

PARConnector : sans détour du Langage R au Cloud Big Data

Application à l'analyse Metagénomique

Ndeye Aram GAYE - Laurent PELLEGRINO



Plan




- Présentation des sociétés
- ProActive Workflows and Scheduling
- Package “PARConnector”
- Cas d'utilisation: analyse métagénomique quantitative

Présentation des sociétés

MetaGenoPolis



Financement des Investissements d'Avenir  19M€ pour 2012-2019

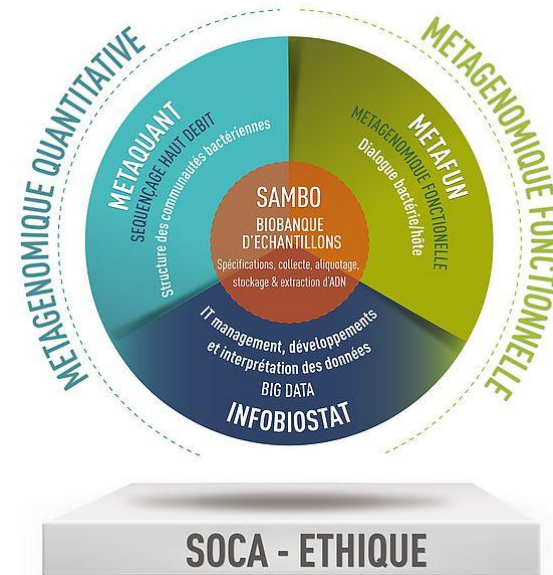
Centre d'excellence en métagénomique humaine : www.mgps.eu

4 plateformes

1 centre de recherche en éthique

Objectifs :

- Explorer le rôle du microbiome intestinal dans la santé et les maladies complexes humaines
- Identification de biomarqueurs et d'outils de diagnostic et de prédiction dans le cadre d'une médecine préventive.
- Ouvrir de nouvelles pistes thérapeutiques centrées sur le microbiome
- Transfert de ces technologies dans le domaine de la nutrition



- Essaimage de l'INRIA, créée en octobre 2007
- 15 collaborateurs (ingénieurs, docteurs, commerciaux)
- Éditeur de logiciel libre, membre de la communauté OW2
- Des solutions applicables à de nombreux domaines :
Ingénierie, Bio Technologies, Finance, Technologies de l'Information, etc.
- Une croissance en France et à l'international

Nos clients



Technologies de l'Information



Ingénierie Energie Aéronautique Espace



Bio Technologies



Medias Distribution



Finance



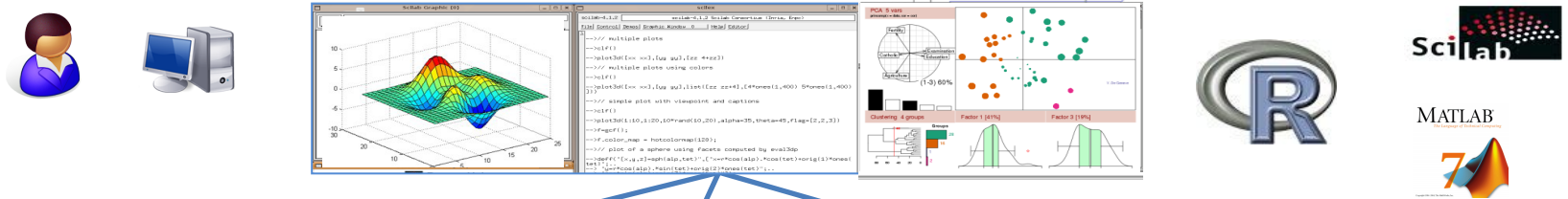
Notre mission



- L'innovation au service de nos clients
- Anticiper les défis technologiques
- Renforcer notre expertise sur des domaines clés
 - Systèmes Distribués et le Cloud Computing
- Être le partenaire privilégié de nos clients
- Être toujours à la pointe de l'innovation grâce à un investissement en R&D important

ProActive Workflows & Scheduling

Workflows & Scheduling (1)



Static Policy
Clusters



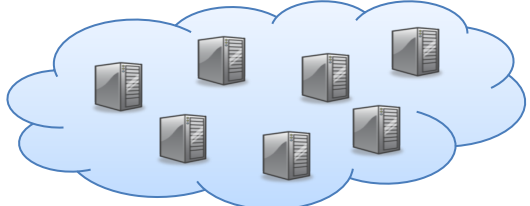
Dedicated resources

Timing Policy
12/24
Desktops



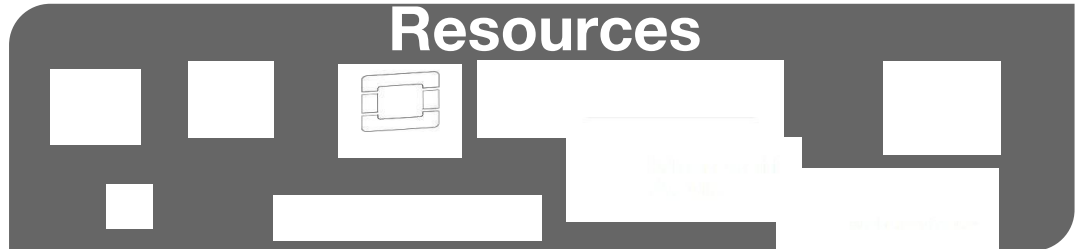
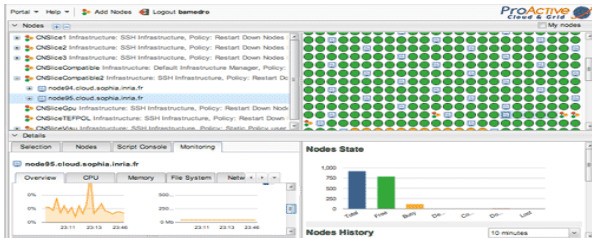
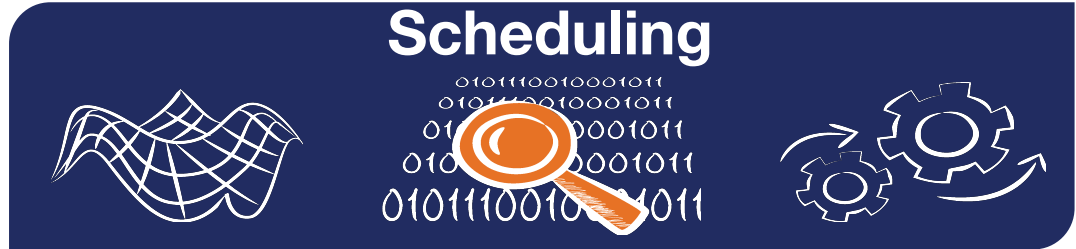
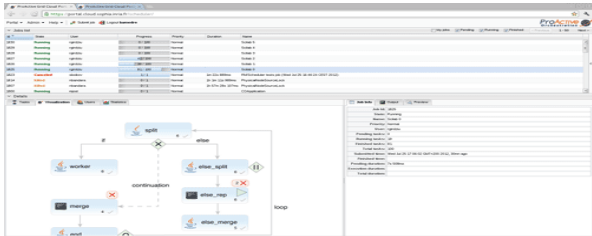
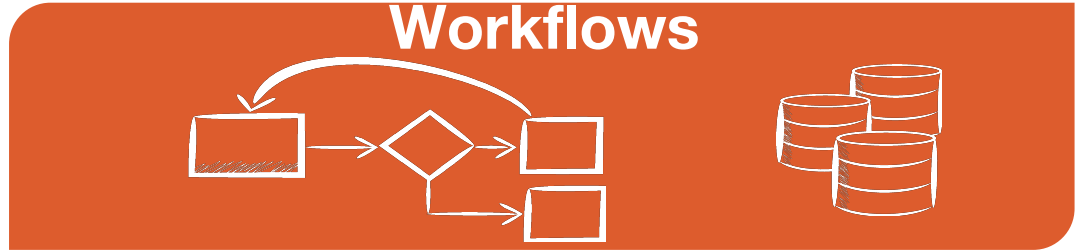
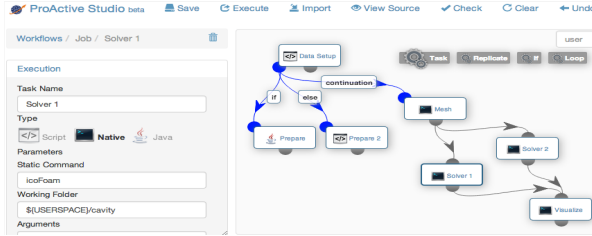
Desktops

Dynamic Workload Policy
EC2, Azure, HP Cloud, ...



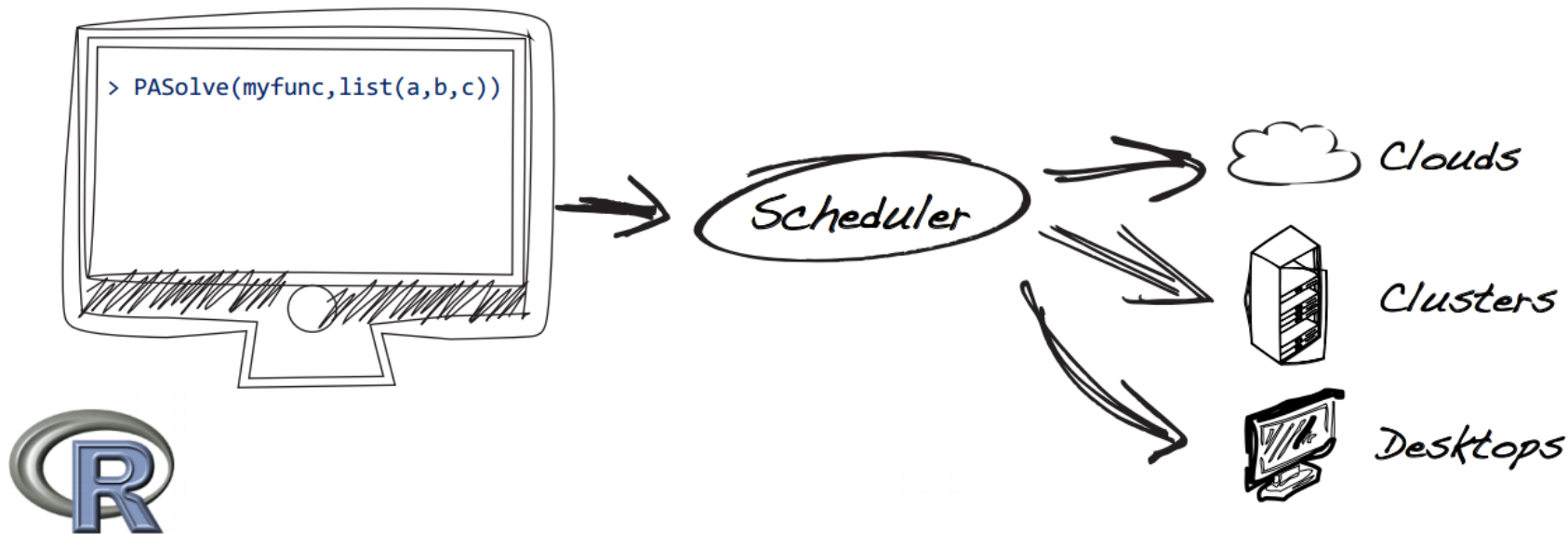
Amazon EC2

Workflows & Scheduling (2)



Package “PARConnector”

PAR Connector



**Intégration transparente
dans votre environnement scientifique**

Fonctions ProActive R



CONNEXION

PAConnect

TRANSFERTS

PAPushFile

PAPullFile

SOUSSION

PASolve

PA

PAM(erge)

PAS(plit)

PAWaitAny

PAWaitFor

MONITORING

PAState

PAJobState

PADebug

Fonctions ProActive R



CONNEXION

PAConnect

TRANSFERTS

PAPushFile

PAPullFile

SOUSSION

PA Solve

PA

PAM(erge)

PAS(plit)

PAWaitAny

PAWaitFor

MONITORING

PAState

PAJobState

PADebug

Soumission d'un calcul



- La soumission d'un calcul s'effectue à l'aide de `PASolve`
- `PASolve` permet des invocations paramétriques
 - exécuter une même fonction avec différentes valeurs en entrée
 - chaque exécution est une tâche indépendante
- `PASolve` a une syntaxe similaire à *map*

```
> res = PASolve(func, arg1, arg2, ..., argn)
```

`func` peut être:

- une fonction définie par fonction (x, y, ...)
- le nom d'une fonction

L'évaluation de `PASolve` est asynchrone

Exemple PASolve



- Invocation paramétrique simple :

```
> job = PASolve('cos',1:5)
```

- Crée 1 job et 5 tâches
- Exécute $\cos(1)$, $\cos(2)$, ... , $\cos(5)$ sur 5 noeuds différents
- Le nombre de tâches est défini par la taille max des paramètres
- Les appels à PASolve ne bloquent pas la session R
- `job` est un élément de substitution qui informe sur l'avancement des tâches soumises à l'exécution
 - `job` est mis à jour dynamiquement au fur et à mesure que les résultats sont reçus

Récupérer des résultats



- `PAWaitFor`: permet d'attendre la fin de l'exécution des calculs distants et de récupérer les résultats

```
> res = PAWaitFor(job, timeout=1000)
```

- Arguments optionnels :
 - `timeout` définit la période maximum d'attente
 - `callback` fonction appelée lorsque les résultats sont disponibles

Transférer des fichiers



- Directement avec PAsolve:
 - Arguments optionnels
 - `input.files`: liste de fichiers à transférer de la machine locale vers les noeuds distants
 - `output.files`: liste de fichiers à transférer des noeuds distants vers la machine locale
- En utilisant des fonctions dédiées:
 - `PAPushFile`
 - `PAPullFile`

Cas d'utilisation: analyse métagénomique quantitative

Pipeline de traitement



Prélèvement de l'échantillon

Séquençage de l'ADN

Alignement sur référence

Matrice de comptage

Analyses bio-informatiques, statistiques}

Échantillons de fèces

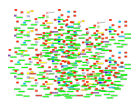
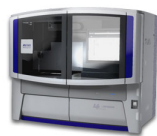
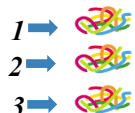
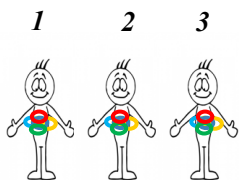
Extraction d'ADN

Séquençage de l'ADN

Courtes séquences
30-50 millions

Alignement des séquences sur un catalogue de gènes

Comptage des gènes



Catalogue de gènes de référence

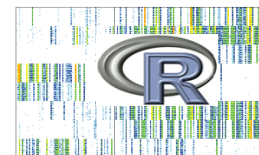
Génomés existants

individuals

		Individuals						
Item		Ind 1	Ind 2	Ind 3	Ind 4	Ind 5	Ind 6	Ind 7
1	0	36	2	0	43	106	1250	
2	0	27	193	0	44	103	8	
3	0	31	0	0	0	0	0	
4	152	59	282	1	0	0	0	
5	115	0	0	1	0	29	2	
6	90	783	26	0	2	0	0	
7	104	1616	0	0	0	0	5	
8	0	82	0	0	0	0	0	
9	2	0	0	0	0	0	0	
10	23	239	1302	10	0	190	0	
11	30	183	900	13	0	172	0	
12	27	228	1120	6	0	324	0	
13	103	0	0	0	0	0	0	
14	0	30	269	0	0	0	0	
15	0	0	0	0	0	95	0	
16	1250	6002	468	607	492	141	8023	
17	0	0	0	0	0	0	0	
18	0	9	108	0	0	55	0	
19	0	0	0	3	0	0	0	
3300000	0	36	2	0	43	106	1250	



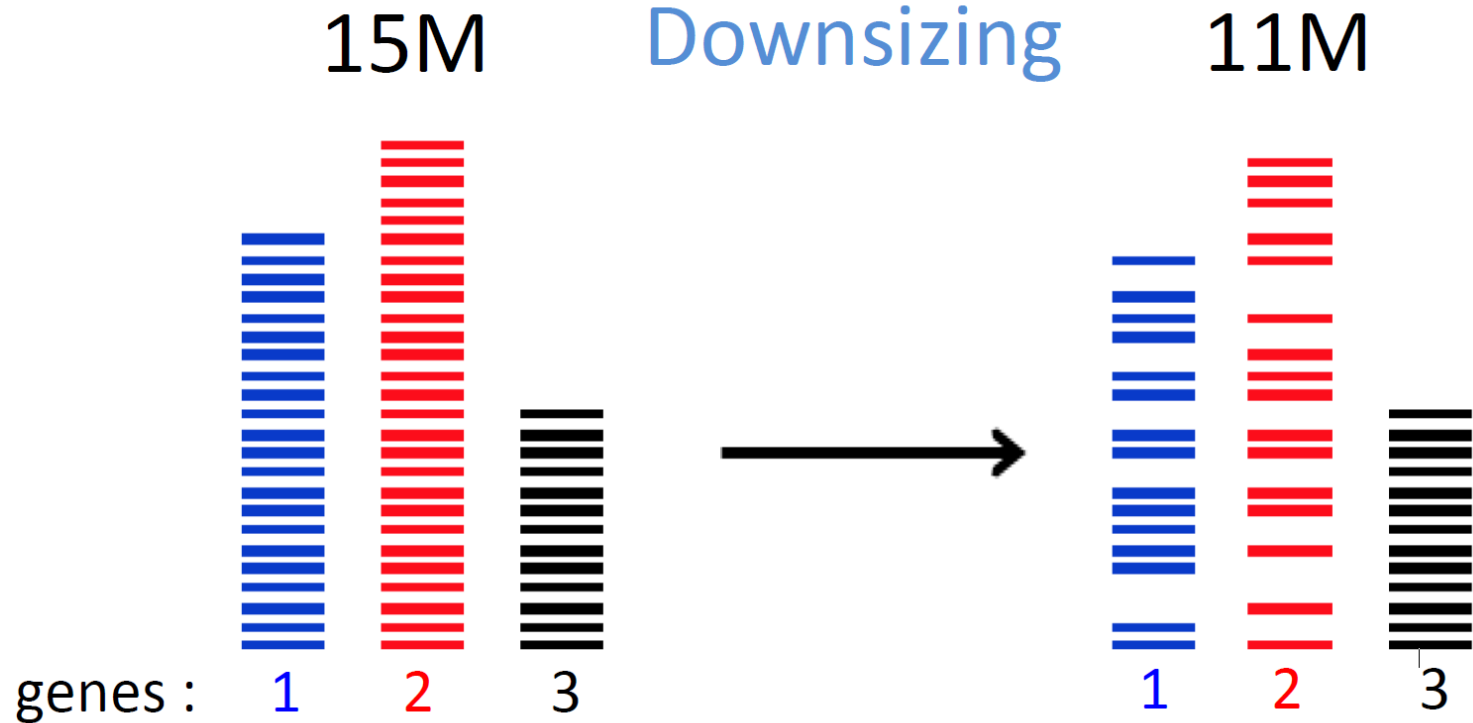
relation with clinical data



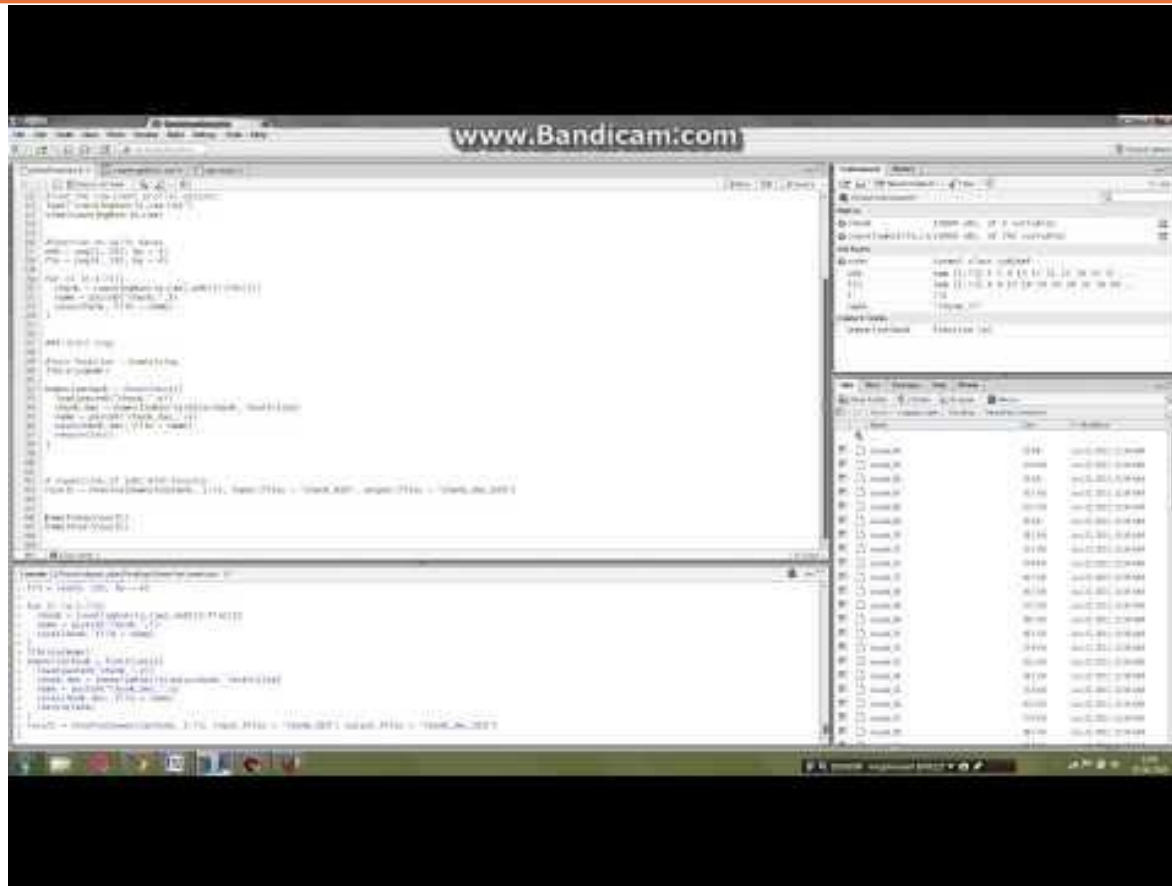
Identify clinically relevant groups

* phases HPC

Démonstration : downsizing matrix



Démo: Downsizing d'une matrice



Conclusion



PAR Connector permet de:

- Distribuer des calculs « facilement » (`PASolve`)
- Rester sur l'environnement R
- Traiter des données de type « Big Data »

Remerciements



MetaGenoPolis

Nicolas Pons

Edi Prifti

Emmanuelle Le Chatelier

Magali Berland

Anne-Sophie Alvarez

Pierre Léonard

Amine Ghozlane

Dusko Ehrlich

Ndeye Aram Gaye

Jean-Michel Batto

ActiveEon

Brian Amedro

Youri Bonnaffé

Denis Caromel

Laurent Pellegrino

Fabien Viale



Thanks for your attention!



Activeeon
SCALE BEYOND LIMITS



metagenopolis
mgps.eu