

Ecole des Mines de Nantes



Self-organization of applications and systems to optimize resources usage in virtualized data centers

Teratec

06/28 2012

Jean-Marc Menaud

Ascola team EMNantes-INRIA, LINA

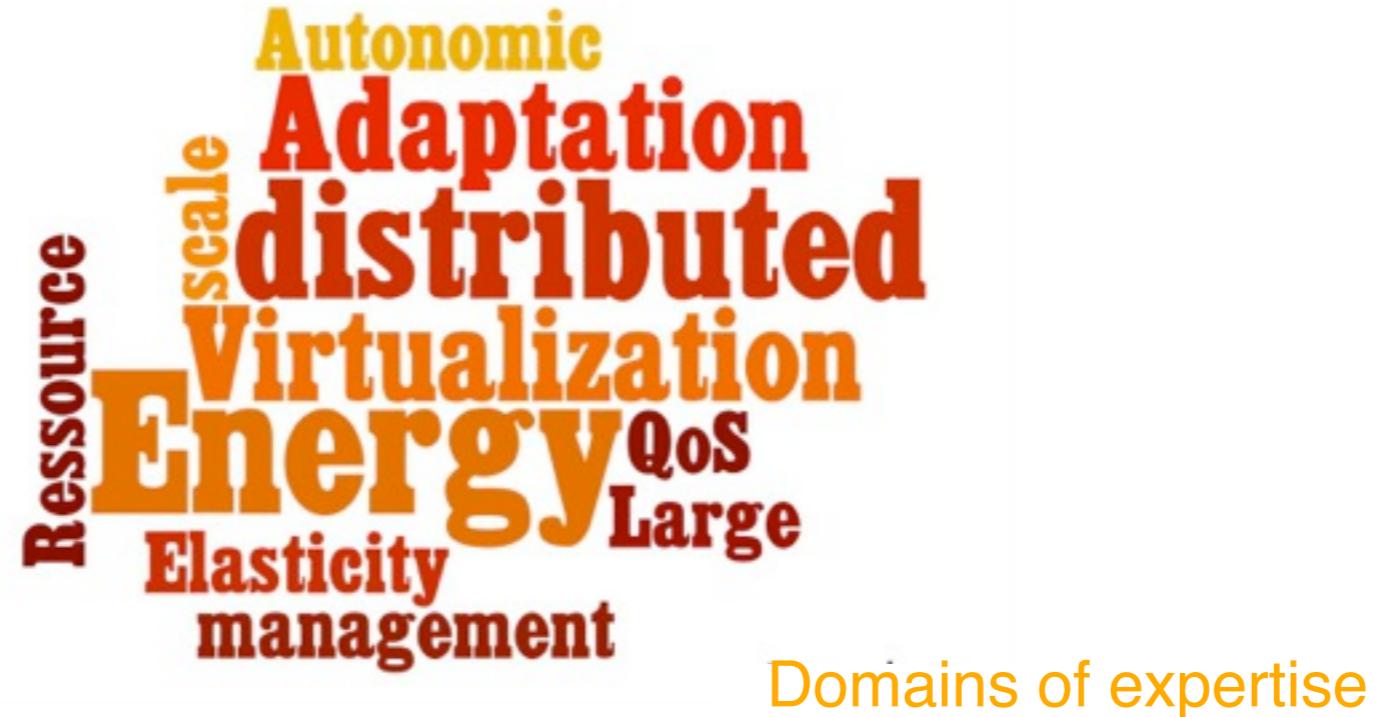




Motivations

- Increasing popularity of Cloud Computing solutions
- Data-centers (DCs) are amazingly growing
- DC providers have to face ressource management and energy consumption concerns

Key characteristics

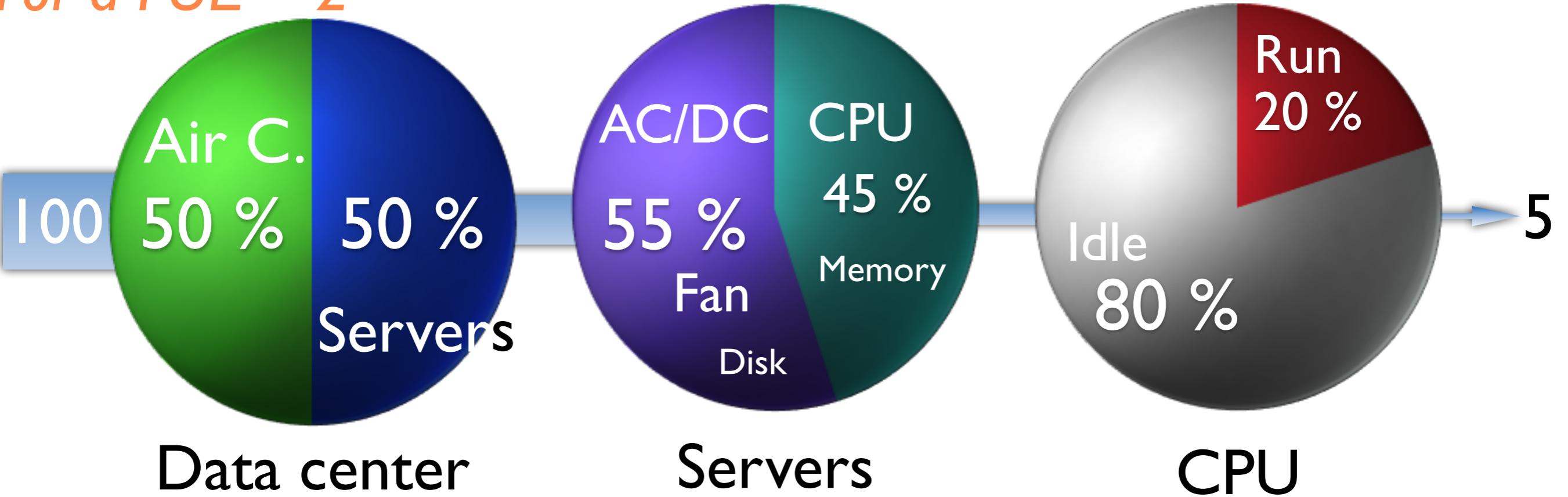


Overall objective:

Capacity planning is the process of planning, analyzing, sizing, managing and optimizing capacity to satisfy demand, swiftly and at a reasonable cost.

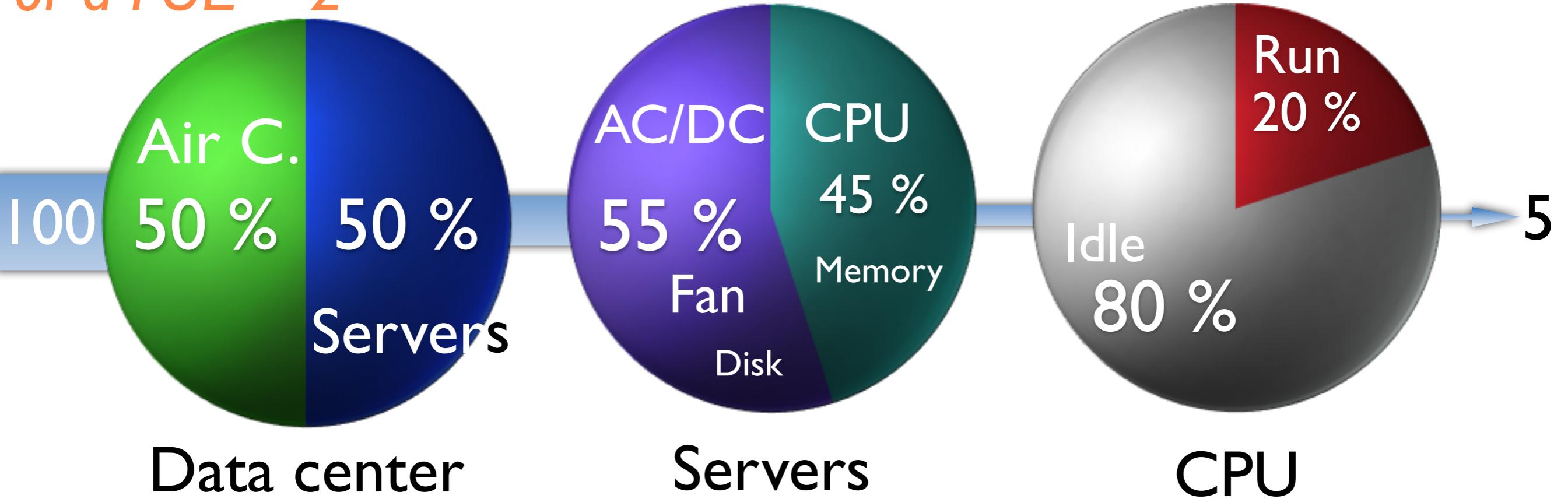
Capacity planning : Energy focus

For a PUE = 2



Capacity planning : Energy focus

For a PUE = 2



Data center

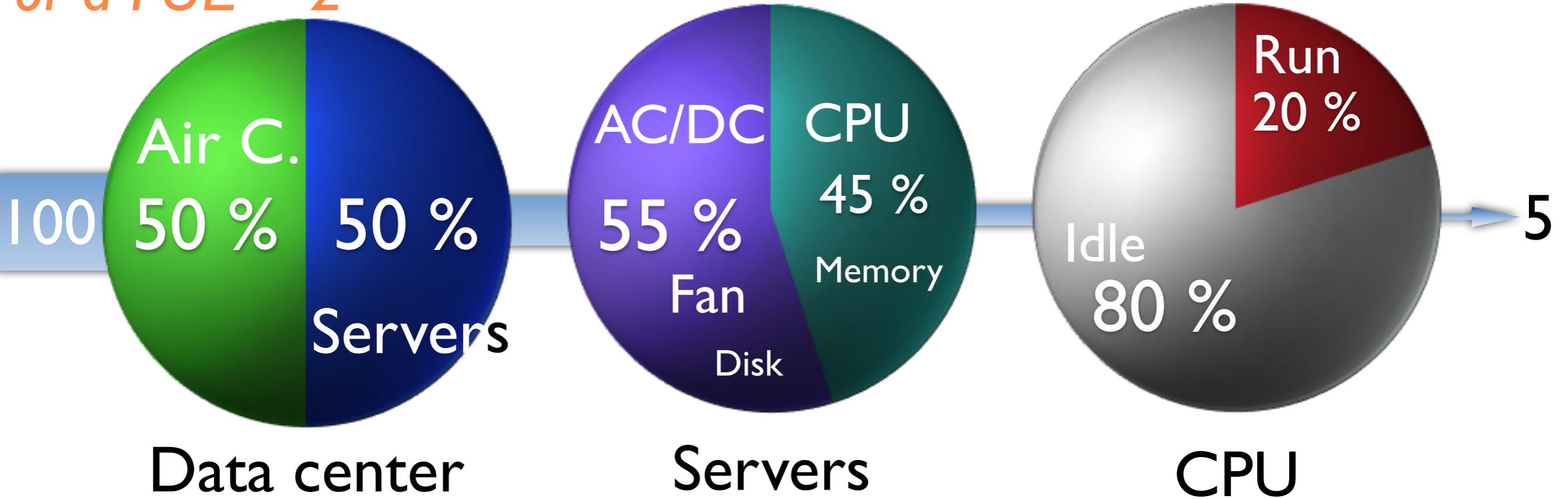
Servers

CPU

- **Analysis of the cost of a 2 MegaWatts DC (5000 nodes, 400w/h)**
 - PUE of 2, 0.06€/kWh => 2 120 886 €
 - A decrease of 5% enables a gain of 110K€
- **Managing DC resources finely becomes a major challenge**

Capacity planning : Energy focus

For a PUE = 2



Data center

Servers

CPU

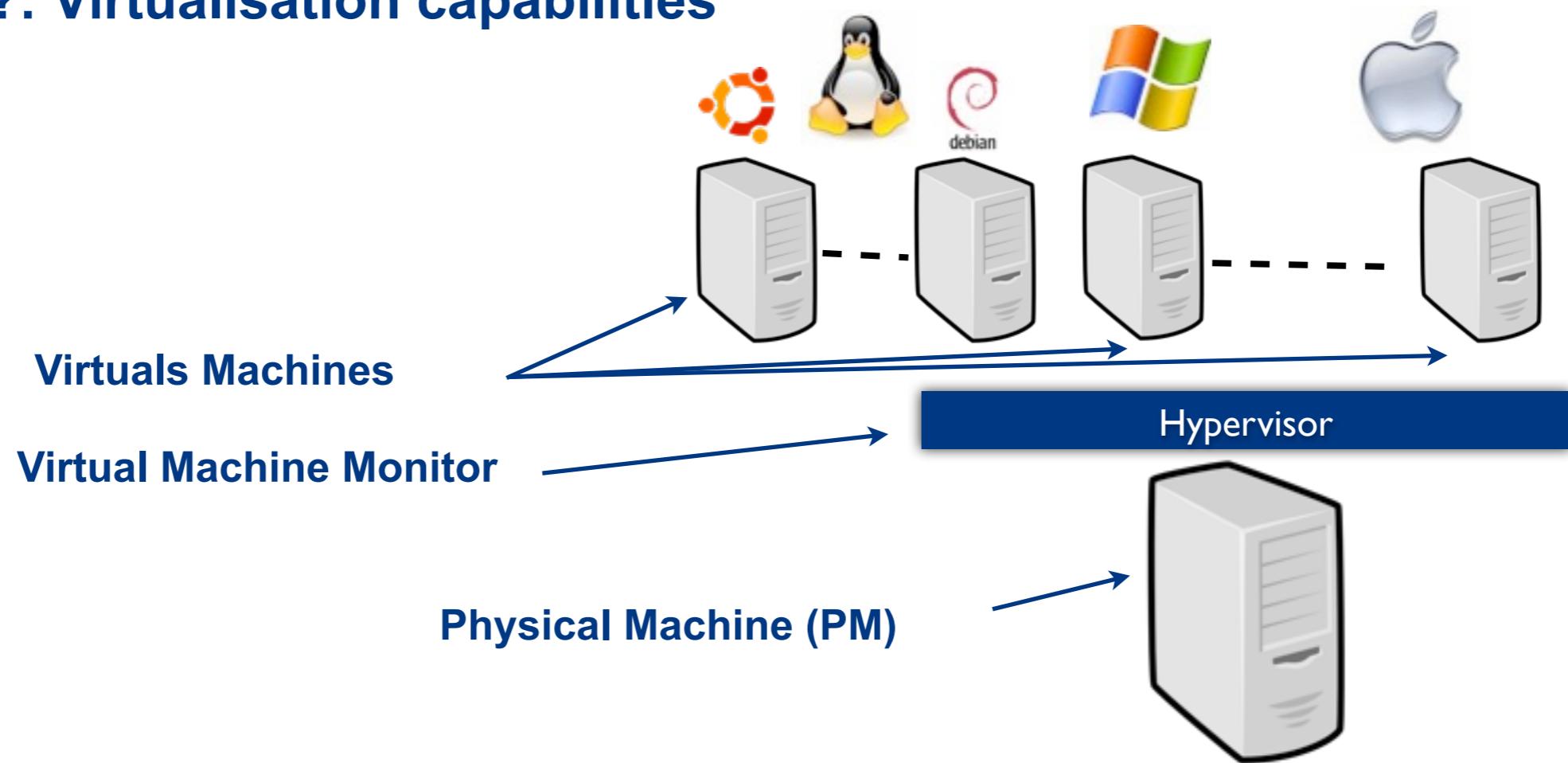
- **Analysis of the cost of a 2 MegaWatts DC (5000 nodes, 400w/h)**

- PUE of 2, 0.06€/kWh => 2 120 886 €
- A decrease of 5% enables a gain of 110K€

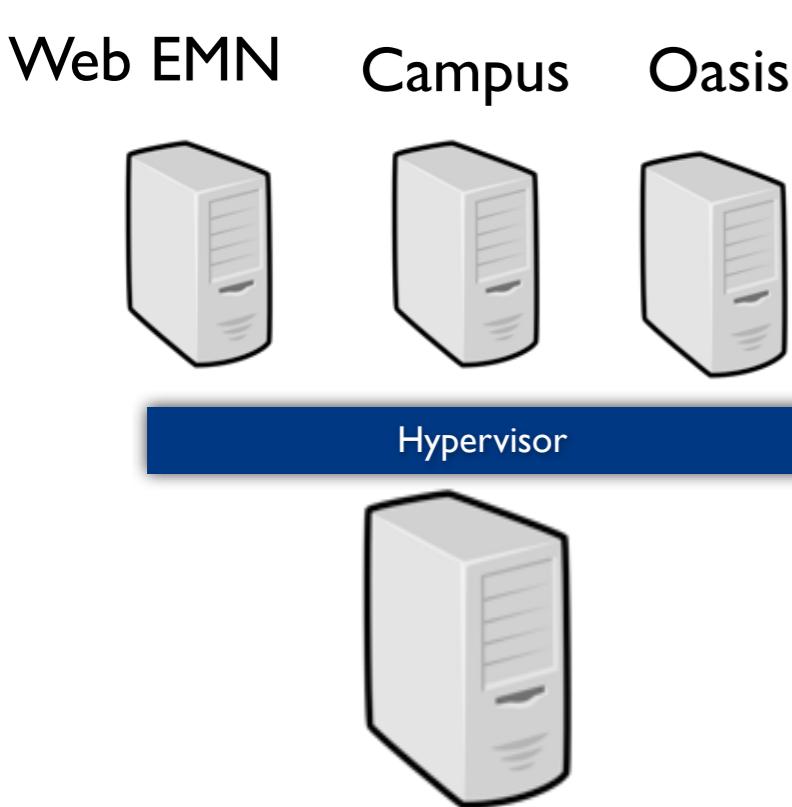
→ **Managing DC resources finely becomes a major challenge**

Consolidation

- **Consolidation :**
 - Consolidating computation reduces the number of running nodes
So energy consumption
 - Reduces hardware costs while providing more efficient node
- **How ?: Virtualisation capabilities**

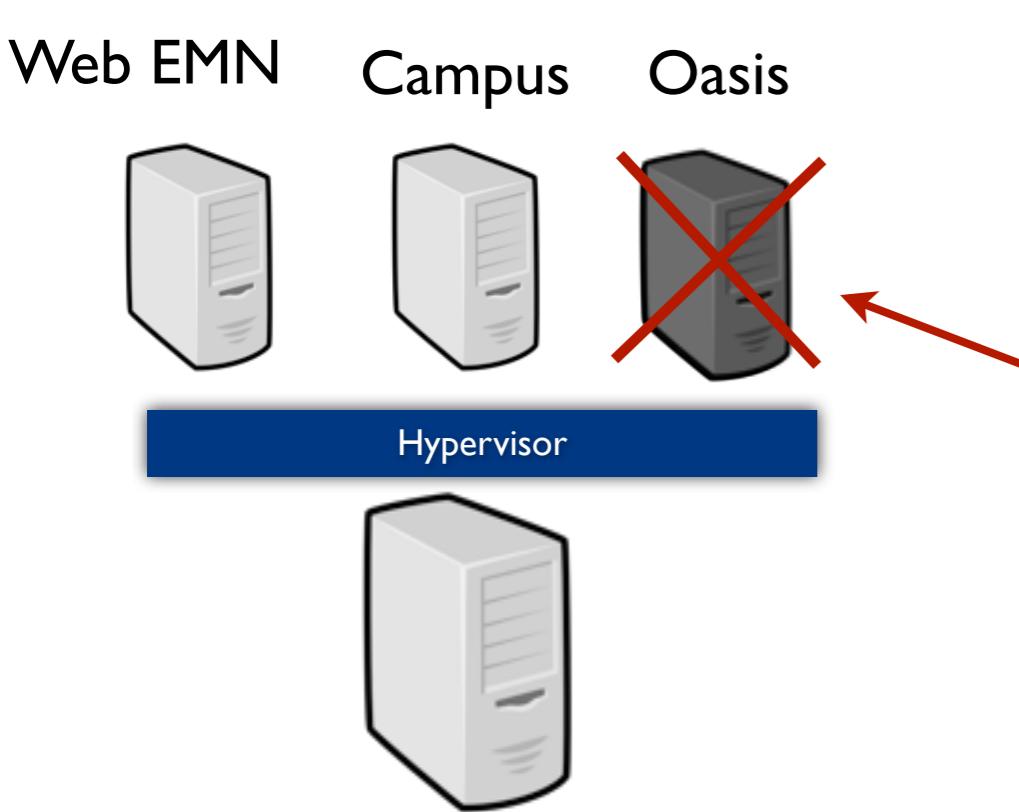


Virtualization capabilities



- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

Virtualization capabilities



- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

Virtualization capabilities

Web EMN

Campus

Oasis



- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

Virus / Invasion / Crash



- **Live migration (load-balancing)**
- **High Availability(downtime ~ 60 ms)**

Web EMN Campus Oasis



Virtualization capabilities

Web EMN

Campus

Oasis



Hypervisor



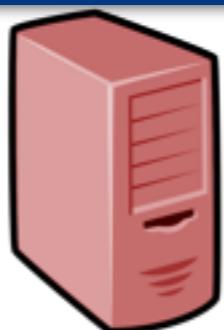
Virus / Invasion / Crash

- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

Web EMN Campus Oasis

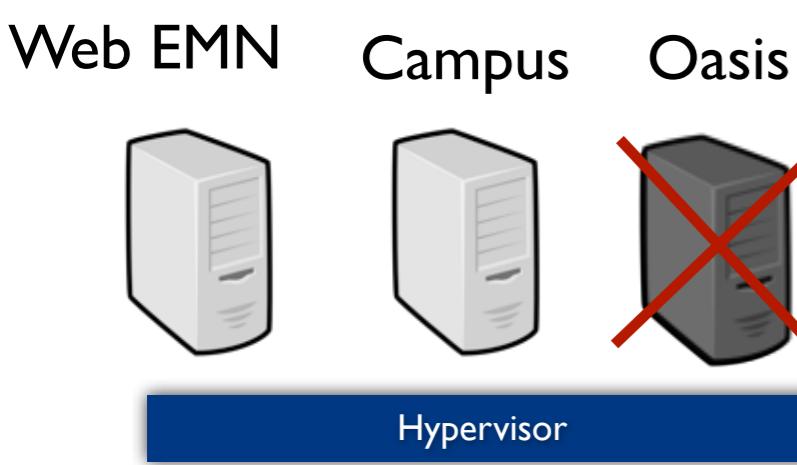


Hypervisor

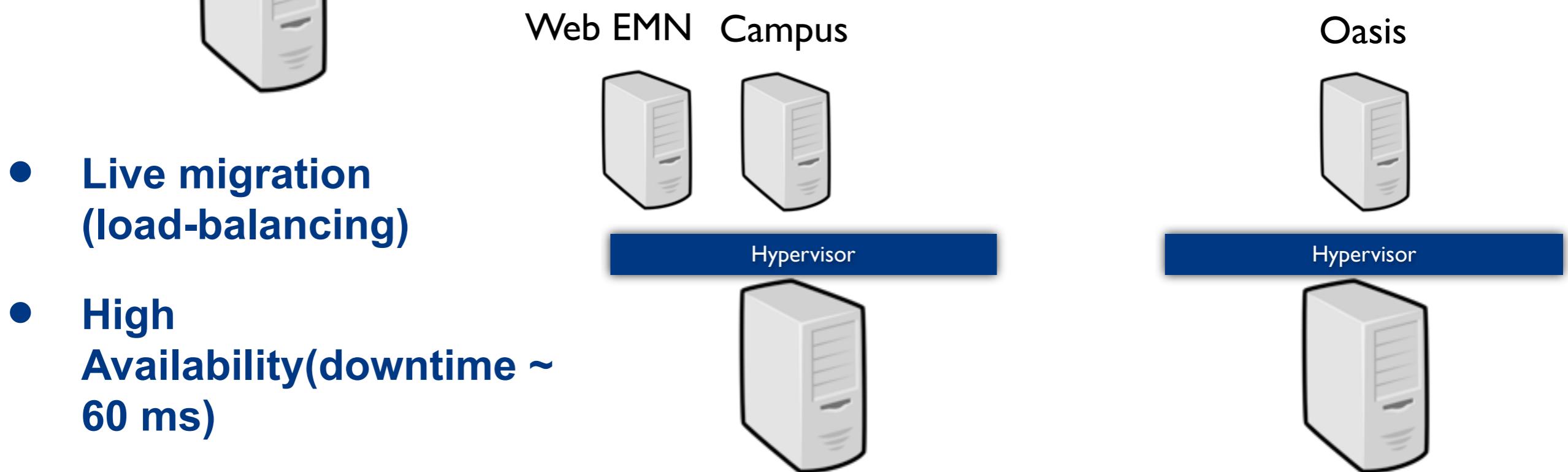


- **Live migration (load-balancing)**
- **High Availability(downtime ~ 60 ms)**

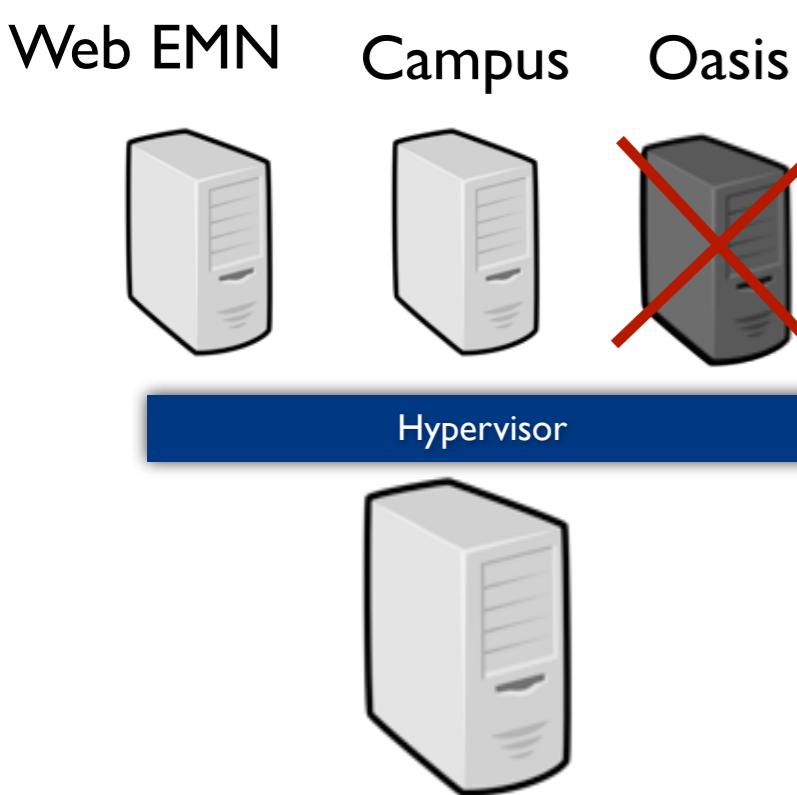
Virtualization capabilities



- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

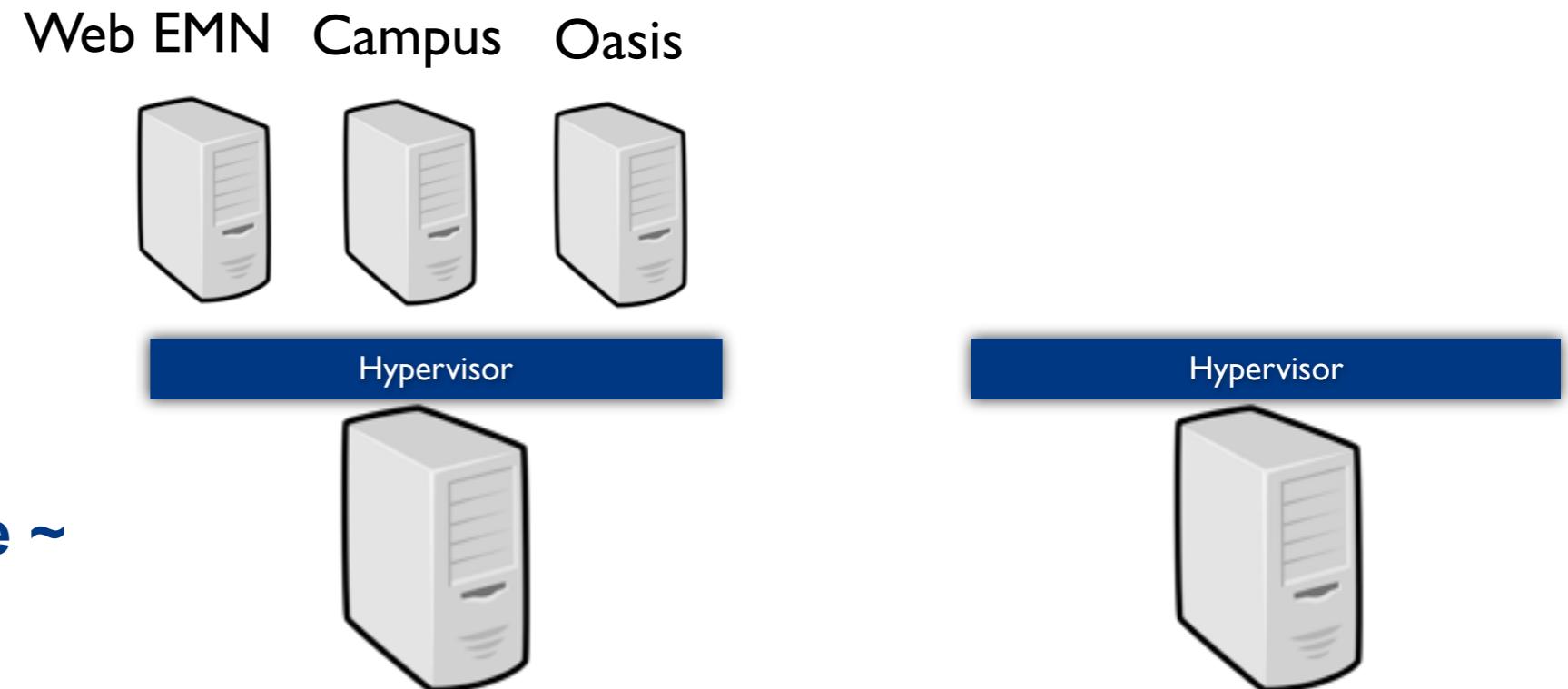


Virtualization capabilities



- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

- **Live migration (load-balancing)**
- **High Availability(downtime ~ 60 ms)**



Virtualization capabilities

Web EMN

Campus

Oasis



Hypervisor

- **Isolation (security between VM)**
- **suspend/resume/reboot (maintenance)**

Virus / Invasion / Crash



Web EMN Campus Oasis



Hypervisor



- **Live migration (load-balancing)**
- **High Availability(downtime ~ 60 ms)**

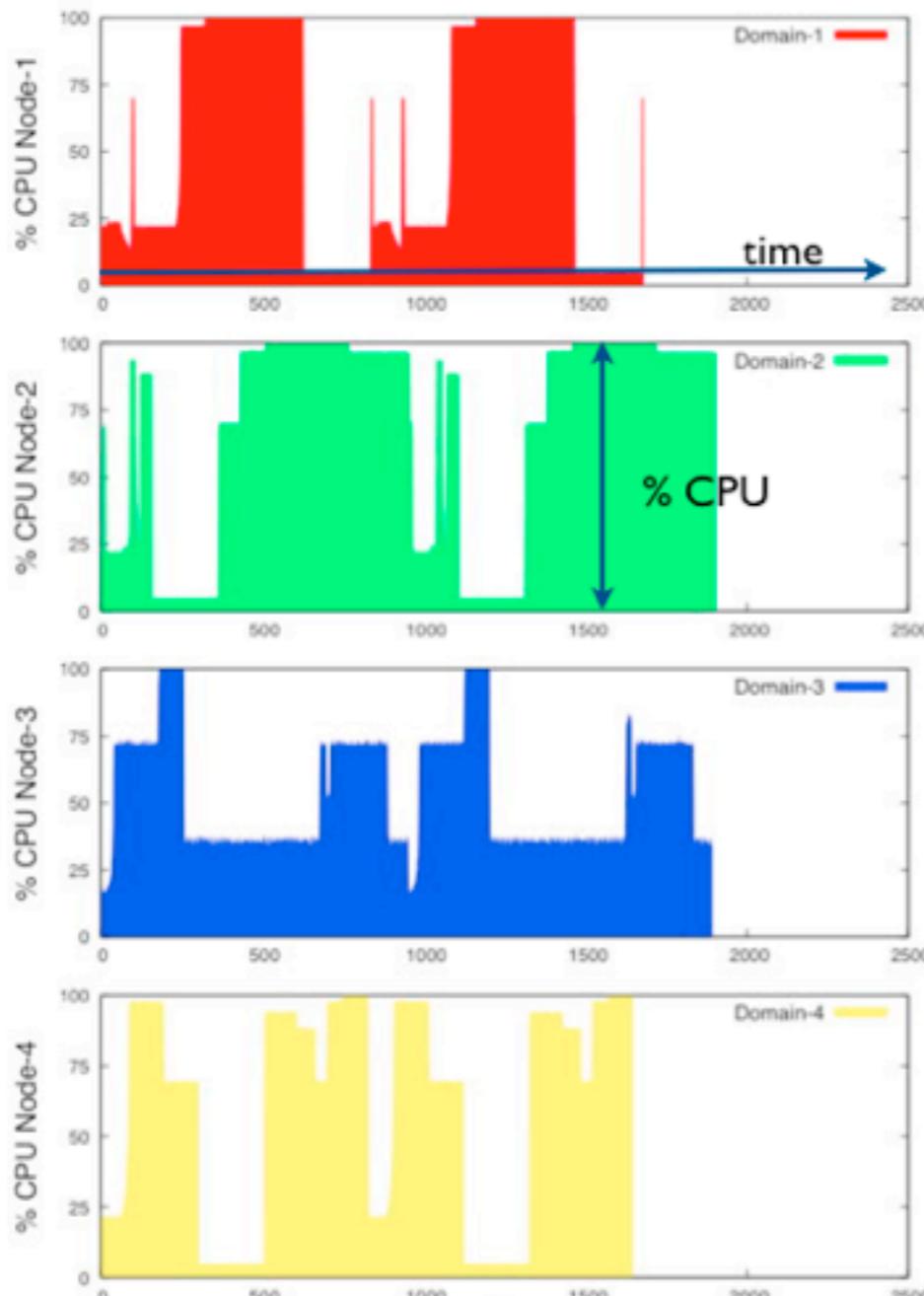


2>Virtual machine manager for clusters

1
2
3
4
5
6
7

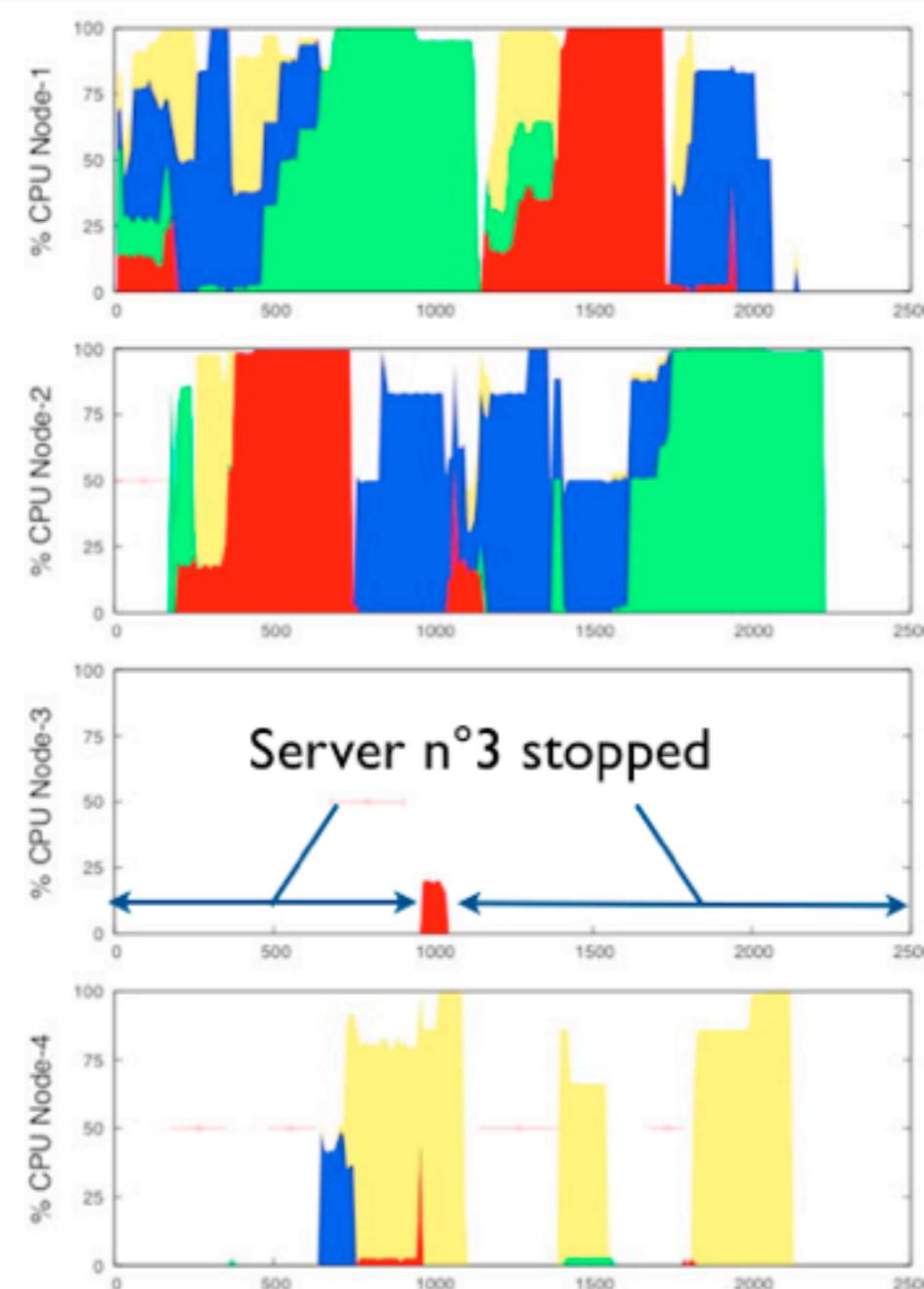
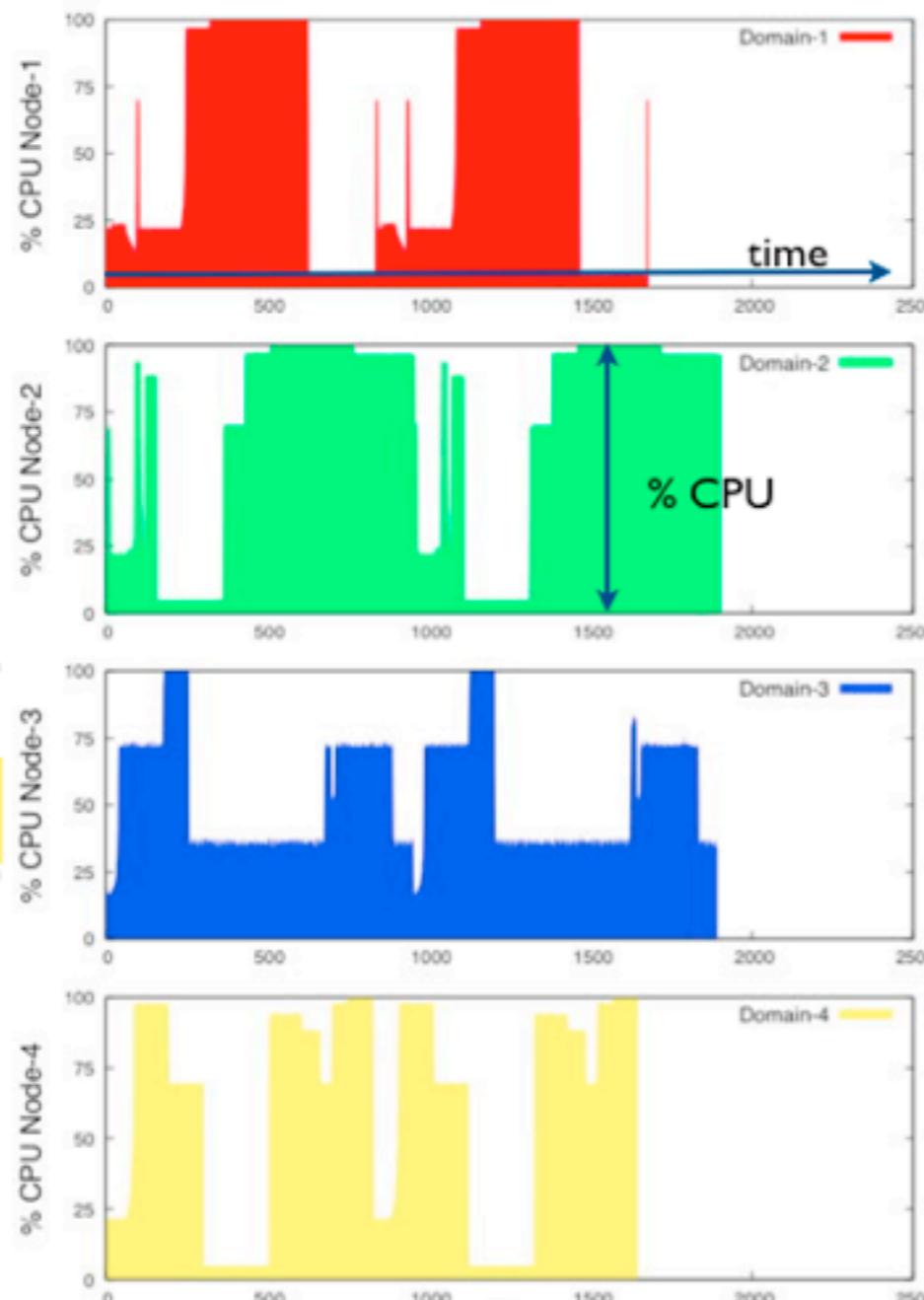
btrPlace: Principles

4 Tasks () 4 servers



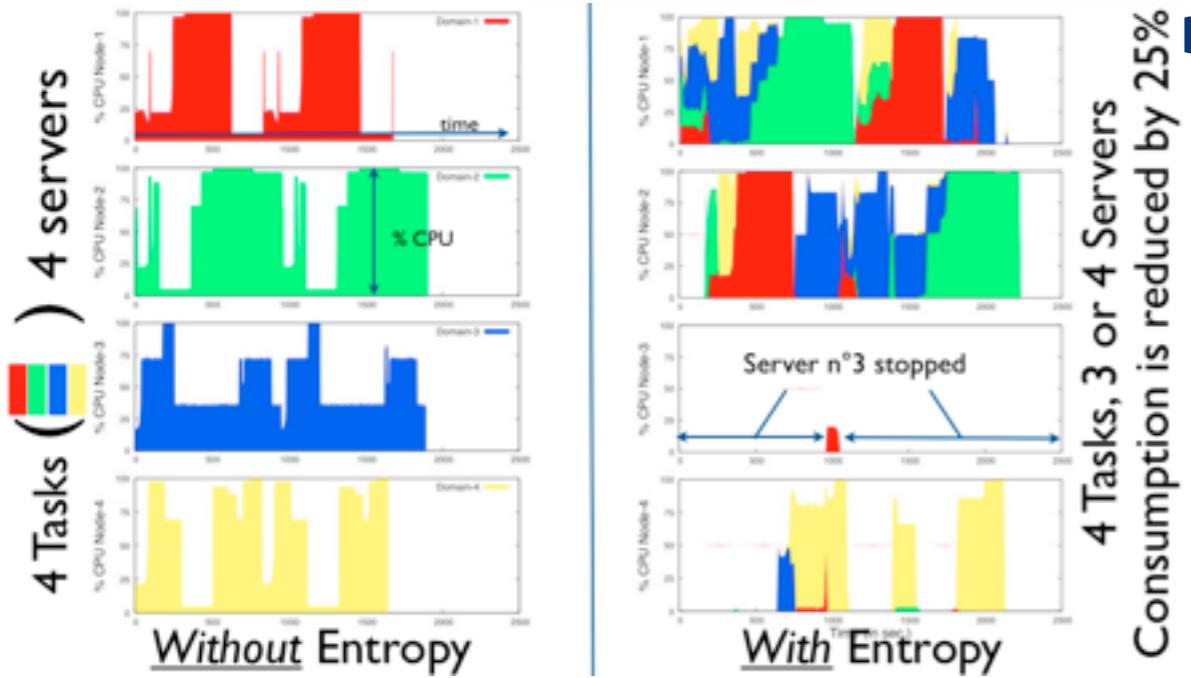
btrPlace: Optimizing the placement of virtual servers

4 Tasks () 4 servers

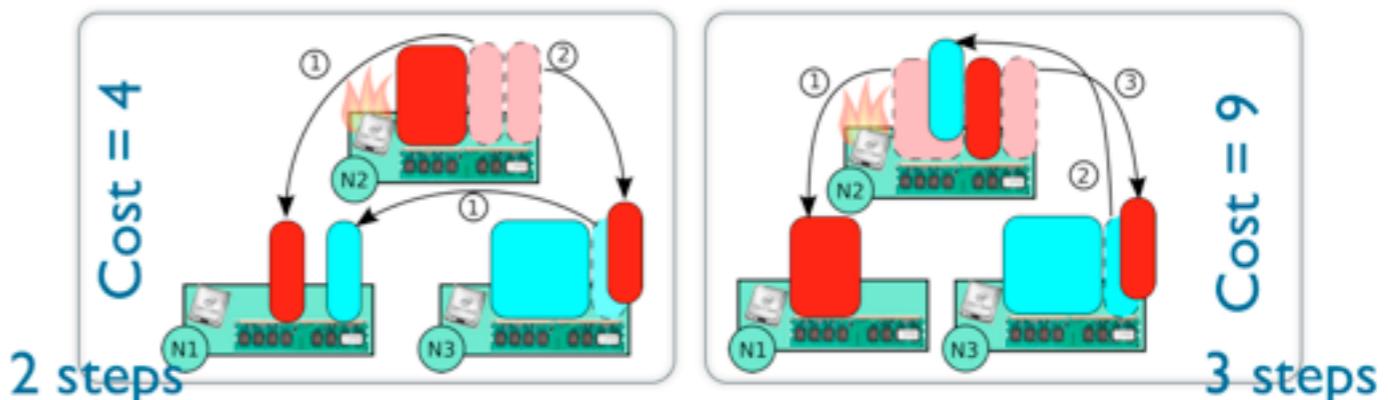


**4 Tasks, 3 or 4 Servers
Consumption is reduced by 25%**

btrPlace: Optimizing the placement of virtual servers



- Determine an efficient reconfiguration plan (thanks to a cost function)



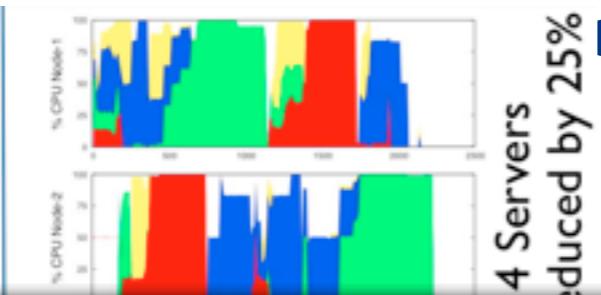
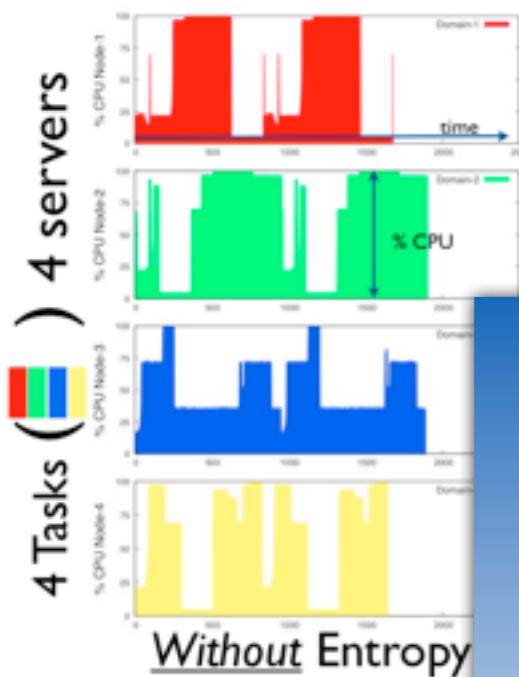
- **Administration and Application placement constraints must be considered**

- **ban({VM1, VM2}, {N1, N2})**
 - prevents a set of VMs from being hosted on a given set of nodes
 - **fence({VM1, VM2}, {N1, N2})**
 - forces a set of VMs to be hosted on a set of nodes
 - **spread({VM1, VM2})**
 - ensures that the specified VMs are never hosted on the same node at the same time
 - **latency({VM1, VM2}, {{N1, N2}, {N3, N4}})**
 - forces a set of VMs to be hosted on a single group of nodes.

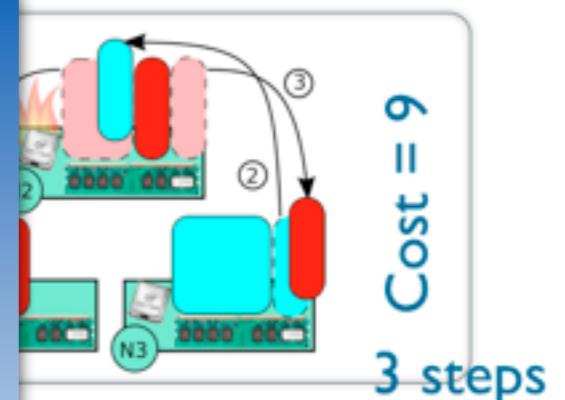
[VEE'09,CFSE'11]

J.M. Menaud,- June 2012 - Ascola 10

btrPlace: Optimizing the placement of virtual servers



- Determine an efficient reconfiguration plan (thanks to a cost function)

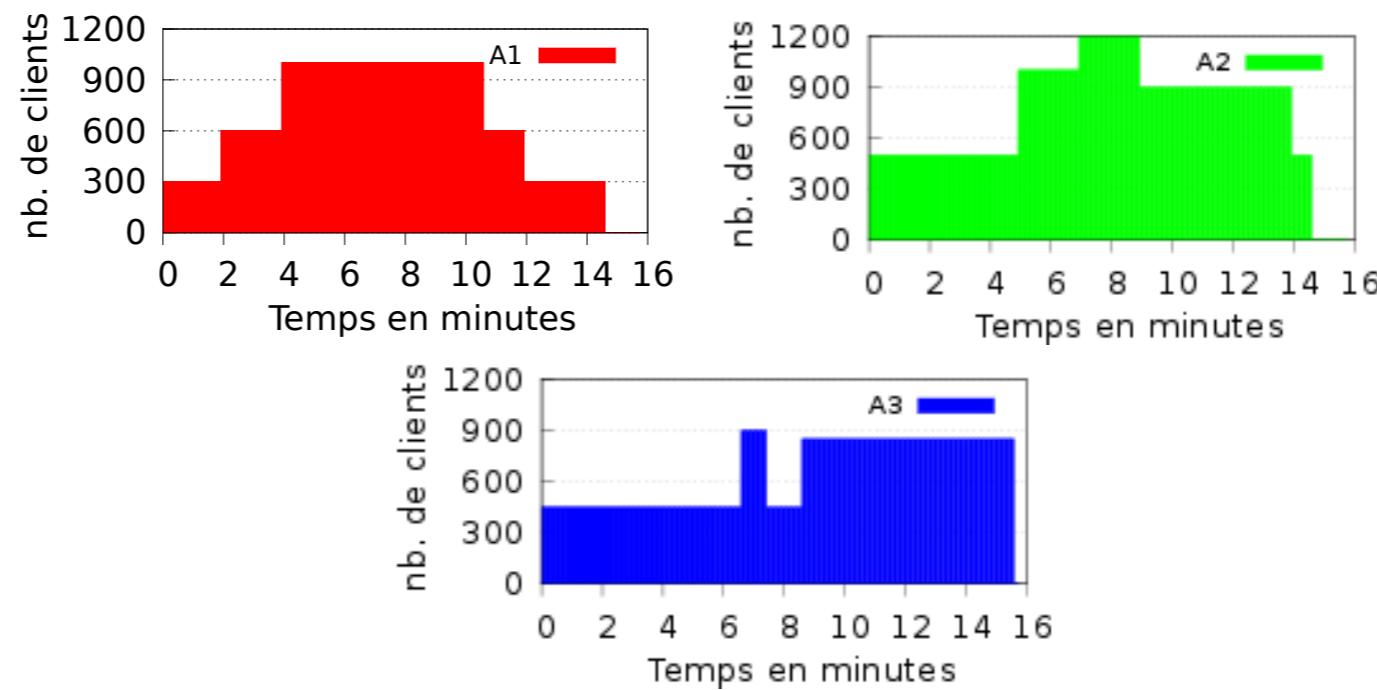


- Administer Application constraints must be considered

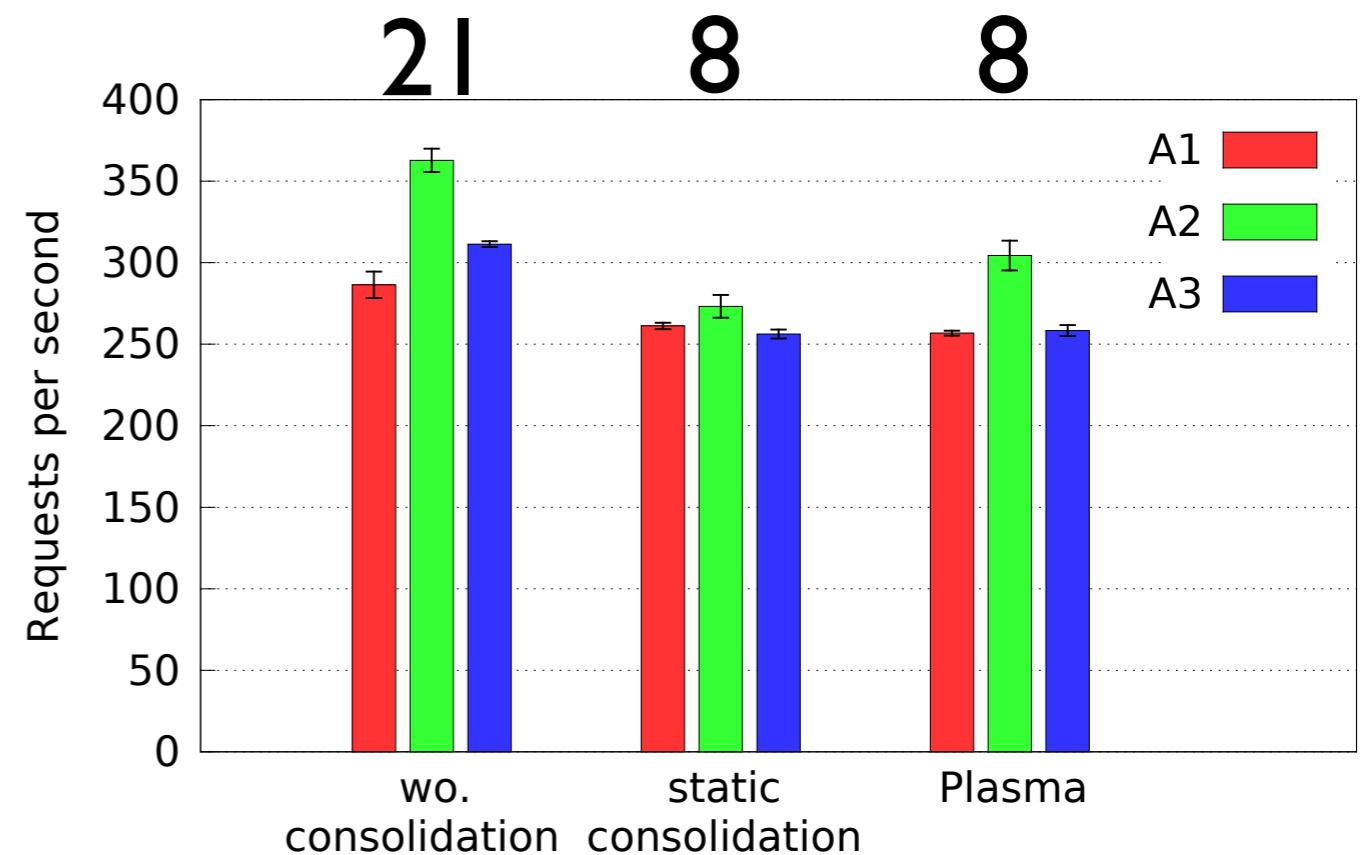
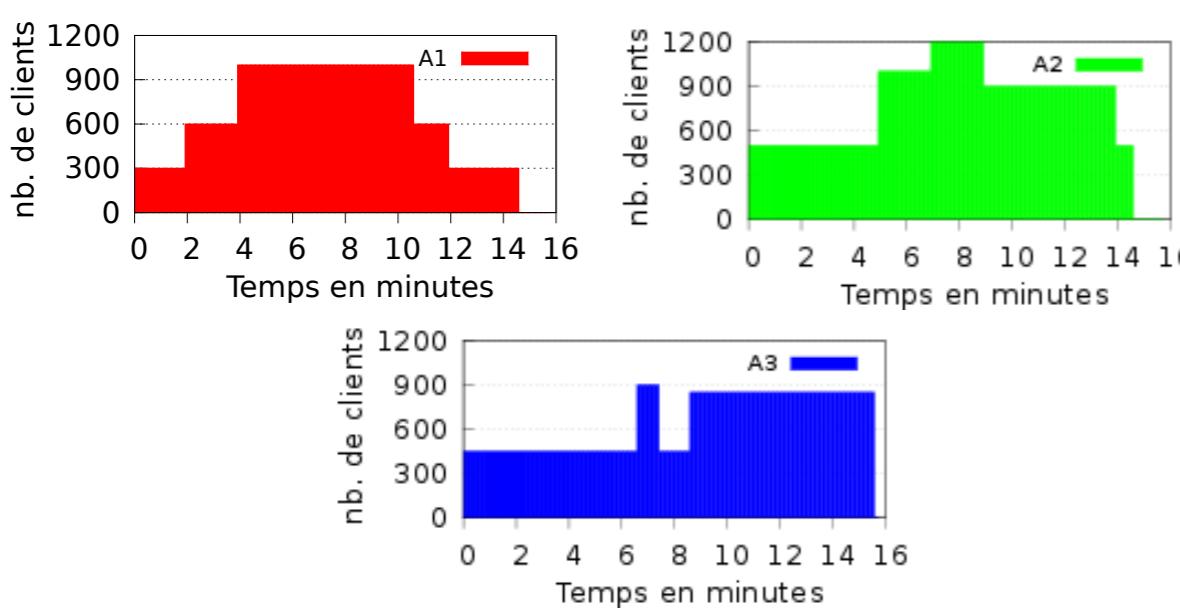
- `spread({VM1, VM2})`
 - ensures that the specified VMs are never hosted on the same node at the same time
- `latency({VM1, VM2}, {{N1, N2}, {N3, N4}})`
 - forces a set of VMs to be hosted on a single group of nodes.

Evaluations

- **RUBiS : The three tiers of each instance of RUBiS are deployed as 7 VMs (2,3,2)**
- **3 applications**
- **21 nodes**



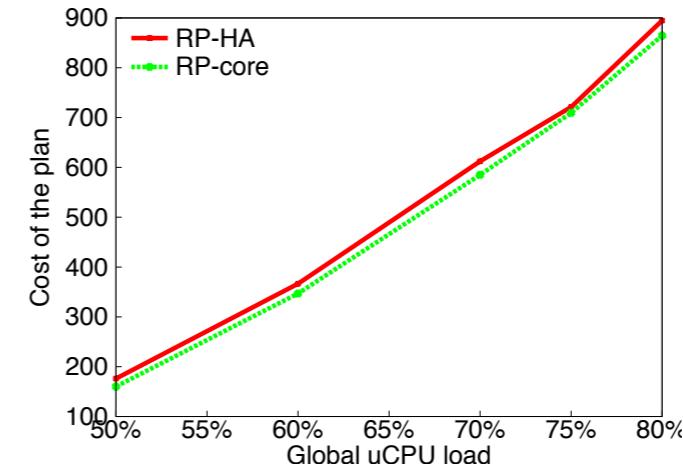
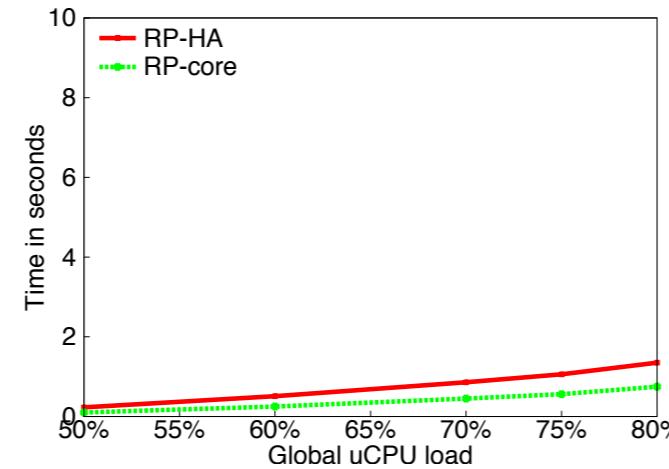
RUBiS Benchmark : Load spikes



- **Improvement wrt. static consolidation (14.7% vs. 17.7%)**
- **About 12 reconfigurations (29 secs) per execution**
- **Longest reconfiguration: 10 migrations in 89 seconds**

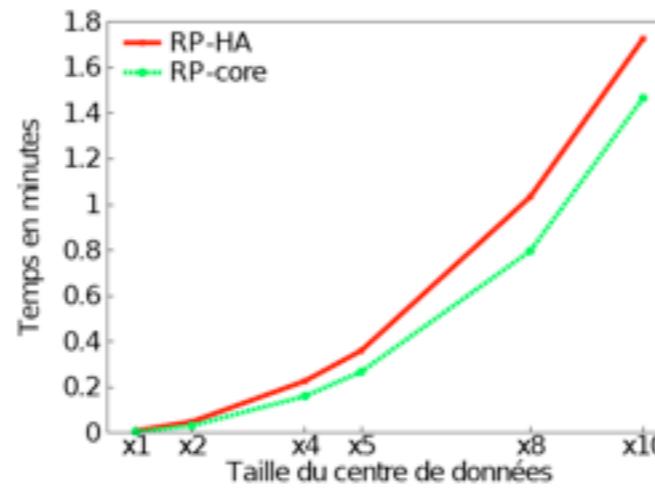
Impact of the global uCPU demand

- Impact of placement constraints is not significant

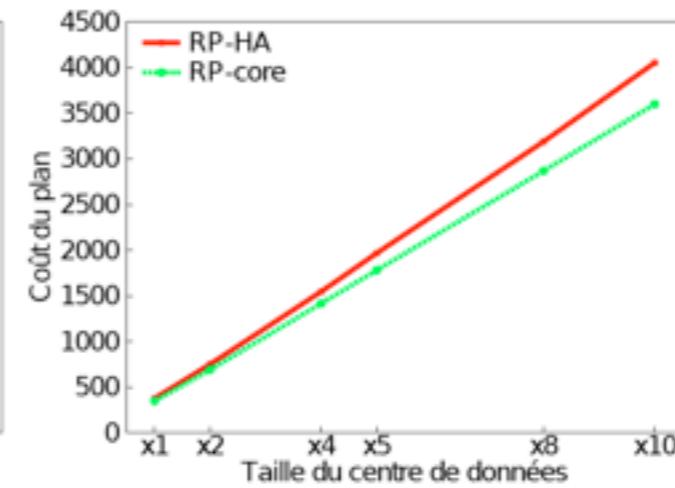


- In practice

- place 1117 candidates VMs on 1980 nodes with 600 spread + 200 latency
- schedule 475 actions



(a) Temps de résolution



(b) Coût des solutions

VM Management : next challenges

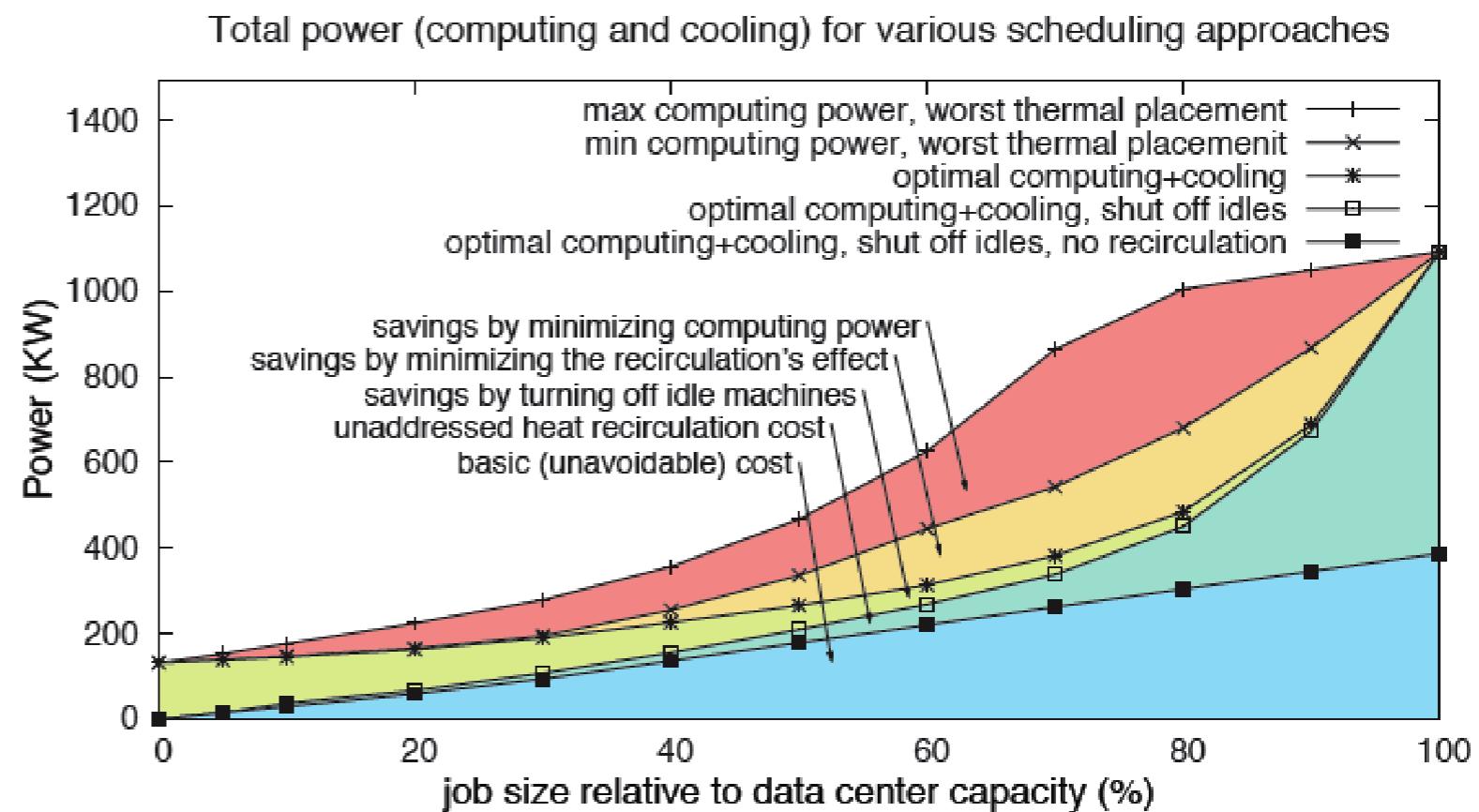
Constraints : hard/soft, flexibility, language ...

Placement : more abstraction, more optimization ...

Energy :

Thermal LoadBalancing

Sustainable Cloud



Distributed model : scalability, reactivity, fault tolerance ...

[CIT'09, DAIS'10, Cloud'10, ICAS'12]

J.M. Menaud, - June 2012 - Ascola 15

VM Management : next challenges

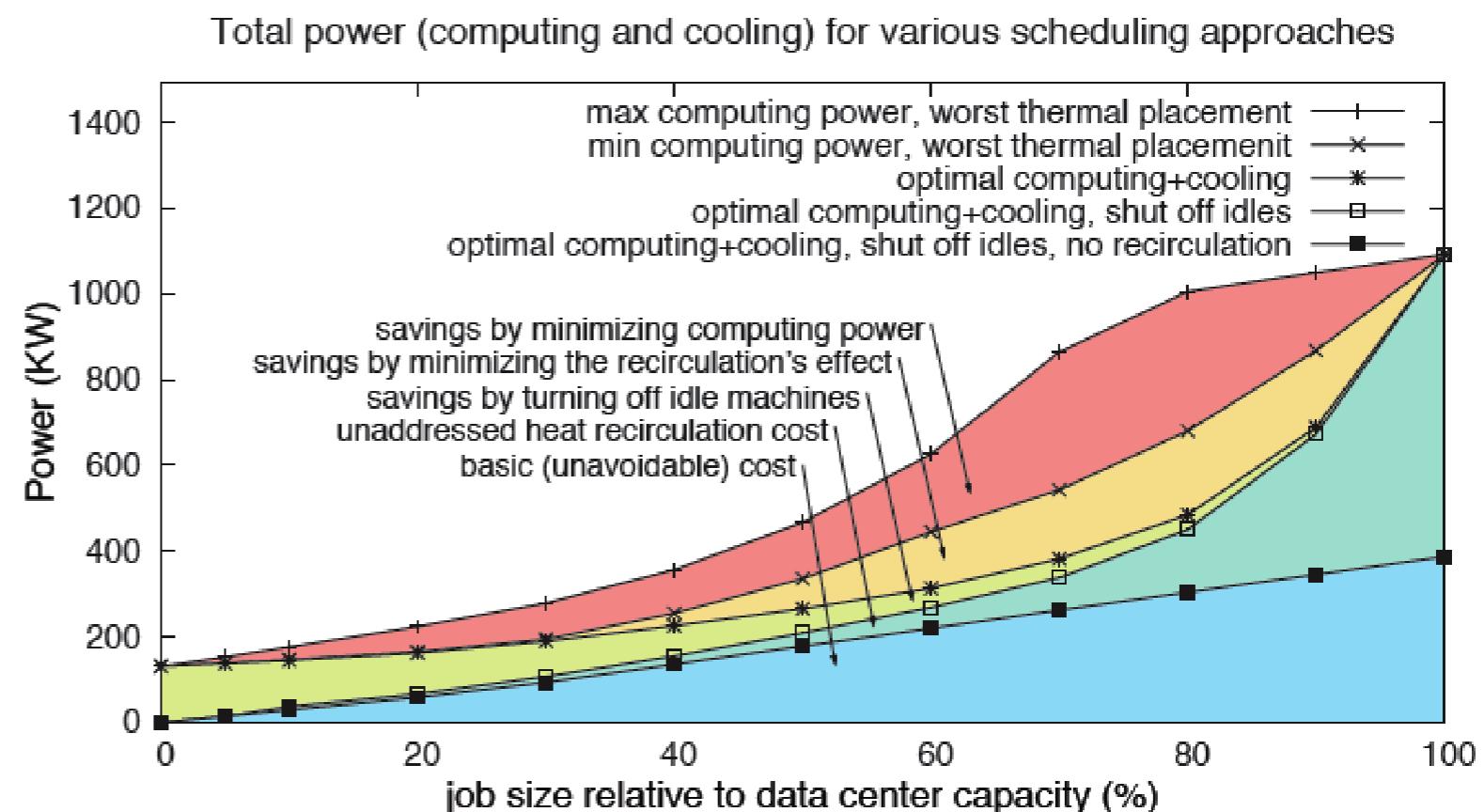
Constraints : hard/soft, flexibility, language ...

Placement : more abstraction, more optimization ...

Energy :

Thermal LoadBalancing

Sustainable Cloud

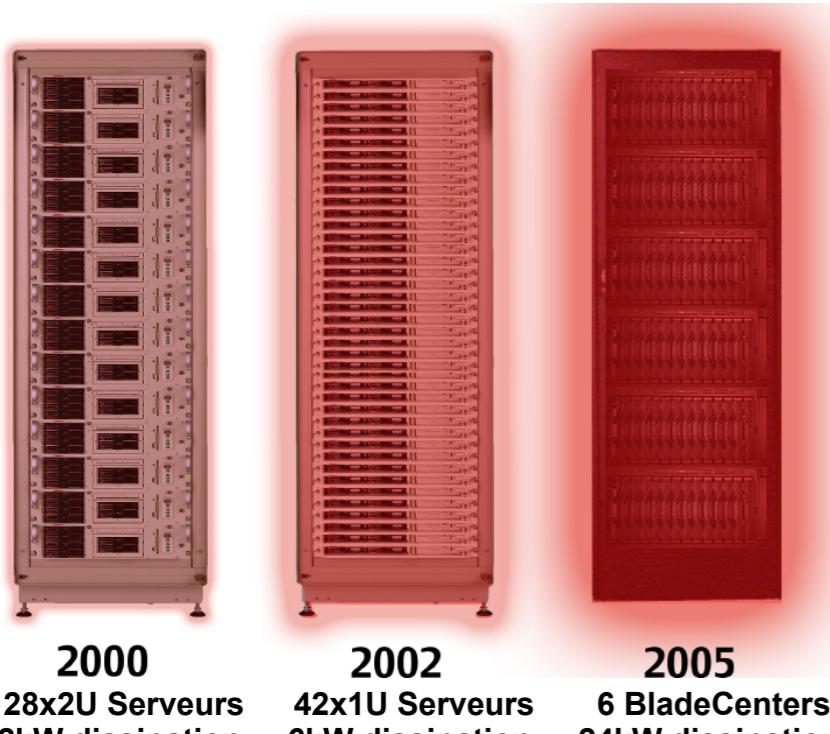


Distributed model : scalability, reactivity, fault tolerance ...

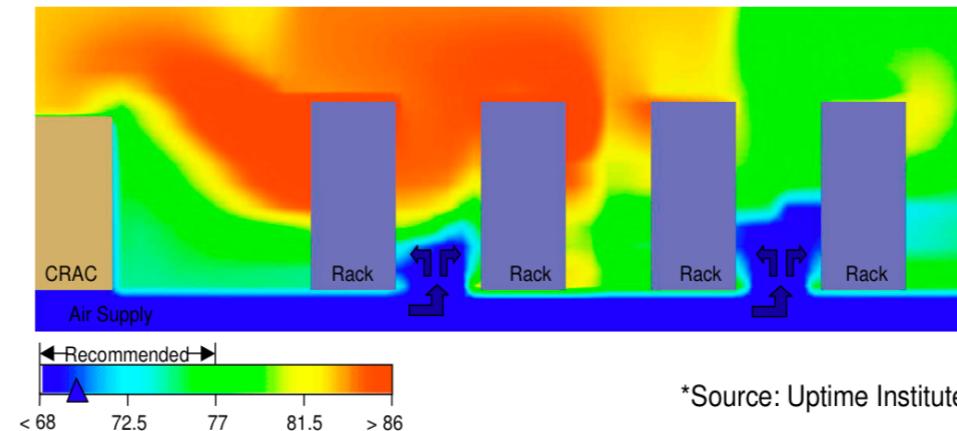
[CIT'09, DAIS'10, Cloud'10, ICAS'12]

J.M. Menaud, - June 2012 - Ascola 15

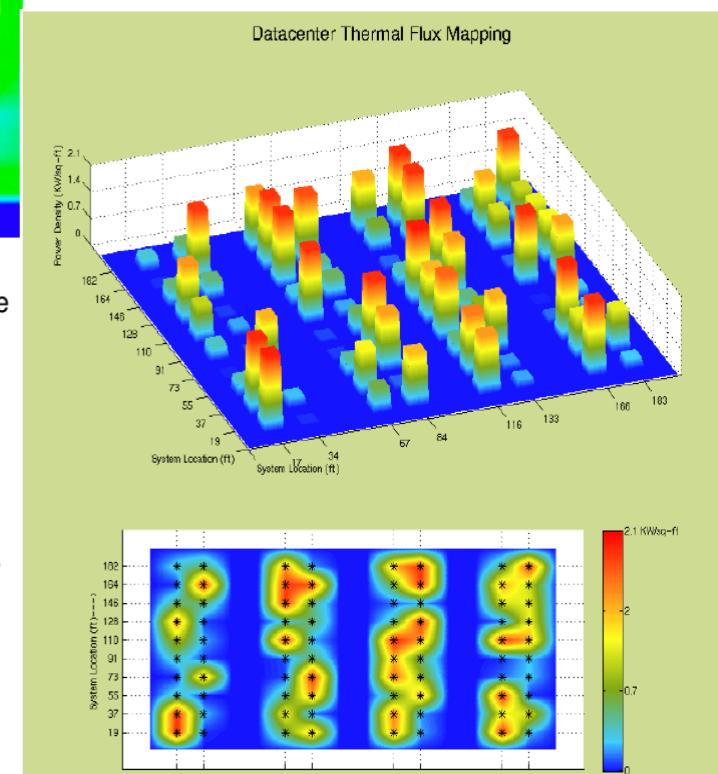
Des solutions ? L'équilibrage de charge thermique



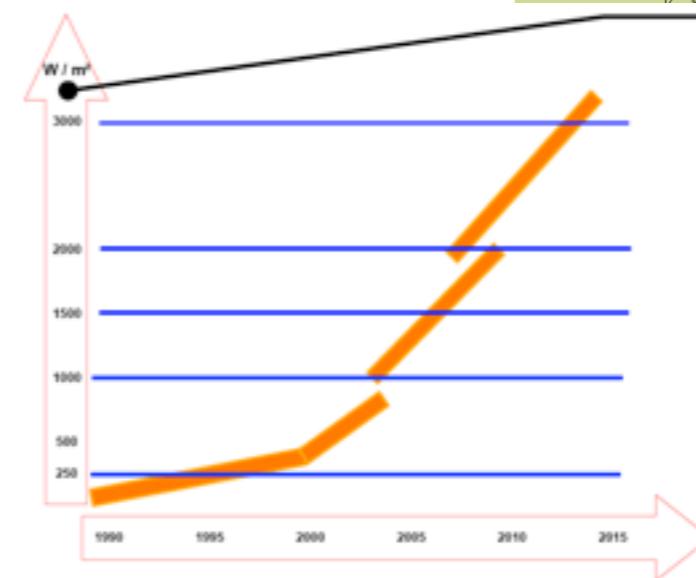
Source: Emerson Network Power/Liebert



*Source: Uptime Institute



Le refroidissement ne suit pas l'augmentation de puissance des équipements



Une densité de 3000 W / m² peut correspondre à une salle équipée en totalité de baies de 10 kw, ou à un mix entre des baies de 4 kw, 10 kw jusqu'à 30 kw

Quand on parle de W / m², il s'agit de la puissance en courant HQ disponible dans la salle, ramenée à la surface de la salle.

Des solutions ? L'équilibrage de charge thermique

Confinement et monter en T°



2000
28x2U Serveurs
2kW dissipation
thermique

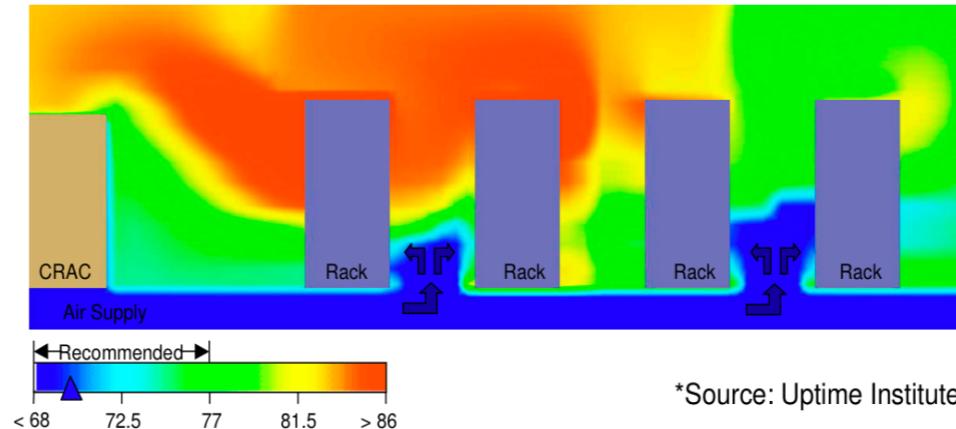
2002
42x1U Serveurs
6kW dissipation
thermique

2005
6 BladeCenters
24kW dissipation
thermique

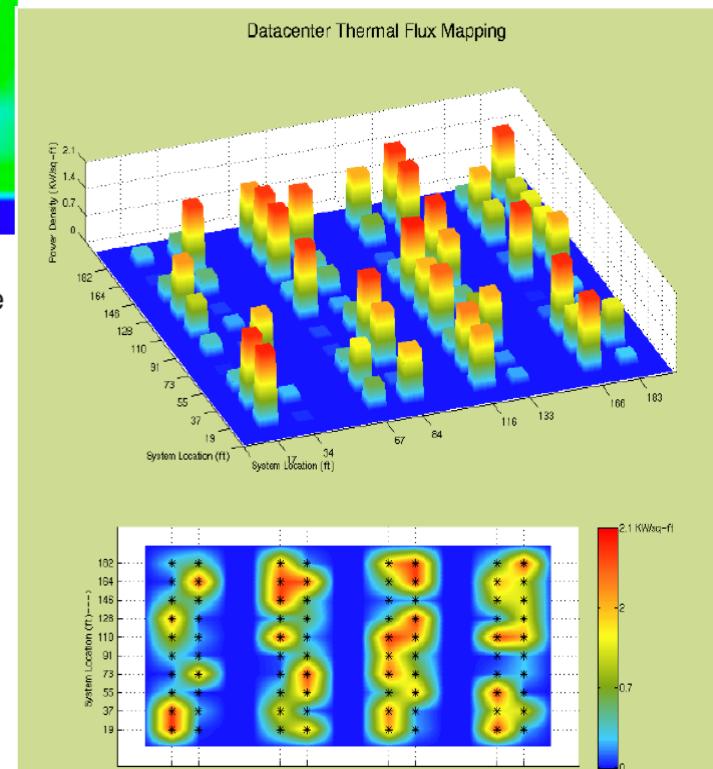
Source: Emerson Network Power/Liebert



2008
6 BladeCenters
30kW dissipation
thermique

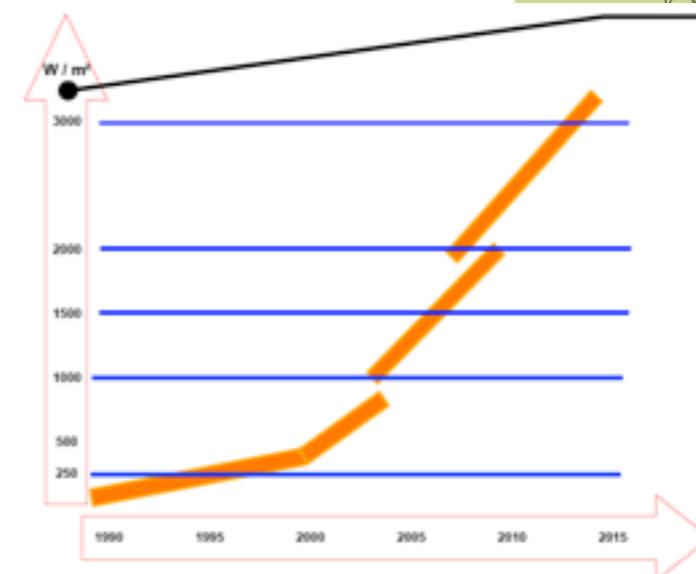


*Source: Uptime Institute



Quand on parle de W / m², il s'agit de la puissance en courant HQ disponible dans la salle, ramenée à la surface de la salle.

Sun



Une densité de 3000 W / m² peut correspondre à une salle équipée en totalité de baies de 10 kw, ou à un mix entre des baies de 4 kw, 10 kw jusqu'à 30 kw

Thermal Load Balancing : Heat and CRAC

- **Each server specifies a maximum temperature of inlet air T_{max} , at a constant rate applied to CARC.**
 $T_{max}(j)$ max inlet temp for server j
- **Heat linear model.**
 - The server i will heat the server j $I(i, j)$ °C by W consumed.
 C_i = power consumption for server i
1) $T_{imp}(j) = \sum_{i} (server\ i) C_i * I(i, j)$
 - The inlet temperature for server j is
2) $T_{in}(j) = T_{crac} + T_{imp}(j) < T_{max}(j)$
 - For the server j the CRAC must be
3) $T_{crac} < T_{max}(j) - \sum_{i} C_i * I(i, j)$
 - For multiple servers, the CRAC must be
4) $T_{crac} = \min((T_{max}(j) - \sum_{i} C_i * I(i, j)) \text{ for all } j)$
 - we need to maximise T_{crac}
- **First results ...**

- **Partners :**
 - **Bull SAS**, Splitted-Desktop Systems (SDS), INRIA, CEA, AVOB, ATRIUM DATA, SINOVIA, WILLELEC, EURODECISION
- **Main Focus**
 - Cooling system, transfer and transport of heat dissipated
Water cooling inside the server
 - Electrical system, integration of uninterruptible power supply (UPS)
By using supra-capacity
 - Collection of energy data
system, power etc.
 - Control of power management
Controlling server (DVFS) and task placement
 - Test and evaluation
- **Control of power management**
 - EuroDecision : Task placement with processor frequency variability

Conclusion

btrCloud

- btrPlace : Configurable consolidation manager
- PPC approach
 - scalable to datacenter with up to 2000 nodes/ 4000 VMs
 - placement constraints do not impact the solving process
- Constraints and optimize
 - Server and CRAC

Futures works

- new placement constraints for new concerns (currently 10 constraints)
- improvement of the scalability using partitioning
- Hard/soft placement constraints.



Questions ?



Jean-Marc.Menaud@mines-nantes.fr