

Energy efficiency in reservation-based large scale distributed systems

Laurent Lefèvre, Anne-Cécile Orgerie

INRIA RESO Project-Team

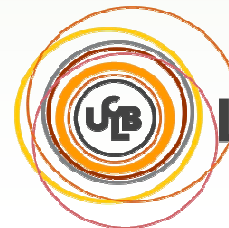
Ecole Normale Supérieure of Lyon - Université de Lyon

laurent.lefevre@inria.fr

HPC Workshop, Cetraro, June 27, 2011



INRIA
RHÔNE-ALPES



Lyon 1



Energy : 1st challenge for large scale systems (datacenter, grids, clouds, internet)?

- Future exascale platforms -> systems from 20 to 100MW (current 4-6 MW – 10 MW for Kei)
- How to build such systems and make them energy sustainable/responsible ?
 - Hardware can help (component by component)
 - Software must be adapted to be scalable but also more energy efficient
 - Usage must be energy aware





Between incentive and reality :

2010 : record year in CO₂ emission : 30.6 Gigatons
(+5% previous record in 2008) IEA

44% coal – 36% petrol – 20% natural gas

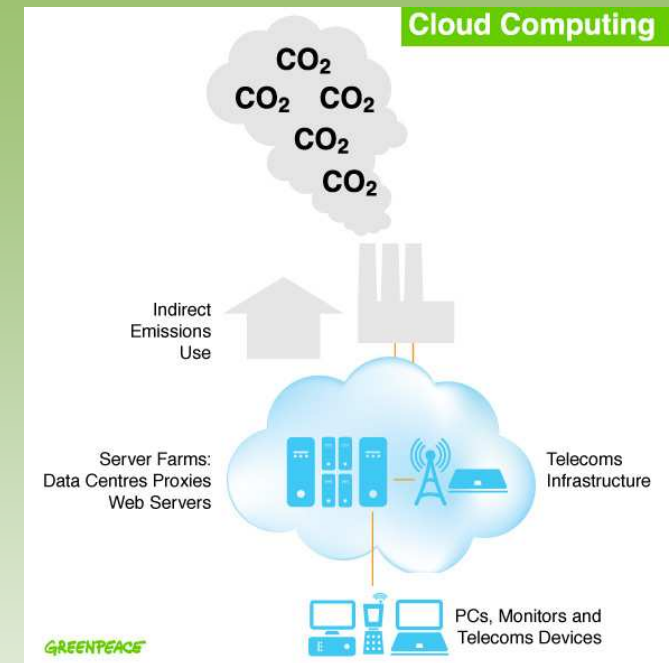
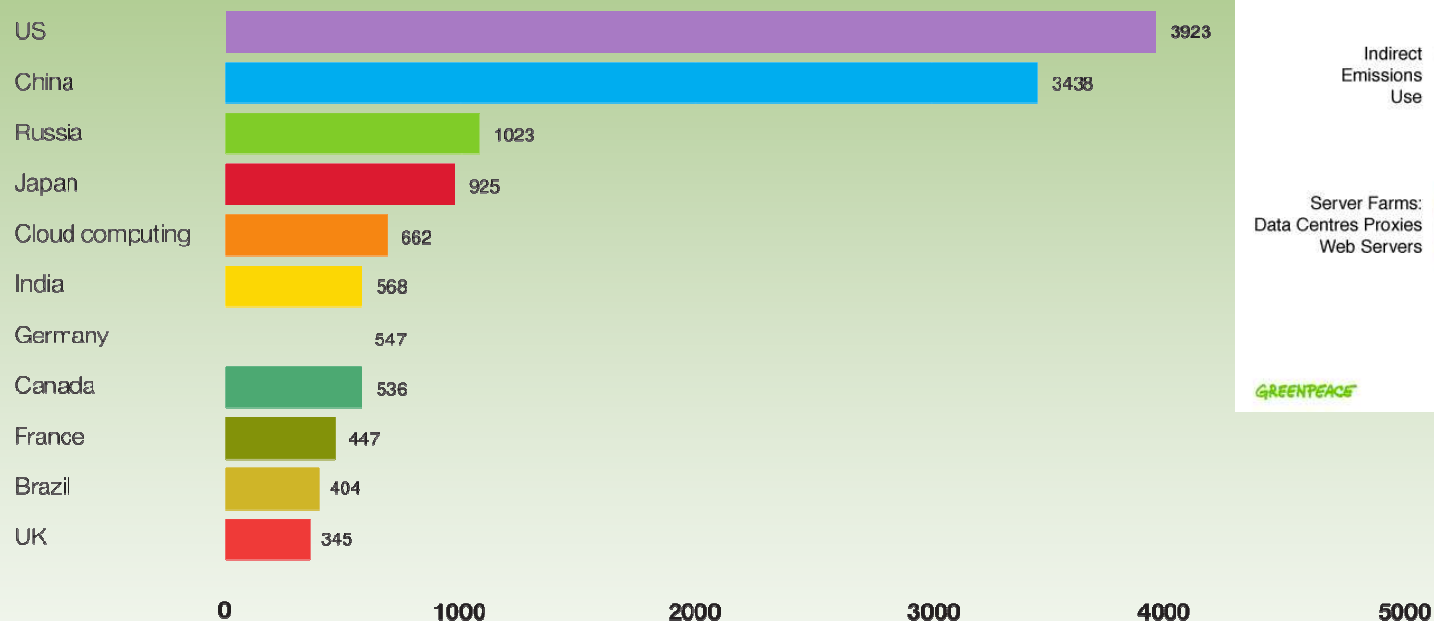
Temperature increasing (2°C – 2100) -> 4°C (50%
chance – 2100)



We are part of climate changing !

- Or at least of enormous electricity usage
- As IT users/designers

2007 electricity consumption. Billion kWh



- « Greenpeace reports 2010-2011 »

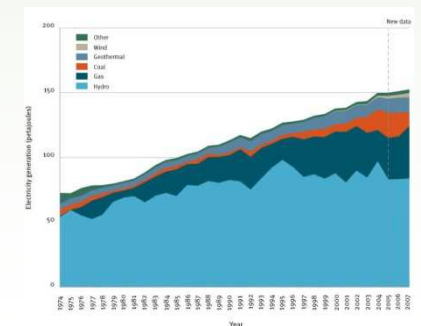
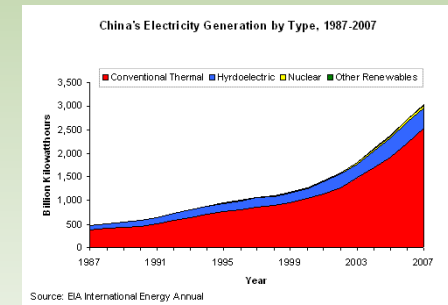
Because it is (intensively) used !



- OK results maybe not very accurate, but estimations...
- Seems exact :
 - 200 M emails / minute
 - 700 000 google request per minute
 - 48 hours of youtube video per minute

Power demand and Green IT

- IT – 2-5% of CO2 emissions
- Green It → reducing electrical consumption of IT equipments - CO2 impact depends on countries
- Focus on usage : fighting un-used/over-provisioned plugged resources



The 4 loops issue

- **Usage loop** : hardware , applications & services energy usage.
What is the impact of my application to energy consumption ?
- **Infrastructure Loop** : embedding into the models the infrastructure energy cost (air/water/free/un cooling). What is the impact of my environment to infrastructure consumption ?
- **Life cycle loop** : from production, usage to destruction (recycling)
: What is the energy cost of my IT during its whole life?
- **Human loop** : add the human factor around the precedent loop.
How to keep the world « happy » while reducing energy usage? 😊
- The holy Grail : a model and frameworks with the 4 loops inside
- This apply to all kind of products : food (local apple vs foreign bananas), cars (prius vs hummer)... etc...
- But with IT it should be easier/quantifiable... hum.. not yet...

Towards Energy Aware Large Scale Systems

How to decrease the energy consumption without impacting the performances?

- How to understand and to analyze the usage and energy consumption of large scale platforms?
- How to monitor lively such usage from pico to large scale views?
- How to design energy aware software frameworks ?
- How to help users to express theirs Green concerns and to express tradeoffs between performance and energy efficiency ?

Green-IT Leverages

- **Shutdown** : reducing the amount of powered unused resources
- **Slowdown** : adapting the speed of resources to real usage
- **Optimizing** : improving hardware and software for energy reduction purpose (i.e. energy aware libraries). Adapt software to green hardware.
- **Coordinating** : using large scale approaches to enhance green leverages



Explosion of (uncoordinated) initiatives

For each domain

- Data centers/HPC : Green500, EU CoC, The Green Grid
- Storage : SNIA
- Networks : Green Touch / EEE



Our Methodology



- Proposing a generic model able to be derivated onto different scenario (Grids, Clouds, Networks)
- Designing software solutions for infrastructures
- Simulating and Validating at medium and large scale

Reservation-based systems

Every usage is based on a reservation (resources, duration, deadline):

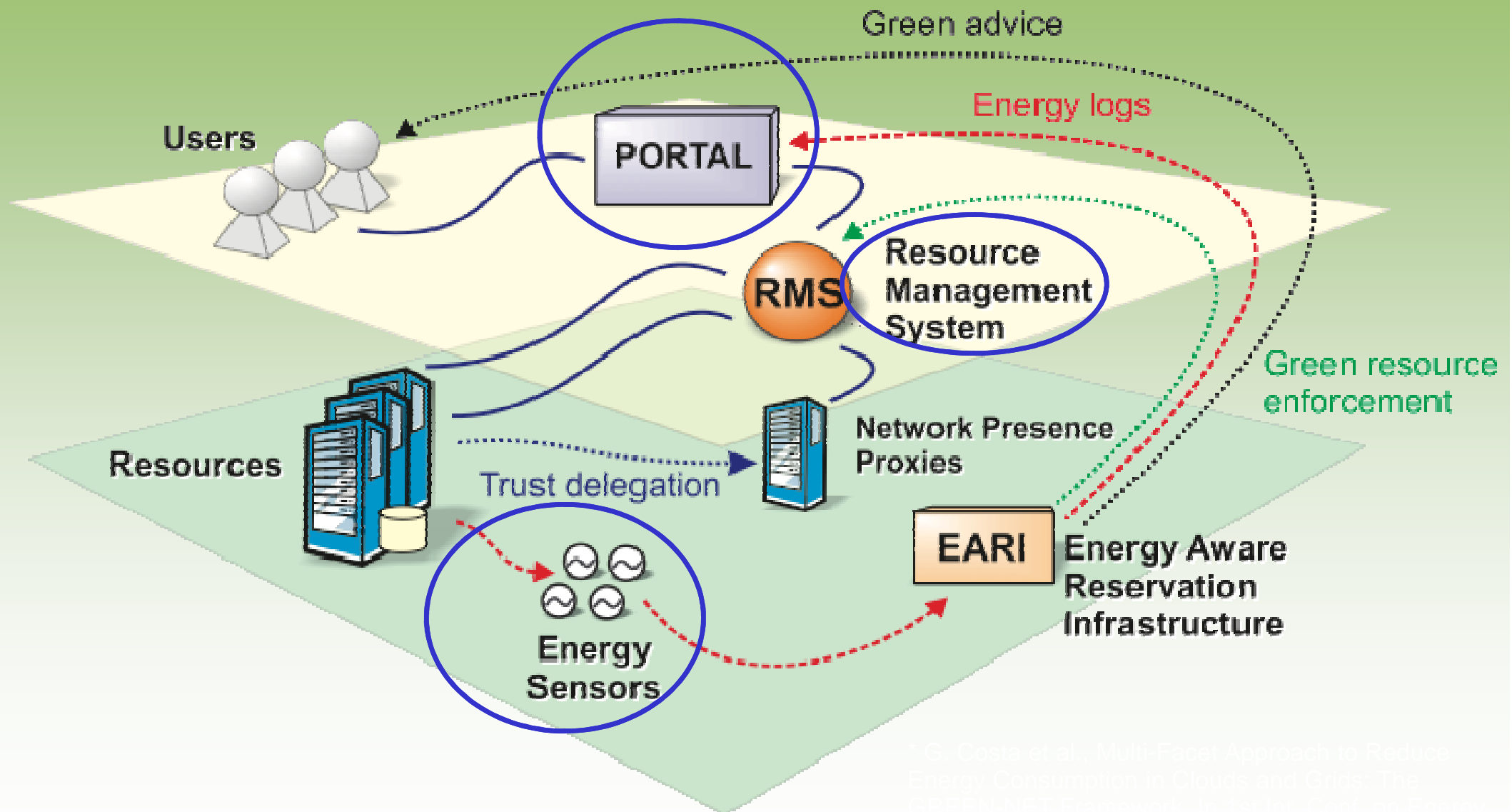
- Reserving cpu in HPC and Grids
- Reserving Virtual machines time in Clouds
- Reserving Bandwidth in large transport of data
- Leverages:
 - Finding and powering the optimal number of resources in front of needs of applications
 - HPC and Grids : switching on/off physical components
 - Clouds : switching on/off VMs
 - Networks : lighting or switching off paths, nics, links, routers, LPI
 - Adapting « speed » (and consumption) to the need of applications/users
 - HPC, Grids : dvfs
 - Clouds : tuning, capping
 - Networks : adaptive link rate

The ERIDIS model

Energy-efficient Reservation Infrastructure for large-scale
Distributed Systems

- Collecting and exposing : usage, energy profiling of applications and infrastructures
- Predicting usage of infrastructures
- Expressing and Proposing : to deal with tradeoffs between perf and energy, Green Policies
- Agregating resources reservations and usage in time and space
- Enforcing Green leverages : switch on/off or adapt performances

The ERIDIS Framework

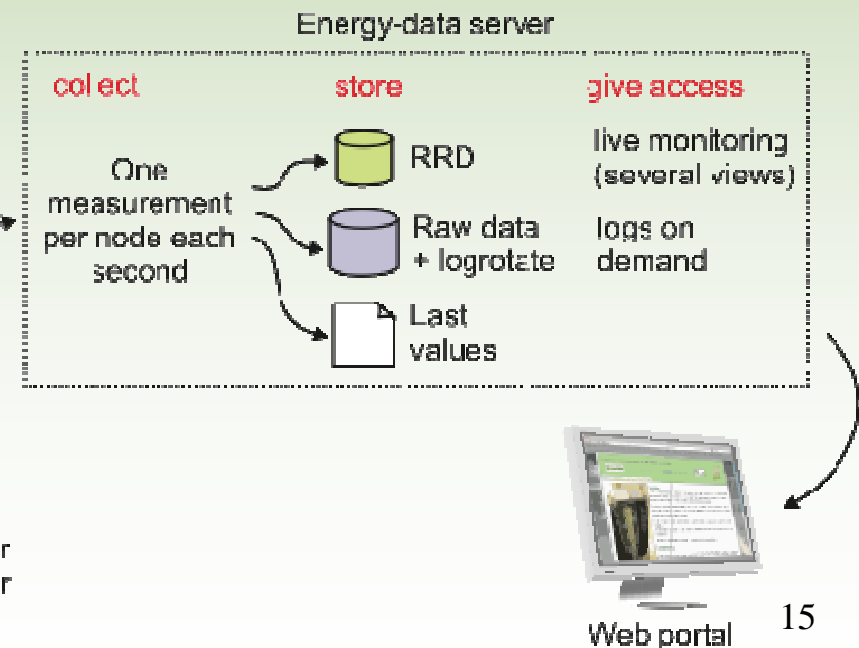
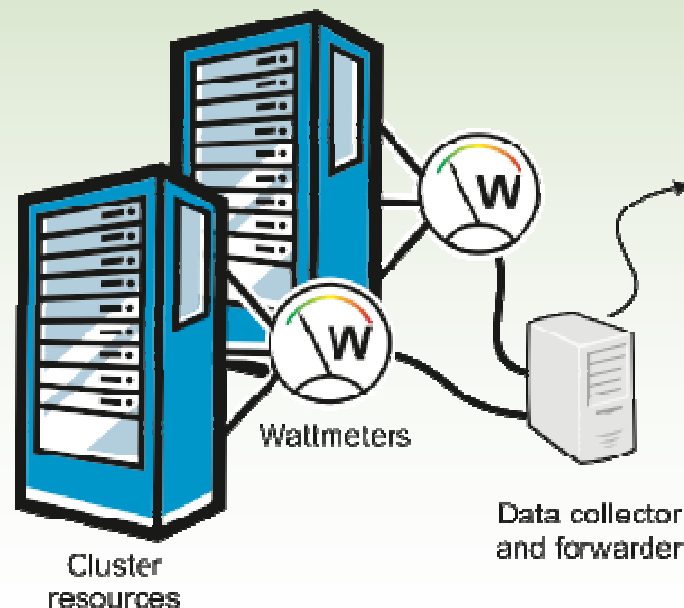
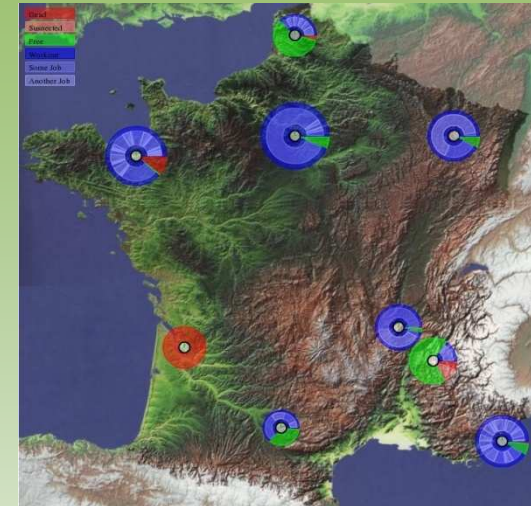


* G. Costa et al. Multi-Facet Approach to Reduce Energy Consumption in Clouds and Grids: The ERIDIS Framework

Collecting and exposing



- (Green) Grid'5000
 - French experimental testbed
 - 7400 cores
 - 10 sites
 - External energy sensors
 - Full site monitoring

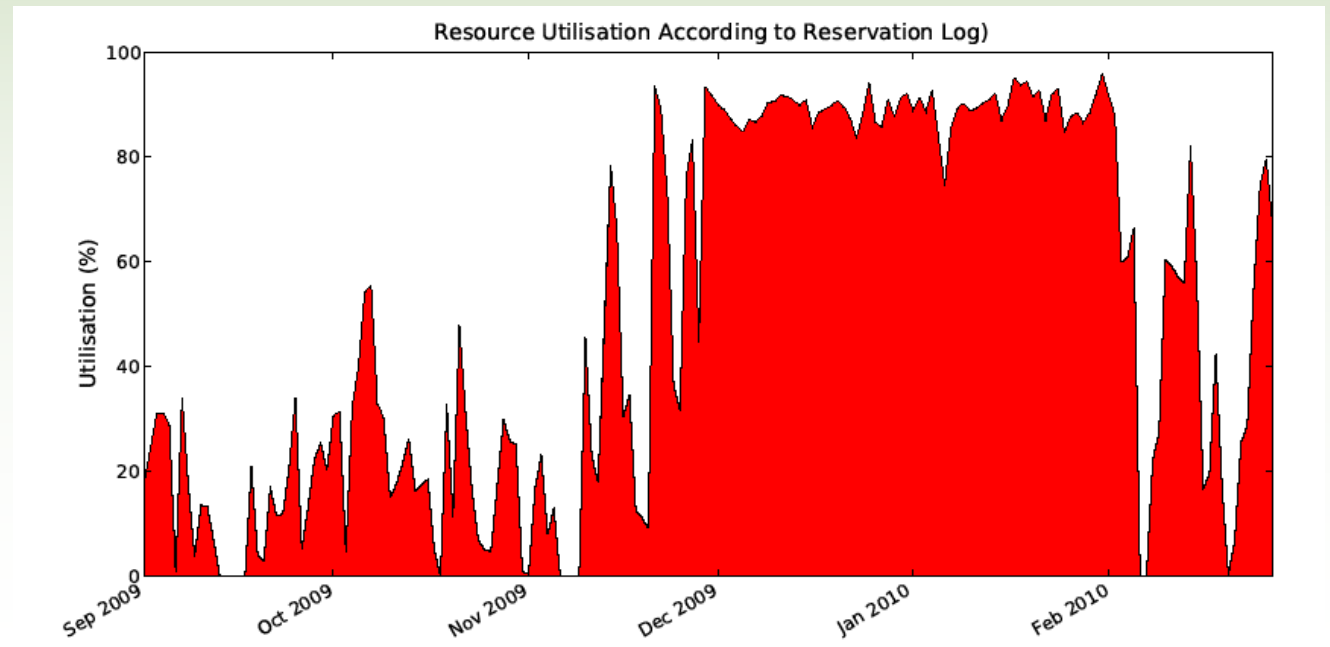
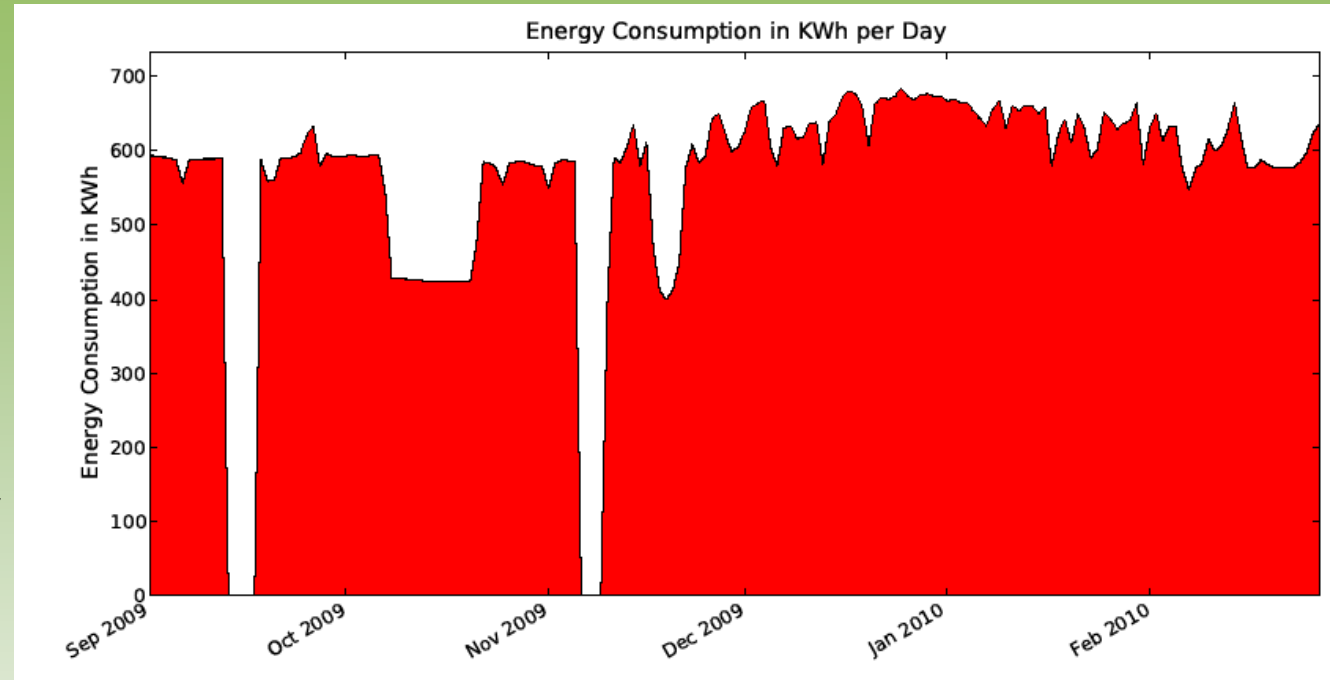


Electrical consumption / Usage

Periodicity of energy measurements:

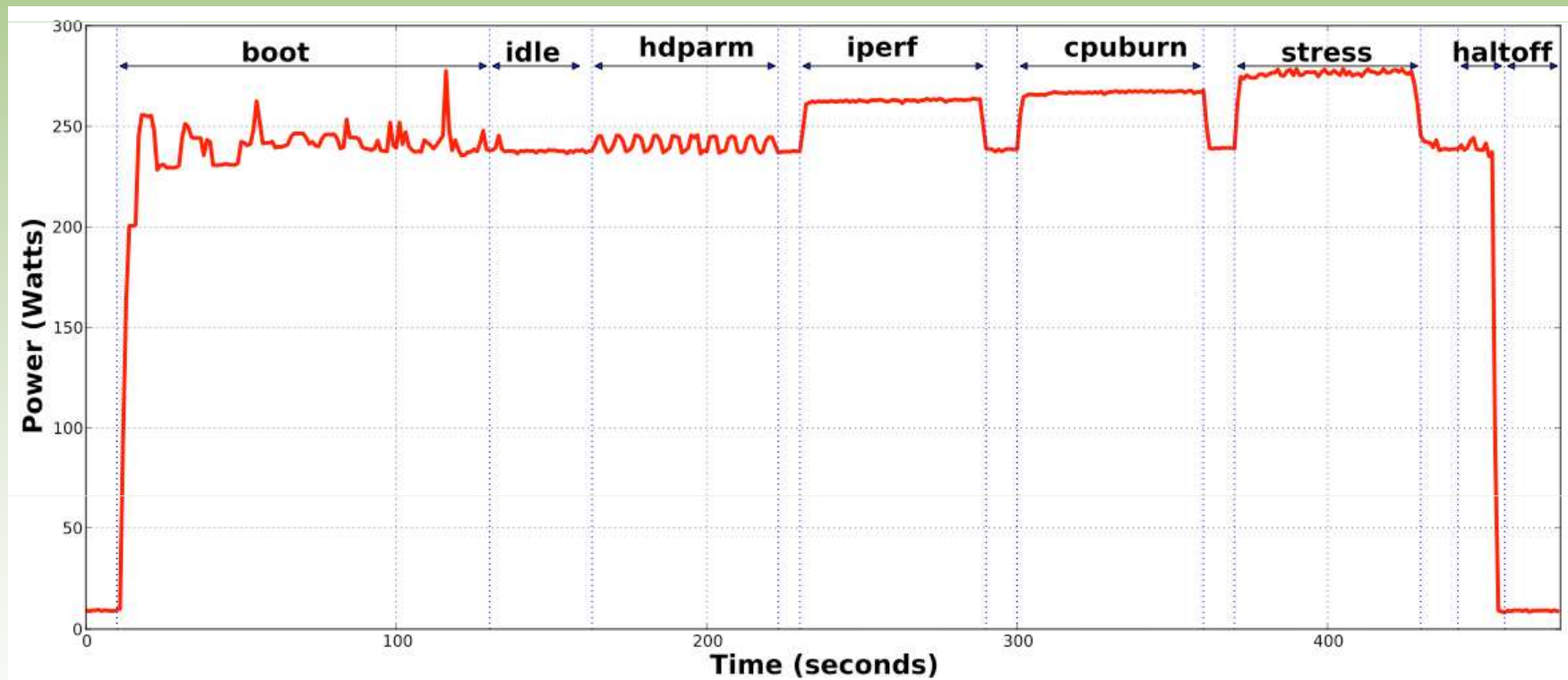
One measurement per **second** for each equipment

*



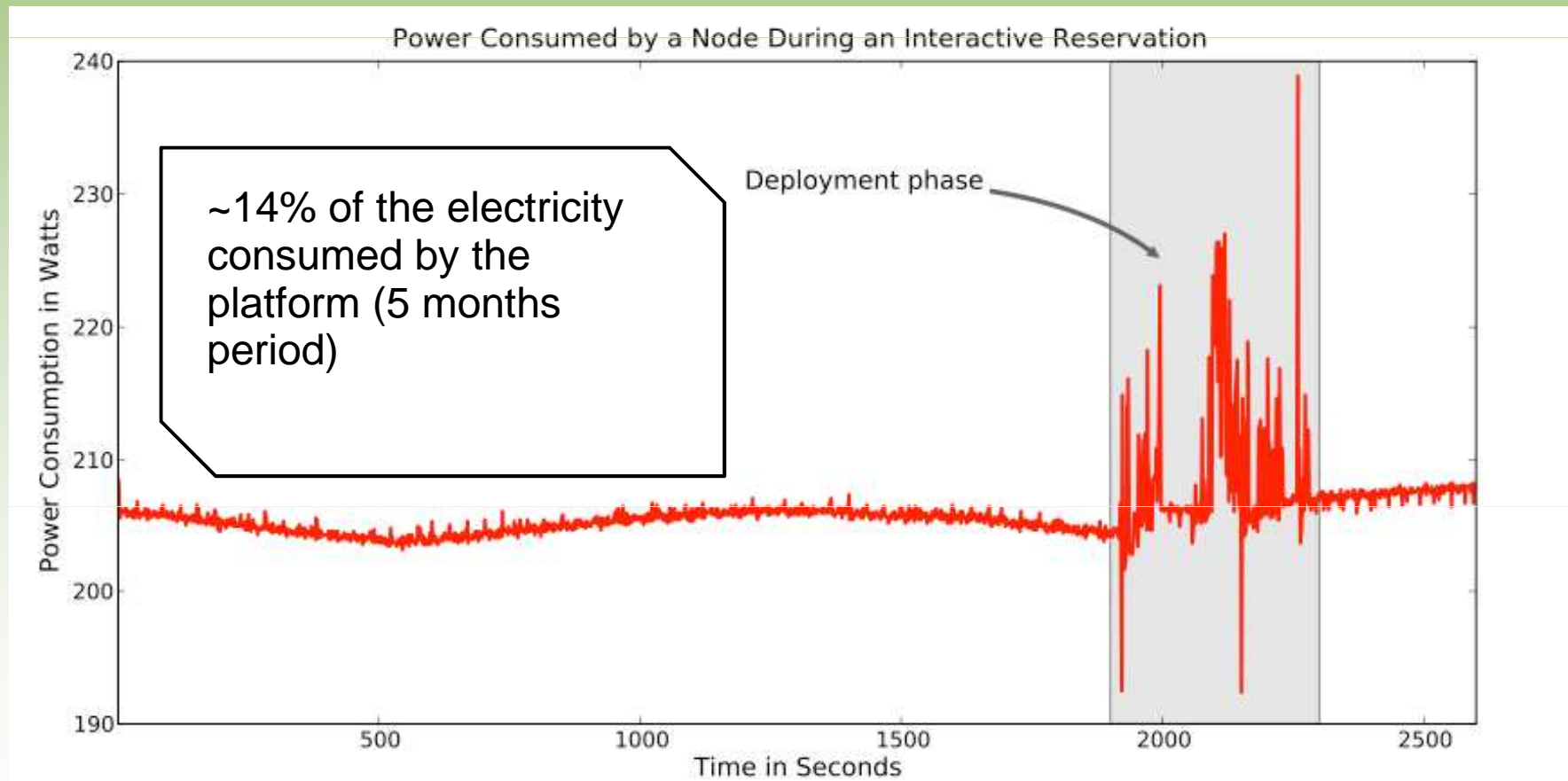
Example I : Profiling applications

Profiling the energy consumption of applications

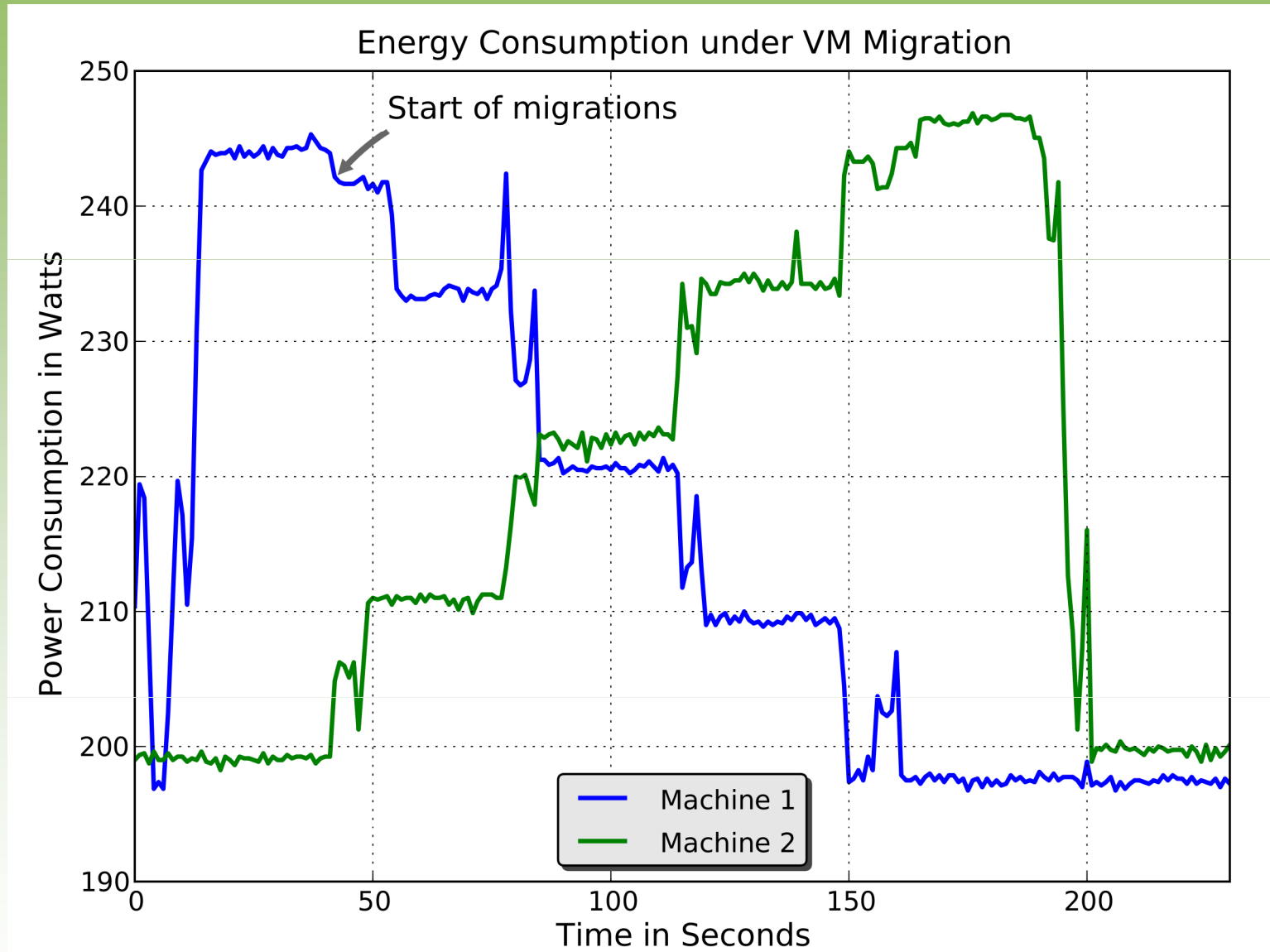


Example II : detecting anomalies

Improving frameworks/middleware and policies



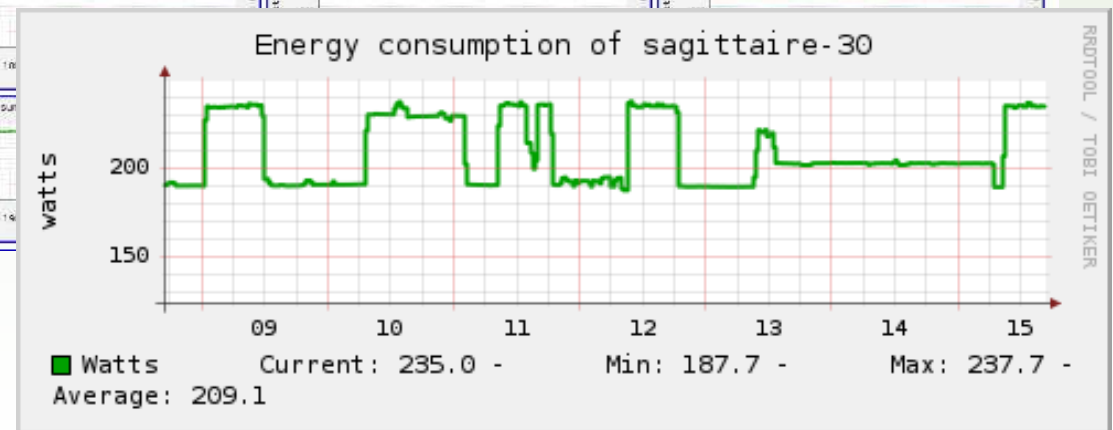
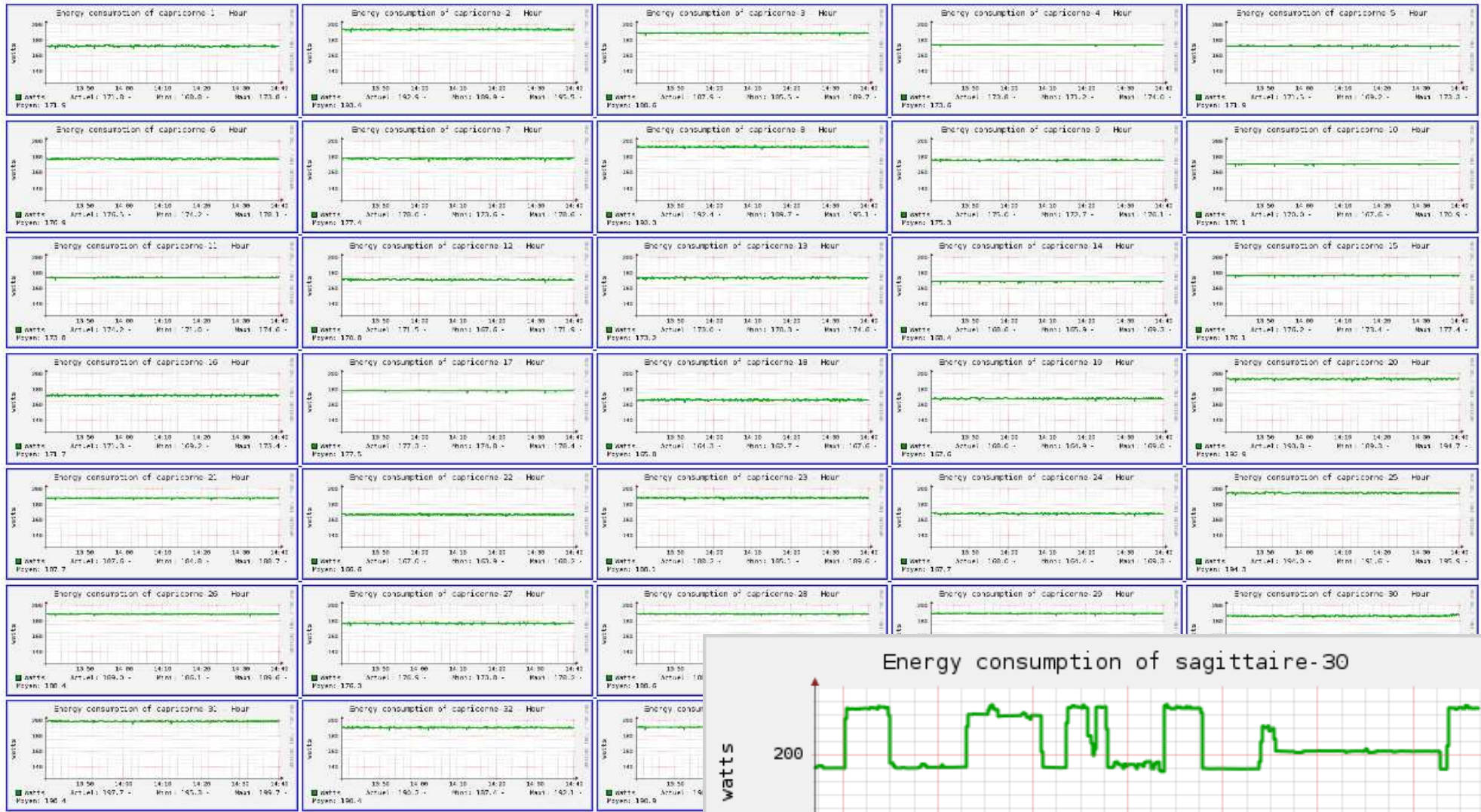
Migration

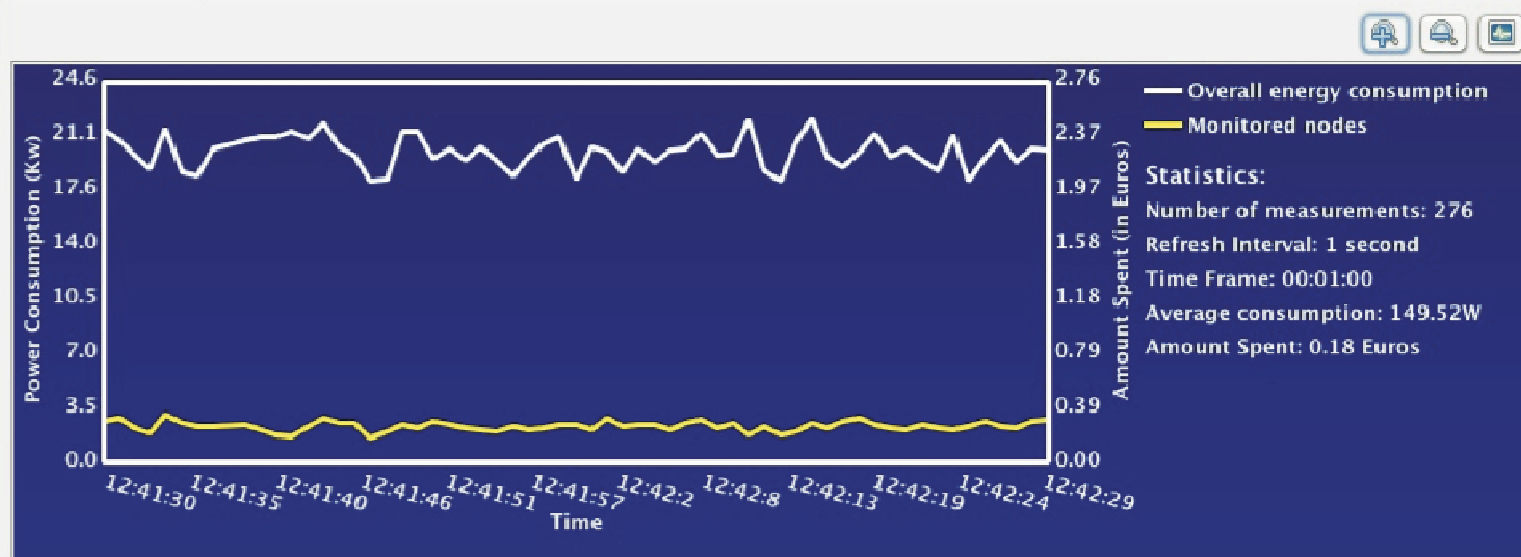


- Bad moment in energy during the migration

Large scale energy exposing

Energy Information of Lyon Grid5000 site

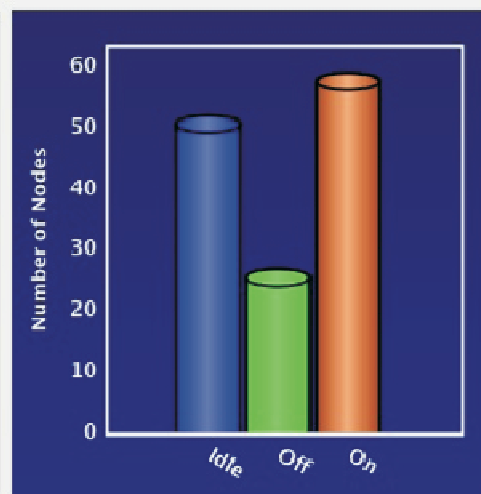




Status of Resources:

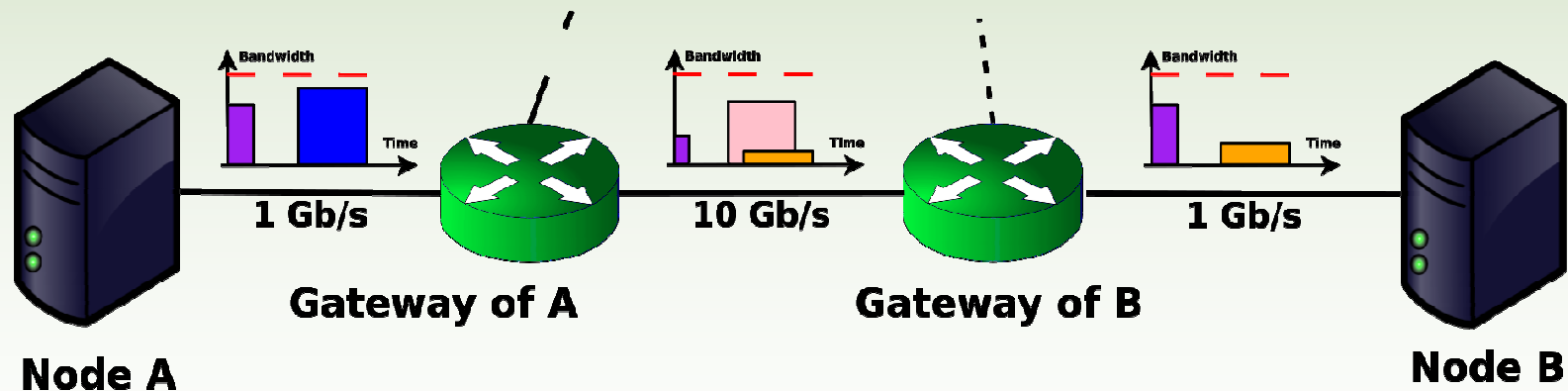
sagit-1 74.81W	sagit-11 294.94W	sagit-21 221.42W	sagit-31 163.69W	sagit-41 43.65W	sagit-51 193.71W	sagit-61 236.40W	sagit-71 64.54W	capric-2 241.65W	capric-12 171.55W	capric-32 83.97W	capric-42 180.02W
sagit-2 162.28W	sagit-12 276.10W	sagit-22 19.56W	sagit-32 274.28W	sagit-42 55.37W	sagit-52 73.74W	sagit-62 189.81W	sagit-72 203.15W	capric-3 192.85W	capric-13 186.97W	capric-43 130.27W	capric-53 226.64W
sagit-3 253.17W	sagit-13 257.72W	sagit-23 74.62W	sagit-33 10.06W	sagit-43 118.46W	sagit-53 220.34W	sagit-63 214.84W	sagit-73 133.10W	capric-4 72.71W	capric-14 52.98W	capric-44 21.68W	capric-54 40.37W
sagit-4 290.73W	sagit-14 32.88W	sagit-24 203.23W	sagit-34 225.22W	sagit-44 8.775W	sagit-54 245.74W	sagit-64 199.51W	sagit-74 234.59W	capric-5 177.33W	capric-15 14.16W	capric-45 2.61W	capric-55 43.12W
sagit-5 11.05W	sagit-15 84.01W	sagit-25 40.13W	sagit-35 298.92W	sagit-45 89.05W	sagit-55 245.91W	sagit-65 36.89W	sagit-75 29.49W	capric-6 22.13W	capric-16 261.25W	capric-46 1.7W	capric-56 171.48W
sagit-6 199.85W	sagit-16 87.00W	sagit-26 121.88W	sagit-36 166.51W	sagit-46 142.07W	sagit-56 69.71W	sagit-66 142.63W	sagit-76 55.75W	capric-7 43.11W	capric-17 12.22W	capric-47 2.11W	capric-57 1.1W
sagit-7 167.38W	sagit-17 103.75W	sagit-27 259.07W	sagit-37 285.37W	sagit-47 214.58W	sagit-57 289.71W	sagit-67 95.29W	sagit-77 287.10W	capric-8 180.11W	capric-18 180.11W	capric-48 9.40W	capric-58 1.56W
sagit-8 12.01W	sagit-18 221.81W	sagit-28 36.93W	sagit-38 213.72W	sagit-48 12.82W	sagit-58 47.50W	sagit-68 244.97W	sagit-78 150.37W	capric-9 203.09W	capric-19 13.43W	capric-49 1.56W	capric-59 1.56W
sagit-9 153.28W	sagit-19 69.04W	sagit-29 201.03W	sagit-39 77.61W	sagit-49 2.38W	sagit-59 298.60W	sagit-69 25.05W	sagit-79 37.01W	capric-10 113.88W	capric-20 113.88W	capric-50 1.56W	capric-60 1.56W
sagit-10 137.56W	sagit-20 216.04W	sagit-30 207.96W	sagit-40 129.01W	sagit-50 223.91W	sagit-60 244.97W	sagit-70 14.47W	capric-1 86.08W	capric-11 215.51W	capric-21 173.91W	capric-51 119.56W	capric-61 1.56W

Resource on Resource idle Resource off Resource monitored

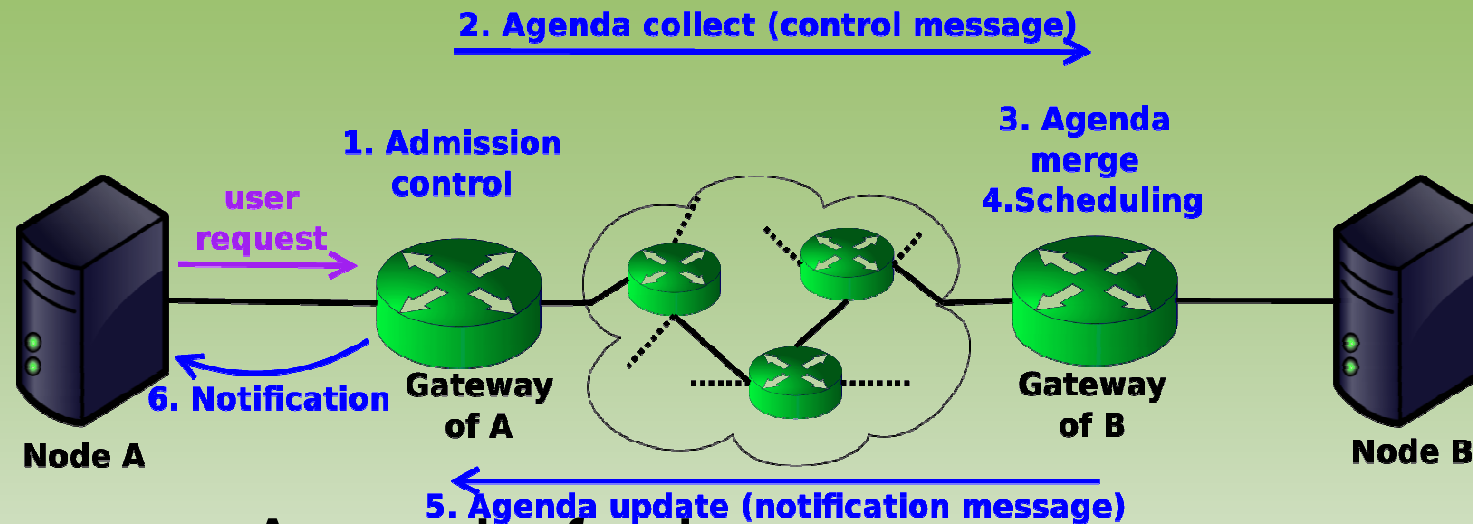


HERMES : High-level Energy-awaRe Model for bandwidth reservation in End-to- end networkS

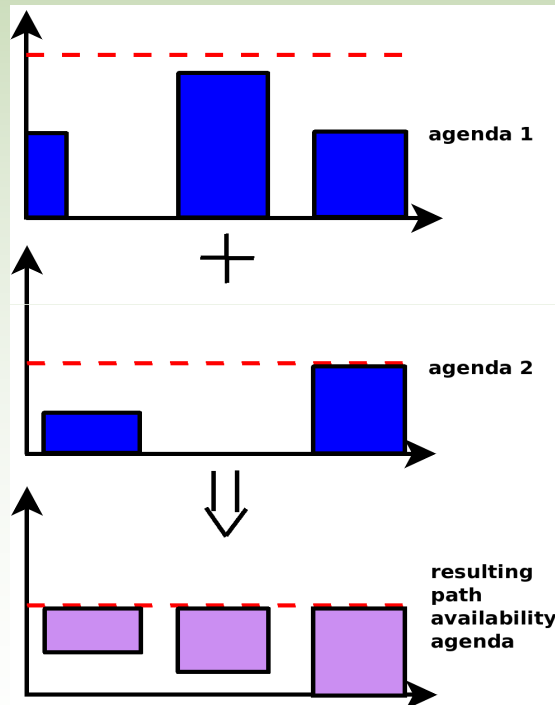
- Switching off unused parts of the network : NIC, routers, links
- Distributed network management
- Energy-efficient scheduling with reservation aggregation
- Usage prediction to avoid on/off cycles
- Minimization of the management messages
- Usage of DTN (Disruptive-Tolerant Network) for network management purpose



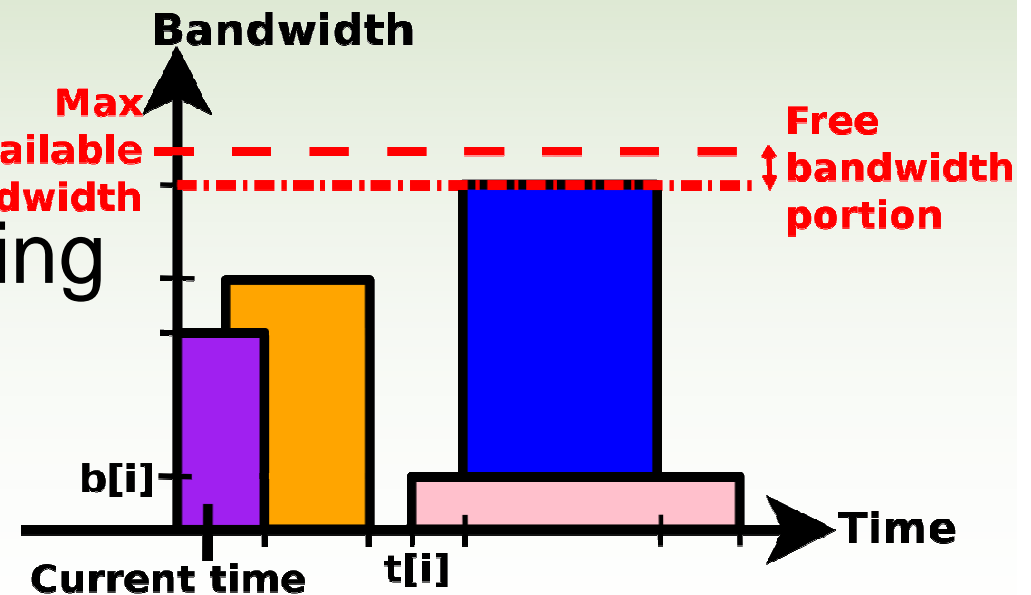
Hermes



Agenda fusion



ABR scheduling



Hermes results

- Network simulated: 500 nodes, 2 462 links.
- Random Network (Molloy & Reed method)
- All the nodes can be sources and destinations.
- Time to boot: 30 s.; time to shutdown: 1 s.
- 1 Gbps per port routers

Component	State	Power
Chassis	ON	150 W
	OFF	10 W
Port	1 Gbps	5 W
	100 Mbps	3 W
	idle, 10 Mbps	1 W

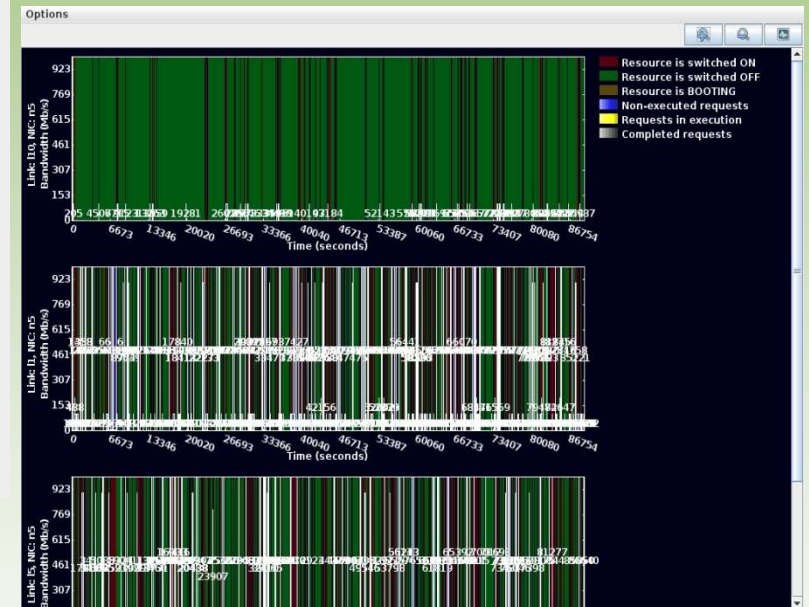
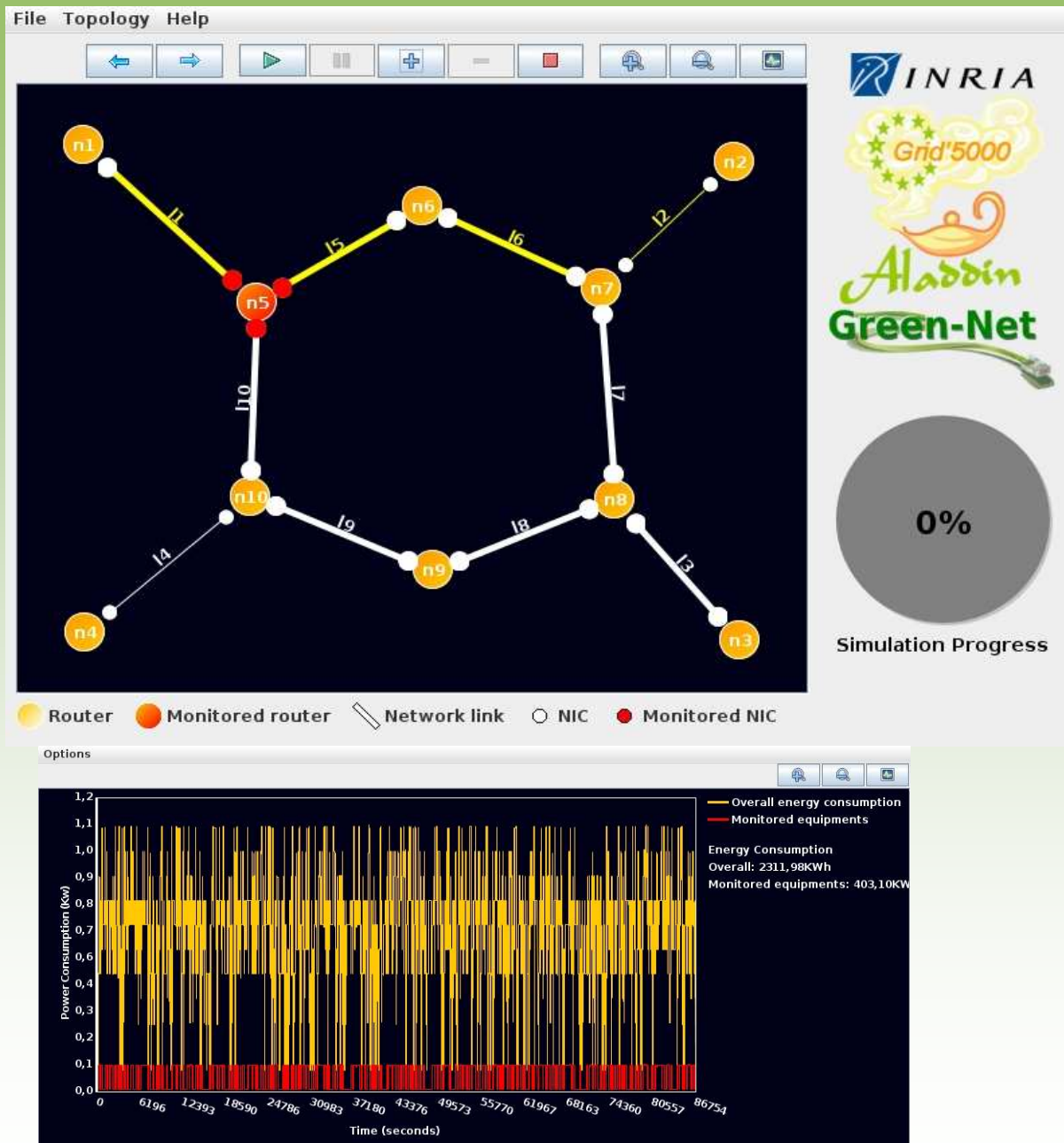
- 31% workload : Energy consumption in Wh

Scheduling	No off	First	First green	Last	Last green	Green
Average	412 306	205 270	203 844	204 949	196 260	203 342
Standard deviation	2 685	2 477	1 938	2 375	2 695	2 145
Accepted volume (Tb)	2 148	2 148	2 128	2 014	1 853	2 149
Cost in Wh per Tb	191.92	95.55	95.78	101.74	105.92	94.60

- Cost in Wh per Tb
- Compared to current case (no-off), HERMES could save 51%, 46% and 43% of the energy consumed depending on the workload

Workload	No off	First	First green	Last	Last green	Green
31%	191.92	95.55	95.78	101.74	105.92	94.60
46%	149.84	81.61	81.95	87.74	92.40	80.63
61%	130.45	74.73	74.91	80.09	84.63	73.79

Replayer



- 2010 SuperComputing demo, Marcos Dias de Assunção

Conclusions



- Big role for IT: Green IT and IT for Green
- Challenge : design energy proportional equipments and frameworks (computing, memory or network usage)
- Need to take out energy efficient models from the lab and put them in operationnal conditions
- Adress all levels from hardware to software – production and usage (4 loops)
- Explore EE best effort environments

Thanks to Anne-Cécile Orgerie

Questions?

Laurent Lefèvre
Laurent.lefevre@inria.fr

