

# Collecte de ressources libres dans une grille en préservant le système de fichiers: une approche autonome

COMPAS 2021, Lyon

Quentin GUILLOTEAU,\* Olivier RICHARD,\* Bogdan ROBU,\*\*  
Éric RUTTEN\*

\*Université Grenoble Alpes, Inria, CNRS, Grenoble INP, LIG

\*\*Université Grenoble Alpes, CNRS, Grenoble INP, GIPSA-lab

2021-07-09

# Contexte

## HPC

Systèmes HPC de + en + complexes:

- imprévisible en performances
- imprévisible en consommation d'énergie
- imprévisible en temps d'accès mémoire.

## Problèmes Potentiels

- Surcharge
- Surchauffe

↔ nécessité d'**une régulation en ligne** pour assurer les performances

## Une solution possible

Informatique Autonومية

# Informatique Autonome

## Définition (IBM 2000s)

Systèmes pouvant **s'auto-réguler** étant donnés des **objectifs de haut-niveau** par les administrateurs

## Outil Principal: La Boucle MAPE-K

### Actionneur(s) et Capteur(s)

(Ex: réguler  $T^{\circ}$  CPU avec sa **freq**)

- 1 **Monitor** (Capte la  $T^{\circ}$ )
- 2 **Analyse** (Erreur signifiante ?)
- 3 **Plan** (Quelle freq choisir ?)
- 4 **Execute** (Applique nouv. freq.)
- 5 **Knowledge** (Freq. min/max)

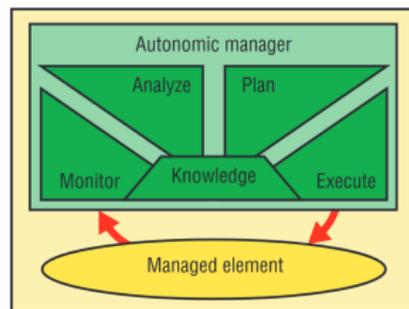


Figure: La Boucle MAPE-K

# Théorie du Contrôle

## Théorie du Contrôle

Régule le comportement de systèmes dynamiques

↔ Interprétation de la Boucle MAPE-K

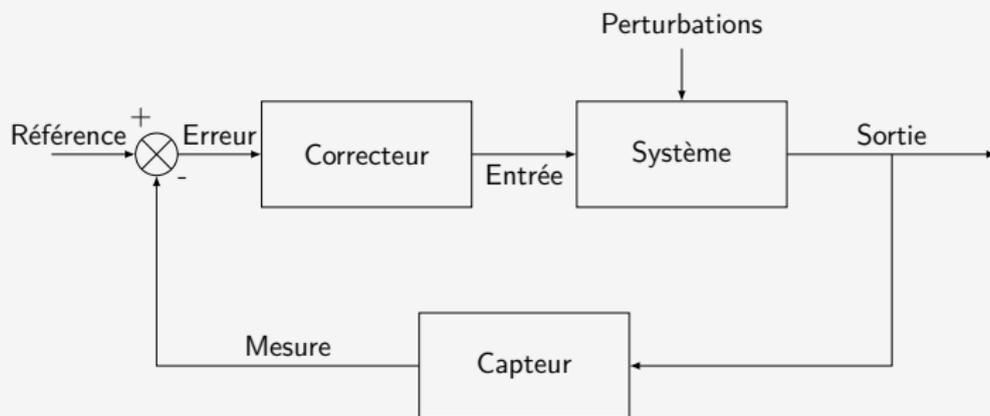


Figure: Boucle de contrôle

1 Introduction & Contexte

2 Collecte de ressources libres & l'intergiciel *CiGri*

3 Solution Proposée

4 Résultats Expérimentaux

5 Conclusion & Perspectives

# Collecte de ressources libres

## Ordonnancement laisse des "trous"

- spécifications
- mauvaise estimation du temps d'exécution
- ...

## Objectif

Utiliser les ressources libres pour exécuter des tâches moins prioritaires

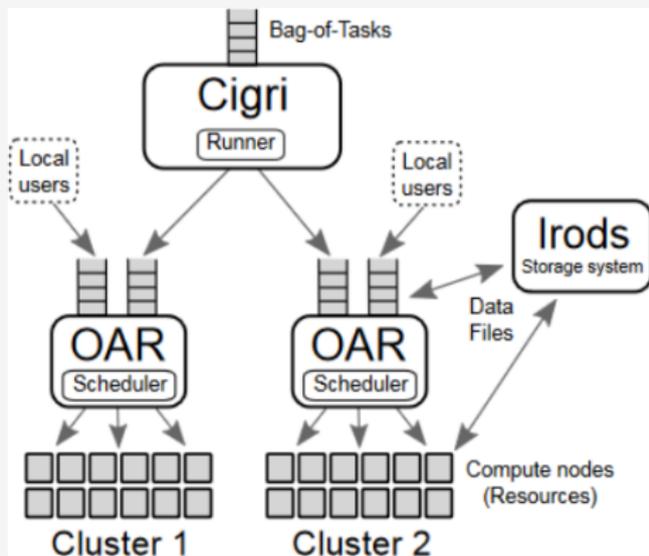
## État de l'art

- BOINC/Condor: récupère des cycles CPU de machines inactives
- OurGrid: partage de machines multi-labos
- BeBiDa: tâches BigData sur machines libres HPC

# CiGri: Présentation

## CiGri (CIMENT Grid)

- Intergiciel résistant aux fautes pour grille légère
- Interagit avec un ensemble d'ordonnanceurs *OAR*
- **But**: Exploiter les ressources libre d'une grille de manière **non intrusive**
- **bag-of-tasks**: Large ensemble de tâches multi-paramétrique
- **Tâches Best-effort**: Tâches avec priorité la plus basse



# CiGri: Boucle de Soumission (1/2)

---

## Algorithm 1: Boucle de Soumission

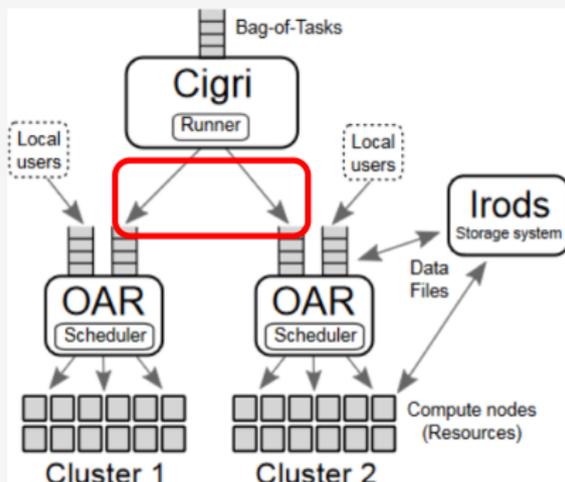
---

```

rate = 3;
increase_factor = 1.5;
while tâches non exécutées do
  if pas de tâches en cours
    d'exécution then
    soumet rate tâches;
    rate = min(rate ×
      increase_factor, 100);
  end
  while nombre de tâches en cours
    d'exécution > 0 do
    dort pendant 30 secondes;
  end
end
end

```

---

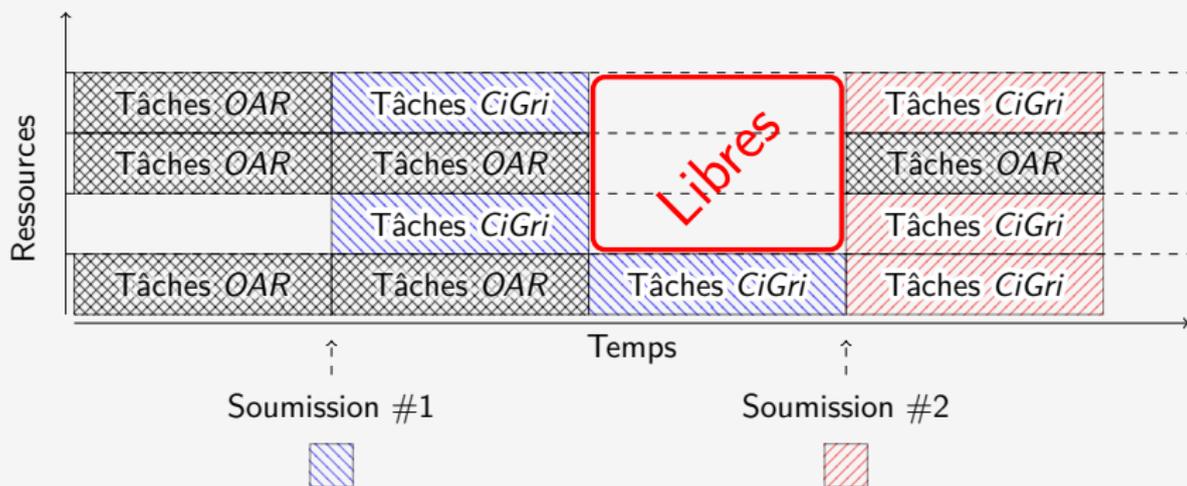


## CiGri: Boucle de Soumission (2/2)

### Le Problème

Doit attendre la fin de la soumission précédente pour soumettre à nouveau

↔ réduit la surcharge du système mais peut introduire une **sous-utilisation** des ressources



## *CiGri*: Le besoin d'amélioration

### Observation

Algo actuel de *CiGri*: **trop protecteur**

↔ peut être amélioré si **prise en compte du l'état de la grille**

### L'idée générale

Appliquer l'informatique autonome à *CiGri*

Réguler le nombre de tâches soumisses à *OAR* en fonction du nombre de ressources libres de la grille

+ ajout de contraintes (e.g. charge du système de fichiers distribué)

# Système de Fichiers Distribu  et son capteur

## Syst me de Fichiers Distribu  (SFD)

stocke les fichiers des utilisateurs de la grappe

## Probl me Potentiel

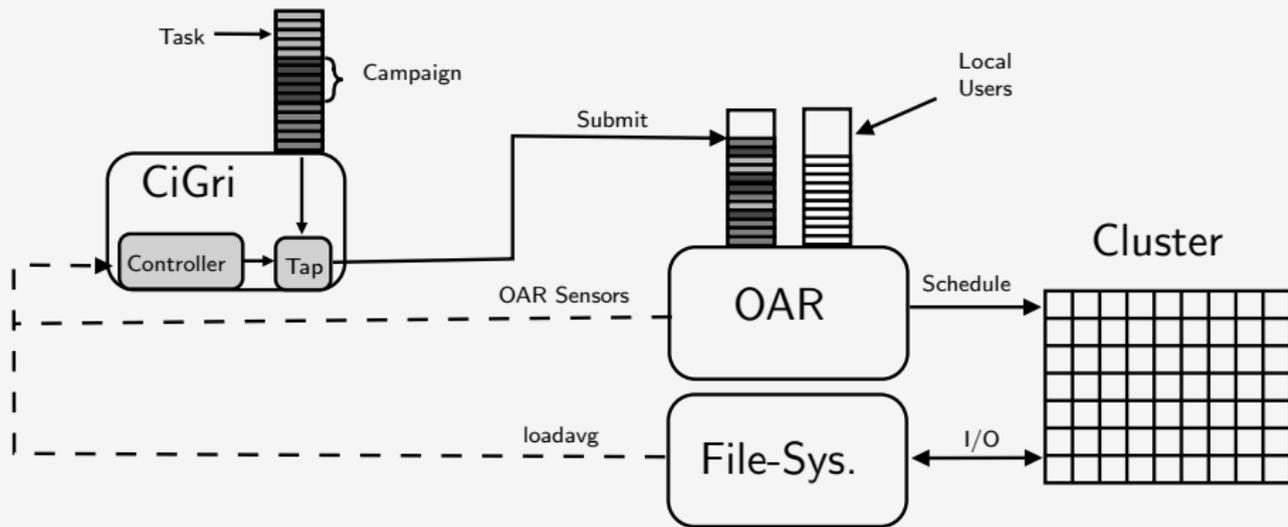
Trop de lectures/ critures simultan es  $\implies \nearrow$  temps lecture/ criture  
 $\hookrightarrow$  perturbation des t ches utilisateurs prioritaires

$\hookrightarrow$  Comment capter la surcharge du SFD ?

## Notre capteur: `/proc/loadavg`

- Repr sente le nombre de processus CPU en cours d'ex cution ou en attente pour le disque
- Lissage exponentiel  $\rightarrow$  Inertie  $\rightarrow$  cool pour contr le

# CiGri: ses boucles de rétro-action, capteurs et actionneurs



- 1 Introduction & Contexte
- 2 Collecte de ressources libres & l'intergiciel *CiGri*
- 3 Solution Proposée**
- 4 Résultats Expérimentaux
- 5 Conclusion & Perspectives

# Définition du Problème et Travaux Précédents

## Objectif

**Utiliser le plus de ressources libres** de la grappe **sans surcharger** le serveur de fichiers distribué

## Travaux précédents

- Correcteur PI pour réguler la quantité de tâches soumises à *OAR*:
  - meilleure utilisation des ressources vs. solution originale
  - serveur fichiers distribué
- Correcteur MPC avec prise en compte du SFD:
  - Modèle du système (trop) simpliste
  - Ne passait pas à l'échelle

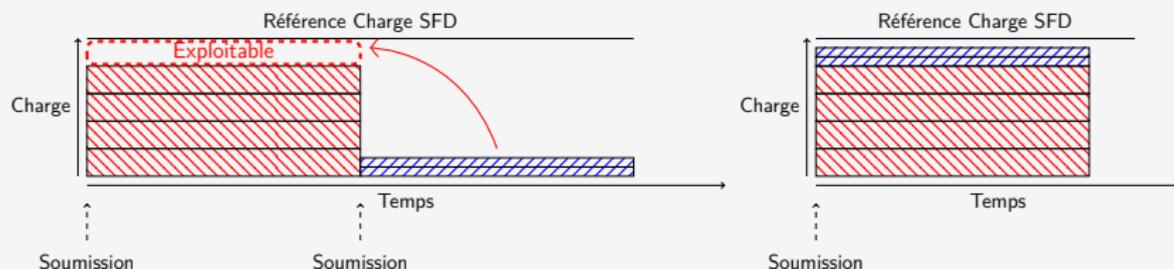
# Observation

## Problème

Originellement: soumissions de *CiGri* composées uniquement de tâches provenant de la **même campagne**

↔  $\simeq$  même comportement (temps exec. + E/S)

↔ peut amener à une sous utilisation de la grappe



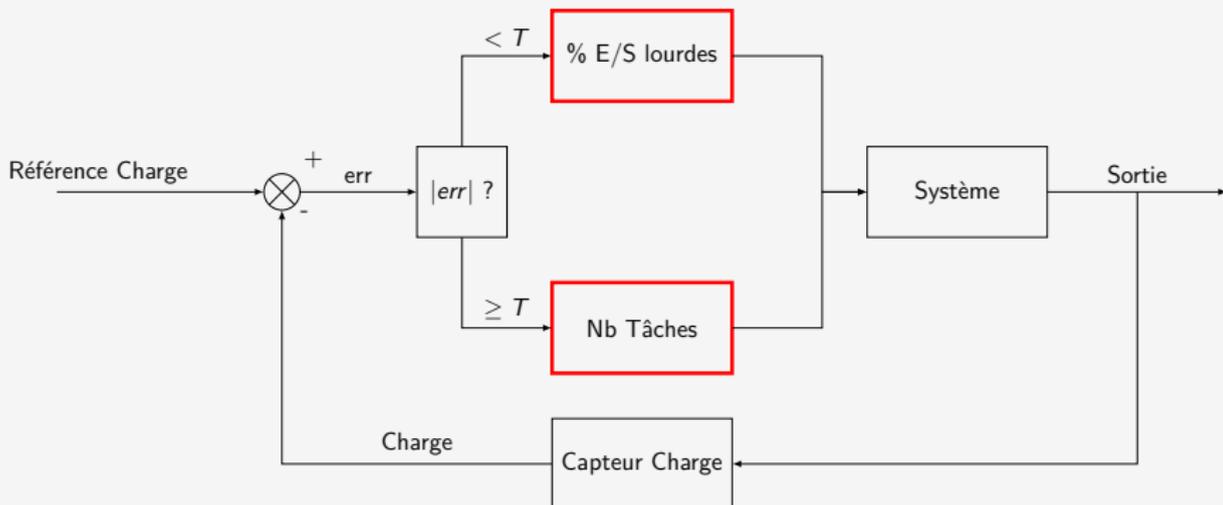
## Stratégie

Soumettre des tâches de 2 campagnes avec des charges E/S  $\neq$

# Régulation Proposée: Approche Biphasique

## Deux modes/phases de contrôle

- 1 Le **nombre total de tâches** soumises à OAR ("big step")
- 2 Le **pourcentage de tâches E/S lourdes** soumises ("small step")



## Quelques détails supplémentaires

### Ajout d'un garde fou

- Correcteur sur la file d'attente *OAR*
- But: éviter de "perdre la main"
- Nb tâches = min des sortie des correcteurs

### Correcteur Intégral: ( $\Delta$ réponse $\propto$ erreur)

$$\begin{aligned}
 U_k &= U_{k-1} + K \times Erreur_k \\
 &= U_{k-1} + K \times (Ref - Charge_k) = U_0 + K \sum_{i=0}^k Erreur_i
 \end{aligned}$$

- 1 Introduction & Contexte
- 2 Collecte de ressources libres & l'intergiciel *CiGri*
- 3 Solution Proposée
- 4 Résultats Expérimentaux**
- 5 Conclusion & Perspectives

# Dispositif Expérimental

## Architecture

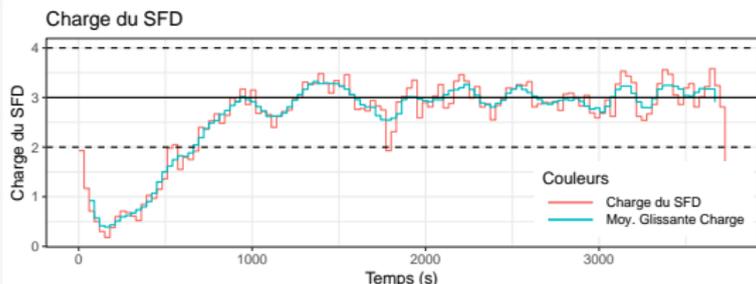
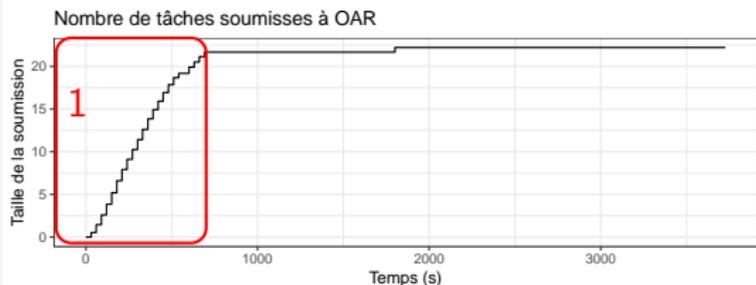
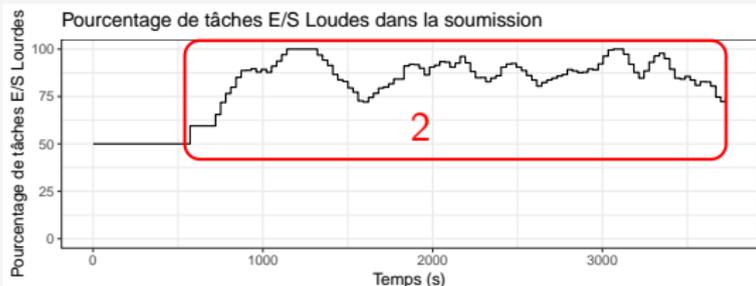
4 noeuds de Grid'5000 (Grisou)

- 1 Serveur *OAR* (v3)
- 1 Serveur *CiGri*
- 1 Serveur de fichiers distribué (NFS)
- émulation d'une grappe de 100 ressources

## Expérience

- Campagne de 1000 tâches: `sleep 30s` + écrit 100Mo
- Campagne de 1000 tâches: `sleep 30s` + écrit 10Mo
- Référence de Charge = 3, Threshold = 1
- Utilisateurs Prioritaires: comportement synthétique

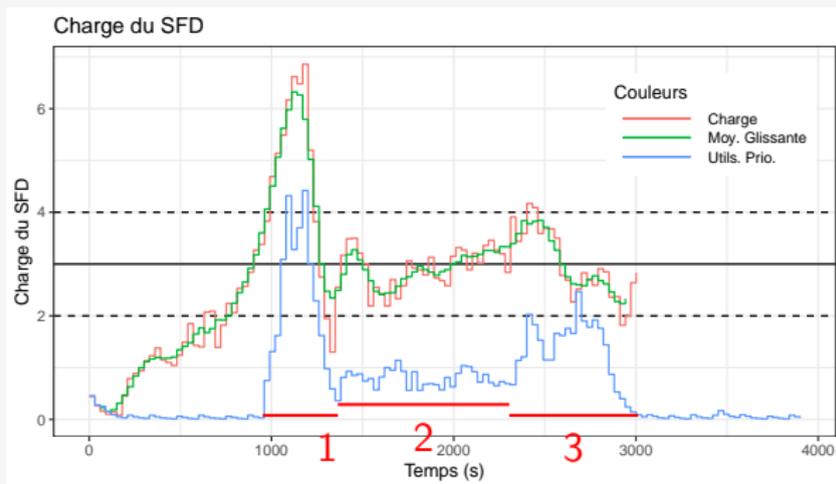
# Résultats sans utilisateurs prioritaires



## Remarques

- 1 Phase de montée ("big steps")
- 2 Puis "small steps"
- 3 Garde la charge dans l'intervalle

# Résultats avec utilisateurs prioritaires



## Remarques

- Charge synthétique (3 Phases)
  - Réaction lente aux variations trop brusques
  - + Sinon, parvient à rester dans l'intervalle

Figure: Charge du SFD avec Utilisateurs Prioritaires

# Comparaison

## Rappel de l'objectif

**Utiliser le + de ressources libres de la grappe sans surcharger le SFD**

Stratégie	Rétro-action	Proportion du temps avec le SFD surchargé (%)			Util. grappe (%)	
		Charge > Ref	Charge > Ref + 33%	Charge > Ref + 66%	absence	présence
Originelle	✗	80.49%	75.61%	60.98%	63.7	73.8
Scan	✓	17.86%	7.14%	3.57%	13.8	16.1
MPC	✓	89.02%	84.15%	59.76%	28.2	44.5
Contrôle Simple	✓	25.93%	11.11%	7.40%	17.4	23.0
Contrôle Biphase	✓	37%	10%	7%	19.6	25.8

## Remarques

- Bonne régulation de la charge pour Biphase (Top 3)
- Meilleure utilisation des ressources dans ce Top 3
- Meilleure utilisation // au contrôle simple → montre l'intérêt de l'approche biphase

- 1 Introduction & Contexte
- 2 Collecte de ressources libres & l'intergiciel *CiGri*
- 3 Solution Proposée
- 4 Résultats Expérimentaux
- 5 Conclusion & Perspectives**

# Conclusion & Perspectives

## Conclusion

- Solution permettant de réguler la charge du SFD autour d'une valeur désirée
- tout en récoltant davantage de ressources

## Perspectives

- Lien loadavg et impact temps lecture/écriture. Autre/combinaison capteurs
- Détection du type de campagne (E/S Lourde ou E/S Légère)
- Charges réelles/Benchmarks
- Identification des gains des correcteurs
- autres formes de régulation (PID, RST, etc)