

# Autonomic Approach to the Runtime Management of HPC Clusters Resources

Quentin GUILLOTEAU, Olivier RICHARD, Eric RUTTEN

Univ. Grenoble Alpes, Inria, CNRS, LIG, F-38000 Grenoble France

## Context & Applicative Problem

### Gricad



- Grenoble Mesocenter
- 169 TFLOPS
- 5700 CPUs on several clusters
- > 1200 users

Idle HPC Resources  $\Rightarrow$  Lost Computing Power  
 $\hookrightarrow$  How to Harvest?

## Perturbations from Harvesting

Cannot simply submit huge amount of jobs.  
 $\nearrow$  Harvesting  $\Rightarrow$   $\nearrow$  Perturbations (e.g., I/O)  
 $\hookrightarrow$  Trade-off to exploit

## Instrumentation

- **Actuator:** #jobs to submit from CiGri to OAR
- **Sensors:** OAR Waiting Queue length  
File-System Load (`loadavg`)

## Controllers

**Controller** = function between objective (or **Reference**) error and system input.

$$Error_k = Reference - Sensor_k$$

Different types of **Controllers** (PID, RST, MFC)

$$\text{Actuator}_{k+1} = K_p \times Error + K_i \times \sum_k Error_k + K_d \times (Error_k - Error_{k-1})$$

## Response to Perturbations

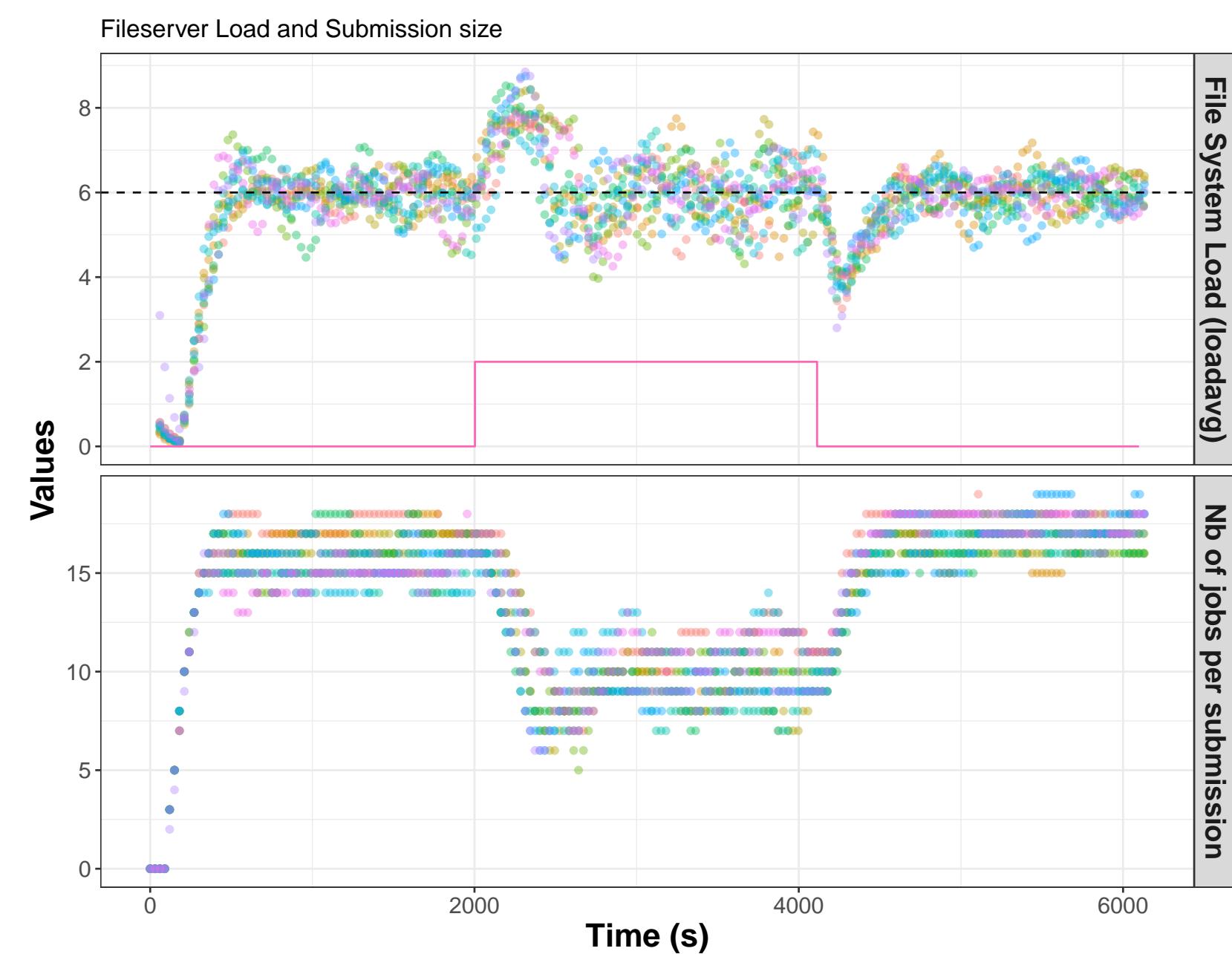
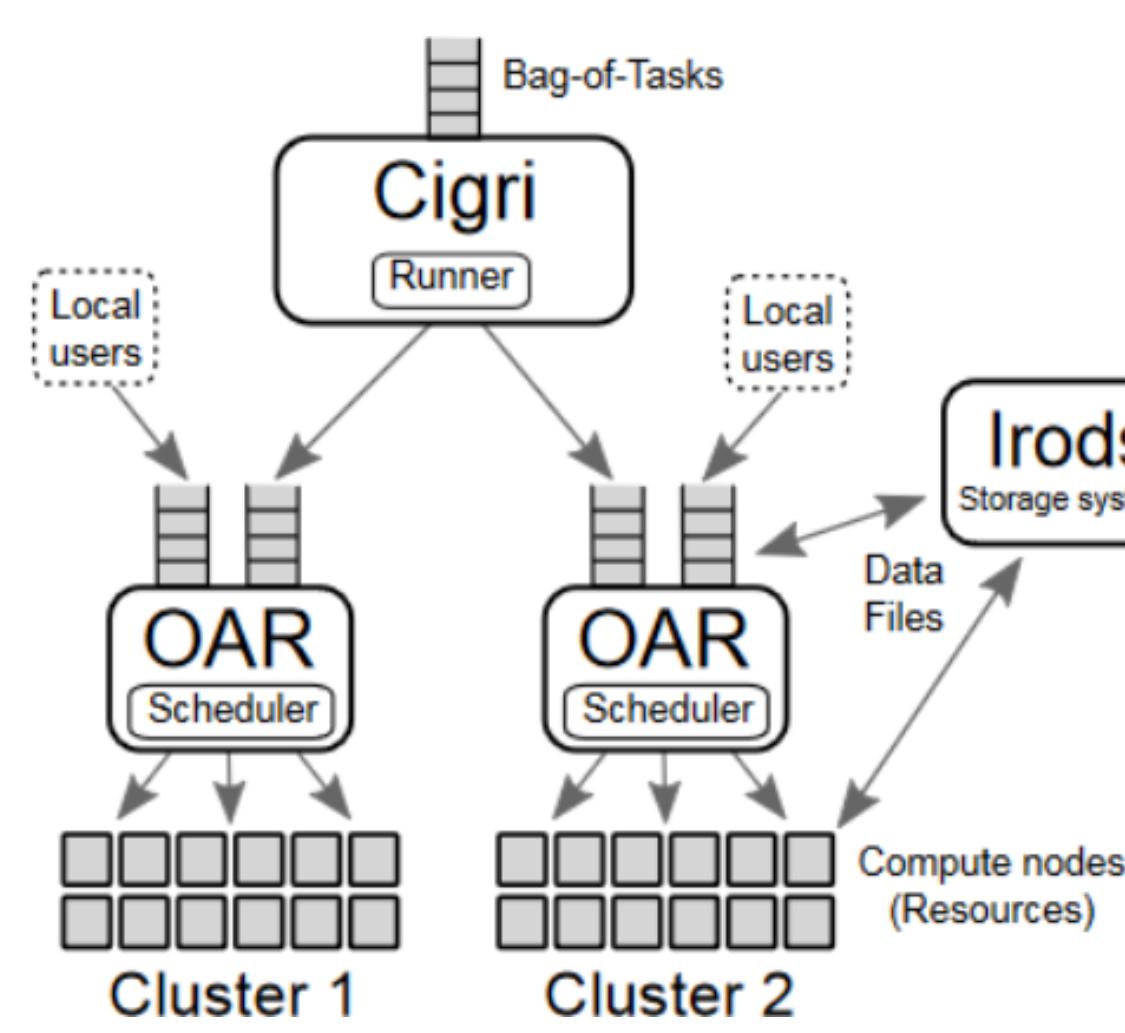


Figure: Response of the closed-loop system (top) to a synthetic step perturbation to keep the load of the File-System around the value 6. Bottom is the number of jobs submitted by CiGri.

## A Solution: CiGri

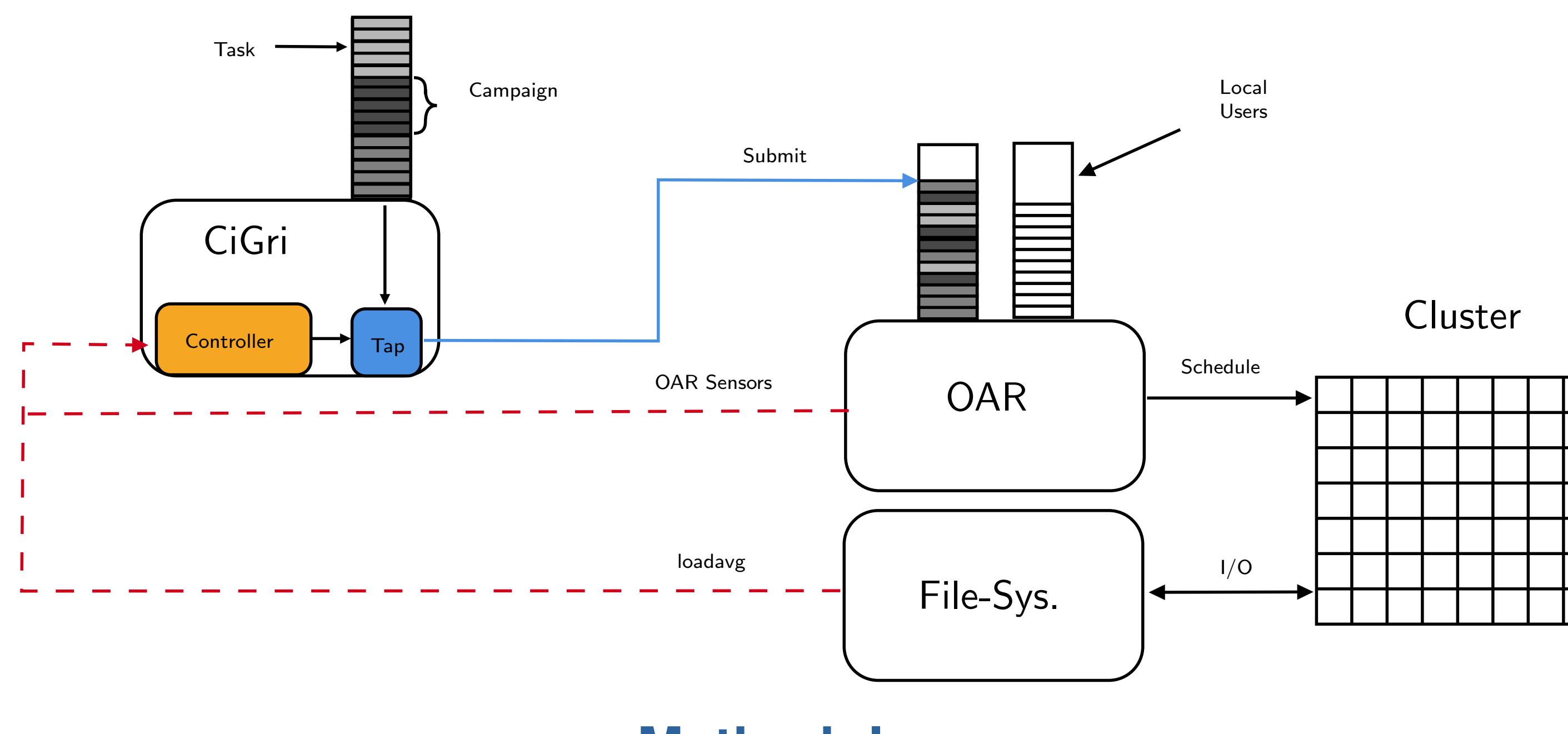
- **bag-of-tasks:** many, multi-parametric jobs
- **Best-effort Jobs:** Lowest priority
- **Objective:** Collect grid idle resources in a non-intrusive way for the users



Various Job Behaviours  
 $\hookrightarrow$  Variability  $\Rightarrow$  Requires **Runtime Management**

## Research Problem

Use Control Theory tools to dynamically **harvest idle HPC resources** in a **non-intrusive** way for the users



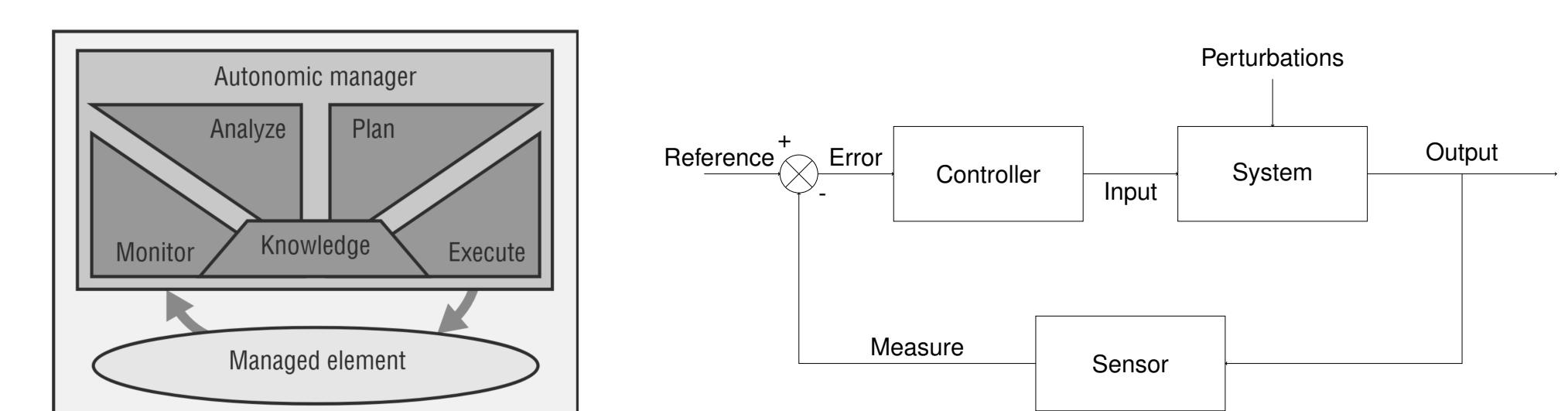
## Methodology

Open-Loop Experiments  $\Rightarrow$  Find Model (1st order)  
 $y(k+1) = ay(k) + bu(k)$   $\Rightarrow$  Choice of Closed-Loop Behaviour  $\Rightarrow$  Deduce Controller Gains  $K_p, K_i, K_d$ ,

## MAPE-K Loop & Control Theory

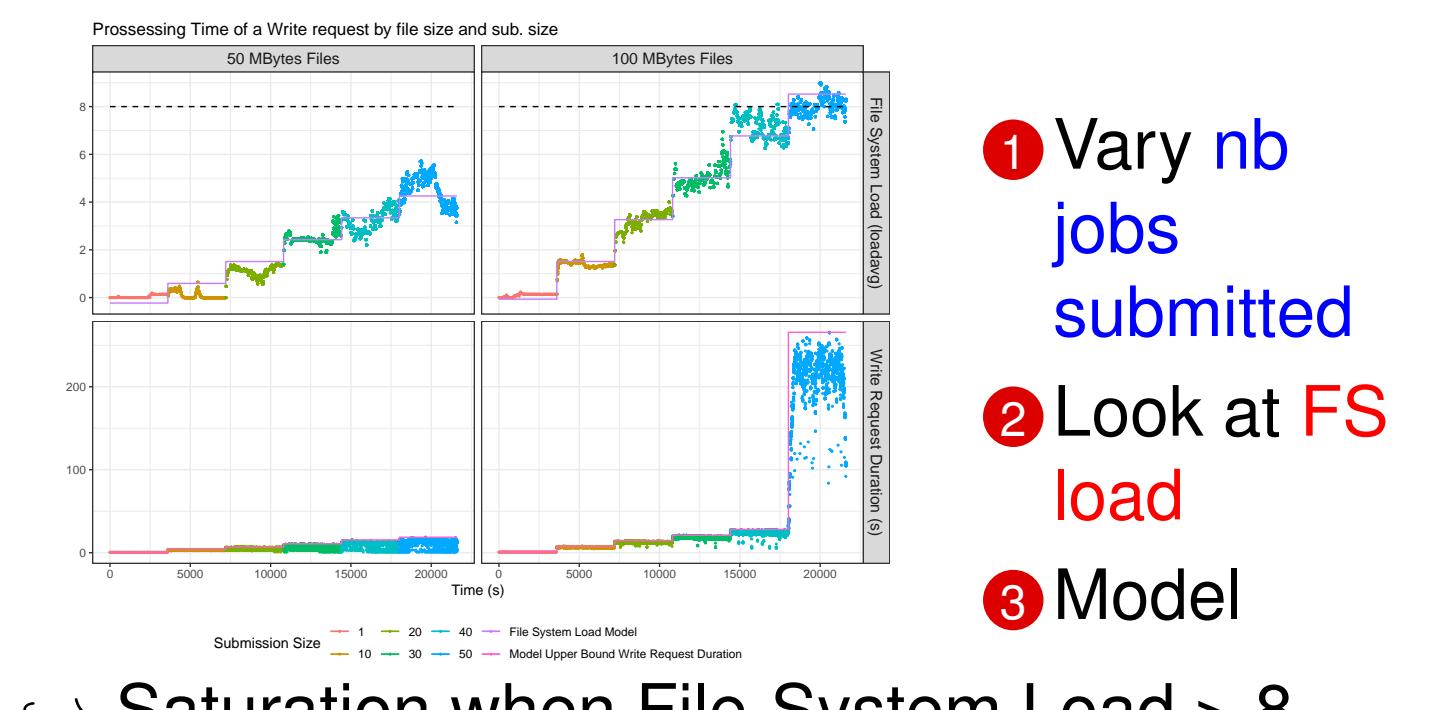
Auto-regulating Dynamical Systems given **high-level objectives**

$\hookrightarrow$  Control Theory as an Interpretation of the MAPE-K Loop



**Control Theory** drives system to a desired state by adapting inputs while **minimizing delay, overshoot, or steady-state error** and ensuring control **stability**

## Open-Loop & Identification

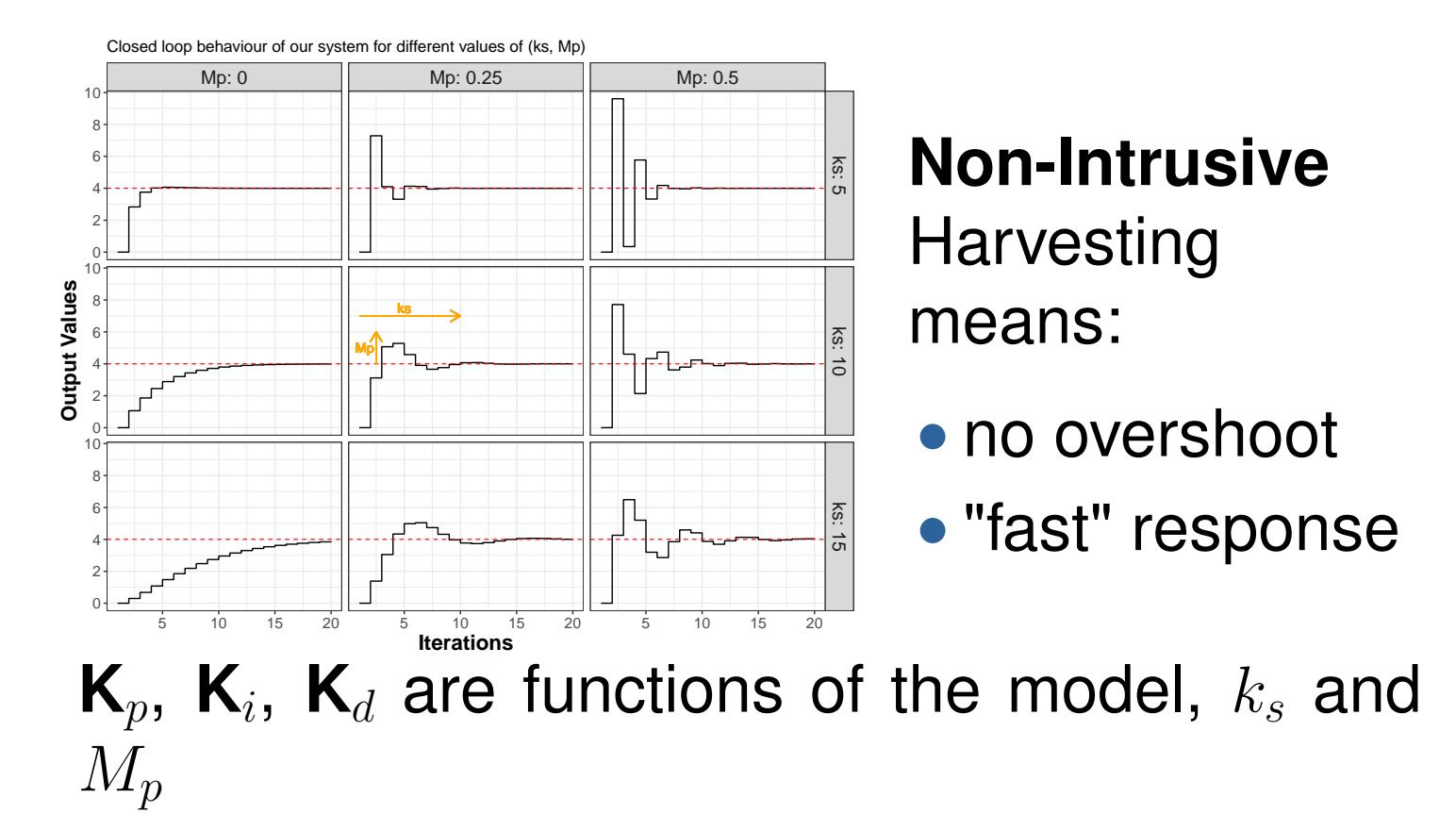


$\hookrightarrow$  Saturation when File-System Load > 8

## Closed-Loop Behaviour

Can pick the **desired system behaviour**:

- $k_s$ : max **time** to steady state
- $M_p$ : max **overshoot** allowed



**Non-Intrusive**  
Harvesting means:

- no overshoot
- "fast" response

## Perspectives

- Adaptive Controllers to deal with different I/O Profiles
- Interaction with Scheduler to anticipate the arrival of new Priority jobs
- Minimize the amount of Best-effort jobs killed (also represents lost computing power)
- Reproducible Experiments with Nix and NixOS

## References

Q. Guilloteau, O. Richard, B. Robu and E. Ruttent. *Controlling the Injection of Best-Effort Tasks to Harvest Idle Computing Grid Resources*, ICSTCC 2021, hal-03363709

## Exploiting the Trade-Off

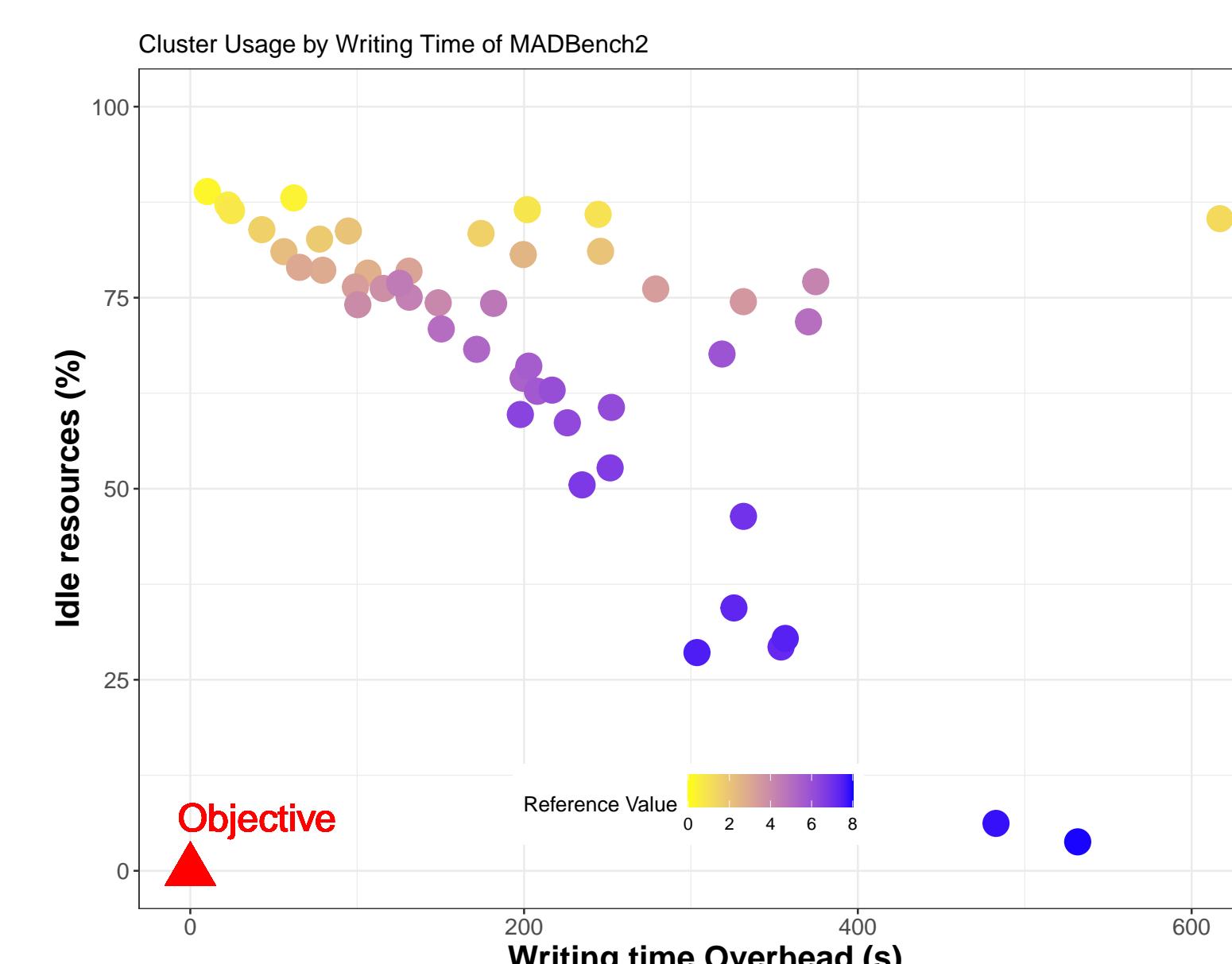


Figure: Writing Overhead of an application (MADBench2 Benchmark) with CiGri based on Reference Value

$\hookrightarrow$  We have **Control over the Degradation and Harvesting** with the choice of Reference Value

Acknowledgements: Experiments presented in this poster were carried out using the Grid'5000 testbed, supported by a scientific interest group hosted by Inria and including CNRS, RENATER and several Universities as well as other organizations (see <https://www.grid5000.fr>).



# CTRL-A

Control for  
safe Autonomic computing systems

