



Contributions à la flexibilité et à l'efficacité énergétique des systèmes distribués à grande échelle

Soutenance d'Habilitation à Diriger les Recherches

Laurent Lefèvre

*Inria - Equipe Avalon – Laboratoire LIP - Ecole Normale Supérieure de Lyon
laurent.lefevre@inria.fr*

Plan

- Introduction
 - μ CV
 - Contributions/ Projets
- Flexibilité
 - Réseaux dynamiques et programmables : L'aventure Tamanoir
 - Nouveaux équipements et services : 1 focus
- Efficacité énergétique
 - Mesurer et comprendre
 - Vers des ordonnanceurs et des nuages verts
 - Améliorer l'efficacité énergétique du HPC avec ou sans connaissance des services et des applications
- Conclusions et perspectives

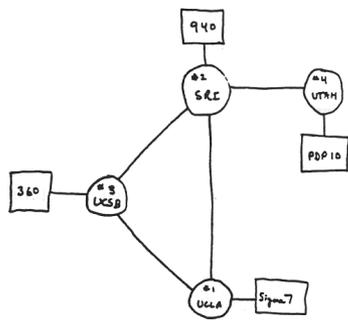
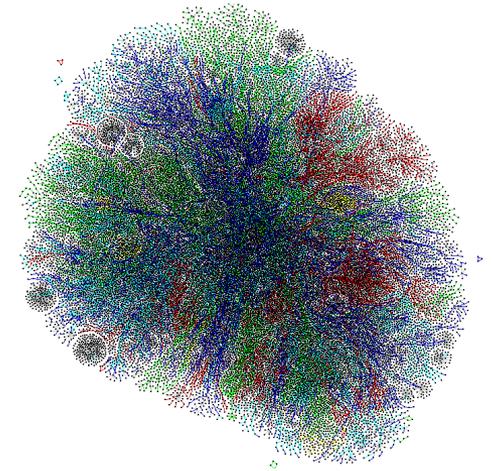
Introduction

Systemes distribués à grande échelle

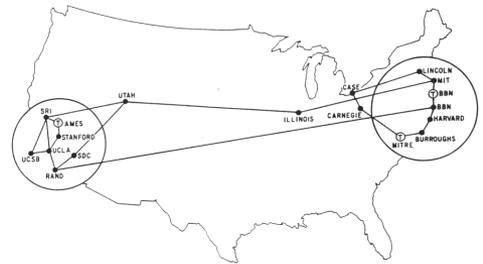
- Hétérogénéité des équipements terminaux
- Des tuyaux, et au bout de grosses infrastructures
 - DataCentres /HPC
 - Clouds / Grilles
 - Réseaux
 - Pas mobile
- Plus imposant, plus de services, plus d'utilisateurs

Systemes distribués à grande échelle

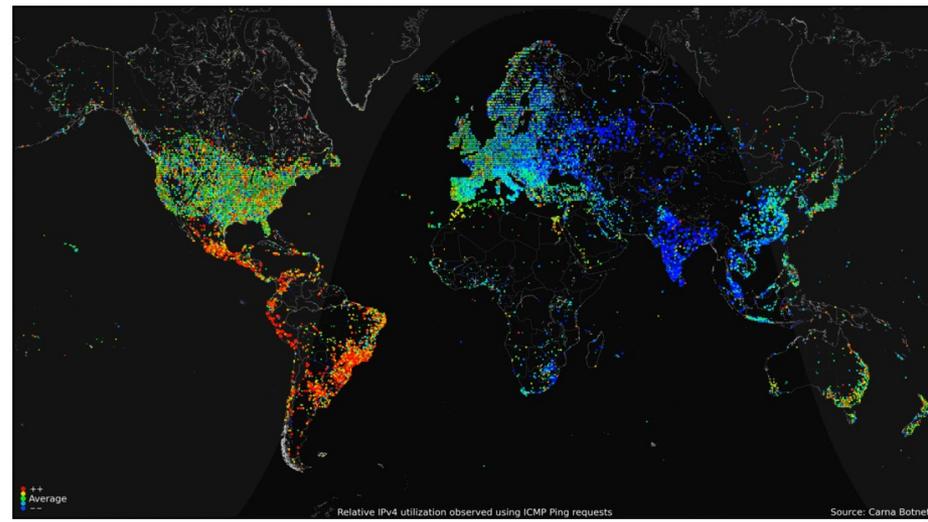
- Réseaux /Internet
- Arpanet (1969) -> Internet (2013)
- 4 sites (69) -> 213 (81) -> 5K (86) -> 1 Milliard de machines connectées (2013)



THE ARPA NETWORK
DEC 1969
4 NODES



MAP 4 September 1971



Relative IPv4 utilization observed using ICMP Ping requests
Source: Carna Botnet

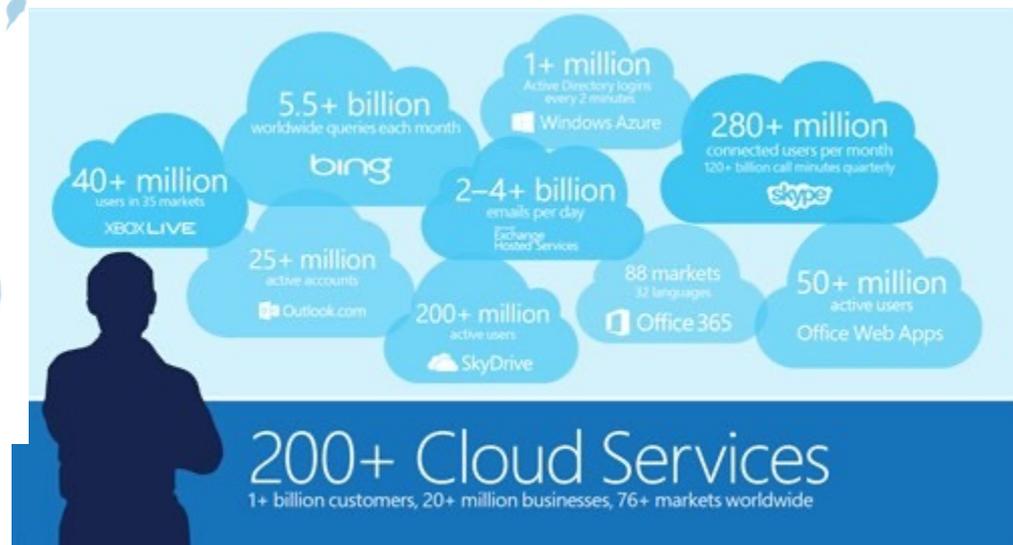
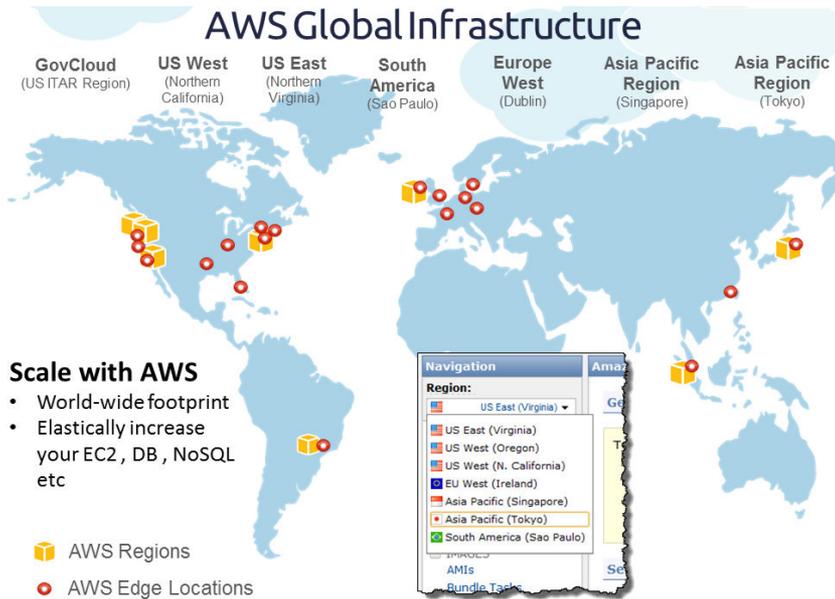
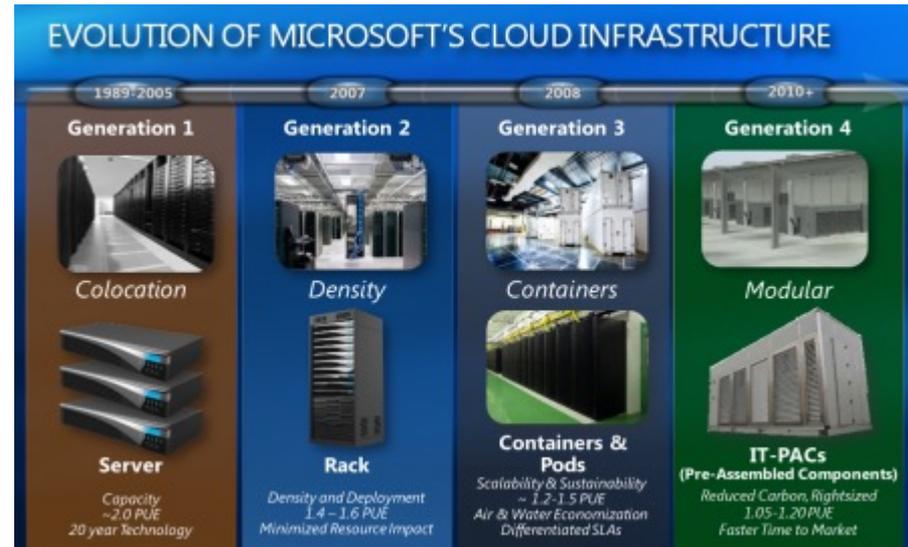
Systemes distribués à grande échelle

- DataCenters
- Google (1997 <10) -
> Google (2000 - 25K) -> Google
(2013 - 1M)



Systemes distribués à grande échelle

- Clouds : des ensembles de datacenters
- Amazon, Azure...
- IaaS, PaaS, HPCaaS...

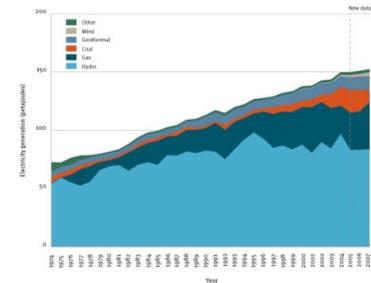
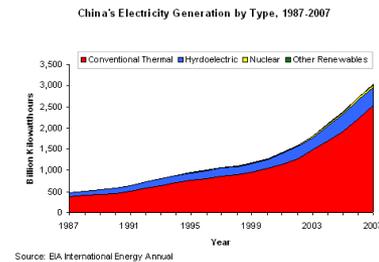
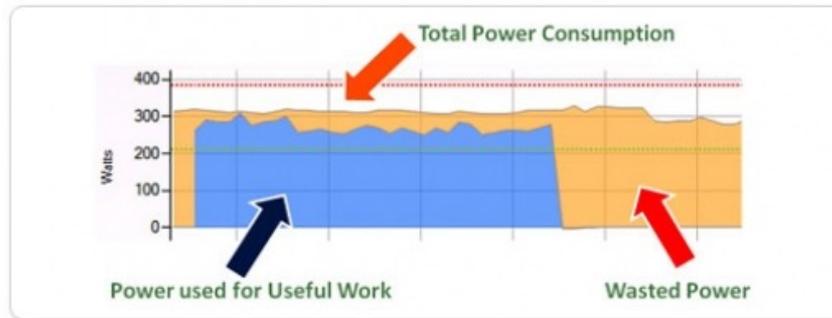


Besoins de nouveaux Services

- Certains services et protocoles (terminaux) voient le jour
- D'autres restent dans les cartons : multicast, QoS, sécurité...
- Temps vers standardisation très long
- Exemple : IPv6 (proposé en 95, finalisé en 98)
- Besoin de flexibilité et intelligence dans le réseau pour :
validation, déploiement expérimental, déploiement
opérationnel, acceptation

Consommation énergétique

- TICs représenteraient équivalent aviation en génération de Co2
- 10% électricité mondiale – 10% augmentation par an
- GreenIT : réduire la consommation électrique des infrastructures – CO2 dépend de la production
- Focalisation sur la phase d'usage



Questions abordées dans cette habilitation

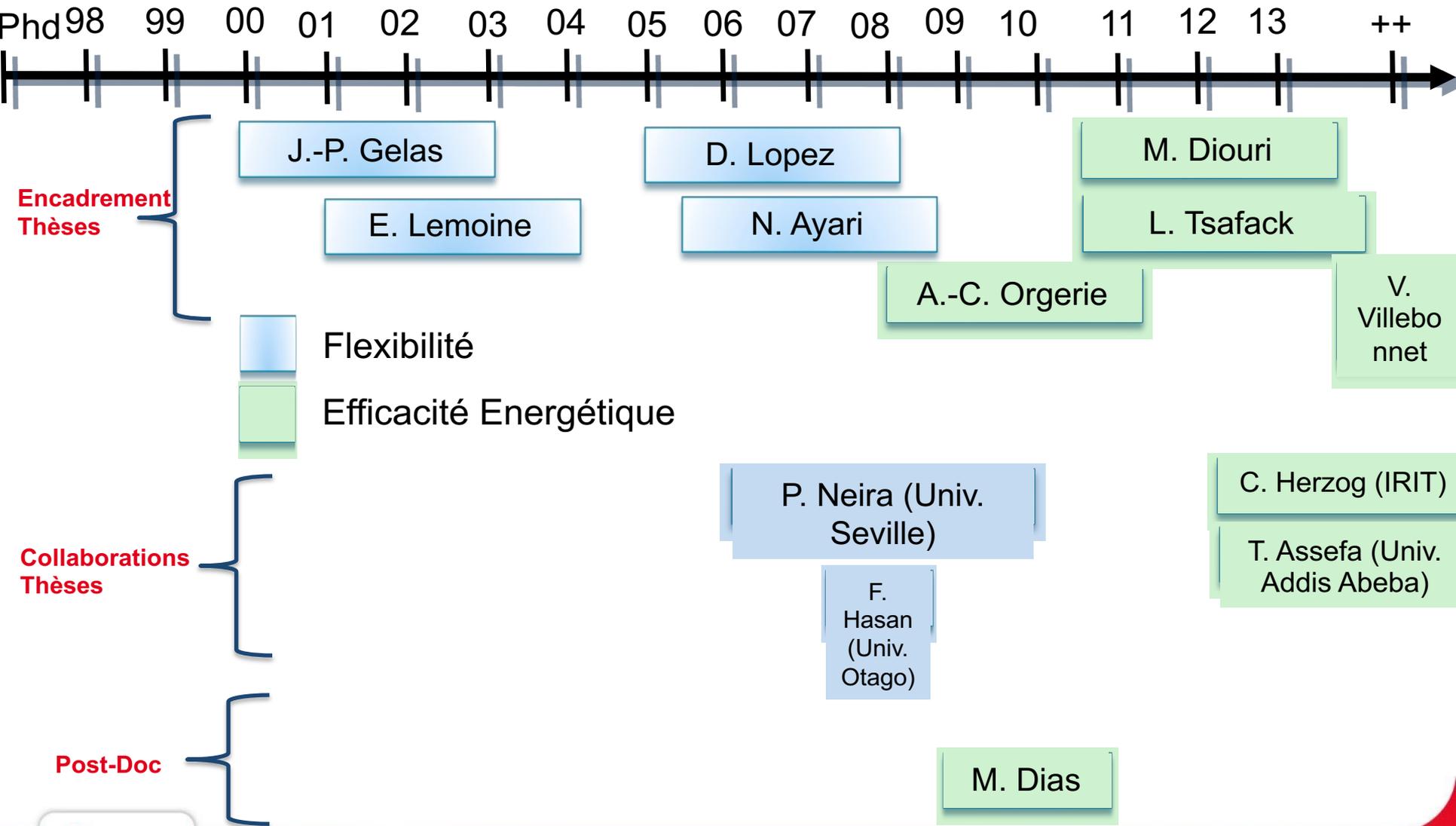
- Comment ajouter de la flexibilité dans les réseaux pour évaluer, expérimenter et valider de nouveaux services et équipements ?
- Comment proposer des modèles et solutions logicielles efficaces pour réduire la consommation énergétique des infrastructures distribuées à grande échelle ?

μCV

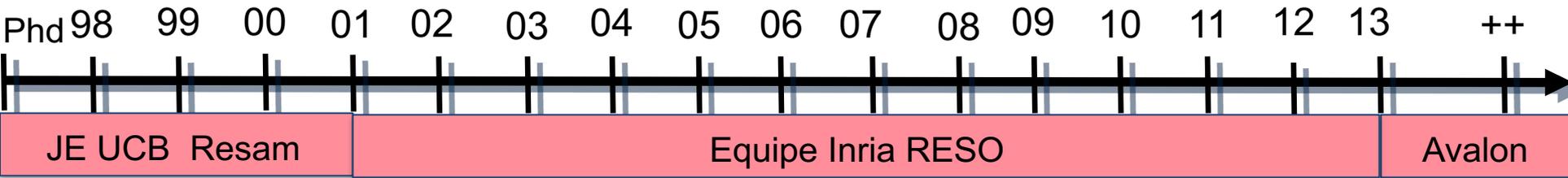
- [1993-1997] Thèse sur les Systèmes à Mémoire Distribuée Partagée
- [1997] Postdoc Rice University (Texas)
- [1997-2001] Maître de Conférences Université Claude Bernard Lyon1
- [2001-now] Chargé de Recherches Inria – Laboratoire LIP - ENS Lyon

- Co-Encadrement : 7 thèses, 1 Postdoc, 10 Ingénieurs, 13 M2, 20 M1
- Projets Nationaux, Européens, Internationaux

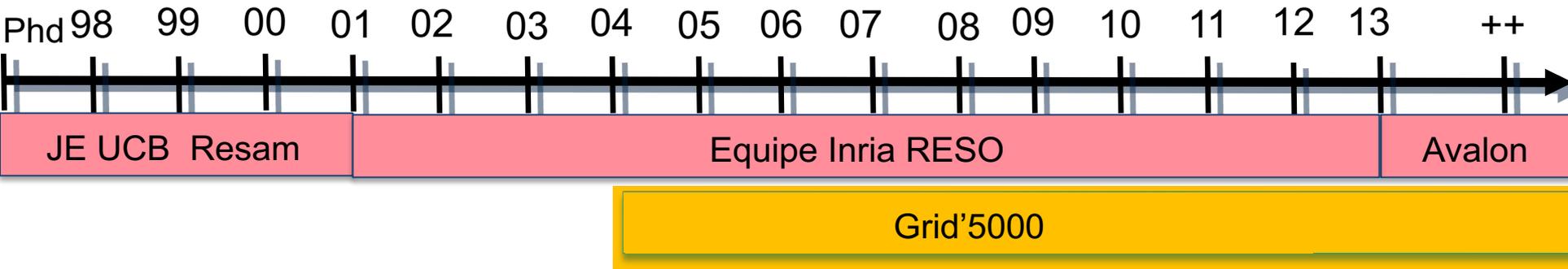
(Co) Encadrements



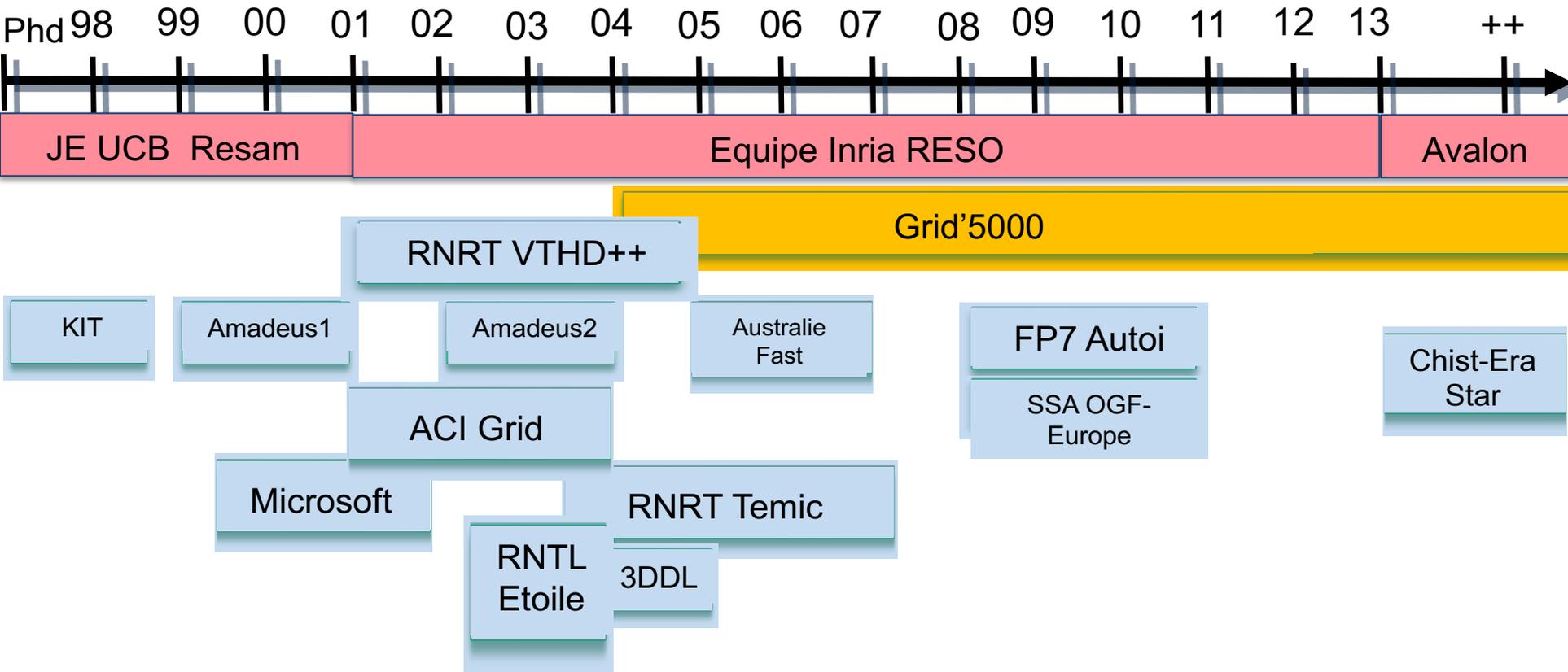
Projets de recherche



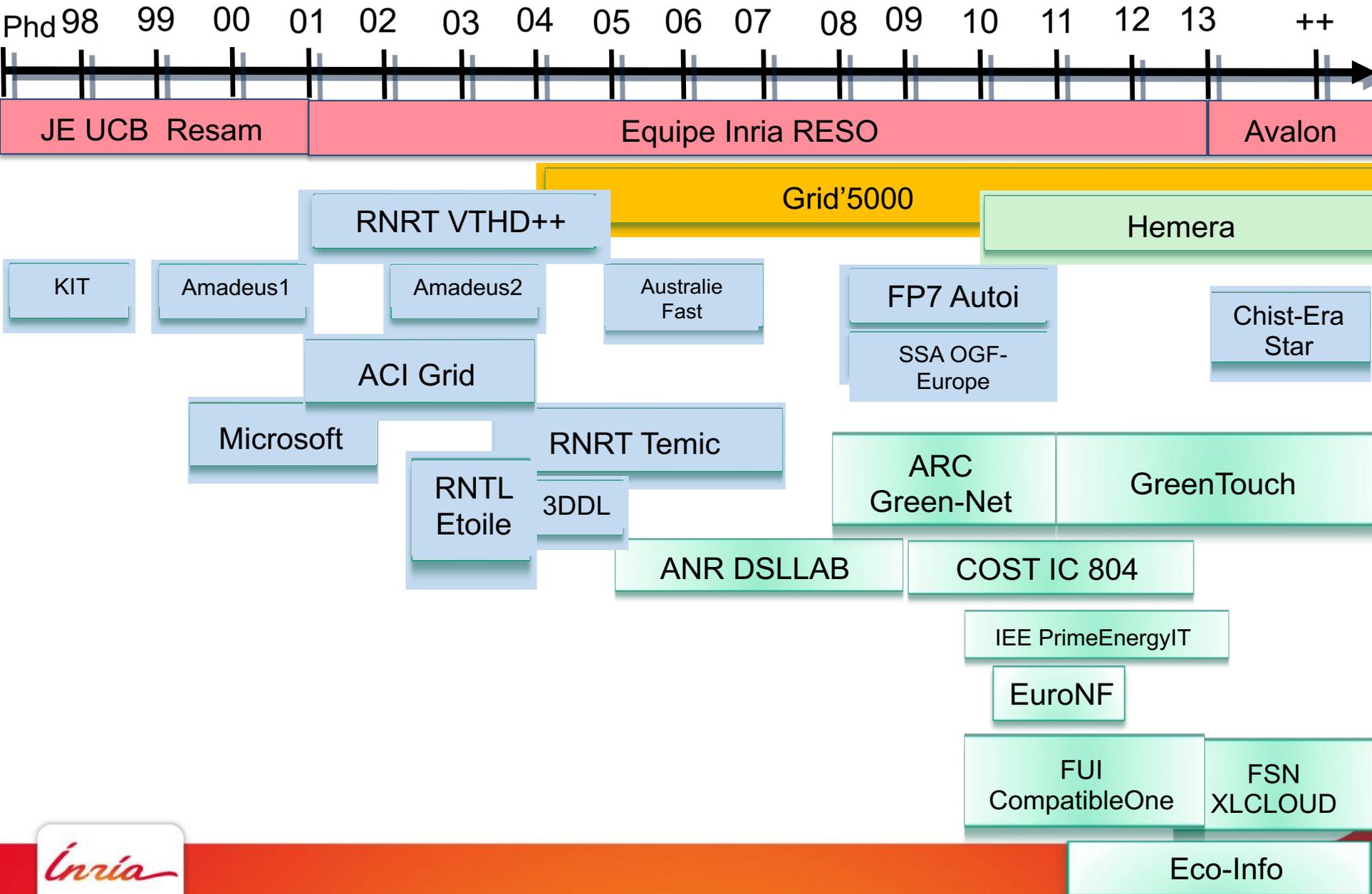
Projets de recherche (G5K)



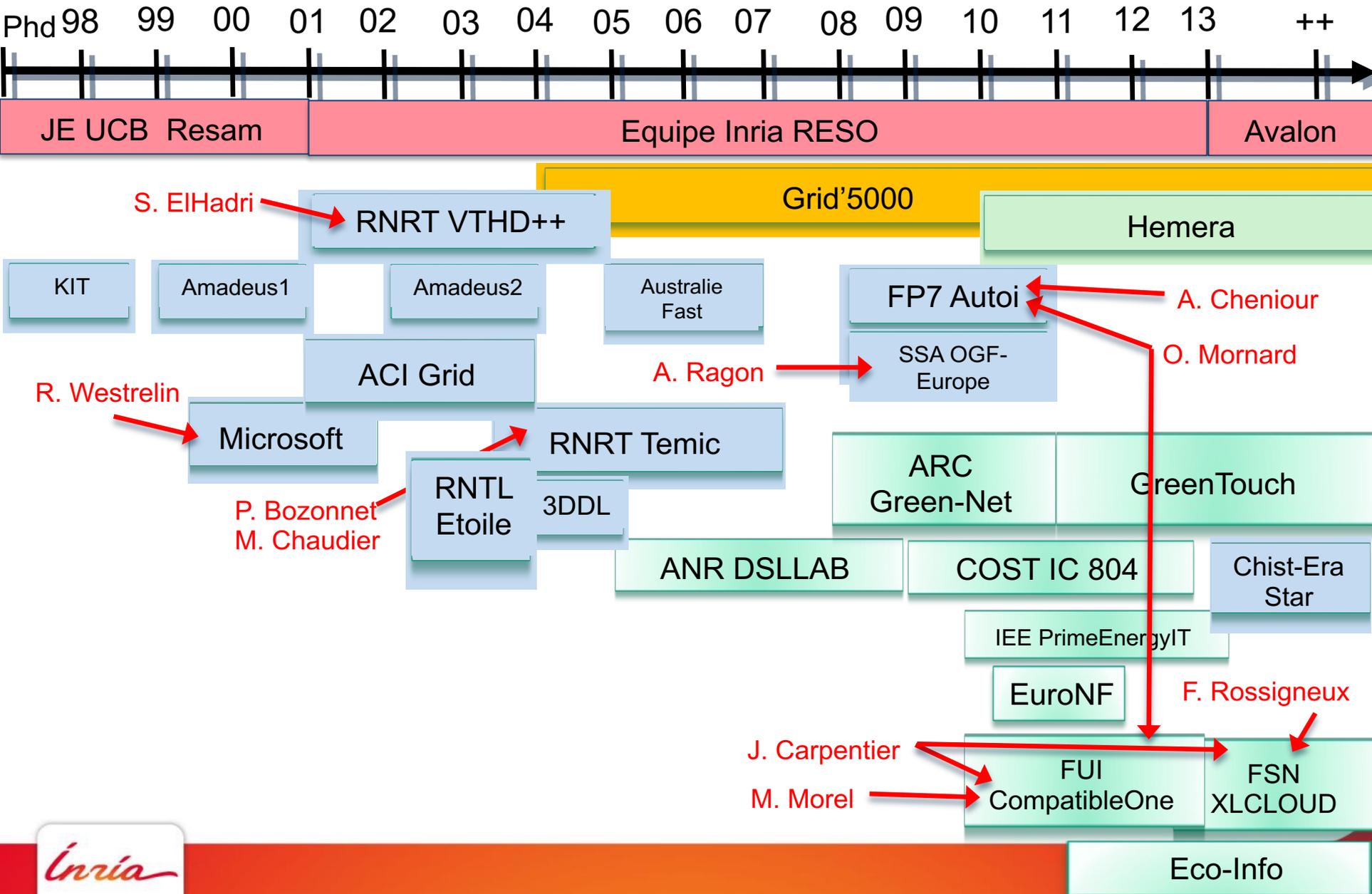
Projets de recherche (flexibilité)



Projets de recherche (EE)



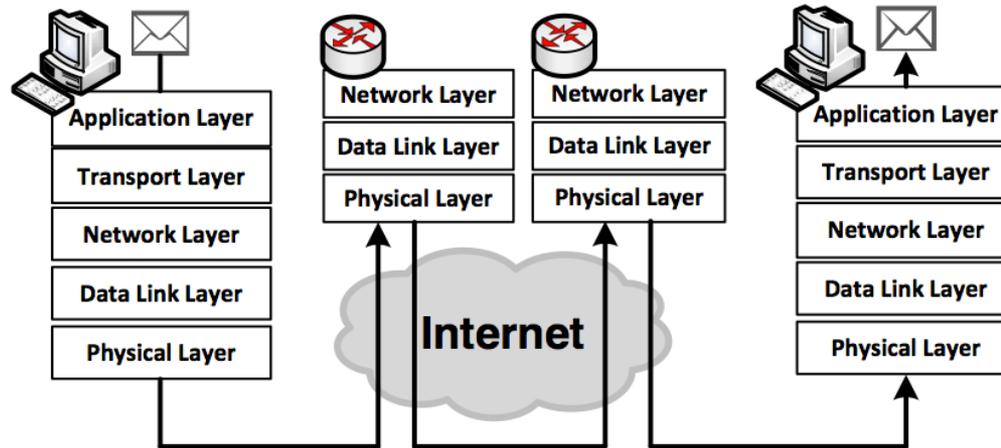
Projets de recherche (ingénieurs)



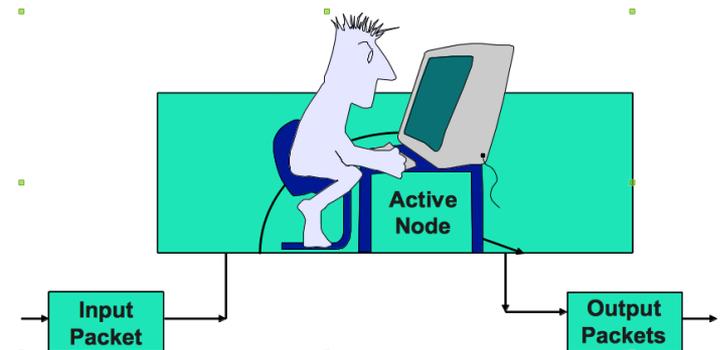
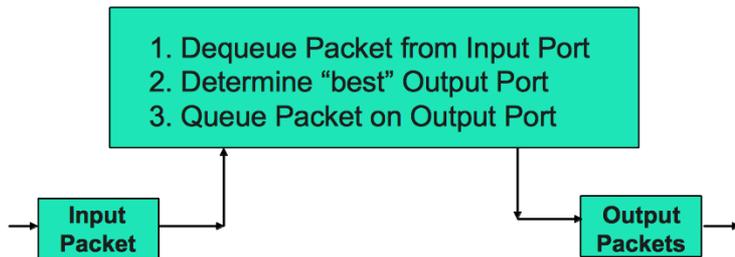
Flexibilité dans les systèmes distribués à grande échelle

Internet

- Le succès d'un réseau bête
- Le principe du bout en bout : *End2End (Saltzer84)* -> on urbanise l'intelligence aux extrémités : des réseaux simples pour des applications intelligentes



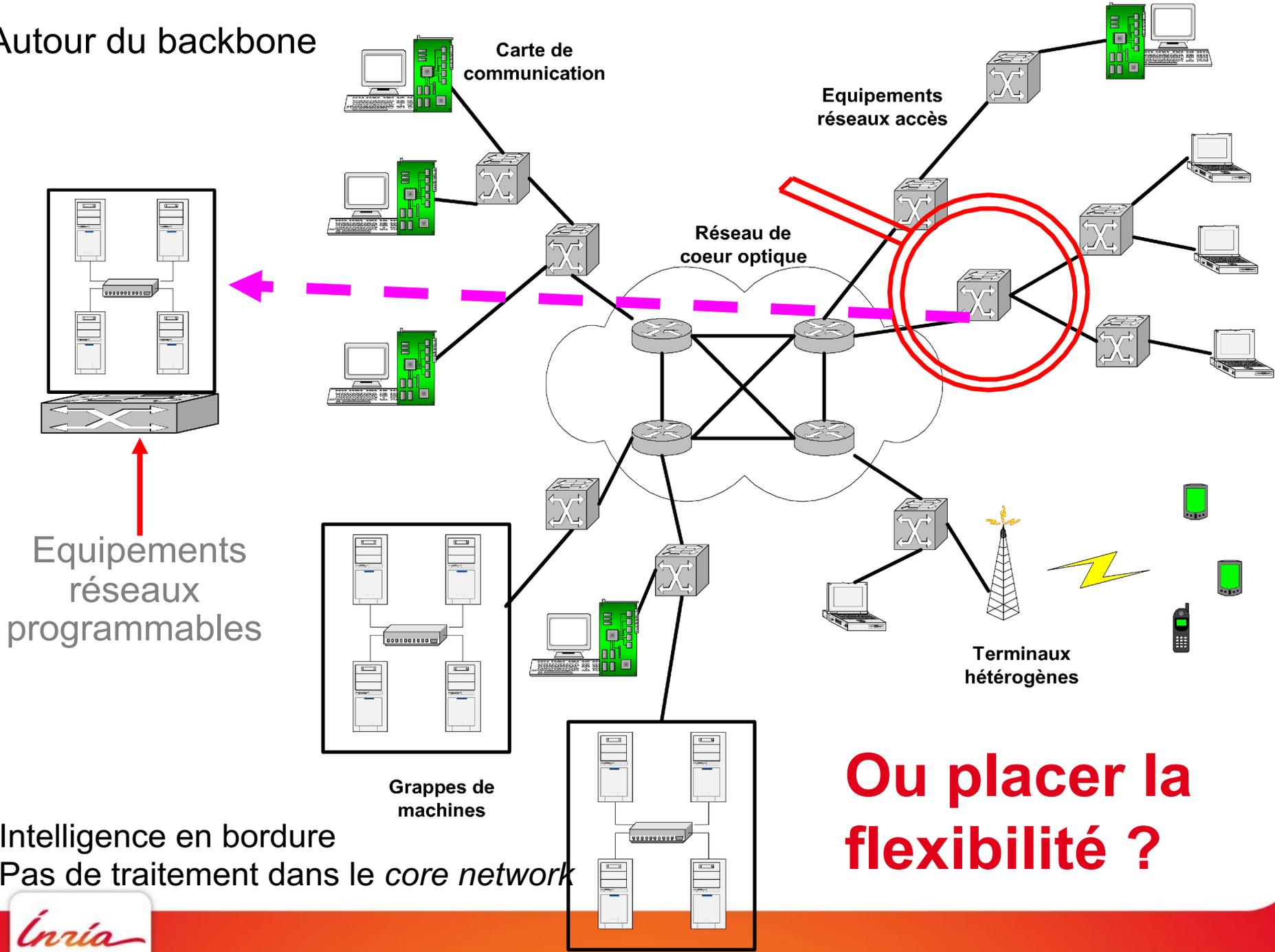
- De *route & route* à *route & process & route* - explorer les limites du E2E et proposer des solutions alternatives



4 vagues de flexibilité

- *Réseaux pour tous* : réseaux actifs et programmables
- *Réseaux sans intervention humaine* : réseaux autonomes
- *Infrastructures à valeur ajoutée* : réseaux virtuels
- *Réseaux flexibles pour maîtriser l'efficacité énergétique*

Autour du backbone



- Intelligence en bordure
- Pas de traitement dans le *core network*

Ou placer la flexibilité ?

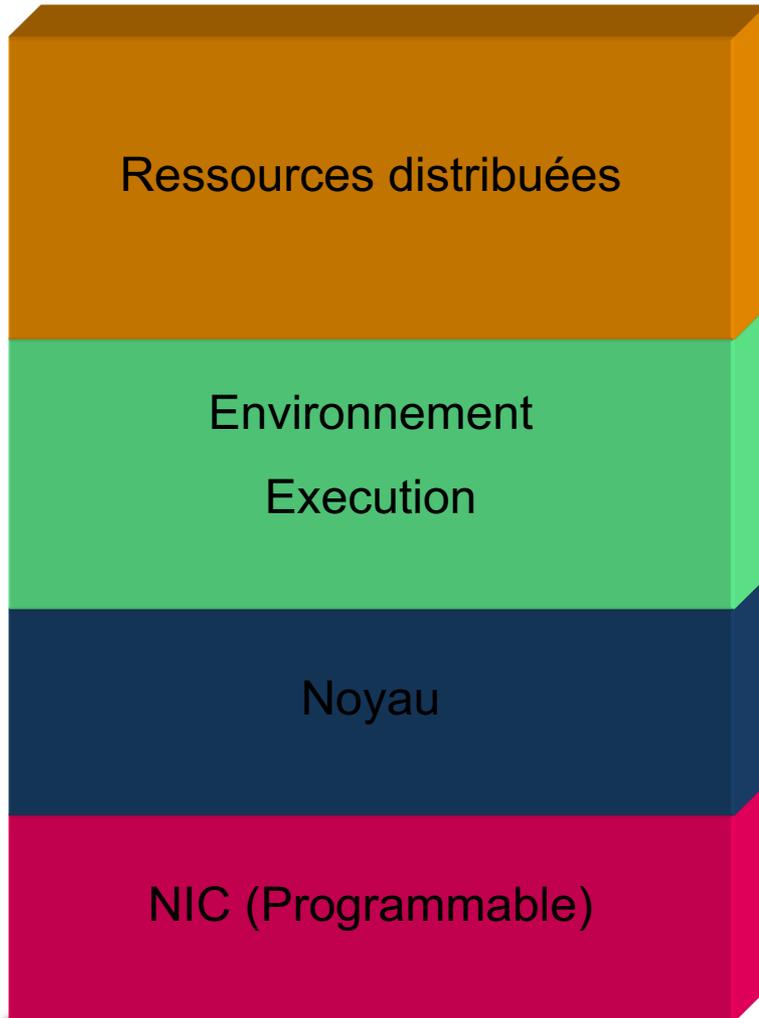
Objectifs

- 1^{er} système de réseaux actifs : ANTS (D. Tenenhouse) 
- Démonstration de la programmation par paquets
- Performances limitées (qq Mbits)
- Explorer les différentes alternatives pour fournir de la haute performance dans les Réseaux Actifs
- Fournir un EE portable et simple (Java) avec des couches adaptées aux besoins des services
- Développer / expérimenter une gamme variée de services actifs (multicast fiable, QoS, adaptation dynamique de flux..)



L'aventure Tamanoir : une architecture adaptée pour des services hétérogènes

Thèse
Jean-Patrick
Gelas



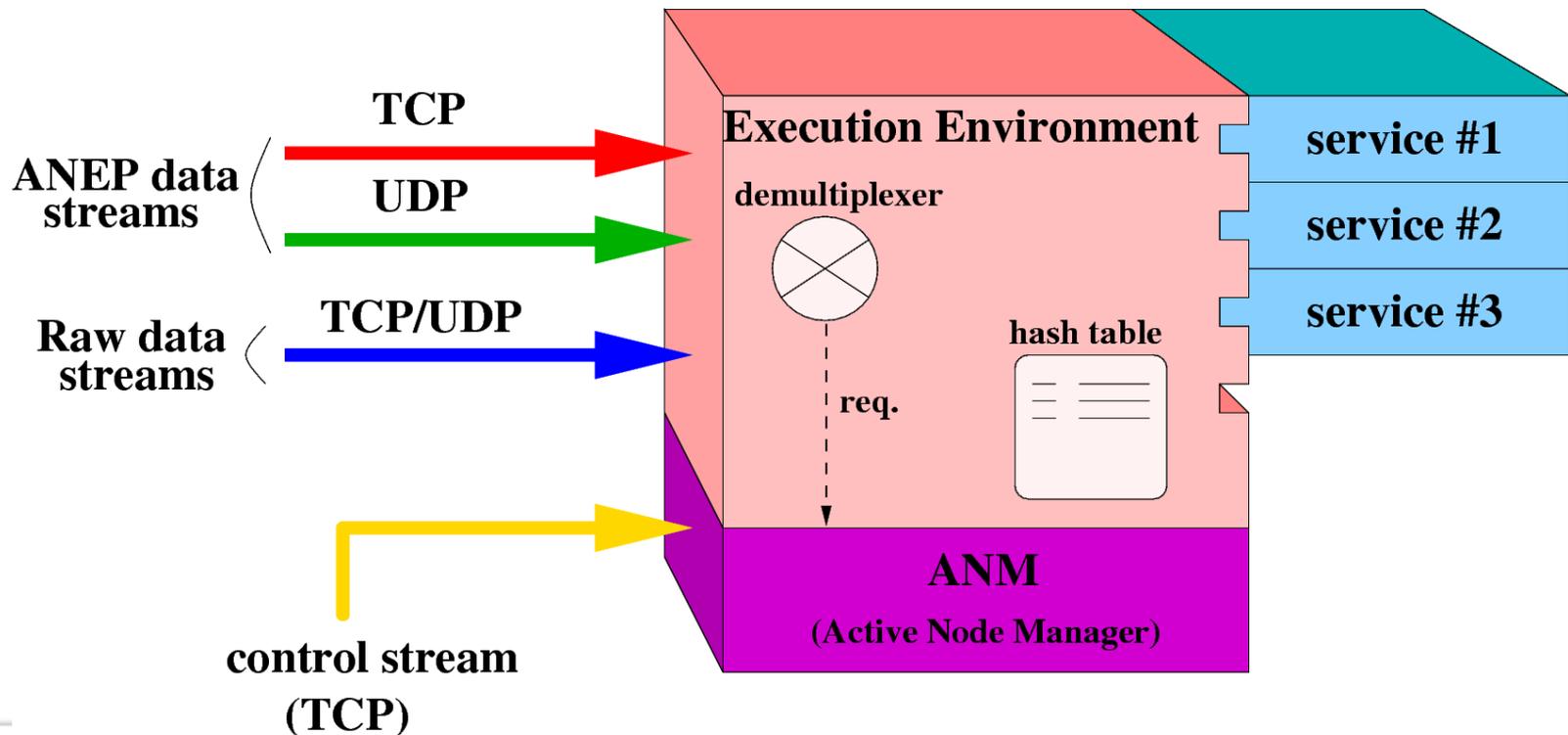
- **Poids lourd** : transcodage multimédia, adaptation d'applications, compression
- **Poids moyen** : caches web, multicast fiable
- **Poids léger** : répartition de charge, pare-feux
- **Poids plume** : support serveurs HP, protocoles de transport

Thèse
Eric
Lemoine



Un nœud actif Tamanoir : TAN

- EE / services en Java (compilation / JIT...)
- Traitement optimisé (limite les copies)
- Format de données ANEP
- Multi services, multi threadé
- Multi protocoles (tcp,udp)
- Déploiement dynamique de services



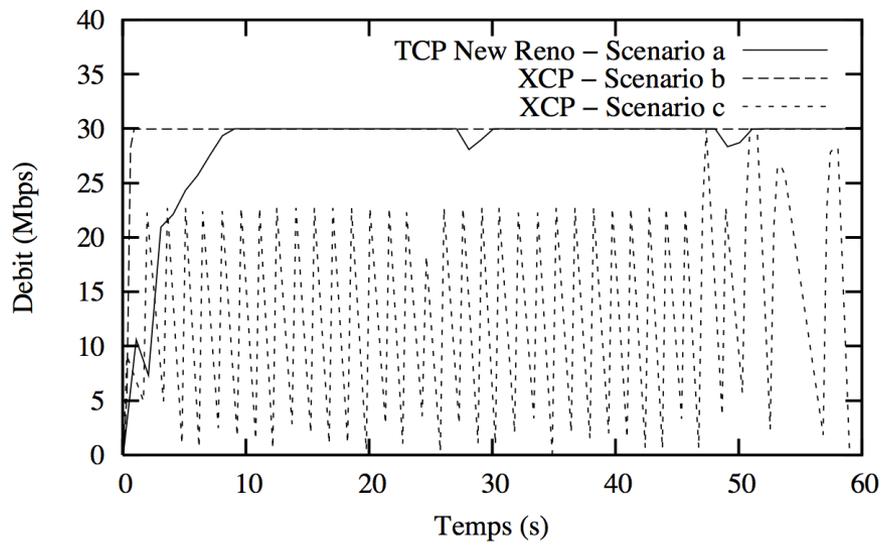
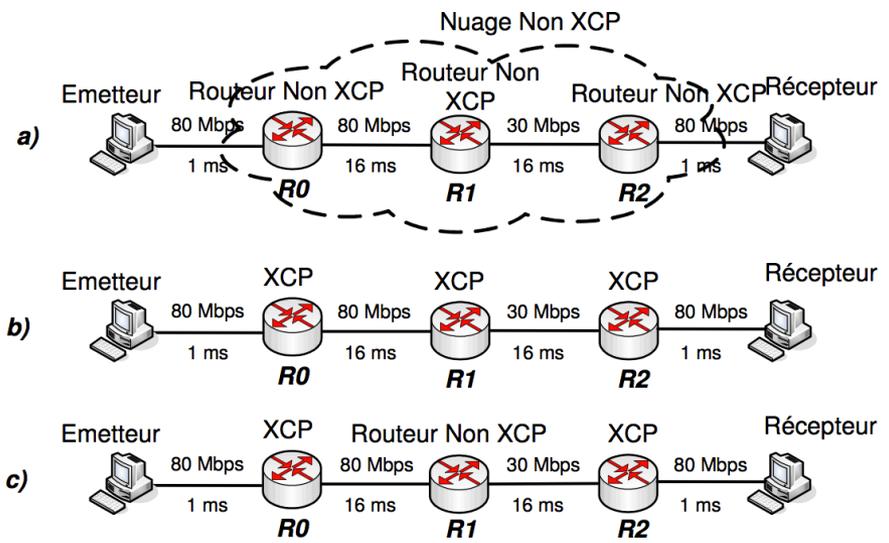
Flexibilité: brique de base pour d'autres projets

- Environnement actif Tamanoir : seul routeur logiciel actif apte à supporter des débits Gbits
- Déploiement et validation sur plate-forme locale et longue distance
- Support de différents projets (RNRT VTHD++, RNTL Etoile, RNRT Temic)
- Déployé et utilisé par partenaires académiques (LAAS, Univ. Vannes, INSA) et contexte "industriel" (3DDL)
- Flexibilité ré-utilisée dans différents équipements : haute disponibilité, pare-feux à états
- 1 focus : protocoles de transport

Thèse
Narjess
Ayari

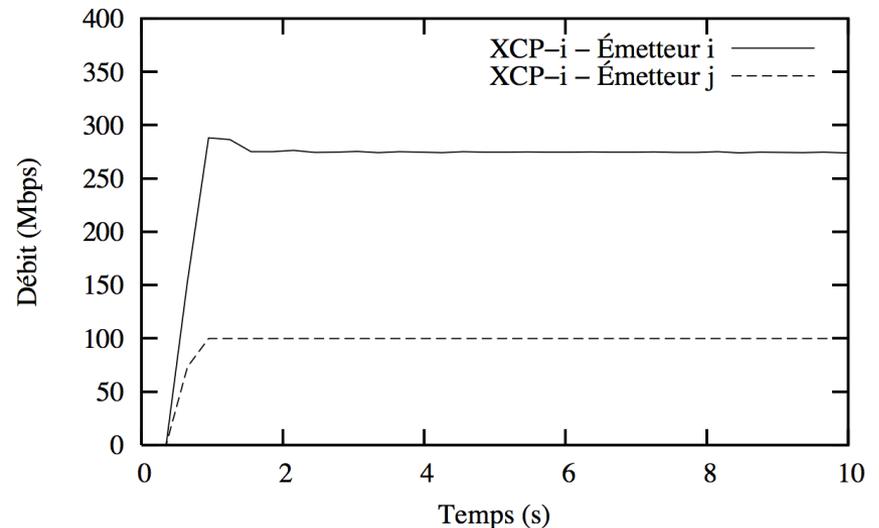
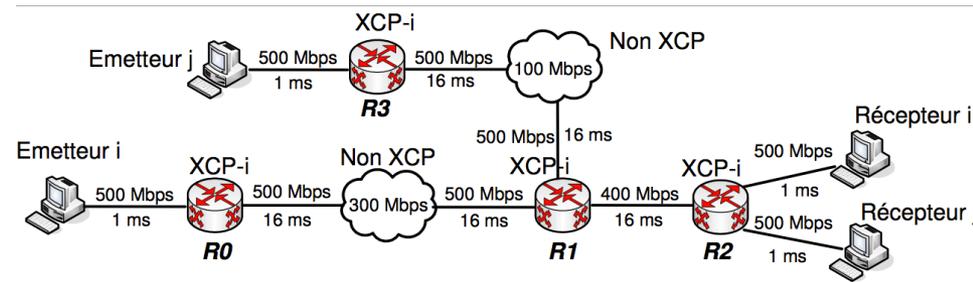
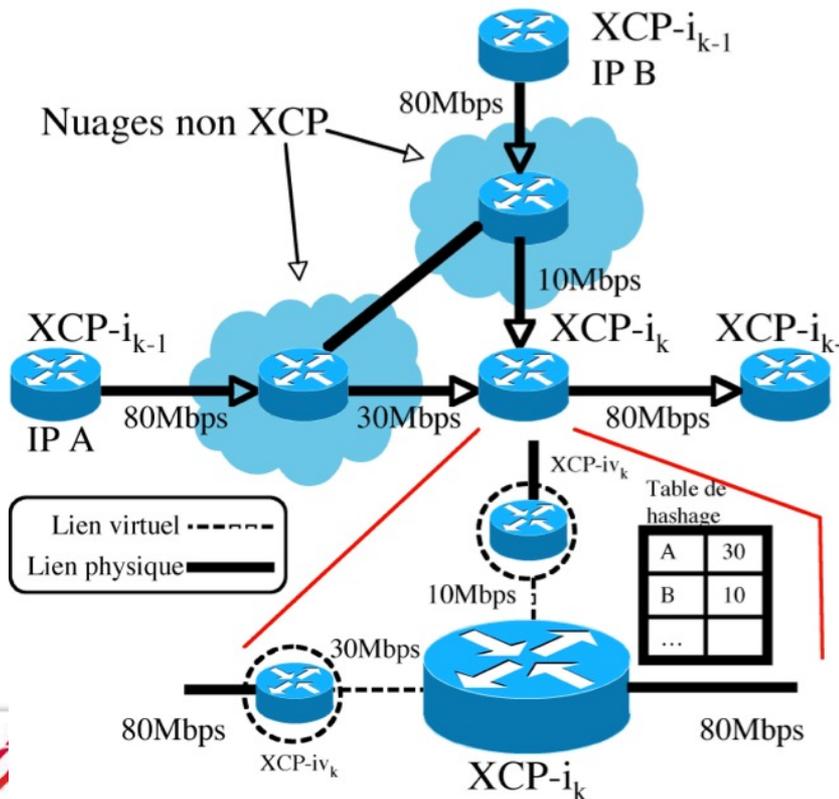
Un service interopérable : XCP-i

- TCP ne récupère pas toute la bande passante
- Proposition de XCP (eXtensible Control Protocol) (Katabi02)
- Protocole de transport à assistance de routeurs => renvoie une info sur l'état du chemin à l'émetteur
- Flexibilité légère dans les équipements réseaux



Un service interopérable : XCP-i

- Sortir ce protocoles du laboratoire?
- XCP-i : remplacer un nuage non-XCP par un routeur virtuel situé aux abords du nuage sur le chemin des données
- Gérer l'interopérabilité avec des équipements non-XCP et l'équité entre flux



Effacité énergétique dans les systèmes distribués à grande échelle

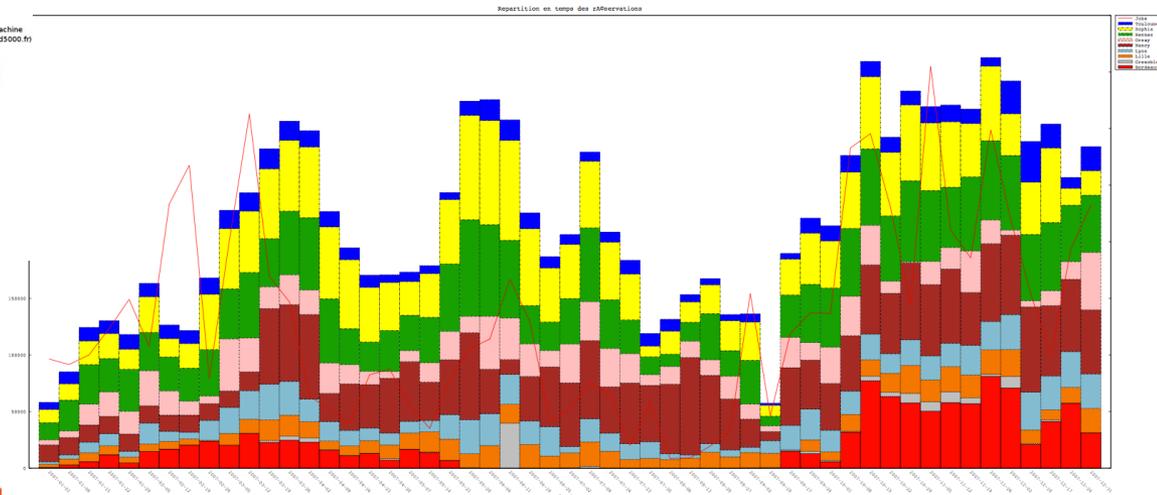
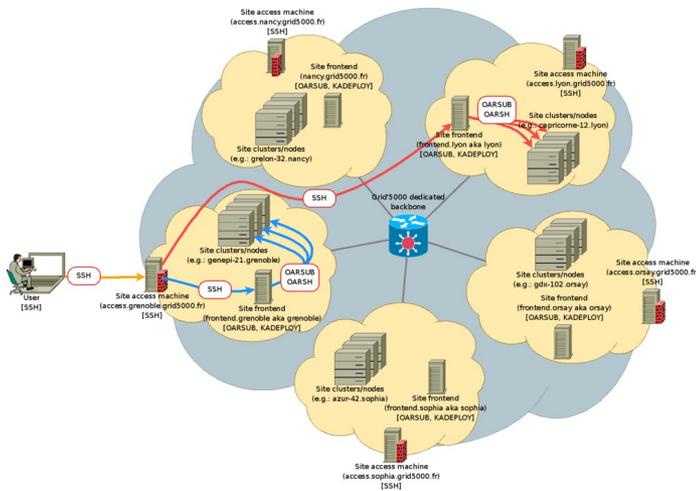
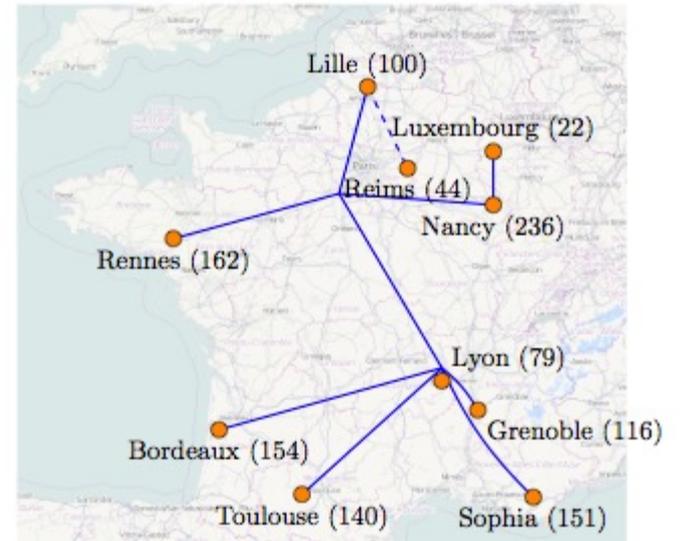
Vers des systèmes distribués efficaces en énergie

Comment diminuer la consommation énergétique sans impacter les performances ?

- Comprendre l'usage des systèmes et la consommation énergétique associée
- Etre capable de mesurer et de proposer des métriques utiles
- Concevoir des modèles et des solutions logicielles efficaces en consommation énergétique
- Proposer des composants logiciels génériques dérivables sur différents scenario (Grids, Clouds, Networks)
- Simuler et valider expérimentalement à grande échelle
- Aider les utilisateurs à faire les bons choix

Notre objet d'étude : Grid5000

- 7000 cœurs de calcul / 10 sites
- Réseau 10 Gbits
- Usage exclusif avec réservation / mode best effort
- Usage pics/creux



Mesurer et comprendre la consommation électrique des TICs

Mesurer des TICs ?

Wh, co2, joules ?

C'est quoi une tonne de co2 ?



©Dave Ames, Cohasset High School



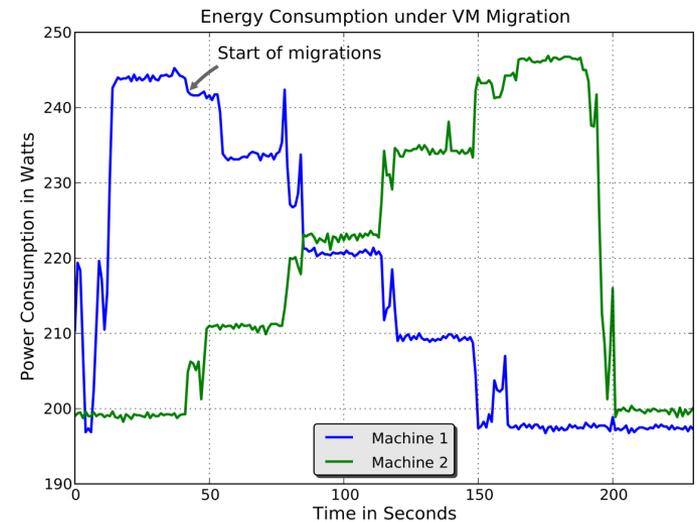
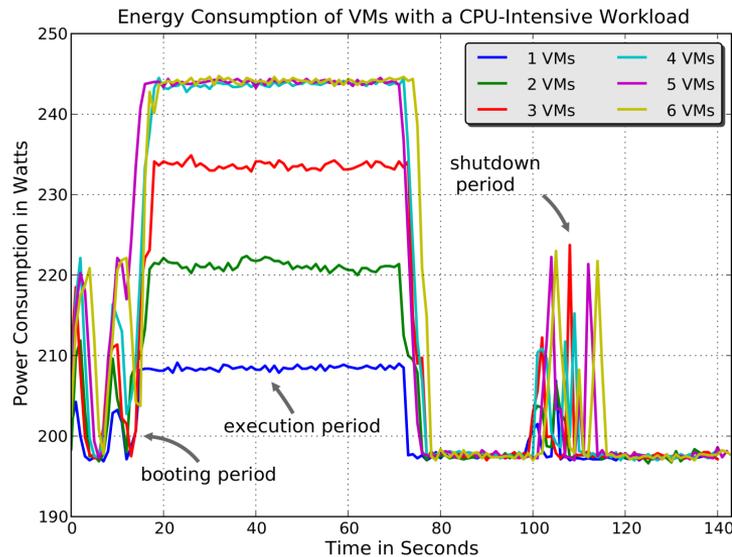
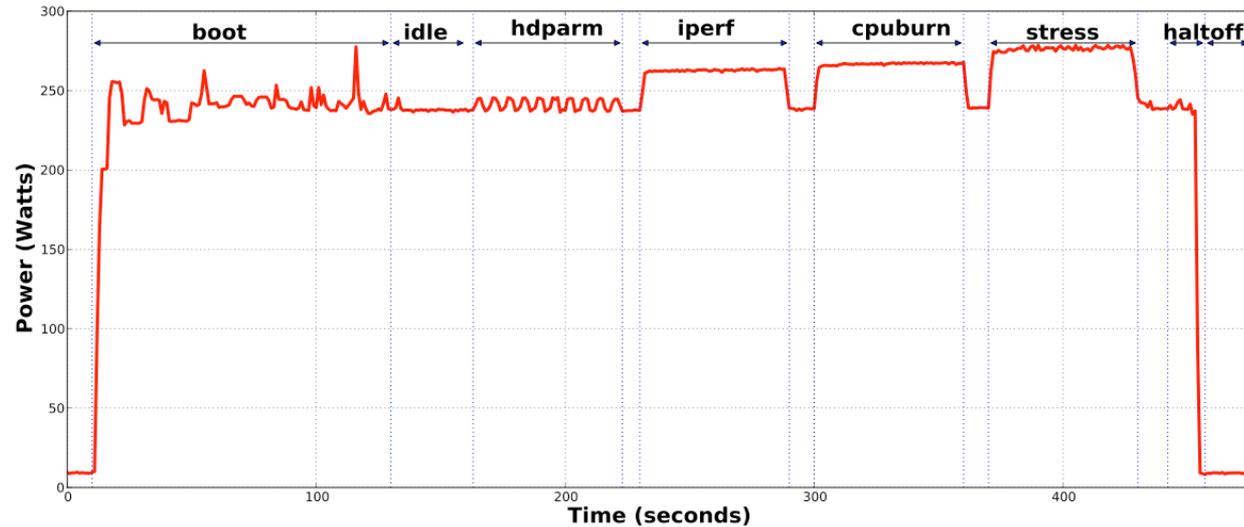
Mesurer ce que l'on comprend : des watts !



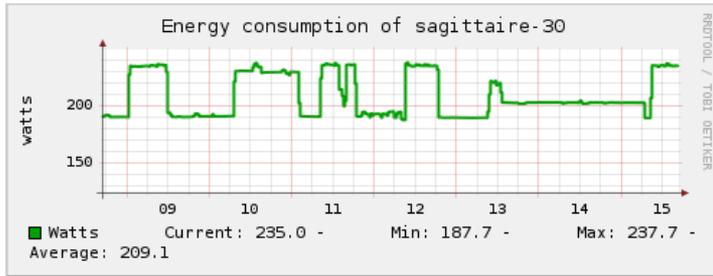
**Capteurs externes, internes, intégrés,
compteurs de performances**



Profiler des infrastructures physiques et virtuelles



Donner des infos (utiles) à l'utilisateur



Energy Information of Lyon Grid5000 site



Green-Net Demo 0.1b

File View Help

Power Consumption (Kw)

Amount Spent (in Euros)

Statistics:

- Number of measurements: 276
- Refresh Interval: 1 second
- Time Frame: 00:01:00
- Average consumption: 149.52W
- Amount Spent: 0.18 Euros

INRIA

Grid'5000

Aladdin

Green-Net

Status of Resources:

sagit-1 74.81W	sagit-11 294.94W	sagit-21 221.42W	sagit-31 163.69W	sagit-41 83.53W	sagit-51 193.71W	sagit-61 236.40W	sagit-71 64.54W	capric-2 117.51W	capric-3 117.51W	capric-4 117.51W	capric-5 117.51W	capric-6 117.51W	capric-7 117.51W	capric-8 117.51W	capric-9 117.51W	capric-10 117.51W	capric-11 117.51W	capric-12 117.51W	capric-13 117.51W	capric-14 117.51W	capric-15 186.97W	capric-16 52.98W	capric-17 71.33W	capric-18 12.22W	capric-19 12.46W	capric-20 12.46W	capric-21 12.46W	capric-22 12.46W	capric-23 12.46W	capric-24 12.46W	capric-25 12.46W	capric-26 261.25W	capric-27 261.25W	capric-28 246.45W	capric-29 246.45W	capric-30 246.45W	capric-31 246.45W	capric-32 246.45W	capric-33 246.45W	capric-34 246.45W	capric-35 246.45W	capric-36 246.45W	capric-37 246.45W	capric-38 246.45W	capric-39 246.45W	capric-40 246.45W	capric-41 246.45W	capric-42 83.97W	capric-43 130.27W	capric-44 130.27W	capric-45 130.27W	capric-46 130.27W	capric-47 130.27W	capric-48 130.27W	capric-49 130.27W	capric-50 130.27W	capric-51 130.27W	capric-52 180.02W	capric-53 226.64W	capric-54 40.37W	capric-55 41.12W	capric-56 171.48W	capric-57 171.48W	capric-58 171.48W	capric-59 171.48W	capric-60 171.48W	capric-61 171.48W	capric-62 171.48W	capric-63 171.48W	capric-64 171.48W	capric-65 171.48W	capric-66 171.48W	capric-67 171.48W	capric-68 171.48W	capric-69 171.48W	capric-70 171.48W	capric-71 171.48W	capric-72 171.48W	capric-73 171.48W	capric-74 171.48W	capric-75 171.48W	capric-76 171.48W	capric-77 171.48W	capric-78 171.48W	capric-79 171.48W	capric-80 171.48W	capric-81 171.48W	capric-82 171.48W	capric-83 171.48W	capric-84 171.48W	capric-85 171.48W	capric-86 171.48W	capric-87 171.48W	capric-88 171.48W	capric-89 171.48W	capric-90 171.48W	capric-91 171.48W	capric-92 171.48W	capric-93 171.48W	capric-94 171.48W	capric-95 171.48W	capric-96 171.48W	capric-97 171.48W	capric-98 171.48W	capric-99 171.48W	capric-100 171.48W
-------------------	---------------------	---------------------	---------------------	--------------------	---------------------	---------------------	--------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	---------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	---------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	---------------------	---------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	----------------------	-----------------------

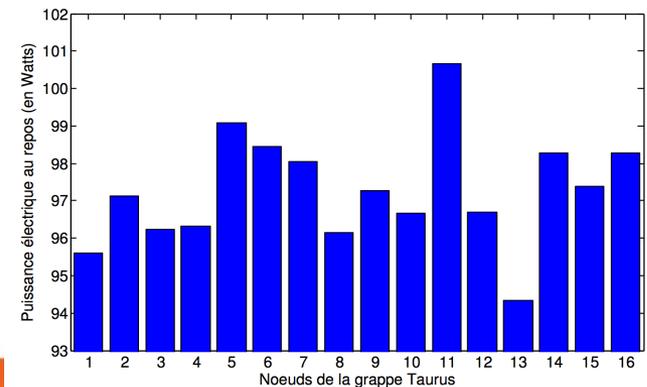
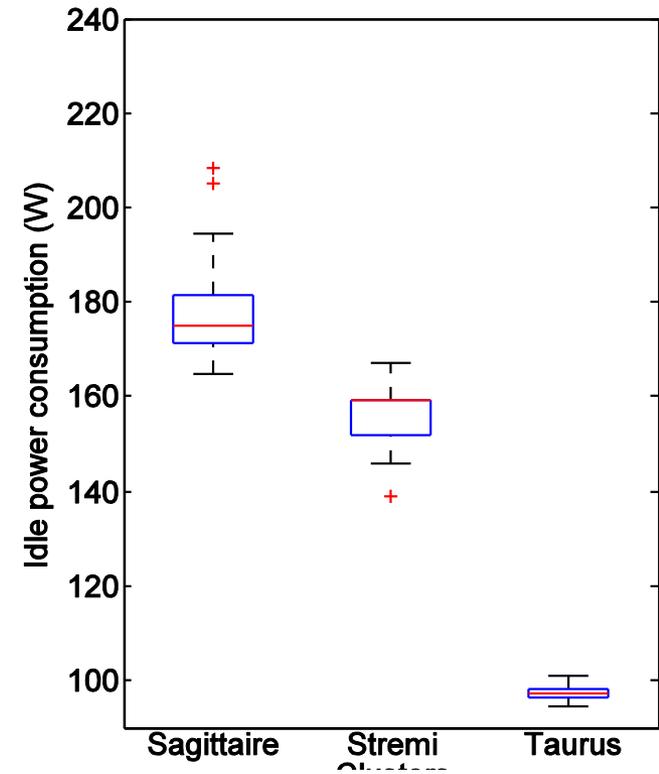
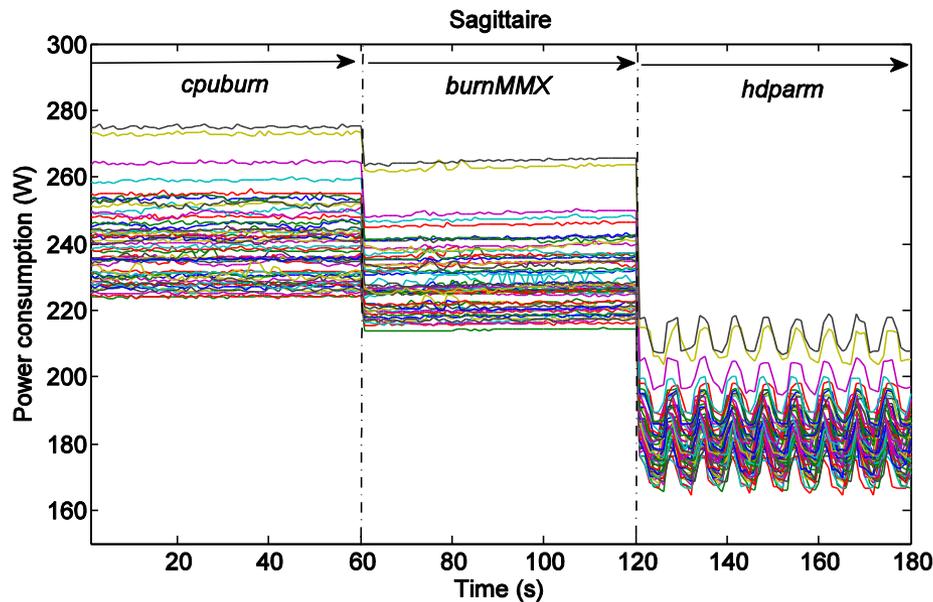
Number of Nodes

Idle Off On

Resource on Resource idle Resource off Resource monitored

Chasse aux mythes : l'homogénéité énergétique

- Mêmes Flops mais flops par watt différents
- Coût Idle
- CPU : responsable principal
- A prendre en compte pour faire des ordonnanceurs verts !

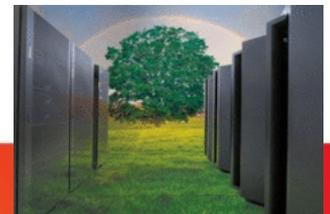


Systemes à base de réservation et nuages verts

Favoriser des systèmes à base de réservations

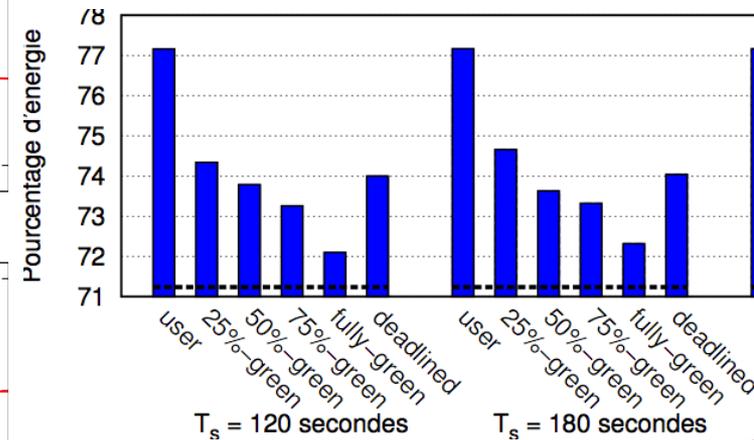
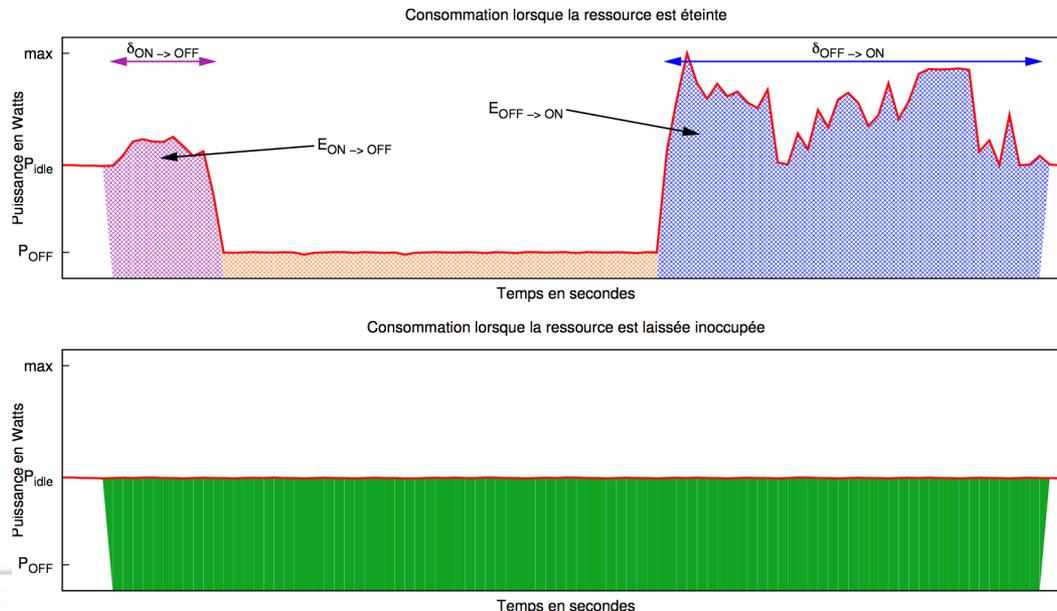
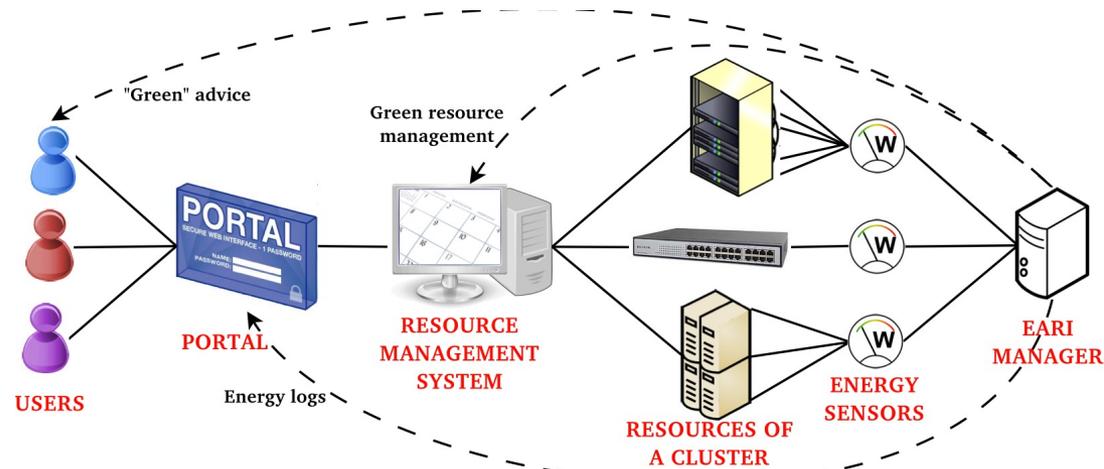
Chaque usage est basé sur une réservation (ressources, durée, temps limite) :

- Réserver des CPUs en HPC
- Réserver des machines physiques ou virtuelles dans les Clouds
- Réserver de la bande passante dans le transport de données
- Leviers:
 - Trouver et alimenter le nombre optimal de ressources :
 - HPC et Grilles : allumer / éteindre des ressources physiques (serveurs)
 - Clouds : allumer/éteindre des VMs
 - Réseaux : allumer/éteindre des ports réseaux, liens, routeurs...
 - Adapter la vitesse et la consommation aux besoins des applications, services et des utilisateurs :
 - HPC: dvfs
 - Clouds : tuning, capping
 - Réseaux : adaptation de bande passante (ALR)
- Favoriser des agrégations : dans le temps, dans l'espace



ERIDIS/EARI

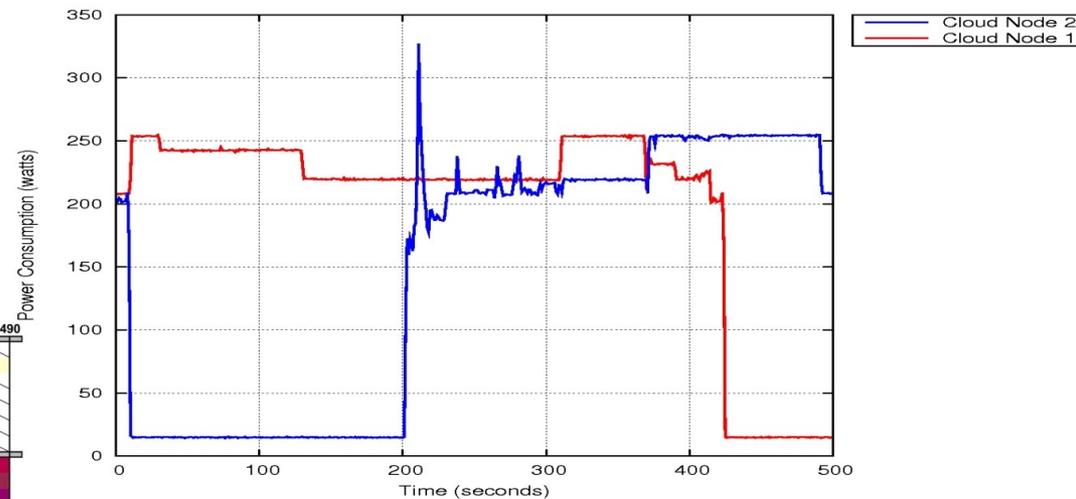
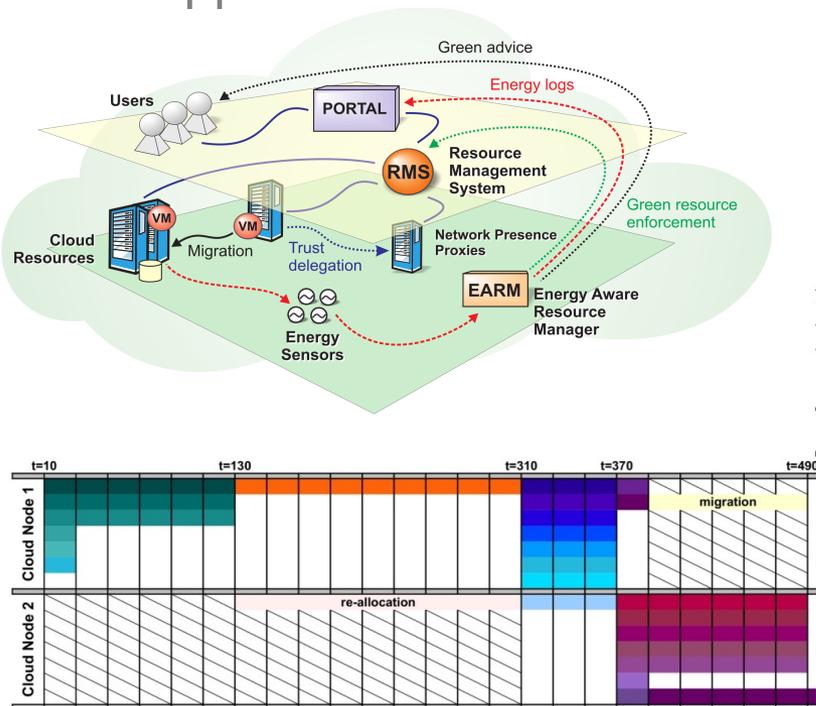
- Extinction de ressources inutilisées
- Prédiction de réservations
- Placement de tâches et de réservations



Nuage Vert / Green Cloud



- Mesurer l'usage énergétique des infrastructures de Cloud
- Prendre en compte la consommation pour proposer des composants logiciels économes
- Proposition de l'architecture Green Open Cloud
- Aider les utilisateurs à exprimer leurs contraintes Green
- Aider les gestionnaires à intégrer les coûts environnementaux
- Supporter de nouveaux leviers : migration, *tuning*, *capping*



Effacité énergétique dans les infrastructures HPC : avec ou sans connaissance des services et des applications

Vers des infrastructures HPC Exascale

Exascale à l'horizon 2020

1^{er} Top500: Tianeh2 : 3M cœurs, 30-50 Pflops, 17 MW

Eviter le mur des 100 MW → 20 MW

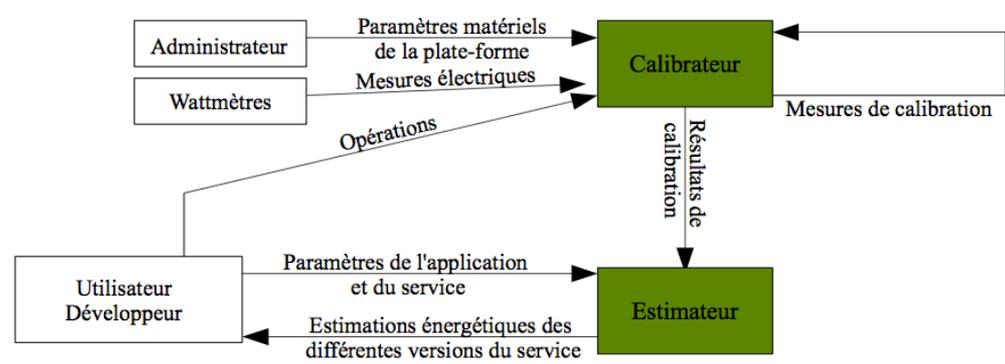
Green500 : meilleurs : 3 Gflops par watt (Tianhe2 1.9Gf/W)

Pour exaflops : 50 Gflops /W

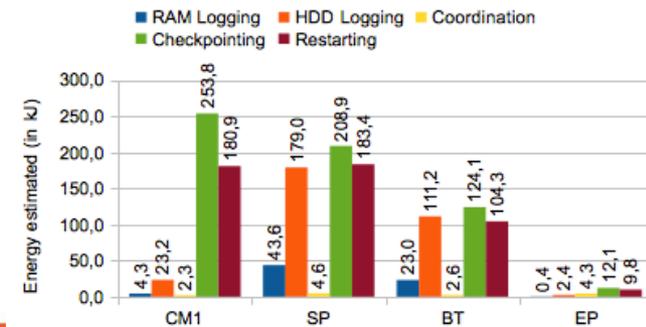
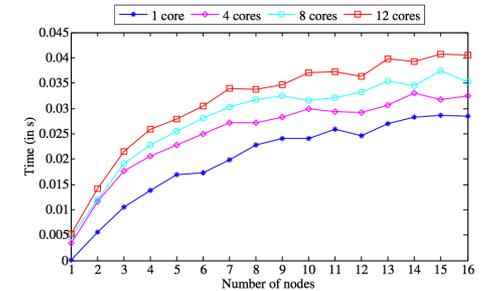
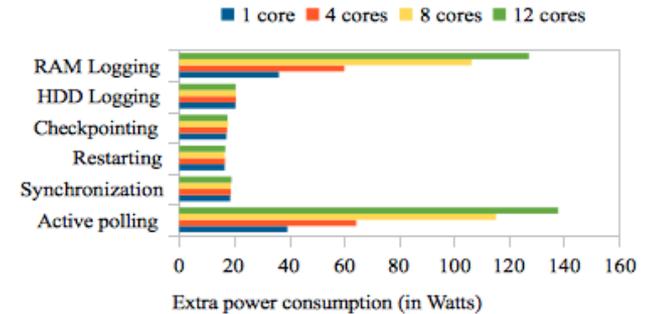
2 approches matérielles : accélérateurs (GPUs) et LPP
(low power processors) (ex: projet MontBlanc)

Focus sur applications et services pour exascale

EE HPC avec connaissance des services



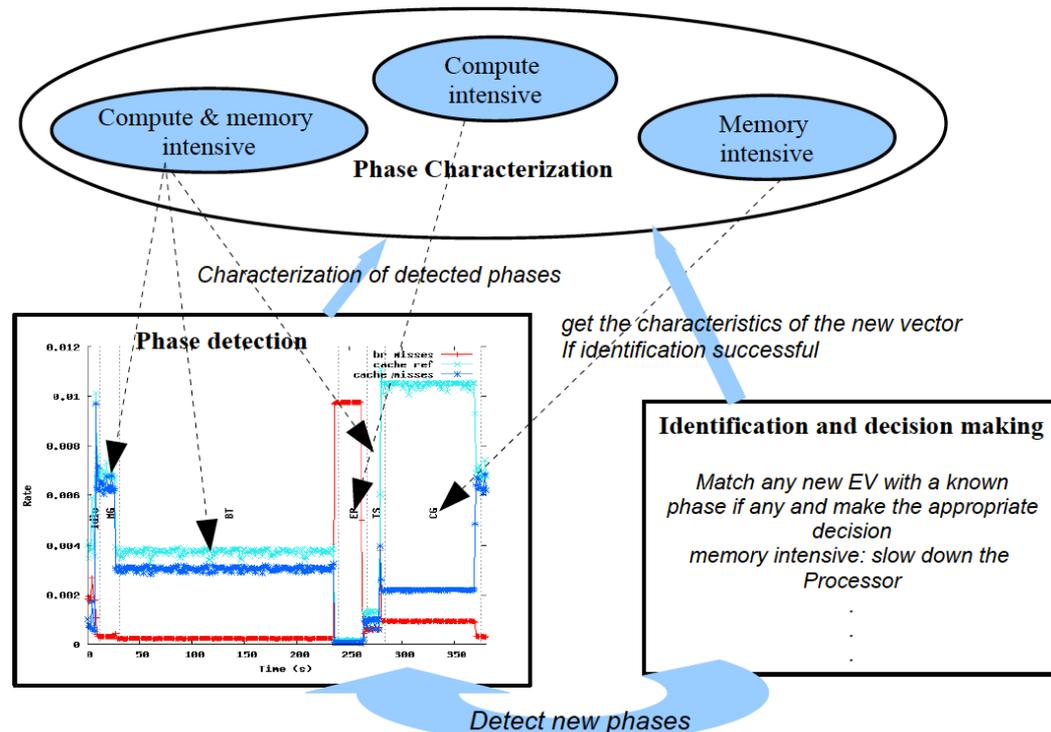
- Exemple services exascale : résilience et diffusion
- 4 étapes
- **Découpage** du service en opérations
- **Calibration** de la consommation énergétique des opérations
- **Estimation** de la consommation électrique
- **Aider** utilisateur à choisir le bon service en fonction des contraintes des applications



Et si le logiciel est trop complexe ?



- Applications ont une utilisation non régulière des ressources
- Analyse ADN live de l'exécution des applications et services
- Détection des phases
- Caractérisation des applications
- Applications de leviers verts



Phase label	Possible reconfiguration decisions
compute intensive	switch off memory banks; send disks to sleep; scale the processor up; put NICs into LPI mode
memory intensive	scale the processor down; decrease disks or send them to sleep; switch on memory banks
mixed	switch on memory banks; scale the processor up send disks to sleep; put NICs into LPI mode
communication intensive	switch off memory banks; scale the processor down switch on disks
I/O intensive	switch on memory banks; scale the processor down; increase disks, increase disks (if needed)

Conclusion et Perspectives

Conclusions

- Œuvre collective : merci !
- Modèles, algorithmes, simulations, validations expérimentales
- Proposition logicielles

- Flexibilité
 - Participation aux 4 vagues de flexibilité des réseaux
 - Proposition de plates-formes expérimentales

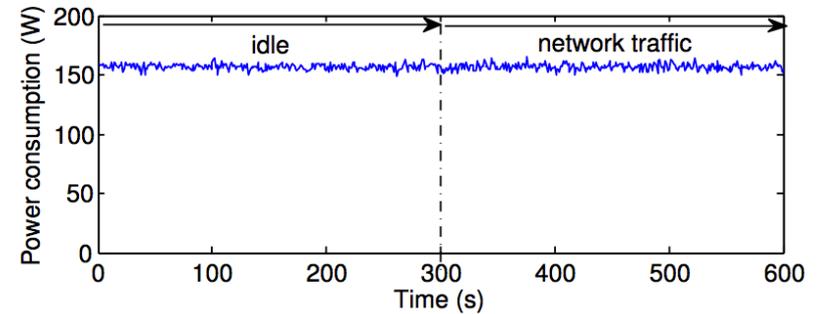
- Efficacité Énergétique
 - Modéliser et mesurer la composante énergétique
 - Proposer des systèmes logiciels d'adaptation et éco-efficaces
 - Animation de la communauté GreenIT

Perspectives

- Contribuer au facteur 1000 dans les réseaux
- Consommer mieux dans les infrastructures distribuées à grande échelle
- Proportionnalité énergétique
- Vers le développement durable

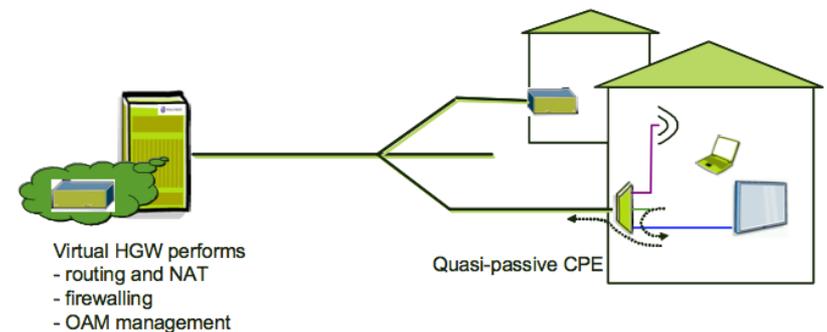
Réseaux flexibles et efficaces en consommation énergétique

- Réseaux : utilisation énergétique plate – non proportionnelle à l'usage
- GreenTouch : réduire la consommation électrique des réseaux d'un facteur 1000 à l'horizon 2015 tout en supportant l'explosion du trafic et la QoS
- Contributions à tous les niveaux: combinaison de matériel et logiciel – sans fil, optique, cœur et extrémités
- *Focus sur les équipements terminaux (10W*box) – virtualisation des services*
- *Virtual Home Gateway*



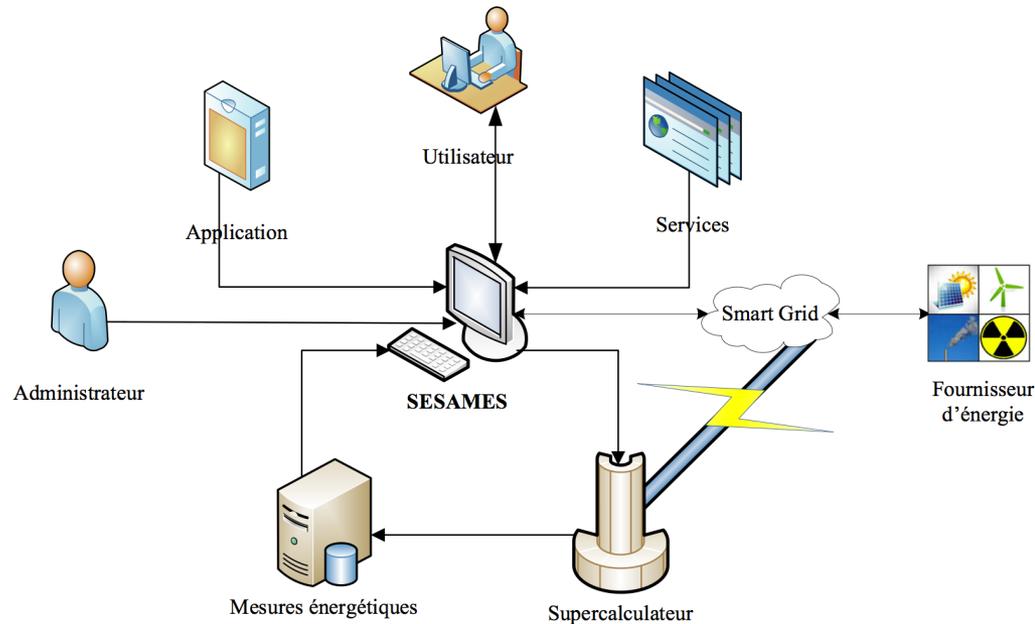
(A) 2015-2020 NETWORK FORECAST: DEVICE DENSITY AND ENERGY REQUIREMENTS IN THE BUSINESS-AS-USUAL CASE (BAU). EXAMPLE BASED ON THE ITALIAN NETWORK.

	<i>power consumption</i> [W]	<i>number of devices</i> [#]	<i>overall consumption</i> [GWh/year]
<i>Home</i>	10	17,500,000	1,533
<i>Access</i>	1,280	27,344	307
<i>Metro/Transport</i>	6,000	1,750	92
<i>Core</i>	10,000	175	15
<i>Overall network consumption</i>			<i>1,947</i>



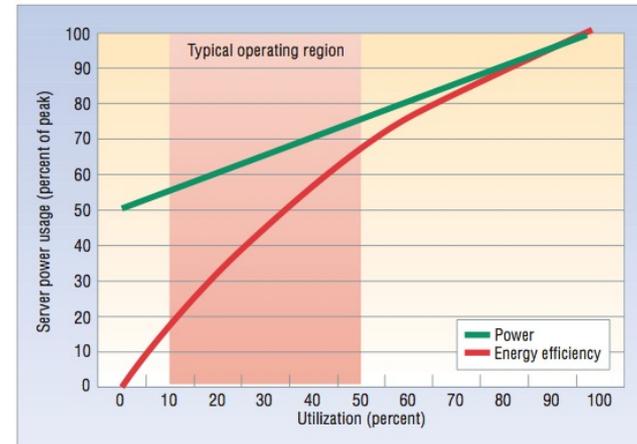
Consommer mieux

- DataCentre : un client « électrique » à part
- Premières expériences avec SESAMES
- Favoriser le dialogue entre fournisseur électrique, utilisateur, admin
- Prise en compte d'autres contraintes/métriques : production d'énergie, disponibilité, cout financier, pollution...

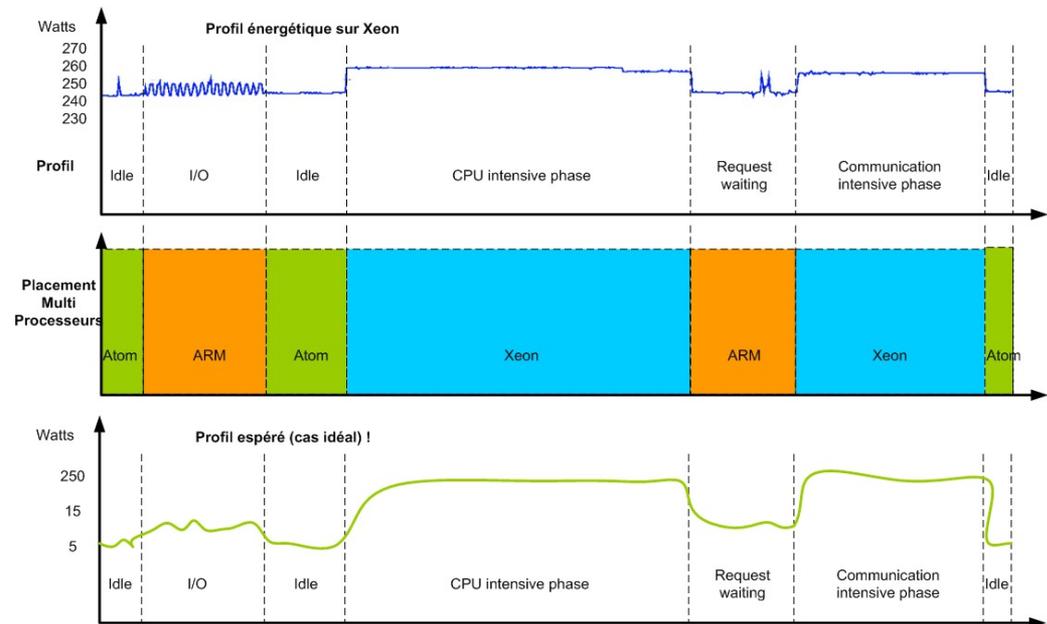
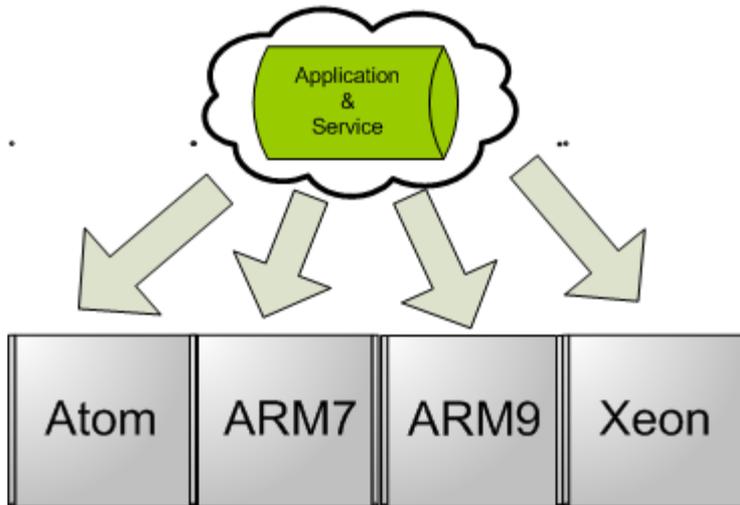


A la recherche de la proportionnalité énergétique

- Exploiter les périodes creuses dans la vie du système
- Combiner des architectures CPU hétérogènes adaptées aux besoins des applications

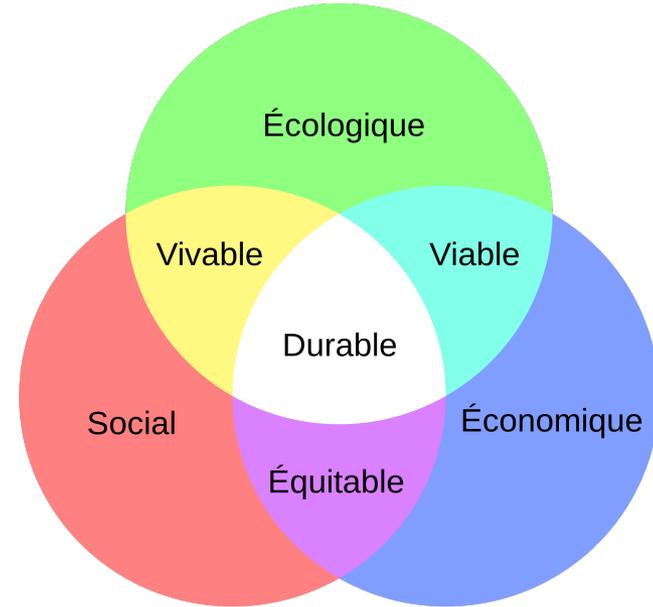


Multi-platform container



IncurSION dans le développement durable

- Supporter un autre usage des systèmes distribués à grande échelle
 - Vers la fin des ressources infinies toujours disponibles ?
- Prendre en compte le cycle de vie des systèmes pour des modèles plus réalistes : production, transport, usage, recyclage



The screenshot shows the EcoInfo website with the following content:

- Header:** EcoInfo logo and the question 'Que faites-vous de vos e-déchets ?'.
- Logos:** SFR, Orange, Inria, and others.
- Main Text:** 'Déchets d'équipements électriques et électroniques : Enquête nationale sur les pratiques actuelles, Bonnes pratiques'. It mentions 'Enquête: 100 réponses dont 41 qui ont permis de recueillir 23 types de pratiques effectuées et de décrire pas de 500. Le mettre en œuvre pour qu'il soit le meilleur possible'.
- Map:** A map of France showing collection points for e-waste.
- Text:** 'Peut-on faire mieux ?' and 'Prix payés par tonne en 2011-2012: De 100€ à 1500€/tonne en moyenne. 400€ pour ceux qui paient le moins cher'.
- Logos:** SFR, Orange, Inria, and others at the bottom.



Remerciements

- Jean-Patrick Gelas
- Olivier Gluck
- Christian Perez
- Eddy Caron
- Frédéric Desprez
- Gilles Fedak
- Paulo Goncalves
- Thomas Begin
- Isabelle Guérin Lassous
- Bernard Tourancheau
- Cong-Duc Pham
- Marcos Dias de Asuncao
- Christina Herzog
- Teferi Assefa
- Mehdi Diouri
- Ghislain Landry Tsafack Chetsa
- Anne-Cécile Orgerie
- Narjess Ayari
- Dino Lopez Pacheco
- Eric Lemoine
- Pablo Neira Ayuso
- Edgar Magana
- Jean-Christophe Mignot
- Francois Rossigneux
- Julien Carpentier
- Maxime Morel
- Abderhaman Cheniour
- Olivier Mornard
- Augustin Ragon
- Pierre Bozonnet
- Martine Chaudier
- Saad ElHadri
- Roland Westrelin
- Roya Golchay
- Sidali Guebli
- Alice Bonhomme
- Walid el Dahabi
- Pablo Pazos
- Chien Jon Soon
- Grégoire Locqueneux
- Aweni Saroukou