

Volet data-centers de SILECS (A.K.A. Grid'5000)

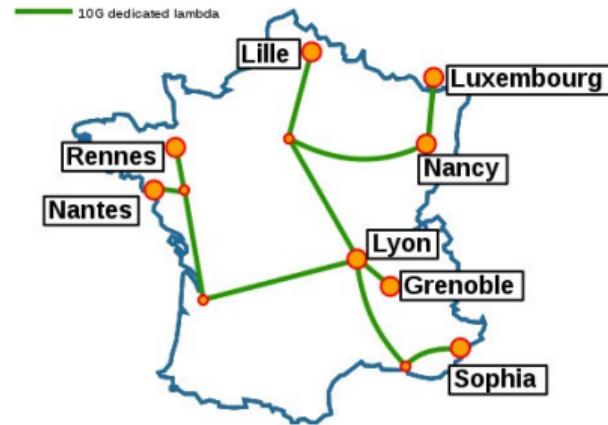
Présentation et exemples d'expériences

Frédéric Desprez & Lucas Nussbaum
Grid'5000 Scientific & Technical Directors

Visite du comité TGIR du CNRS
2019-04-19

The Grid'5000 testbed

- ▶ A large-scale testbed for distributed computing
 - ◆ 8 sites, 31 clusters, 828 nodes, 12328 cores
 - ◆ Dedicated 10-Gbps backbone network
 - ◆ 550 users and 120 publications per year



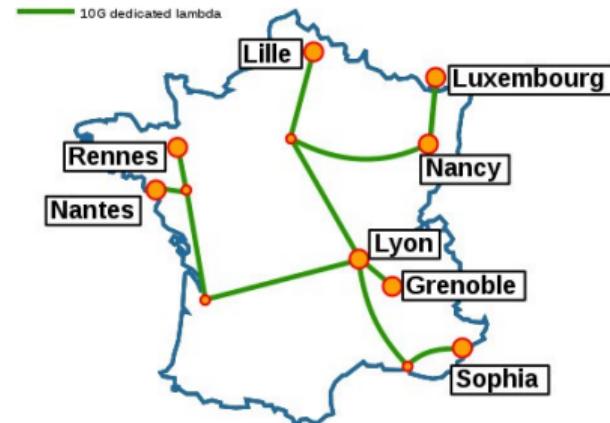
The Grid'5000 testbed

- ▶ A large-scale testbed for distributed computing

- ◆ 8 sites, 31 clusters, 828 nodes, 12328 cores
- ◆ Dedicated 10-Gbps backbone network
- ◆ 550 users and 120 publications per year

- ▶ A meta-cloud, meta-cluster, meta-data-center

- ◆ Used by CS researchers in HPC, Clouds, Big Data, Networking, AI
- ◆ To experiment in a fully controllable and observable environment
- ◆ Similar problem space as Chameleon and Cloudlab (US)
- ◆ Design goals
 - ★ Support high-quality, reproducible experiments
 - ★ On a large-scale, distributed, shared infrastructure



Landscape – cloud & experimentation¹

- ▶ **Public cloud infrastructures** (AWS, Azure, Google Cloud Platform, etc.)
 - ⌚ No information/guarantees on placement, multi-tenancy, real performance
- ▶ **Private clouds:** Shared observable infrastructures
 - ⌚ Monitoring & measurement
 - ⌚ No control over infrastructure settings
 - ↝ Ability to **understand** experiment results
- ▶ **Bare-metal as a service, fully reconfigurable infrastructure** (Grid'5000)
 - ⌚ Control/alter all layers (virtualization technology, OS, networking)
 - ↝ *In vitro* Cloud

And the same applies to all other environments (e.g. HPC)

¹Inspired from a slide by Kate Keahey (Argonne Nat. Lab.)

Some recent results from Grid'5000 users

- ▶ Portable Online Prediction of Network Utilization (Inria Bdx + US)
- ▶ Energy proportionality on hybrid architectures (LIP/IRISA/Inria)
- ▶ Maximally Informative Itemset Mining (Miki) (LIRM/Inria)
- ▶ Damaris (Inria)
- ▶ BeBida: Mixing HPC and BigData Workloads (LIG)
- ▶ HPC: In Situ Analytics (LIG/Inria)
- ▶ Addressing the HPC/Big-Data/IA Convergence
- ▶ An Orchestration Syst. for IoT Applications in Fog Environment (LIG/Inria)
- ▶ Toward a resource management system for Fog/Edge infrastructures
- ▶ Distributed Storage for Fog/Edge infrastructures (LINA)
- ▶ From Network Traffic Measurements to QoE for Internet Video (Inria)

Portable Online Prediction of Network Utilization

► Problem

- ◆ Predict network utilization in near future to enable optimal utilization of spare bandwidth for low-priority asynchronous jobs co-located with an HPC application

► Goals

- ◆ High accuracy, low compute overhead, learn on-the-fly without previous knowledge

► Proposed solution

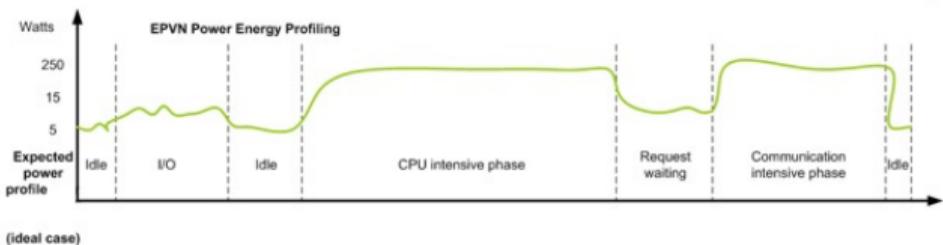
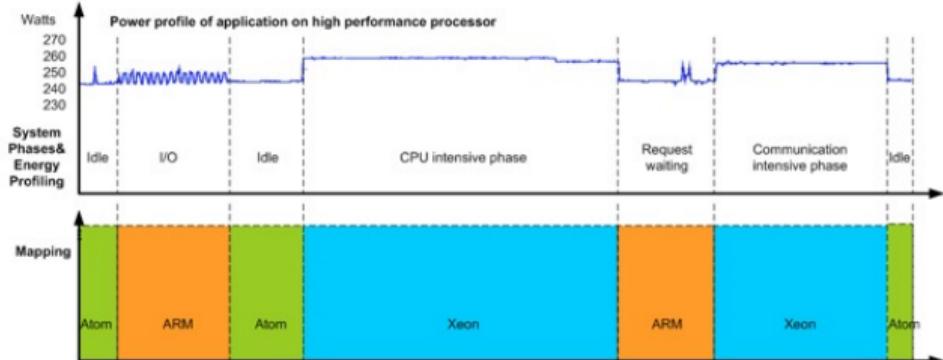
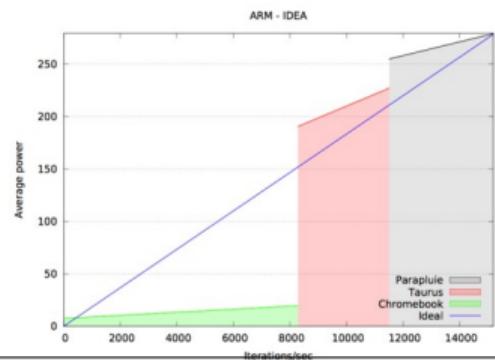
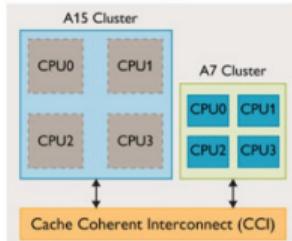
- ◆ Dynamic sequence-to-sequence recurrent neural networks that learn using a sliding window approach over recent history
- ◆ Evaluate the gain of a tree-based meta-data management
- ◆ INRIA, The Univ. of Tennessee, Exascale Comp. Proj., UC Irvine, Argonne Nat. Lab.

► Grid'5000 experiments

- ◆ Monitor and predict network utilization for two HPC applications at small scale (30 nodes)
- ◆ Easy customization of environment for rapid prototyping and validation of ideas (in particular, custom MPI version with monitoring support)
- ◆ Impact: Early results facilitated by Grid'5000 are promising and motivate larger scale experiments on leadership class machines (Theta@Argonne)

Energy proportionality on hybrid architectures²

- ▶ Hybrid computing architectures : low power processors, co processors, GPUs...
- ▶ Supporting a “Big, Medium, Little” approach : the right processor at the right time



²V. Villebonnet, G. Da Costa, L. Lefèvre, J.-M. Pierson and P. Stolf. "Big, Medium, Little" : Reaching Energy Proportionality with Heterogeneous Computing Scheduler", Parallel Processing Letters, 25 (3), Sep. 2015

Maximally Informative Itemset Mining (Miki)³

Extracting knowledge from data

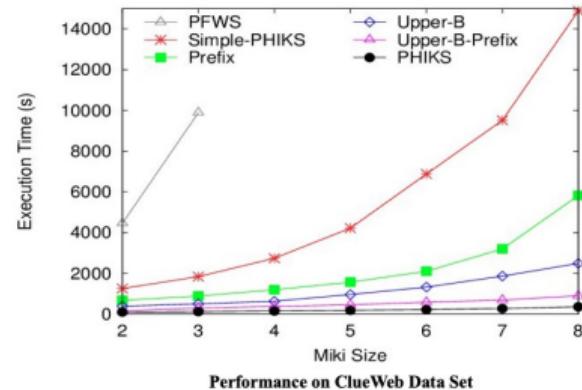
Miki: measures the quantity of information (e.g., based on joint entropy measure) delivered by the itemsets of size k in a database (i.e., k denotes the number of items in the itemset)

► PHIKS, a parallel algorithm for mining of maximally informative k-itemsets

- ◆ Very efficient for parallel miki discovery
- ◆ High scalability with very large amounts of data and high size of the itemsets
- ◆ Includes several optimization techniques
- ◆ Communication cost reduction using entropy bound filtering
- ◆ Incremental entropy computation
- ◆ Prefix/Suffix technique for reducing response time

► Experiments on Grid'5000

- ◆ Hadoop/Map Reduce on 16 and 48 nodes
- ◆ Datasets of 49 Gb (English Wikipedia, 5 millions articles),
1 Tb (ClueWeb, 632 millions articles)
- ◆ Metrics: Response time, communication cost, energy consumption



³S.Salah, R. Akbarinia, F. Masseglia. A Highly Scalable Parallel Algorithm for Maximally Informative k-Itemset Mining. Knowledge and Information Systems (KAIS), Springer, 2017, 50 (1)

Damaris

Scalable, asynchronous data storage for large-scale simulations using the HDF5 format

► Traditional approach

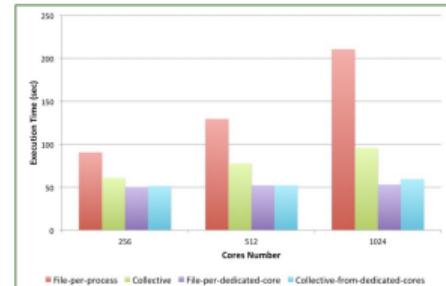
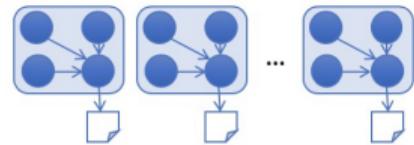
- ◆ All simulation processes (10K+) write on disk at the same time synchronously
- ◆ Problems: 1) I/O jitter, 2) long I/O phase, 3) Blocked simulation during data writing

► Solution

- ◆ Aggregate data in dedicated cores using shared memory and write asynchronously

► Grid'5000 used as a testbed

- ◆ Access to many (1024) homogeneous cores
- ◆ Customizable environment and tools
- ◆ Repeat the experiments later with the same environment saved as an image
- ◆ The results show that Damaris can provide a jitter-free and wait-free data storage mechanism
- ◆ G5K helped prepare Damaris for deployment on top supercomputers (Titan, Pangea (Total), Jaguar, Kraken, etc.)
- ◆ <https://project.inria.fr/damaris/>



BeBida: Mixing HPC and BigData Workloads

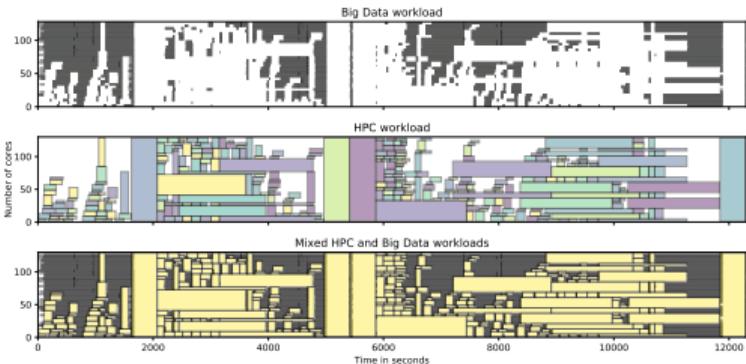
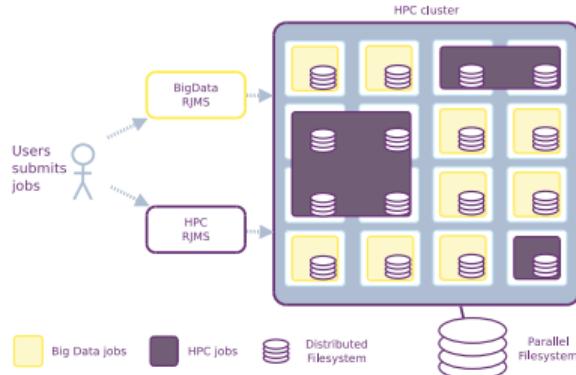
Objective: Use idle HPC resources for BigData workloads

► Simple approach

- ◆ HPC jobs have priority
- ◆ BigData Framework: Spark/Yarn, HDFS
- ◆ Evaluating costs of starting/stopping tasks (Spark/Yarn) and data transferts (HDFS)

► Results

- ◆ It increases cluster utilisation
- ◆ Disturbance of HPC jobs is small
- ◆ Big Data execution time varies (WIP)



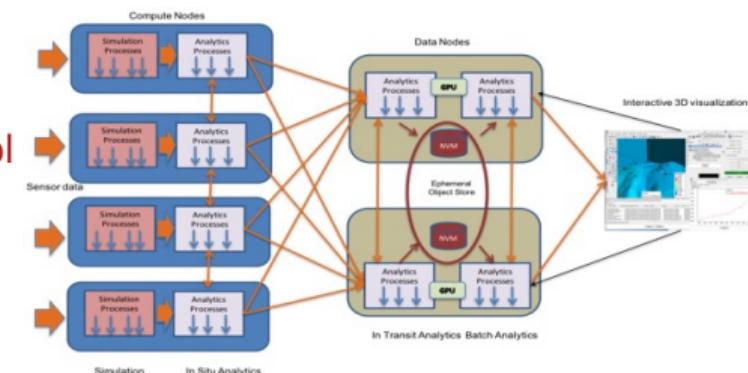
HPC: In Situ Analytics

Goal: improve organization of simulation and data analysis phases

- ▶ Simulate on a cluster; move data; post-mortem analysis
 - ◆ Unsuitable for Exascale (data volume, time)
- ▶ Solution: analyze on nodes, during simulation
 - ◆ Between or during simulation phases?
dedicated core? node?

Grid'5000 used for development and test, because control

- ▶ of the software environment (MPI stacks),
- ▶ of CPU performance settings (Hyperthreading),
- ▶ of networking settings (Infiniband QoS).



Then evaluation at a larger scale on the Froggy supercomputer (CIMENT center/GRICAD, Grenoble)

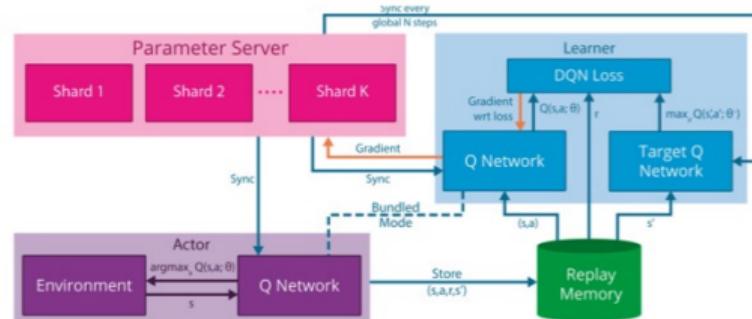
Addressing the HPC/Big-Data/IA Convergence⁴

Gathering teams from HPC, Big Data, and Machine Learning to work on the convergence of

- ▶ Smart Infrastructure and resource management
- ▶ HPC acceleration for AI and Big Data
- ▶ AI/Big Data analytics for large scale scientific simulations

Current work

- ▶ Molecular dynamics trajectory analysis with deep learning
 - ◆ Dimension reduction through DL, accelerating MD simulation coupling HPC simulation and DL
- ▶ Flink/Spark stream processing for in-transit on-line analysis of parallel simulation outputs
- ▶ Shallow Learning
 - ◆ Accelerating Scikit-Learn with task-based programming
(Dask, StarPU)
- ▶ Deep Learning
 - ◆ TensorFlow graph scheduling for efficient parallel executions
 - ◆ Linear algebra and tensors for large scale machine learning
 - ◆ Large scale parallel deep reinforcement learning



⁴<https://project.inria.fr/hpcbigdata/>

An Orchestration Syst. for IoT Applications in Fog Environment

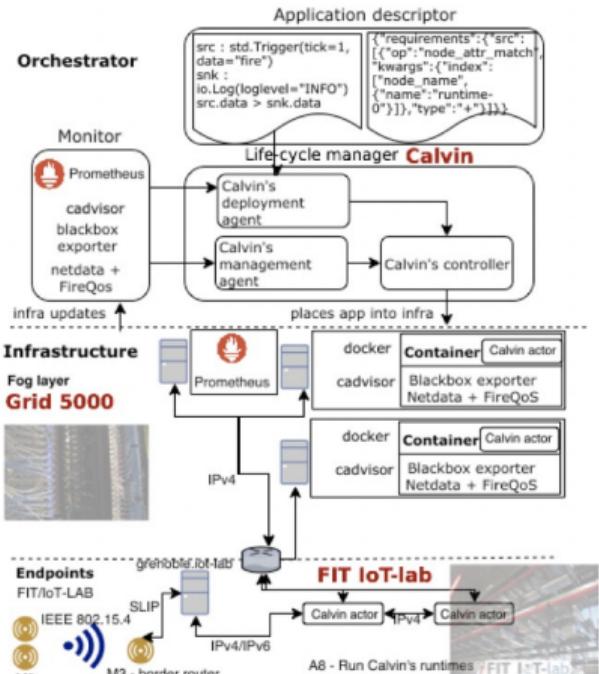
Objective: Design a Optimized Fog Service Provisioning strategy (O-FSP) and validate it on a real infrastructure

► Contributions

- ◆ Design and implementation of FITOR, an orchestration framework for the automation of the deployment, the scalability management, and migration of micro-service based IoT applications
- ◆ Design of a provisioning solution for IoT applications that optimizes the placement and the composition of IoT components, while dealing with the heterogeneity of the underlying Fog infrastructure

► Experiments

- ◆ Fog layer = 20 servers from Grid5000 which are part of the genepi cluster, Mist layer = 50 A8 nodes from IOTLab
- ◆ Use of a software stack made of open-source components (Calvin, Prometheus, Cadvisor, Blackbox exporter, Netdata)
- ◆ Experiments show that the O-FSP strategy makes the provisioning more effective and outperforms classical strategies in terms of: i) acceptance rate, ii) provisioning cost, and iii) resource usage



Toward a resource management system for Fog/Edge infras.

Inria Project Lab: Discovery

- ◆ Design a resource management system (a.k.a. a cloudkit) for Fog/Edge infrastructures
- ◆ A four year project started in 2015 with Inria, Orange (and initially Renater)
- ◆ Designing from scratch such a system cannot be envisioned (OpenStack 13 Millions of LOCs)

Contributions

- ◆ Implementation of a complete workflow to evaluate OpenStack WANWide scenarios
- ◆ Evaluate OpenStack up to 1000 compute nodes (Grid'5000, oct 2016)
- ◆ Evaluate OpenStack WANWide (impact of latency and throughput constraints) (oct 2017)
- ◆ Evaluation of communication bus for Fog/Edge scenarios (May 2018)
- ◆ Evaluation of database backends (NewSQL, NoSQL, etc. (May 2018)

Multi-Level Elasticity for Data Stream Processing

Vanja Maragozova-Martin, Noëll de Palma and Ahmed El Rheddiene
Univ. Grenoble Alpes, CNRS, LIG, F-38000 Grenoble France
E-mail: first.name.secondName@imag.fr

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TPDS.2018.2867966, IEEE Transactions on Parallel and Distributed Systems

SUBMISSION TO TPDS

would need to be deployed in containers with different capacities which in turn call for multi-dimensional-bin-packing-oriented scheduling [45].

[18] "CoMD," <https://gpugreen.com/compute-product/comd/>.

[19] Y. Wu and K. L. Tan, "ChronoStream: Elastic Stateful Stream Computation," in *2015 IEEE International Conference on Data Engineering*, April 2015, pp. 723–734.

[20] V. Gultiano, R. Jimenez-Perez, M. Patino-Martinez, C. Soriente, and P. Valduriez, "StreamCloud: An Elastic and Scalable Data Streaming System," *IEEE Trans. Parallel Distrib. Syst.*, vol. 23, no. 12, pp. 2351–2365, Dec. 2012.

[21] L. Neuemyer, B. Robbins, A. Nair, and A. Kesari, "S4: Distributed stream computing platform," in *Data Mining Workshops (ICDMW), 2010 IEEE International Conference on*, IEEE, 2010, pp. 170–177.

[22] OpenStack, <https://www.openstack.org/>

[23] Grid'5000, <http://www.grid5000.fr/>

[24] J. Chernousov, D. Pettin, A. Simonet, A. Lebre, and M. Simonin, "Toward a Holistic Framework for Conducting Scientific Evalu-

The experimental work presented in this paper would not have been possible without the existence of the Grid'5000 platform and the help of the supporting teams. The authors would also like to thank the *enos* team who made the OpenStack deployment process a child's play.

ACKNOWLEDGEMENTS



Last update: Mon Apr 16 14:49:02 CEST 2018

Deploy a micro DC on each Network Point of Presence

Toward a resource management system for Fog/Edge infras.

► Inria Project Lab: Discovery (contd)

- ◆ The creation of a dedicated working group within the OpenStack community that deals with Fog/Edge challenges (now managed by the foundation with key actors such as ATT, Verizon, CISCO, China mobile etc.)
- ◆ Several presentations / publications (see the DISCOVERY website)
- ◆ France has the main academic actor in the worldwide community (Inria/IMT Atlantique) thanks to the G5K testbed in particular.
 - ★ A leadership position
 - ★ A strong expertise for experiments related to performance, scalability of OpenStack components (concrete actions with RedHat, ongoing actions with Huawei, etc.)

OPENSTACK COLLABORATES WITH OTHER EDGE GROUPS

openEDGE computing ETSI BeyondTheCloud.github.io

Get more involved with the
Fog Edge Massively Distributed Computing SIG

F. Desprez - SILECS - Frédéric.Desprez@inria.fr



Open Infrastructure summit
Vancouver May 2018
(3000 participants)

Distributed Storage for Fog/Edge infrastructures

► Objective

- ◆ Design of a storage system taking locality of edge resources into account
- ◆ “Must-have” Properties: data locality, network containment, mobility support, disconnected mode, scalability

► Contributions

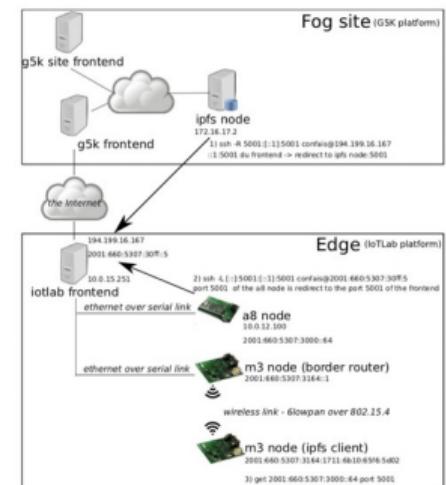
- ◆ Improving data locality by interconnecting Fog Scale-Out NAS systems with IPFS
- ◆ Improving meta-data locality thanks to a tree based approach inspired by the DNS

► Grid'5000 based experiments

- ◆ Evaluate the gain of using IPFS and scale out NAS systems for a 10 fog site infrastructure emulated on Grid'5000 (clients are deployed within Grid'5000).
- ◆ Evaluate the gain of a tree-based meta-data management
- ◆ ICFEC'2017 and GLOBECOM 2018

► Grid'5000-FIT experiments

- ◆ Evaluate the penalties/side effects of using representative Fog clients (Fog servers are deployed on Grid'5000 and clients on the IoTLab platform)
- ◆ Enabled us to identify several limitations (experiments using IoTLab and Grid'5000 are (currently) not easy to perform)



From Network Traffic Measurements to QoE for Internet Video

Problem solved: Estimation of QoE from encrypted video traces using network level measurements only)

- ▶ Play out a wide range of videos under realistic network conditions to build ML models (classification and regression) that predict the subjective MOS (Mean Opinion Score) based on the ITU P.1203 model along with the QoE metrics of startup delay, quality (spatial resolution) of playout and quality variations using only the underlying network Quality of Service (QoS) features

▶ A diverse QoS-QoE dataset

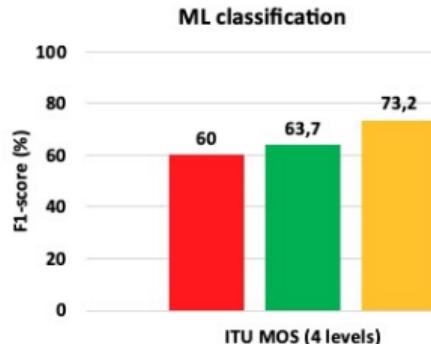
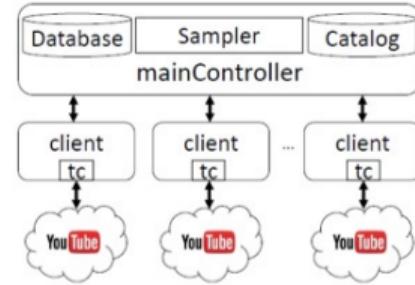
- ◆ Around 100k unique video playouts from geographically distributed locations (Sophia Antipolis, Grenoble, Rennes, Nancy, Nantes) using compute resources from AWS, Grid5000, and R2lab platforms

▶ Input features for ML:

- ◆ Network QoS (outband,inband, inband+chunks)

▶ Output labels:

- ◆ App QoS (startup delay, resolution, quality switches) and ITU P.1203 MOS



An experiment's outline

- ① Discovering resources and selecting resources
- ② Reconfiguring the resources to meet experimental needs
- ③ Monitoring experiments, extracting and analyzing data
- ④ Controlling experiments \leadsto automation, reproducible research

Discovering and selecting resources

► Describing resources ~ understand results

- ◆ Covering nodes and network infrastructure
- ◆ Machine-parsable format ~ scripts
- ◆ Human-readable description on the web⁵
- ◆ Archived (*State of testbed 6 months ago?*)
- ◆ Verified
 - ★ Avoid inaccuracies/errors ~ wrong results
 - ★ Self-checking by nodes before each reservation

► Selecting resources

- ◆ Complex queries using resource manager

```
oarsub -p "wattmeter='YES' and gpu='YES'"
```

```
oarsub -l "{cluster='a'}/nodes=1+
```

```
{cluster='b' and eth10g='Y'}/nodes=2"
```

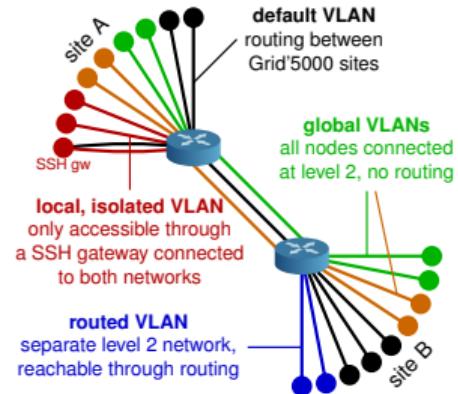
⁵<https://www.grid5000.fr/w/Hardware>

#	Cluster	Queue	Date of arrival	Nodes	CPU	Core	Memory	Storage	Network	Accelerators
Grid5000	clustu	default	2014-03-07	32	2 x Intel Xeon E5-2620 v3	36 cores/CPUs	192 GB	240 GB SSD + 480 GB SSD + 4.8 TB HDD	10 Gbps + 200 Gbps Direct Path	
Grid5000	jerli	default	2014-03-05	4	2 x Intel Xeon E5-2620 v3	32 cores/CPUs	192 GB	480 GB SSD + 2 x 1.8 TB SSD + 3 x 2.8 TB HDD	10 Gbps + 200 Gbps Direct Path	
Life	cheeset	default	2014-12-01	15	2 x Intel Xeon E5-2620 v4	36 cores/CPUs	256 GB	2 x 360 GB SSD	2 x 10 Gbps	
Life	chicken	default	2018-08-05	8	2 x AMD EPYC 7371	35 cores/CPUs	128 GB	480 GB SSD + 2 x 4 TB HDD	2 x 25 Gbps	
Life	coffee	default	2018-12-01	8	2 x Intel Xeon E5-2620 v4	34 cores/CPUs	196 GB	2 x 480 GB SSD + 2 x 4 TB HDD	2 x 10 Gbps	
Life	coffee	default	2018-08-05	9	2 x Intel Xeon Gold E5-2620 v3	32 cores/CPUs	192 GB	2 x 480 GB SSD + 4 x 4.8 TB HDD	2 x 25 Gbps	
Lyon	lyonhpc	prodshare	2012-12-01	11	2 x Intel Xeon E5-2620 v3	40 cores/CPUs	198 GB	544 GB SSD	2 x 10 Gbps + 100 Gbps	
Lyon	lyonhpc	prodshare	2013-09-30	15	2 x Intel Xeon E5-2620 v3	40 cores/CPUs	198 GB	256 GB SSD	2 x 10 Gbps	
Lyon	hercule	default	2012-10-02	4	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	32 GB	3 x 2.0 TB HDD	10 Gbps	
Lyon	rome	default	2016-12-01	23	2 x Intel Xeon E5-2620 v4	8 cores/CPUs	64 GB	596 GB SSD	10 Gbps	
Lyon	colin	default	2012-09-14	4	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	32 GB	596 GB SSD	10 Gbps	
Lyon	cogito	default	2008-07-01	32	2 x AMD Opteron 250	32 cores/CPUs	2 GB	73.6 GB HDD	1 Gbps	
Lyon	louis	default	2012-09-14	14	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	32 GB	596 GB SSD	10 Gbps	
Nancy	grapheo	production	2018-01-04	18	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	128 GB	2 x 800 GB SSD	10 Gbps + 50 Gbps Infiniband	
Nancy	graphique	production	2018-05-12	6	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	64 GB	256 GB SSD	10 Gbps + 50 Gbps Infiniband	1/2 x Nvidia Tesla P40
Nancy	graphix	default	2013-12-05	4	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	256 GB	2 x 300 GB SSD	10 Gbps + 50 Gbps Infiniband	Intel Xeon Phi 7210P
Nancy	goings	production	2013-04-09	48	2 x Intel Xeon E5-2620 v3	48 cores/CPUs	64 GB	1 TB HDD	1 Gbps + 56 Gbps Infiniband	
Nancy	goole	production	2017-06-26	14	2 x Intel Xeon E5-2620 v4	32 cores/CPUs	128 GB	2 x 299 GB SSD	10 Gbps + 200 Gbps Direct Path	2 x Nvidia GTX 1080 Ti
Nancy	gorenix	production	2018-08-03	6	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	64 GB	1 TB HDD	10 Gbps + 300 Gbps Direct Path	2 x Nvidia Tesla K40M
Nancy	gramene	default	2018-01-02	8	2 x Intel Xeon E5-2620 v3	8 cores/CPUs	128 GB	200 GB SSD + 5 x 600 GB HDD	4 x 10 Gbps + 96 Gbps Infiniband	

Reconfiguring resources

- ▶ Operating System reconfiguration with **Kadeploy**:
 - ◆ Provides a *Hardware-as-a-Service* cloud infrastructure
 - ◆ Enable users to deploy their own software stack & get *root* access
 - ◆ **Scalable, efficient, reliable and flexible:**
200 nodes deployed in ~5 minutes
- ▶ Customize **networking** environment with **KaVLAN**
 - ◆ Protect the testbed from experiments (Grid/Cloud middlewares)
 - ◆ Avoid network pollution
 - ◆ Create custom topologies
 - ◆ By reconfiguring VLANS \leadsto almost no overhead

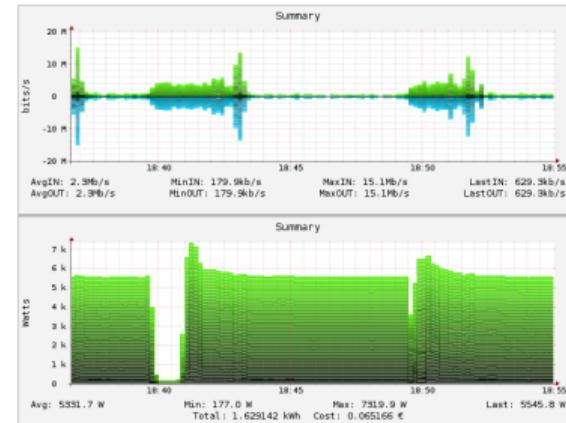
KADEPLOY



Monitoring experiments

Goal: enable users to understand what happens during their experiment

- ▶ System-level probes (usage of CPU, memory, disk, with Ganglia)
- ▶ Infrastructure-level probes: Kwapi
 - ◆ Network, power consumption
 - ◆ Captured at high frequency (≈ 1 Hz)
 - ◆ Live visualization
 - ◆ REST API
 - ◆ Long-term storage



Controlling experiments

- ▶ Legacy way of performing experiments: shell commands
 - ⌚ time-consuming
 - ⌚ error-prone
 - ⌚ details tend to be forgotten over time
- ▶ Promising solution: **automation of experiments**
 - ~ Executable description of experiments
 - ~ Reproducible research
- ▶ Support from the testbed: Grid'5000 RESTful API
(Resource selection, reservation, deployment, monitoring)
- ▶ Several higher-level tools to help automate experiments
Execo, Python-Grid5000 (Python), Ruby-cute (Ruby)
<https://www.grid5000.fr/w/Grid5000:Software>



Users and publications

Year	2010	2011	2012	2013	2014	2015	2016	2017	2018
Active users	564	553	592	514	528	458	573	600	564
Publications	154	141	101	134	106	143	122	151	127
PhD & HDR	14	20	9	27	24	30	27	23	22
Usage rate	50%	56%	58%	63%	63%	63%	55%	53%	70%

- ▶ 1313 active users over the last 3 years
- ▶ 3769 active users since 2003
- ▶ 2007 publications that benefited from Grid'5000 in our **HAL collection**⁶
 - ◆ Computer Science: 96%, Mathematics: 2.4%, Physics: 2.4%
 - ◆ Since 2015: LORIA: 23%, IRISA: 23%, LIG: 19%, LIP: 13%, LS2N: 13%, CRISTAL: 5%, LIRMM: 5%, LIP6: 3%

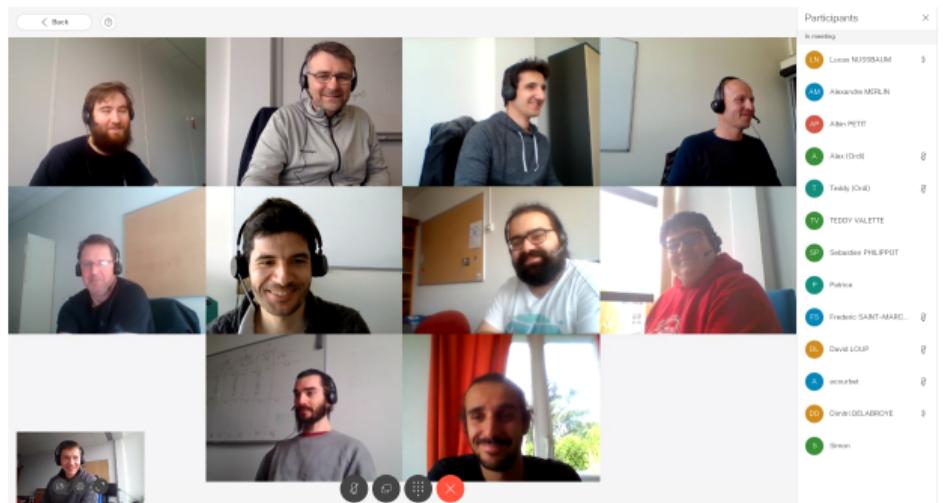
⁶<https://hal.archives-ouvertes.fr/GRID5000>

Organization and governance

- ▶ **Director** – Frédéric Desprez
 - ▶ **Bureau** (6 members: FD, LN, Christian Perez, Adrien Lebre, Laurent Lefevre, David Margery)
 - ▶ **Comité des responsables de sites**
 - ▶ **Technical Director** – Lucas Nussbaum
 - ◆ Technical team
 - ▶ **Architects committee** (6 members)
 - ▶ **Conseil de groupement**
 - ◆ Inria, CNRS, RENATER, CEA, CPU, CDEFI, IMT
(≈ Allistène + RENATER)
 - ▶ **Conseil scientifique**
 - ◆ 10 members
- 
- institutional and scientific steering
- technical steering
- advisory and evaluation bodies

Technical organization

- ▶ Distributed infrastructure, but managed by a single distributed team
 - ◆ Strong coherence and coordination between sites
- ▶ Current composition: 8.13 full-time engineers
 - ◆ Inria: 5.91 (perm: 0.86, CDD: 5.1), CNRS: 1.02 (perm: 1.02), U. Rennes: 0.6 (perm: 0.6), IMT Atlantique: 0.4 (CDD: 0.4), U. Lorraine: 0.2 (perm: 0.2)
 - ◆ Including Pierre Neyron: Médaille de Cristal du CNRS 2019



Positionnement: France Grilles & EGI

- ▶ Confusion entre Grid'5000 et France Grilles, à cause du nom (« grilles »)
- ▶ Positionnements très différents:
 - ◆ Grid'5000:
 - ★ Plate-forme expérimentale très homogène, gérée de manière uniforme
 - ★ Communauté utilisatrice: surtout recherche en informatique
 - ◆ France Grilles & EGI:
 - ★ Plate-forme mutualisant des ressources de calcul et de traitement, avec un modèle de fédération basé sur des VO (Virtual Organizations)
 - ★ Communauté utilisatrice: tous les domaines scientifiques (peu de recherche en informatique)

Positionnement: Supercalculateur IA/HPC Jean Zay

- ▶ Positionnement difficile à définir à ce stade
- ▶ Supercalculateur IA/HPC Jean Zay (GENCI/IDRIS) :
 - ◆ Usages production, *batch*
 - ◆ Grosses campagnes de calcul
 - ◆ Capacité de calcul très importante
 - ◆ Volonté d'une flexibilité accrue par rapport au mode de fonctionnement habituel à l'IDRIS
- ▶ Grid'5000 :
 - ◆ Ressources flexibles, agiles, de proximité, facilement accessibles
 - ◆ Usages développement / test / expérimentation, souvent en interactif
 - ◆ Reconfiguration – *HPC on steroids*
- ▶ Besoin de développer les interactions et de faciliter les passages entre Grid'5000 et Jean Zay

Conclusions

- ▶ An advanced and established infrastructure for the *data-center* facets of Computer Science
 - ◆ Large-scale, distributed
 - ◆ Shared (many involved laboratories and institutions)
 - ◆ Designed for reconfigurability, observability, reproducible research
- ▶ Looking forward to work with FIT in the context of SILECS
 - ◆ We share the same design goals, principles, but focus on different objects
 - ★ *Core of the Internet vs Edge of the Internet*
 - ★ *Internet of servers vs Internet of clients*
 - ◆ Strong needs of joint experiments

Backup slides

Reconfiguring the testbed

- ▶ Typical needs:
 - ◆ Install specific software
 - ◆ Modify the kernel
 - ◆ Run custom distributed middlewares (Cloud, HPC, Grid)
 - ◆ Keep a stable (over time) software environment

Reconfiguring the testbed

- ▶ Typical needs:
 - ◆ Install specific software
 - ◆ Modify the kernel
 - ◆ Run custom distributed middlewares (Cloud, HPC, Grid)
 - ◆ Keep a stable (over time) software environment
- ▶ Likely answer on any production facility: **you can't**
- ▶ Or:
 - ◆ Install in \$HOME, modules ↵ no root access, handle custom paths
 - ◆ Use virtual machines ↵ experimental bias (performance), limitations
 - ◆ Containers: kernel is shared ↵ various limitations

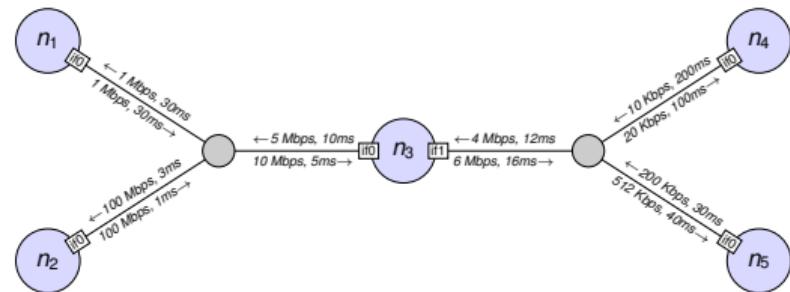
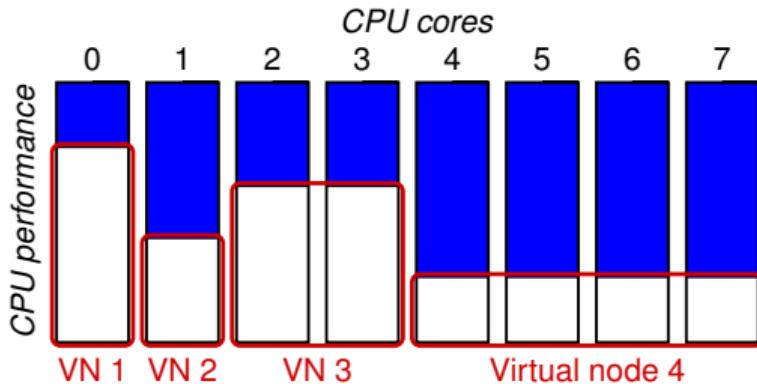
Creating and sharing Kadeploy images

- ▶ When doing manual customization:
 - ◆ Easy to forget some changes
 - ◆ Difficult to describe
 - ◆ The full image must be provided
 - ◆ Cannot really serve as a basis for future experiments
(similar to binary vs source code)
- ▶ Kameleon: Reproducible generation of software appliances
 - ◆ Using *recipes* (high-level description)
 - ◆ Persistent cache to allow re-generation without external resources (Linux distribution mirror) ↵ self-contained archive
 - ◆ Supports Kadeploy images, LXC, Docker, VirtualBox, qemu, etc.

<http://kameleon.imag.fr/>

Changing experimental conditions

- Reconfigure experimental conditions with **Distem**
 - ◆ Introduce heterogeneity in an homogeneous cluster
 - ◆ Emulate complex network topologies



<http://distem.gforge.inria.fr/>

(Including a tutorial about OpenFlow and P4)



Utilisateurs IA actuels sur Grid'5000

- ▶ Au LORIA (Nancy):
 - ◆ Équipes MULTISPEECH et SYNALP (deep learning pour le **traitement automatique des langues et de la parole**)
 - ◆ Équipe BISCUIT et LARSEN (deep learning pour la **robotique**)
- ▶ Dans le centre Inria Lille :
 - ◆ Équipe Bonus (optimisation massivement parallèle assistée par les métamodèles avec des applications dans les domaines de l'**ordonnancement de systèmes complexes et l'engineering design**)
 - ◆ Équipe Magnet (Apprentissage de représentations dans le domaine du **traitement automatique des langues**)
 - ◆ Équipe Sequel (apprentissage profond pour la **vision** et apprentissage par renforcement profond avec notamment des applications dans les **systèmes de dialogues parlés**)
 - ◆ Équipe Spirals (IA et data mining pour le **génie logiciel** en particulier pour des approches empiriques autour de l'**automated software repair**)
- ▶ Au LIG (Grenoble):
 - ◆ Équipe MRIM (Deep Learning pour l'**indexation de contenus multimédia**)

Modes d'utilisation et support de l'IA

Mode d'utilisation / réservation et interface

- ▶ Beaucoup d'usage interactif
 - ◆ Nécessaire pour préparer les expériences, et souvent pour expérimenter
 - ◆ Du coup, utilisation fréquente de réservations à l'avance
 - ◆ Pour les usages plus *production*, file et ressources spécifiques (en mode batch uniquement)
- ▶ Utilisateurs informaticiens: SSH+shell et/ou API REST (pour automatiser) ne posent pas de problèmes
 - ◆ Pour les utilisateurs pour qui ça pose un problème, il vaut mieux abandonner (\leadsto filtre : de toute façon, cela présage qu'ils auront d'autres difficultés insurmontables par la suite)
 - ◆ Un portail web de soumission existe, mais n'est pas maintenu (était plutôt utilisé à des fins de démos)

Utilisation des machines

- ▶ Trois modes d'utilisation:
 - ◆ Utilisation sans accès root
 - ◆ Déploiement d'un environnement personnalisé avec Kadeploy (OS, hyperviseur, ...)
 - ◆ *sudo-g5k*: comme utilisation sans déploiement, mais commande permettant de passer root
 - ★ Permet à l'utilisateur de modifier facilement son environnement d'expérimentation, sans limites
 - ★ Le noeud est nettoyé (voire réparé) en fin de job (redéployé avec Kadeploy)
 - ★ C'est aussi la technique utilisée pour lancer des containers (permet de contourner les problèmes de sécurité posés par les containers)

Containers

- ▶ Containers:
 - ◆ Outil automatisant la configuration de docker
 - ◆ Pour l'instant, pas de Singularity, Charliecloud, et autres (peu de demande utilisateur)
- ▶ Containers pour l'IA: peu de demandes. Plutôt utilisation de pip, virtualenv ou conda

Sécurité

- ▶ Politique plutôt permissive (blacklist plutôt que whitelist)
- ▶ Mais peu (pas?) de données sensibles
- ▶ Nous nous protégeons de l'extérieur (bonnes pratiques habituelles)
- ▶ Nous faisons confiance dans une certaine mesure à nos utilisateurs, et c'est difficile de faire autrement :
 - ◆ L'accès root permet beaucoup de choses
 - ◆ L'accès à Internet depuis les noeuds aussi

Matériel

- ▶ Segmentation de marché NVidia ↗ choix difficile pour l'IA
 - ◆ GPUs Tesla très chers
 - ◆ GPUs GeForce peu chers, mais:
 - ★ Warnings sur la durée en vie en utilisation DC (non supportée par le constructeur)
 - ★ Interdiction de l'utilisation en DC (licence drivers) – tolérance ?
- ▶ Pour l'instant, beaucoup d'usage mono-noeud, ou multi-GPU avec peu de communications (multi-paramétrique)
 - ◆ On achète surtout des GeForce (par ex Titan 1080 Ti) sur Grid'5000
- ▶ Projet Inria Project Lab en démarrage ↗ convergence HPC & IA
 - ◆ Besoin grandissant pour des GPUs Tesla (NVLink) et les technos émergentes pour l'IA (FPGA, etc.) ?

Matériel (2)

Accelerator model	Lille	Lyon	Nancy	Total
Intel Xeon Phi 7120P			4	4
Nvidia GTX 1080 Ti	16		28	44
Nvidia GTX 980			10	10
Nvidia Tesla K40M			12	12
Nvidia Tesla M2075		4		4
Nvidia Tesla P100	12			12
Nvidia Tesla V100	4			4
Nvidia Titan Black			2	2
Total	32	4	56	92

Conclusions – besoins utilisateurs IA

- ▶ Les piles logicielles évoluent très rapidement
 - ◆ Difficile (impossible?) de suivre avec un support par l'infrastructure
 - ◆ Il faut donner de la flexibilité aux utilisateurs pour qu'ils puissent adapter eux-mêmes leur pile logicielle
 - ★ Déploiement bare-metal avec Kadeploy
 - ★ Droits *root* avec sudo-g5k (+ Kadeploy *behind the scenes*)
~~ environnement *HPC on steroids*
- ▶ Accès facile à la plate-forme
 - ◆ Pas beaucoup plus compliqué qu'une machine dans son bureau
 - ★ Quelques heures pour obtenir un compte
- ▶ Utilisation simple
 - ◆ Souvent des utilisateurs IA peu familiers des environnements HPC
 - ◆ Pas beaucoup plus compliqué que travailler sur son portable Linux/Mac
 - ◆ Usage interactif et ressources rapidement disponibles pour la mise au point