# Grid'5000
## a scientific instrument for experiment-driven research on parallel, large-scale and distributed systems

Lucas Nussbaum

`lucas.nussbaum@loria.fr`
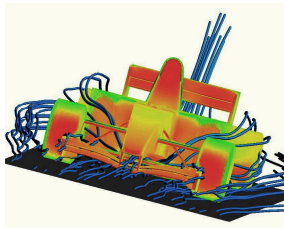
*Grid'5000 executive committee member*
*in charge of following the technical team*

# Experimentation for distributed systems
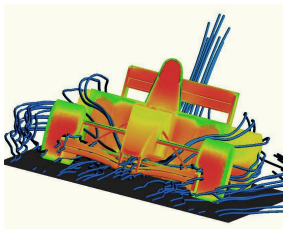
# Experimentation for distributed systems

### Simulation



1. Model application
2. Model environment
3. **Compute** interactions

# Experimentation for distributed systems

## Simulation



1. Model application
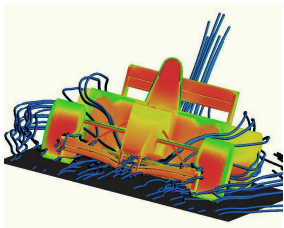2. Model environment
3. **Compute** interactions

## Real-scale experiments



Execute the **real** application on **real** machines

# Experimentation for distributed systems

## Simulation



**1** Model application
**2** Model environment
**3** **Compute** interactions

## Real-scale experiments



Execute the **real** application
on **real** machines

## Complementary solutions:

☺ Work on algorithms
☺ Scalable, more user-friendly

☺ Work on applications
☺ Closer to production use

# Grid'5000

- **Testbed for research on distributed systems**
    - High Performance Computing
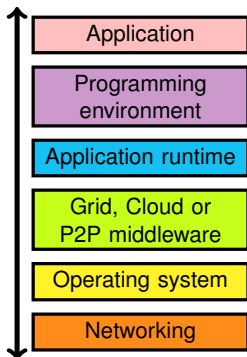    - Grids
    - Peer-to-peer systems
    - Cloud computing
- History:
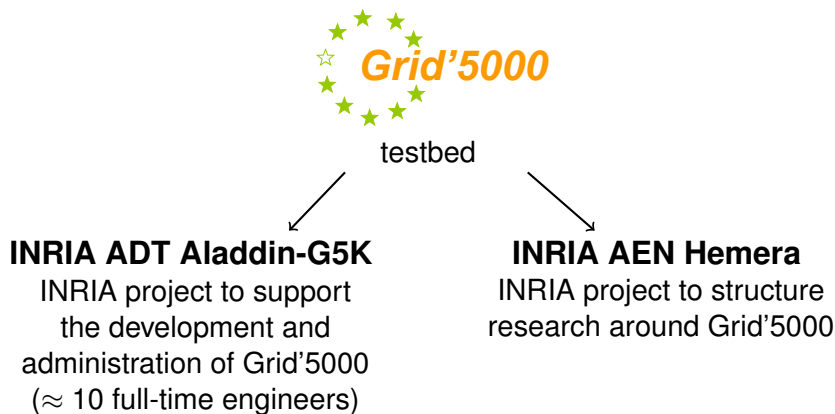    - 2003: Project started (ACI GRID)
    - 2005: Opened to users
- Funding: Inria, CNRS and many local entities (regions, universities)
- Only for research on distributed systems → no production usage
  Litmus test: *are you interested in the result of the computation?*
    - Free nodes during daytime to prepare experiments
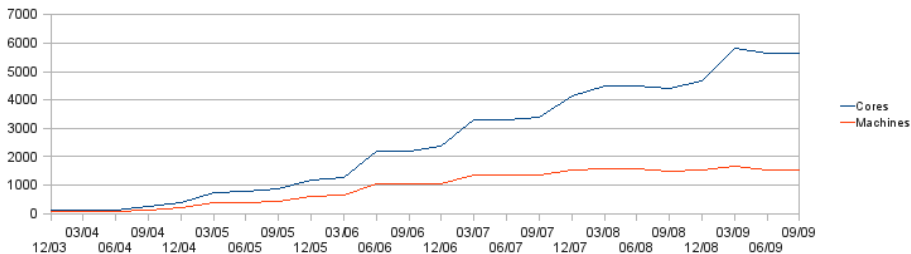    - Large-scale experiments during nights and week-ends

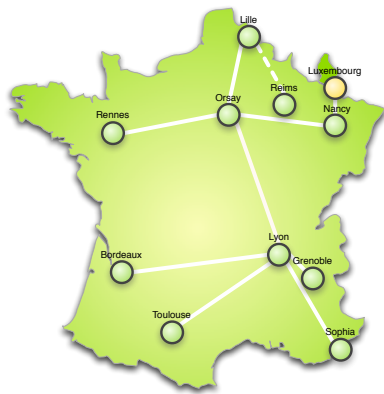| Application |
| Programming environment |
| Application runtime |
| Grid, Cloud or P2P middleware |
| Operating system |
| Networking |

# Organization



testbed

**INRIA ADT Aladdin-G5K**
INRIA project to support
the development and
administration of Grid'5000
($\approx$ 10 full-time engineers)

**INRIA AEN Hemera**
INRIA project to structure
research around Grid'5000

People:

► Scientific director: Frédéric Desprez
► Technical director: David Margery
► Hemera director: Christian Perez

# Current status
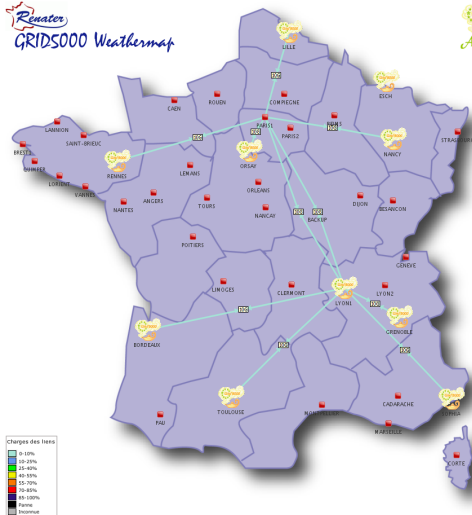
- 11 sites (1 outside France)
- 26 clusters
- 1700 nodes
- 7400 cores
- Diverse technologies:
  - Intel (60%), AMD (40%)
  - CPUs from one to 12 cores
  - Myrinet, Infiniband {S,D,Q}DR
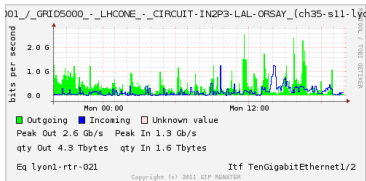  - Two GPU clusters
- **500+ users per year**

# Backbone network

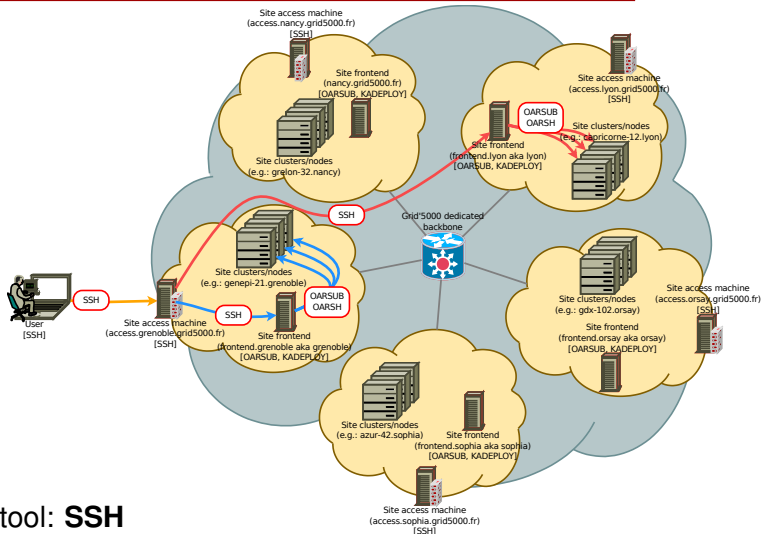Dedicated 10 Gbps backbone provided by RENATER (french NREN)



Work in progress:

- ▶ packet-level and flow-level monitoring
- ▶ bandwidth reservation and limitation

# Using Grid'5000: the user's point of view



- ▶ Key tool: **SSH**
- ▶ Private network: connect through access machines
- ▶ Data storage: **NFS** (one server per Grid'5000 site)

# Grid'5000 software stack

- ▶ Resource management: **OAR**

- ▶ System reconfiguration: **Kadeploy**

- ▶ Network isolation: **KaVLAN**

- ▶ Monitoring: **Ganglia**, **Kaspied**, **Energy**

- ▶ Putting it all together: **Grid'5000 API**

# Resource management: OAR



- ▶ Batch scheduler with specific features
    - ▶ interactive jobs
    - ▶ advance reservations
    - ▶ powerful resource matching
- ▶ Resources hierarchy: cluster / switch / node / cpu / core
- ▶ Properties: memory size, disk type & size, hardware capabilities, network interfaces, . . .
- ▶ Other kind of resources: VLANs, IP ranges for virtualization

*I want 1 core on 2 nodes of the same cluster with*
*4096 GB of memory and Infiniband 10G +*
*1 cpu on 2 nodes of the same switch with dualcore processors*
*for a walltime of 4 hours. . .*

```
oarsub -I -l "{memnode=4096 and
    ib10g='YES'}/cluster=1/nodes=2/core=1
+{cpucore=2}/switch=1/nodes=2/cpu=1,walltime=4:0:0"
```
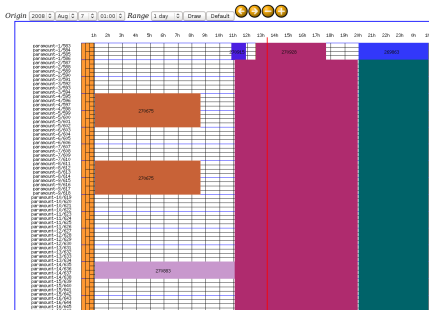
# Resource management: OAR - visualization



Resources status



Gantt chart

# Kadeploy – scalable cluster deployment tool

- ▶ Provides a *Hardware-as-a-Service* Cloud infrastructure
- ▶ Built on top of PXE, DHCP, TFTP
- ▶ **Scalable, efficient, reliable and flexible**:
  - ▶ Chain-based and BitTorrent environment broadcast
  - ▶ **255 nodes deployed in 7 minutes**
- ▶ Support of a **broad range of systems** (Linux, Xen, *BSD, etc.)
- ▶ Command-line interface & asynchronous interface (REST API)



**http://kadeploy3.gforge.inria.fr/**

# Network isolation: KaVLAN

- ▶ Reconfigures switches for the duration of a user experiment to achieve **complete level 2 isolation**:
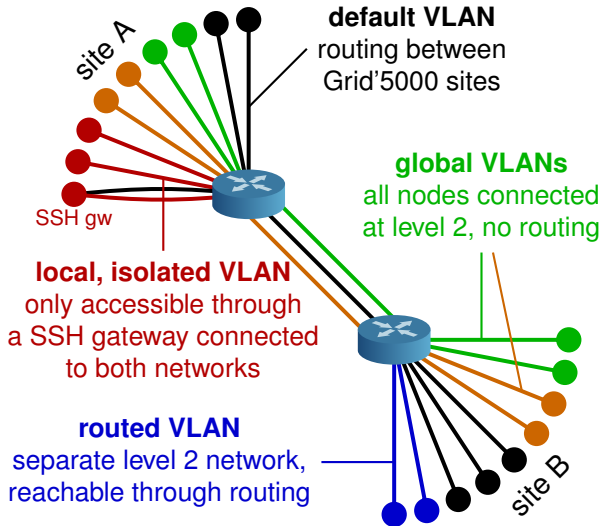    - ▶ Avoid network pollution (broadcast, unsolicited connections)
    - ▶ Enable users to start their own DHCP servers
    - ▶ Experiment on ethernet-based protocols
    - ▶ Interconnect nodes with another testbed without compromising the security of Grid'5000
- ▶ Relies on **802.1q (VLANs)**
- ▶ Compatible with many network equipments
    - ▶ Can use SNMP, SSH or telnet to connect to switches
    - ▶ Supports Cisco, HP, 3Com, Extreme Networks and Brocade
- ▶ Controlled with a command-line client or a REST API

# KaVLAN - different VLAN types



**default VLAN**
routing between
Grid'5000 sites

site A

SSH gw

**global VLANs**
all nodes connected
at level 2, no routing

**local, isolated VLAN**
only accessible through
a SSH gateway connected
to both networks

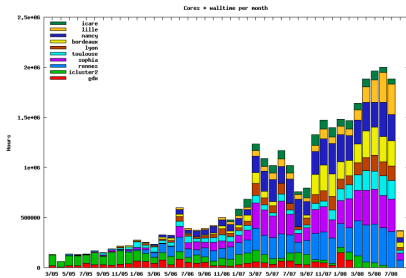**routed VLAN**
separate level 2 network,
reachable through routing

site B

# Monitoring: Ganglia, Kaspied, Energy



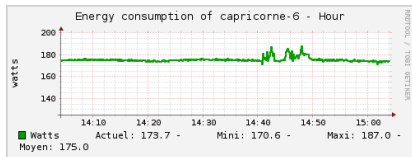Ganglia



Kaspied
(Grid'5000 usage over time)



Power consumption

# Putting it all together: Grid'5000 API

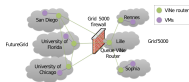- Individual services & command-line interfaces are painful
- REST API for each Grid'5000 service:
  - **Reference API**: versioned description of Grid'5000 resources
  - **Monitoring API**: state of Grid'5000 resources
  - **Metrology API**: Ganglia data
  - **Jobs API**: OAR interface
  - **Deployments API**: Kadeploy interface
  - . . .
- Also some nice Web interfaces on `https://api.grid5000.fr/`

# Leading to results in several fields

Cloud: Sky computing on FutureGrid and Grid'5000

- ► Nimbus cloud deployed on 450+ nodes
- ► Grid'5000 and FutureGrid connected using ViNe



HPC: factorization of RSA-768

- ► Feasibility study: prove that it can be done
- ► Different hardware $\leadsto$ understand the performance characteristics of the algorithms
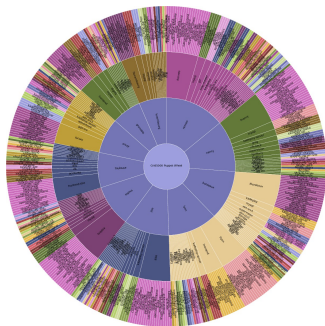


Grid: evaluation of the gLite grid middleware

- ► Fully automated deployment and configuration on 1000 nodes (9 sites, 17 clusters)

# Open challenges

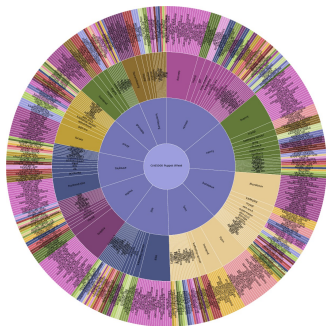Testbeds optimize for experimental capabilities, not performance

- **Access** to the modern architectures / technologies
    - Not necessarily the fastest CPUs
    - But still expensive $\rightsquigarrow$ funding!

- Ability to **trust** results
    - Regular checks of testbed for bugs

- Ability to **understand** results
    - Documentation of the infrastructure
    - Instrumentation & monitoring tools
      *network, energy consumption*
    - Evolution of the testbed
      *maintenance logs, configuration history*

- **Empower** users to perform complex experiments
    - Facilitate access to advanced software tools

# Open challenges

Testbeds optimize for experimental capabilities, not performance

- **Access** to the modern architectures / technologies
    - Not necessarily the fastest CPUs
    - But still expensive $\rightsquigarrow$ funding!

- Ability to **trust** results
    - Regular checks of testbed for bugs

- Ability to **understand** results
    - Documentation of the infrastructure
    - Instrumentation & monitoring tools
      *network, energy consumption*
    - Evolution of the testbed
      *maintenance logs, configuration history*

- **Empower** users to perform complex experiments
    - Facilitate access to advanced software tools ← this afternoon

# Conclusions

- Grid'5000: a testbed for experimentation on distributed systems
- With a unique combination of features
  - *Hardware-as-a-Service* cloud: redeployment of operating system on the bare hardware by users
  - Access to various technologies (CPUs, high performance networks, etc.)
  - Networking: dedicated backbone, monitoring, isolation
  - Programmable through an API

**Interested in trying it? Contact us!**

https://www.grid5000.fr/

lucas.nussbaum@loria.fr