# Monitoring GreenClouds

# Evaluating the trade-off between Performance and Energy Consumption in DAS-4



## System and Networking Engineering

Renato Fontana | Katerina Mparmpopoulou

# Presentation Flow

- Green concepts

- Project objective

- Experimental environment

- Metrics and Workload

- Experiment Results

- Conclusions

- Future work

# Green Concepts

- What does it mean to be green?
- Refers to environmentally sustainable
- Energy becomes a key challenge in large-scale distributed systems
- IT requires more and more power

# Known techniques

- Event-monitoring counters
  - Deducing energy consumption

- On/off algorithms
  - Switch on/off nodes in long idle state

- Load balancing
  - Distribute workload amongst multiple nodes

- Task scheduling
  - Slowdown factors

- Thermal management
  - Monitoring heat generation

# Research Question

- How to evaluate the trade-off between energy and performance in DAS-4?

- How to correlate performance and energy consumption in Cloud Computing Systems?
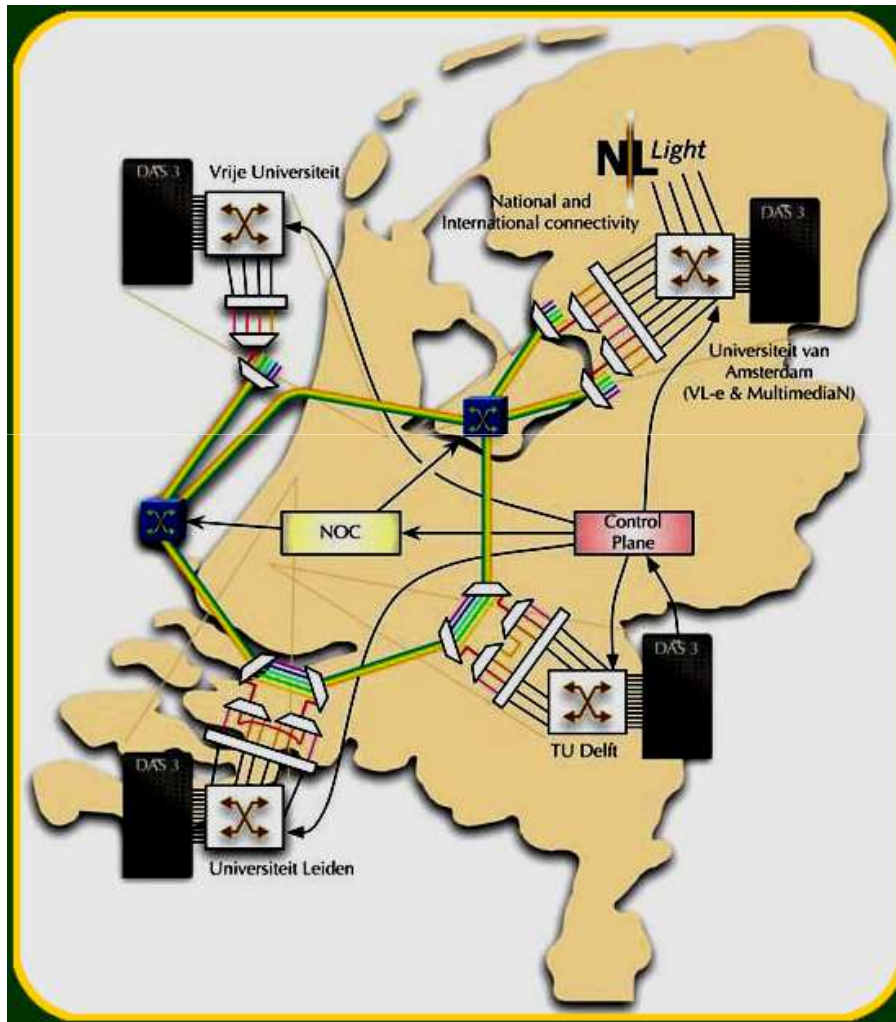
# Approach

- Compare workload with power-monitoring tools
- Estimate energy consumption in nodes
- Correlate main components (CPU, memory)
- CPU load and energy consumed

# Experimental environment

DAS-4 (The Distributed ASCI Supercomputer 4)

- Six-cluster wide-area distributed system
  - UvA and VU nodes (PDU enable)
- Grid Computing
  - DAS-4 mainly composed by cluster nodes
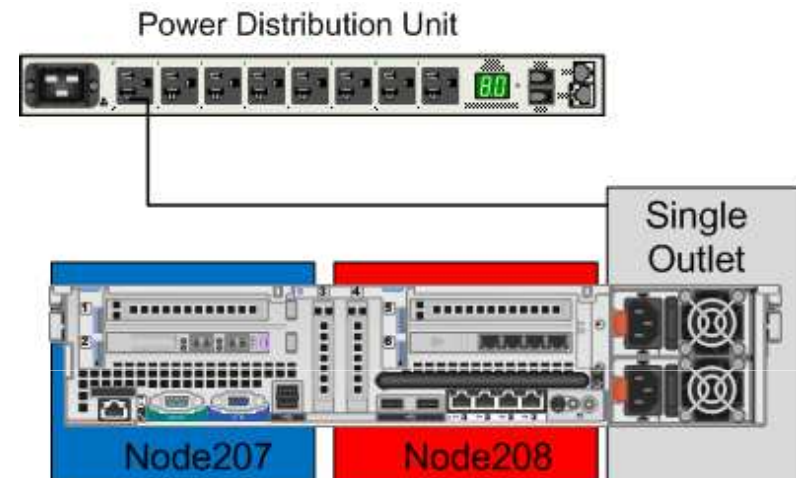- Cloud Computing
  - OpenNebula

# Topology



| Cluster | Head node | Compute nodes |
|---------|-----------|---------------|
| VU | fs0.das4.cs.vu.nl | 001-075 |
| LU | fs1.das4.liacs.nl | 101-116 |
| UvA | fs2.das4.science.uva.nl | 201-218 |
| TUD | fs3.das4.tudelft.nl | 301-332 |
| UvA-MN | fs4.das4.science.uva.nl | 401-436 |
| ASTRON | fs5.das4.astron.nl | 501-523 |

# Current Setup

## Cluster environment
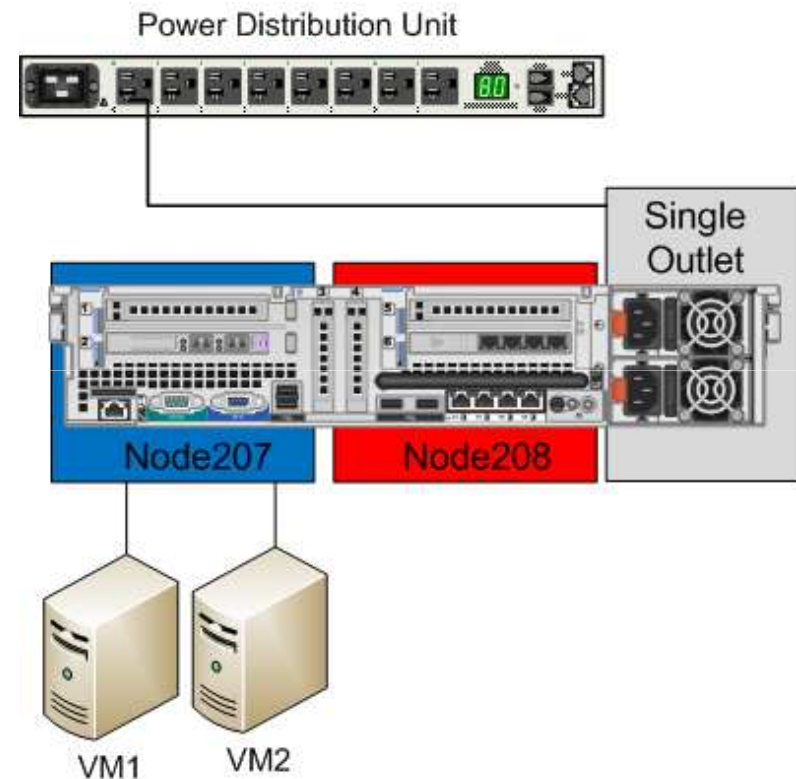
- 2U Twin Server

- Single outlet for the entire server



*Rear View*

## Cloud environment

- Single node with two VMs

- Only one energy source for both VMs

- Why?

  - No monitoring tools;

  - Concurrent resource share;

Workload measurement

- Bright Cluster Manager

Power management

- Racktivity PDUs

Correlation of the two systems

- Workload and energy

# Bright Cluster Manager

# Metrics

| Metric | Extraction Method | Source |
|---|---|---|
| Execution time | As reported by the Job | Job |
| Power Consumption | Python Script | PDU |
| Energy Consumption | Python Script | PDU |
| CPU Load | Python script | Bright Cluster Manager |

# Linpack Vs Polyphase Filter

- Linpack lacks the configuration option to control the amount of resources that it uses

- Polyphase filter is configurable, as regards the number of its runs and the used threads

- We define two different jobs; job1 and job2, so that job1 causes the double workload of job2

- We treat every single job as a unit and measure the power produced by each of them under various rates of CPU utilization
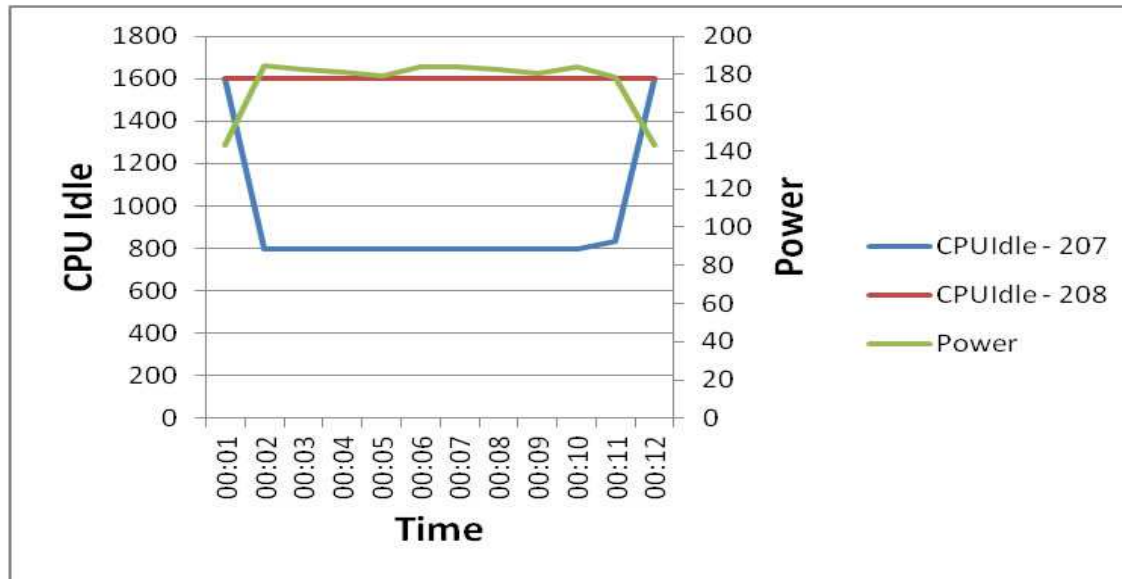
job 1 is running on node-207 and the adjacent node-208 is idle

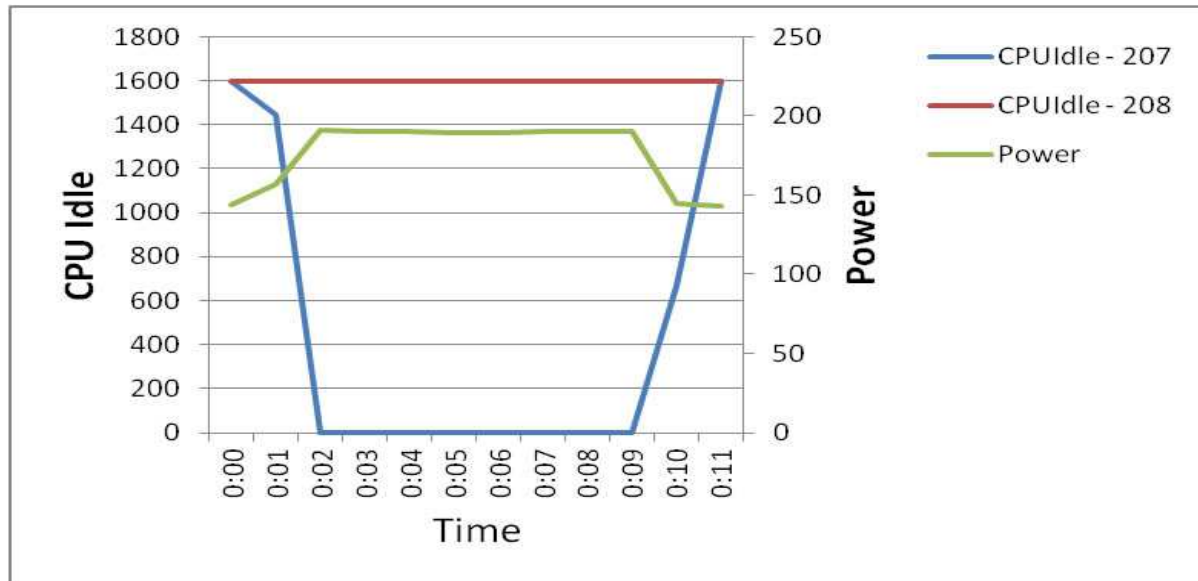| CPU Load Node-207 | CPU Load Node-208 | Peak of Power Consumption | Max Execution Time |
|---|---|---|---|
| 25% | 0% | 165,4 W | 1028 sec |

# Polyphase Filter – 50% workload



job 1 is running on node-207 and the adjacent node-208 is idle

| CPU Load Node-207 | CPU Load Node-208 | Peak of Power Consumption | Max Execution Time |
|---|---|---|---|
| 50% | 0% | 184 W | 587,6 sec |

job 1 is running on node-207 and the adjacent node-208 is idle

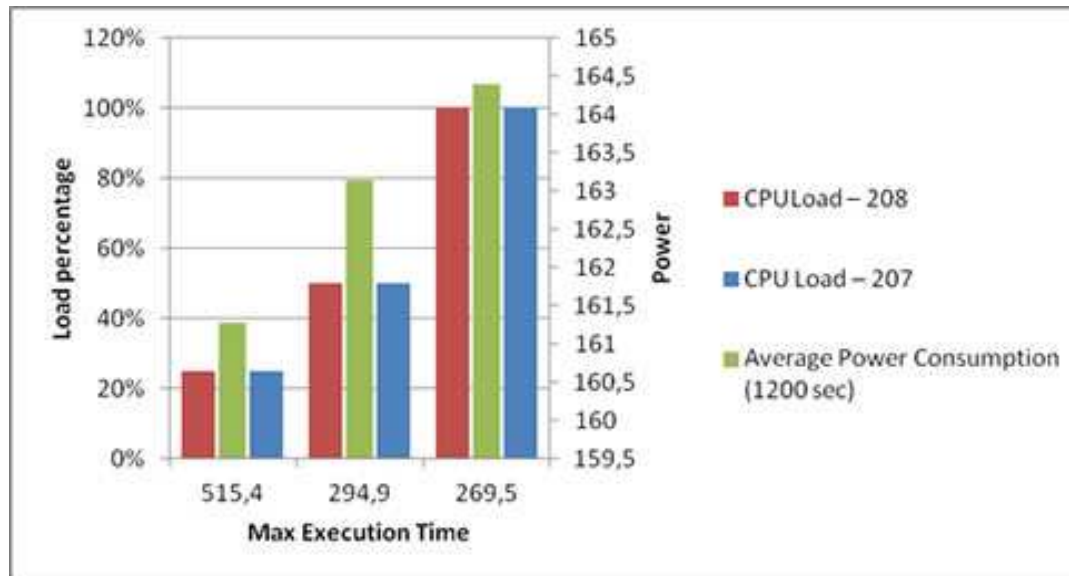| CPU Load Node-207 | CPU Load Node-208 | Peak of Power Consumption | Max Execution Time |
|---|---|---|---|
| 100% | 0% | 190 W | 530,3 sec |

# Results evaluation



To evaluate the trade-off between power consumption and performance for all the above cases, we built a **coupled in time** environment of 1200 sec

| CPU Load Node-207 | CPU Load Node-208 | Average Power Consumption In time interval equal to 1200 sec | Max Execution Time |
|---|---|---|---|
| 25% | 0% | 161,30 W | 1028 sec |
| 50% | 0% | 162,54 W | 587,6 sec |
| 100% | 0% | 162,62 W | 530,3 sec |

# Results evaluation



Finally in a short time interval, approximately equal to the longer execution time, **gains in power saving are almost negligible**.
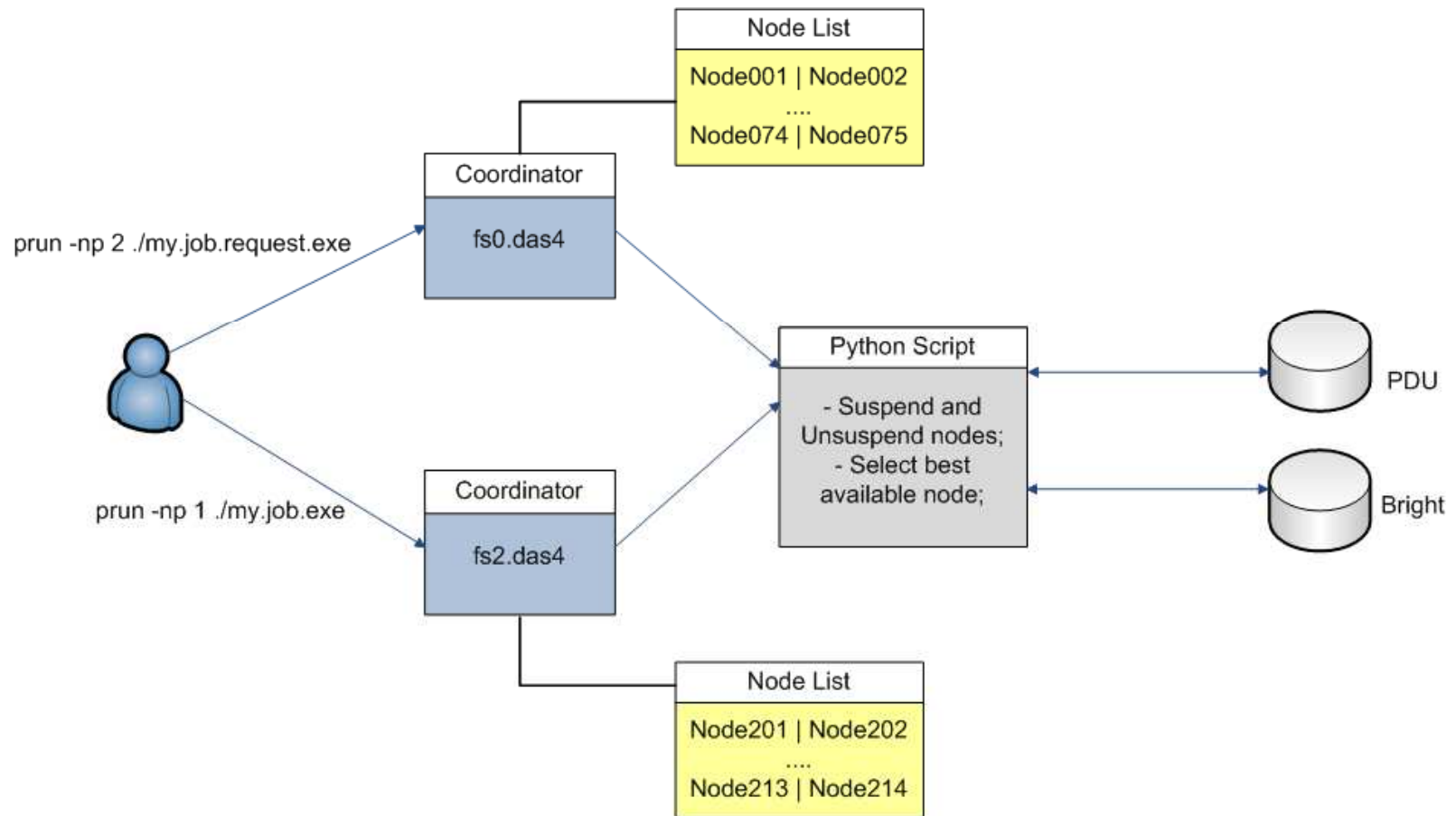
job2 = ½ job1

| CPU Load Node-207 | CPU Load Node-208 | Average Power Consumption In time interval equal to 1200 sec | Max Execution Time |
|---|---|---|---|
| 25% | 25% | 161,27 W | 515,4 sec |
| 50% | 50% | 163,14 W | 294.9 sec |
| 100% | 100% | 164,39 W | 269,5 sec |

# Conclusions

- Definite execution time job
  - Better performance using roughly the same amout of power
  - Grant execution in available nodes which share the same physical server

- In the current cluster implementation, it is impossible to execute more then one job at a time
  - Queue system

**Questions?**