

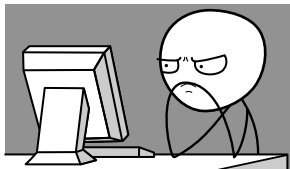
# QUALITY ASSURANCE OF A HPC CLUSTER: TESTING FOR PERFORMANCE NON-REGRESSION

---

Tom Cornebize, Arnaud Legrand  
Laboratoire d'Informatique de Grenoble

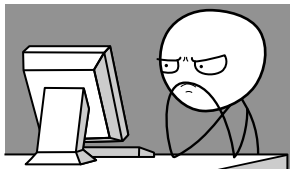
November 5, 2019

## Typical Performance Evaluation Questions (Given my application and a supercomputer)



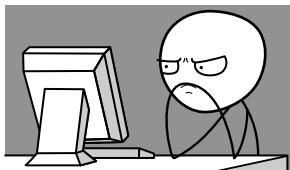
- **Before** running
  - How many nodes?
  - For how long?
  - Which parameters?

## Typical Performance Evaluation Questions (Given my application and a supercomputer)



- **Before** running
  - How many nodes?
  - For how long?
  - Which parameters?
- **After** running
  - Performance as “expected”?
  - Problem in the app or the platform?

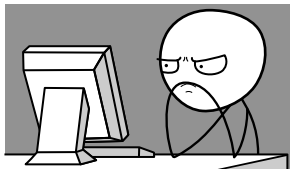
### Typical Performance Evaluation Questions (Given my application and a supercomputer)



- **Before** running
  - How many nodes?
  - For how long?
  - Which parameters?
- **After** running
  - Performance as “expected”?
  - Problem in the app or the platform?

So many large-scale runs, solely to tune performance?!?

### Typical Performance Evaluation Questions (Given my application and a supercomputer)



- **Before** running
  - How many nodes?
  - For how long?
  - Which parameters?
- **After** running
  - Performance as “expected”?
  - Problem in the app or the platform?

So many large-scale runs, solely to tune performance?!?

Holy Grail: Predictive Simulation on a “Laptop”

Building a predictive **model** of the durations:

- Computations (`dgemm`, ...)
- Communications (`MPI_send`, ...)

A lot of measures, with different input sizes



Building a predictive **model** of the durations:

- Computations (`dgemm`, ...)
- Communications (`MPI_send`, ...)

A lot of measures, with different input sizes

Some troubles, **wrong predictions**

⇒ Needed to investigate.



## Possible issues

- Most of the problems are human mistakes (wrong library version, wrong options, ...)
- A lot of transient phenomena: OS scheduler, temperature changes, core frequencies oscillation...
- A measure can have an impact (positive or negative) on the next measure (cache effects...)



## Possible issues

- Most of the problems are human mistakes (wrong library version, wrong options, ...)
- A lot of transient phenomena: OS scheduler, temperature changes, core frequencies oscillation...
- A measure can have an impact (positive or negative) on the next measure (cache effects...)

## Design of experiments

- Randomizing the sequence of measures to reduce bias
- Tools: ad-hoc scripts to generate *experiment files*

### Automating the setup

- Job submission, deployment, software stack installation, experiment execution...
- Tools: OAR, Kadeploy, Peanut

## DOWN THE RABBIT HOLE (2)

### Automating the setup

- Job submission, deployment, software stack installation, experiment execution...
- Tools: OAR, Kadeploy, Peanut

### Automating the metadata collection

- Date, kernel and library versions, output of every command, CPU temperature, core frequencies...
- Tools: Peanut, custom scripts

## DOWN THE RABBIT HOLE (2)

### Automating the setup

- Job submission, deployment, software stack installation, experiment execution...
- Tools: OAR, Kadeploy, Peanut

### Automating the metadata collection

- Date, kernel and library versions, output of every command, CPU temperature, core frequencies...
- Tools: Peanut, custom scripts

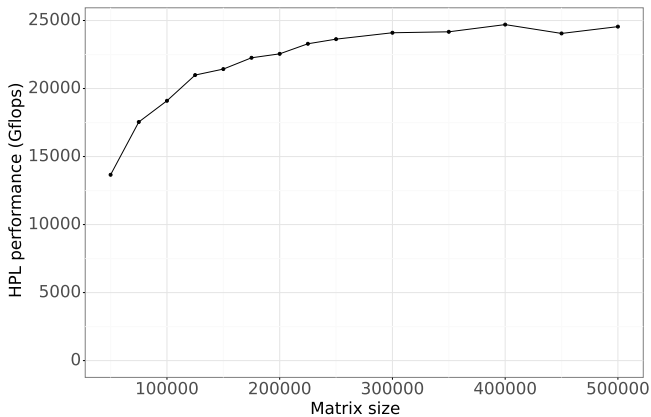
### Data analysis

- Data visualization (correlations, temporal patterns, distributions)
- Statistics (linear regressions, ANOVA)
- Tools: Python & R with Jupyter, ggplot...

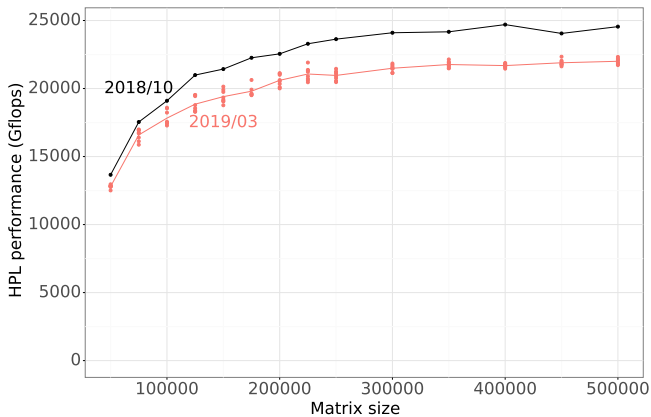
DAHU@GRID'5000

---

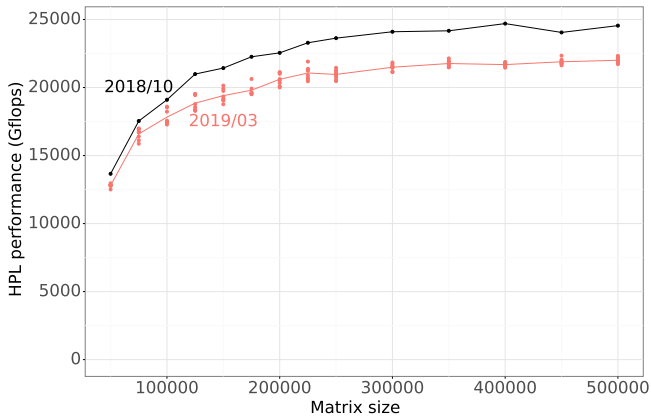
## PERFORMANCE OF THE WHOLE CLUSTER: PARALLEL APPLICATION



# PERFORMANCE OF THE WHOLE CLUSTER: PARALLEL APPLICATION



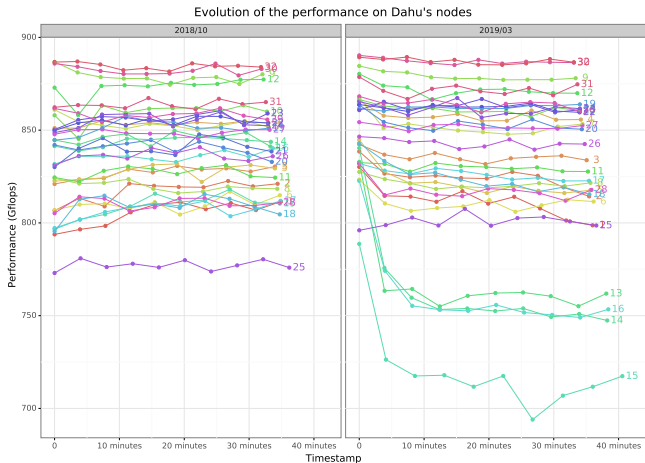
# PERFORMANCE OF THE WHOLE CLUSTER: PARALLEL APPLICATION



Same software, same hardware, 10% performance drop

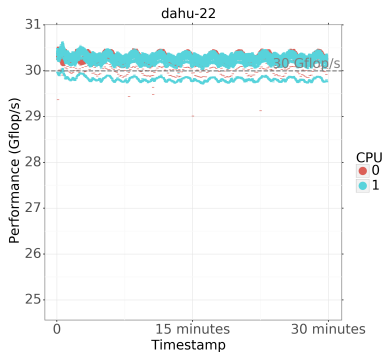


# SINGLE-NODE PERFORMANCE

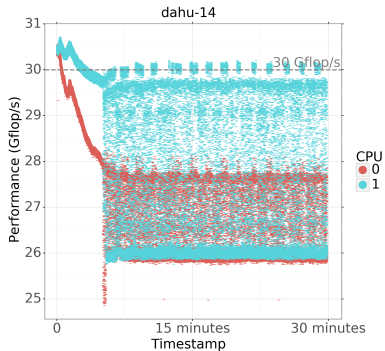
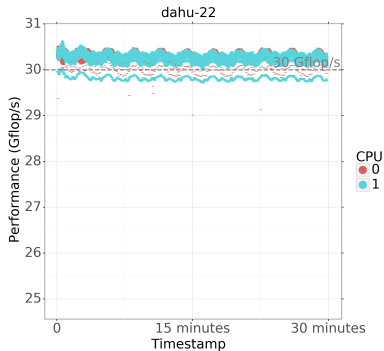


Performance drop for dahu- $\{13..16\}$  after a few minutes

# THE GOOD & THE BAD: PERFORMANCE EVOLUTION

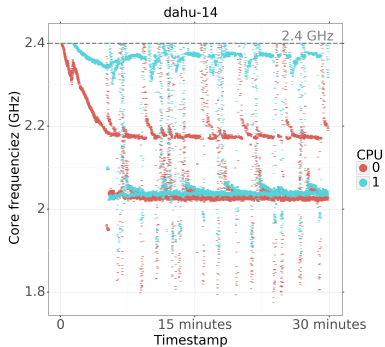
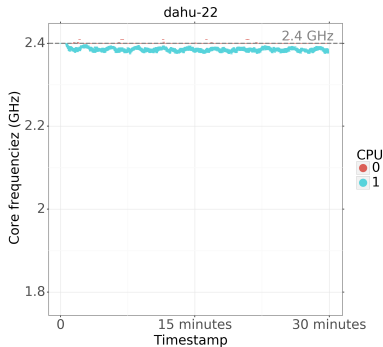


# THE GOOD & THE BAD: PERFORMANCE EVOLUTION



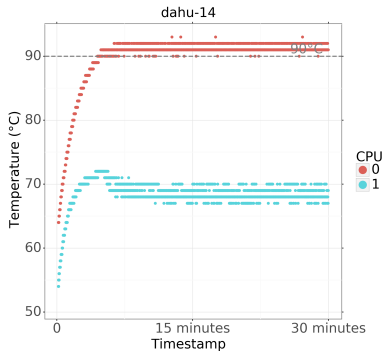
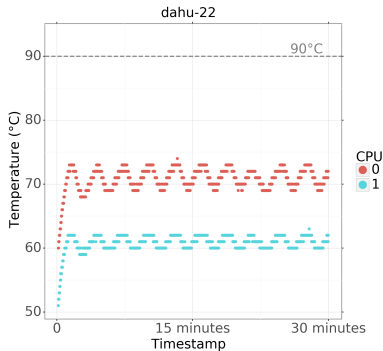
Performance drop, huge variability, CPU n°0 is worse

# THE GOOD & THE BAD: FREQUENCY EVOLUTION



Frequency drop, huge variability, CPU n°0 is worse

# THE GOOD & THE BAD: TEMPERATURE EVOLUTION



Very high temperature  $\Rightarrow$  probably a cooling issue

Fixed by changing the node frames

Several other problems encountered on this cluster:

- Connectivity issue ( $\Rightarrow$  replug the Omnipath cable, then reboot)

Several other problems encountered on this cluster:

- Connectivity issue ( $\Rightarrow$  replug the Omnipath cable, then reboot)
- Memory bandwidth issue ( $\Rightarrow$  change the faulty memory stick)

Several other problems encountered on this cluster:

- Connectivity issue ( $\Rightarrow$  replug the Omnipath cable, then reboot)
- Memory bandwidth issue ( $\Rightarrow$  change the faulty memory stick)
- Important heterogeneity between the nodes (10% difference between the slowest and the fastest, even without cooling problems)



# PERSPECTIVES

---

Objective: ~~finding bugs~~ building models

Side effect: use these models for [statistical tests](#), to automatically detect performance problems

Objective: ~~finding bugs~~ building models

Side effect: use these models for [statistical tests](#), to automatically detect performance problems

Regular and [semi-automated performance measures](#) on Grid'5000:  
[https://gitlab.in2p3.fr/tom.cornebize/g5k\\_data](https://gitlab.in2p3.fr/tom.cornebize/g5k_data)

Several problems detected, some very severe, others more subtle



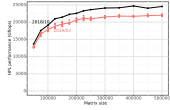
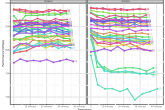
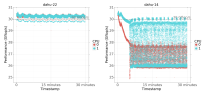
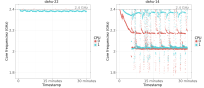
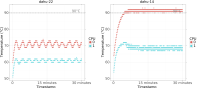
Objective: ~~finding bugs~~ building models

Side effect: use these models for [statistical tests](#), to automatically detect performance problems

Regular and [semi-automated performance measures](#) on Grid'5000:  
[https://gitlab.in2p3.fr/tom.cornebize/g5k\\_data](https://gitlab.in2p3.fr/tom.cornebize/g5k_data)

Several problems detected, some very severe, others more subtle

Dahu@Grid'5000 had troubles. [What about Dahu@Ciment?](#)

CONTEXT	PERFORMANCE MEASURES	DOWN THE RABBIT HOLE (2)
<p>Typical Performance Evaluation Questions (Given my application and a supercomputer)</p>  <ul style="list-style-type: none"> <li>Before running           <ul style="list-style-type: none"> <li>How many nodes?</li> <li>For how long?</li> <li>Which parameters?</li> </ul> </li> <li>After running           <ul style="list-style-type: none"> <li>Performance as "expected"?</li> <li>Problem in the app or the platform?</li> </ul> </li> </ul> <p>So many large-scale runs, solely to tune performance?!</p> <p>Holy Grail: Predictive Simulation on a "Laptop"</p>	<p>Building a predictive <b>model</b> of the durations:</p> <ul style="list-style-type: none"> <li>Computations (<code>gemm...</code>)</li> <li>Communications (<code>MPI_send...</code>)</li> </ul> <p>A lot of measures, with different input sizes</p> <p>Some troubles, <b>wrong predictions</b>        =&gt; Needed to investigate.</p> 	<p>Automating the setup</p> <ul style="list-style-type: none"> <li>Job submission, deployment, software stack installation, experiment execution...</li> <li>Tools: OAR, Kadeploy, Peanut</li> </ul> <p>Automating the metadata collection</p> <ul style="list-style-type: none"> <li>Data, kernel and library versions, output of every command, CPU temperature, core frequencies...</li> <li>Tools: Peanut, custom scripts</li> </ul> <p>Data analysis</p> <ul style="list-style-type: none"> <li>Data visualization (correlations, temporal patterns, distributions)</li> <li>Statistics (linear regressions, ANOVA)</li> <li>Tools: Python &amp; R with jupyter, ggplot...</li> </ul>
<p>PERFORMANCE OF THE WHOLE CLUSTER: PARALLEL APPLICATION</p>  <p>Same software, same hardware, 10% performance drop</p>	<p>SINGLE-NODE PERFORMANCE</p>  <p>Performance drop for dahu-[13..16] after a few minutes</p>	<p>THE GOOD &amp; THE BAD: PERFORMANCE EVOLUTION</p>  <p>Performance drop, huge variability, CPU n°0 is worse</p>
<p>THE GOOD &amp; THE BAD: FREQUENCY EVOLUTION</p>  <p>Frequency drop, huge variability, CPU n°0 is worse</p>	<p>THE GOOD &amp; THE BAD: TEMPERATURE EVOLUTION</p>  <p>Very high temperature =&gt; probably a cooling issue        Fixed by changing the node frames</p>	<p>STEPPING BACK</p> <p>Objective: <b>finding bugs</b> building models</p> <p>Side effect: use these models for <b>statistical tests</b>, to automatically detect performance problems</p> <p>Regular and semi-automated performance measures on Grid'5000:  <a href="https://gitlab.in2p3.fr/tom.cornebize/g5k_data">https://gitlab.in2p3.fr/tom.cornebize/g5k_data</a></p> <p>Several problems detected, some very severe, others more subtle</p> <p>Dahu@Grid'5000 had troubles. <a href="#">What about Dahu@Ciment?</a></p>