green HPC à EDF

Green Days, Toulouse

2 Juillet 2018

**Mehdi Dogguy**

EDF



# 1 Calcul scientifique

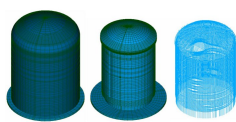
# Calcul scientifique à EDF

- ▶ R&D
  - ▶ Conception
  - ▶ Technologie de l'information
  - ▶ Énergies renouvelables
  - ▶ Réseaux électriques
  - ▶ ...
  
- ▶ Ingénierie
  
- ▶ Le management de l'énergie
  - ▶ Réduire le temps d'arrêt de tranches dans les réacteurs nucléaires
  - ▶ Planification de la consommation et production des semaines en avance

# Calcul scientifique en quelques mots

## ► Modélisation

- Approximer la réalité à l'aide d'un modèle



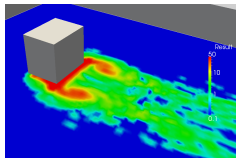
## ► Simulation

- Exécution d'un code qui calcule le comportement d'un système
- Un domaine de recherche en soi
- Besoin d'avoir du matériel adapté pour calculer rapidement sur des nombre à virgule flottante

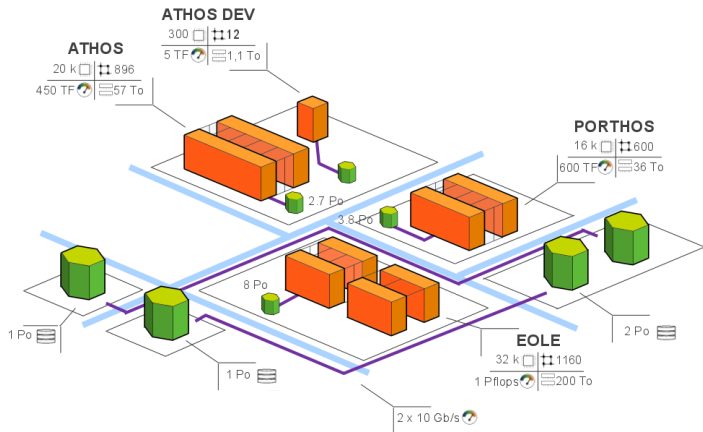


## ► Visualisation

- Exploration des résultats et analyses
- Besoin d'avoir les cartes graphiques et écrans adaptés



# Infrastructures de calcul



Plus de 2 Pflops de puissance de calcul qui consomment beaucoup d'électricité !

## 2 Le Green HPC

## Des datacenters green

Certification ISO 50001<sup>1</sup> des Datacenters d'EDF : Mise en place d'un système de Management de l'énergie qui vise à faire un meilleur usage de l'énergie

L'ISO 50001 en quelques actions :

- ▶ Élaborer une politique pour une utilisation plus efficace de l'énergie avec des objectifs de mise en œuvre
- ▶ Examiner l'efficacité de la politique
- ▶ Améliorer en continu le management de l'énergie.

---

<sup>1</sup>[https://fr.wikipedia.org/wiki/ISO\\_50001](https://fr.wikipedia.org/wiki/ISO_50001)

<sup>2</sup><https://bit.ly/2LkB19j>

## Des datacenters green

Certification ISO 50001<sup>1</sup> des Datacenters d'EDF : Mise en place d'un système de Management de l'énergie qui vise à faire un meilleur usage de l'énergie

L'ISO 50001 en quelques actions :

- ▶ Élaborer une politique pour une utilisation plus efficace de l'énergie avec des objectifs de mise en œuvre
- ▶ Examiner l'efficacité de la politique
- ▶ Améliorer en continu le management de l'énergie.

En Février 2016, EDF a obtenu<sup>2</sup> cette certification environnementale pour la performance énergétique de tous ses datacenters.

---

<sup>1</sup>[https://fr.wikipedia.org/wiki/ISO\\_50001](https://fr.wikipedia.org/wiki/ISO_50001)

<sup>2</sup><https://bit.ly/2LkB19j>



## Le HPC, grand consommateur d'énergie

Chez EDF, la consommation électrique du HPC peut atteindre 60% d'un des deux datacenters.

# Le HPC, grand consommateur d'énergie

Chez EDF, la consommation électrique du HPC peut atteindre 60% d'un des deux datacenters.

Quelques actions identifiées pour réduire la consommation des clusters (de manière directe ou indirecte) :

- ▶ Réduire la fréquence CPU des nœuds, dès que possible
- ▶ Éteindre les nœuds inutilisés
- ▶ Une utilisation plus judicieuse des ressources de calcul
- ▶ Un meilleur rendement sur les clusters (ratio du nombre de jobs passés par électricité consommée)
- ▶ Augmentation de la température en salle d'hébergement
- ▶ ...

# Le HPC, grand consommateur d'énergie

Chez EDF, la consommation électrique du HPC peut atteindre 60% d'un des deux datacenters.

Quelques actions identifiées pour réduire la consommation des clusters (de manière directe ou indirecte) :

- ▶ Réduire la fréquence CPU des nœuds, dès que possible
- ▶ **Éteindre les nœuds inutilisés**
- ▶ Une utilisation plus judicieuse des ressources de calcul
- ▶ Un meilleur rendement sur les clusters (ratio du nombre de jobs passés par électricité consommée)
- ▶ Augmentation de la température en salle d'hébergement
- ▶ ...

## Configurer l'extinction automatique

Dans Slurm, la configuration est assez simple :

```
ResumeProgram=/usr/lib/slurm-pwmgmt/exec/slurm-start-nodes  
SuspendProgram=/usr/lib/slurm-pwmgmt/exec/slurm-stop-nodes  
SuspendExcNodes=atcn[001-100,145],atbm[001-010,025-027]  
SuspendTime=1800
```

Fonctionnement du mécanisme :

- ▶ Au bout de 30 minutes d'inactivité, le nœud est éteint
- ▶ Dès que Slurm le juge utile, il rallume le nœud, l'affecte à un job et met à jour son statut

Cf. <https://github.com/edf-hpc/slurm-llnl-misc-plugins/tree/master/pwmgmt>

# Configurer l'extinction automatique

```
SuspendExcNodes=atcn[001-100,145],atbm[001-010,025-027]
```

Quelques remarques :

- ▶ Exclure quelques nœuds du mécanisme pour qu'il y ait toujours des nœuds disponibles
- ▶ Exclure les nœuds difficile à démarrer
- ▶ Exclure les nœuds ayant un rôle spécial (Par exemple : Quorum GPFS)
- ▶ Souvent des régressions dans le code de Slurm, pour ce mécanisme. Il faut tester les nouvelles versions avant de déployer.
- ▶ Nécessite un système de déploiement pouvant encaisser le démarrage simultané de plusieurs (centaines+) de nœuds et leur intégration de manière automatique

# Déploiement du mécanisme d'économie d'énergie sur les calculateurs

Déploiement sur les clusters :

- ▶ Décembre 2016 : Athos Dev
- ▶ Juin 2017 : Eole
- ▶ Juillet 2017 : Athos



# Déploiement du mécanisme d'économie d'énergie sur les calculateurs

Déploiement sur les clusters :

- ▶ Décembre 2016 : Athos Dev
- ▶ Juin 2017 : Eole
- ▶ Juillet 2017 : Athos

Résultats :

- ▶ Diminution de l'ordre de 20 % sur Athos Dev
- ▶ Diminution de l'ordre de 10 % sur Eole
- ▶ Économie non mesurée pour Athos



## 3 Suivi de la performance



# Outils développés chez EDF

- ▶ UncleBench : automatisation du lancement d'une série de benchmarks sur un cluster et agrégation des résultats dans un rapport  
<https://github.com/edf-hpc/unclebench>
  
- ▶ LpProf : Outil de profiling léger  
<https://github.com/edf-hpc/lpprof>
  
- ▶ JobMetrics : Suivi de l'utilisation des ressources d'un job en temps réel  
<https://github.com/edf-hpc/jobmetrics>

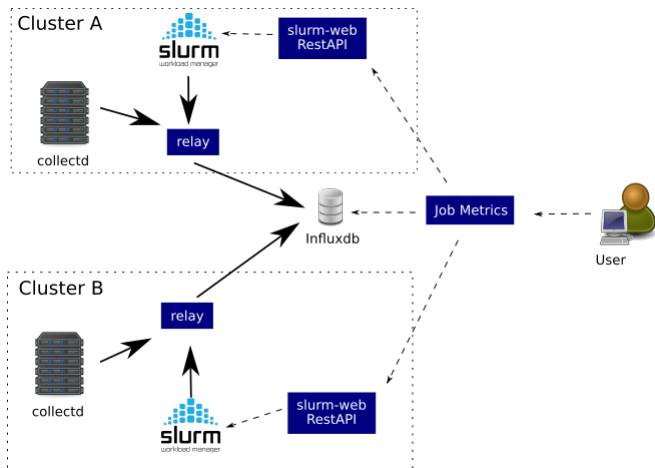
## Outils développés chez EDF

- ▶ UncleBench : automatisation du lancement d'une série de benchmarks sur un cluster et agrégation des résultats dans un rapport  
<https://github.com/edf-hpc/unclebench>
  
- ▶ LpProf : Outil de profiling léger  
<https://github.com/edf-hpc/lpprof>
  
- ▶ **JobMetrics** : Suivi de l'utilisation des ressources d'un job en temps réel  
<https://github.com/edf-hpc/jobmetrics>

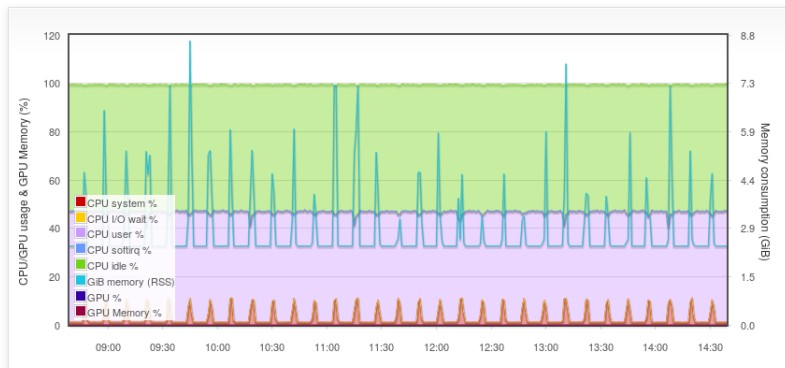
Systèmes de visualisation de la consommation des ressources sur les jobs

- ▶ Graphique
- ▶ Données temporelles
- ▶ Granularité fine
- ▶ Léger
- ▶ Temps Réel

# JobMetrics - Fonctionnement



## HPC metrics: cluster porthos job 248



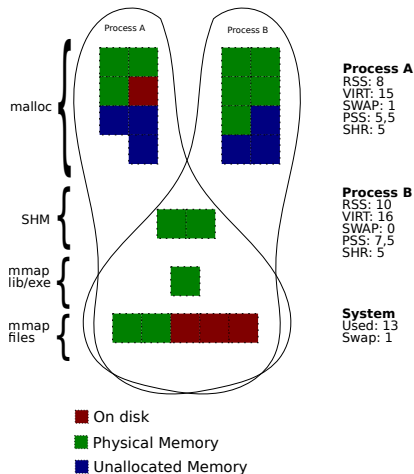
## JobMetrics - CPU

Utilisation des CPU alloués par les process attachés au job, les métriques affichées sont :

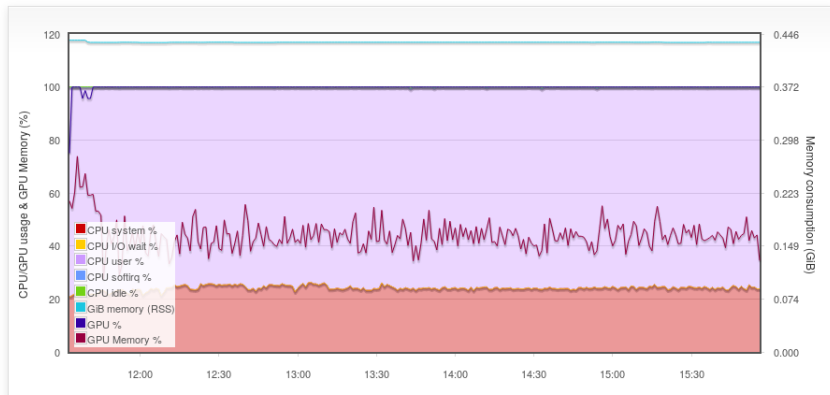
- ▶ User: Utilisation du CPU par du code en espace utilisateur (idéal=100%).
- ▶ System: Utilisation du CPU par du code en espace kernel (idéal=0%).
- ▶ IO Wait: Le CPU est en attente d'une IO (idéal=0%).
- ▶ SoftIRQ: Le CPU traite une interruption logicielle (idéal=0%).
- ▶ Idle: Le CPU n'a rien à faire, aucun processus ne le sollicite (idéal=0%).

# JobMetrics - Mémoire

- ▶ La mémoire est la somme des Resident Set Size (RSS) des processus attachés au job.
- ▶ L'utilisation des RSS fait que la mémoire peut être comptés plusieurs fois dans certains cas



## HPC metrics: cluster eole job 95





# JobMetrics - GPU

Pourcentage d'utilisation des ressources GPU sur les noeuds :

- ▶ GPU: Pourcentage d'utilisation des GPU des noeuds du job.
- ▶ GPU Memory: Pourcentage d'utilisation de la mémoire GPU des noeuds du job.

L'information vient de l'outil `nvidia-smi`. Elle est globale au noeud.

## JobMetrics - Limites

- ▶ Certaines métriques lourdes a capturer (RSS/PSS)
- ▶ Difficultés liées au cgroups (process en dehors, utilisation non associée...)
- ▶ Information pas suffisante pour trouver tous les problèmes (flops, réseau...)

# JobMetrics - Évolutions

- ▶ Archivage des jobs passés
- ▶ Analyse des jobs pour détecter les cas défailants
- ▶ Amélioration de la présentation

# JobMetrics - Évolutions

- ▶ Archivage des jobs passés
- ▶ Analyse des jobs pour détecter les cas défailants
- ▶ Amélioration de la présentation

# Merci !