



# *BigData accelerated computing in R: an application in Metagenomics*

**Ndeye Aram GAYE**  
***[ndeeye-aram.gaye@jouy.inra.fr](mailto:ndeeye-aram.gaye@jouy.inra.fr)***

# What is BIG DATA ?

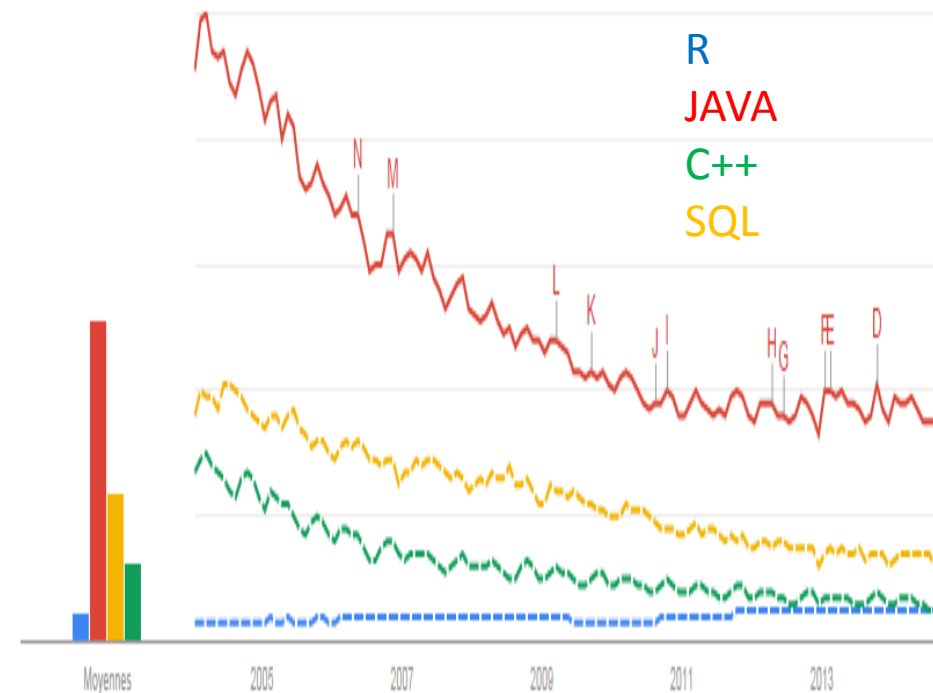
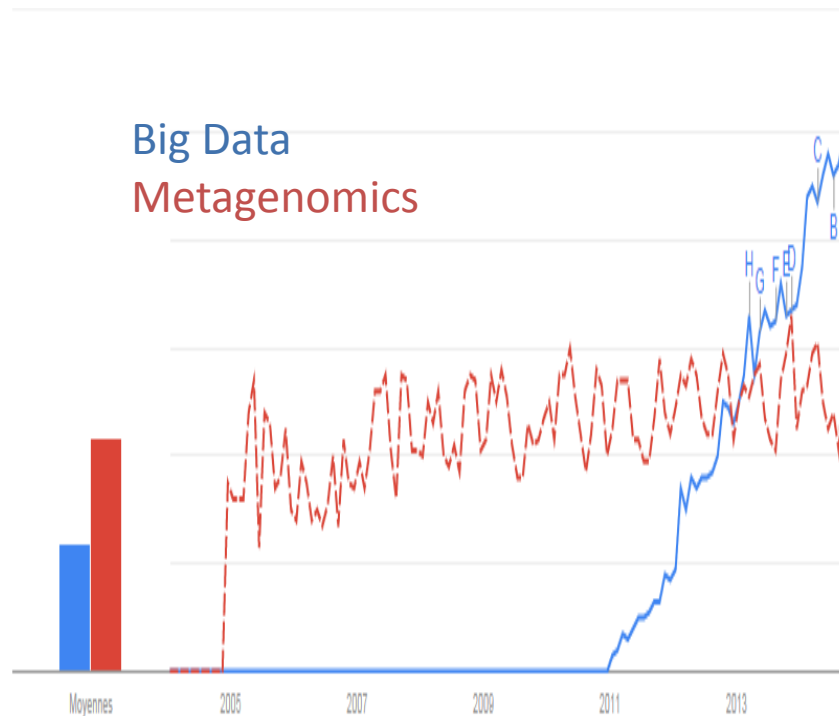
➤ Ambiguous and vague

- **Volume** : data intensity
- **Velocity** : data fast produced and the speed
- **Variety** : the degree of diversity
- **Veracity** : quality and provenance of data
- **Complexity** : data management become complex



➤ What are applications domains ?

# BIG DATA: evolution

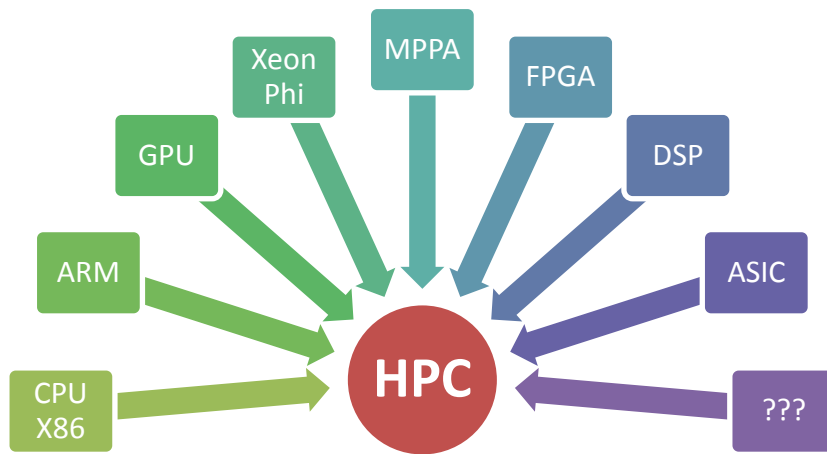


Source: Google Trends

Mckinsey prediction : big data the next frontier for innovation , May 2011

## How to process Big Data ?

### Hardware solutions



### Software solutions



- What is the complexity of these solutions ?
- How introduce such solutions to final users ?



- Language and scientific computing environment
- Free and open source software
- Large community of users
- Interpreted language:

Pros: easy to use, high level language,

Cons: slow, no suitable for big data analysis

# **An application domain: Metagenomics**

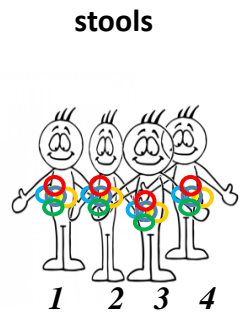
Biology

DNA  
preparation

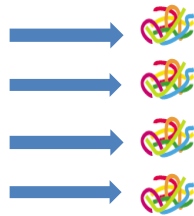
Get  
Sequences

Compare to  
reference

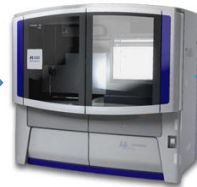
Count



DNA preparation



NGS



Short reads  
30-50M

mapping



Discovery

clustering

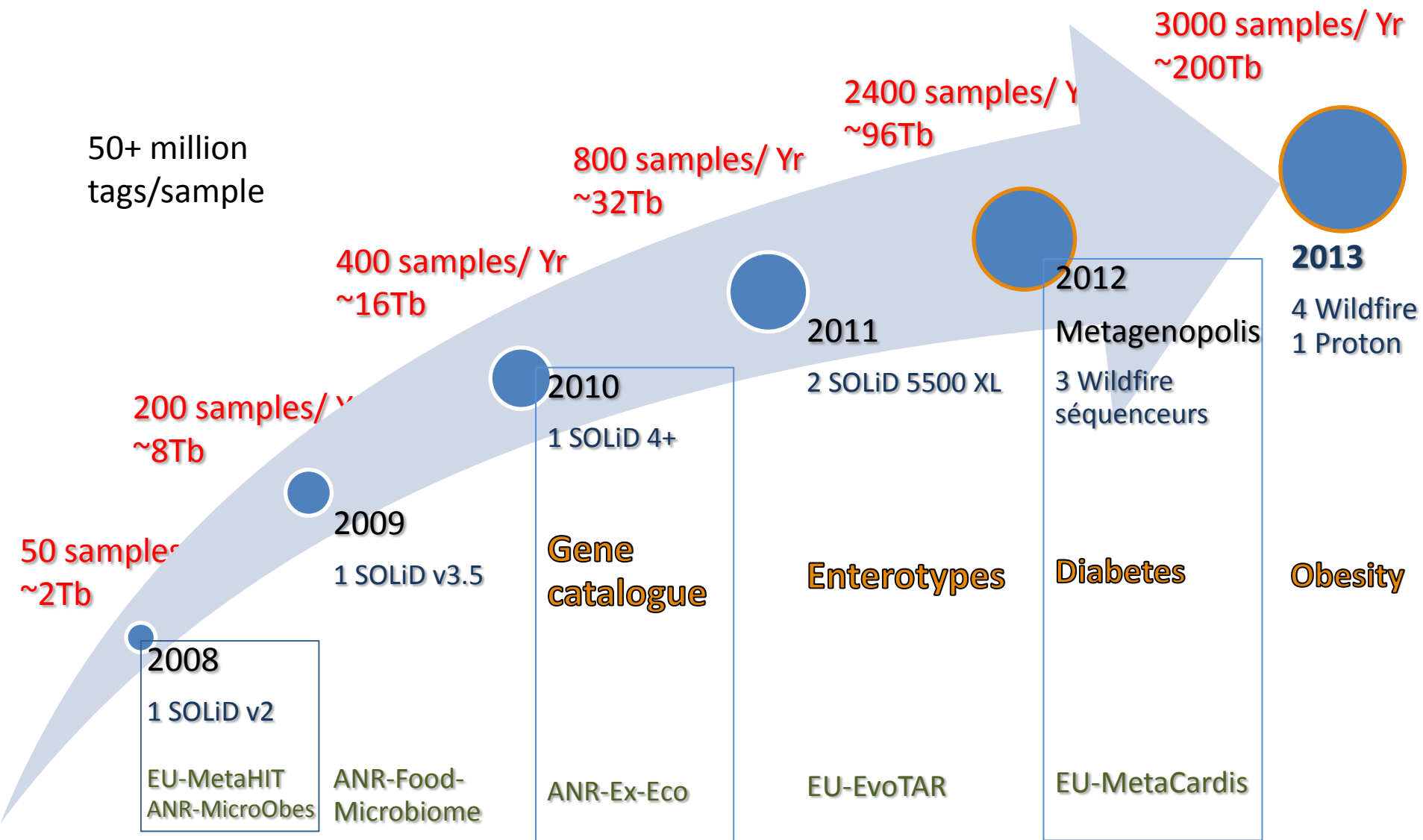
10 M  
Genes

counting *individual*

item	Individuals						
	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



BiG dAtA  
BiG dAtA





✓ data analysis of quantitative metagenomics

## *Preprocessing*

Downsizing

Normalization



## *Processing*

Gene Counting

Sample  
Clustering

Biomarker  
Identification

•  
•  
•

## Wish of Users



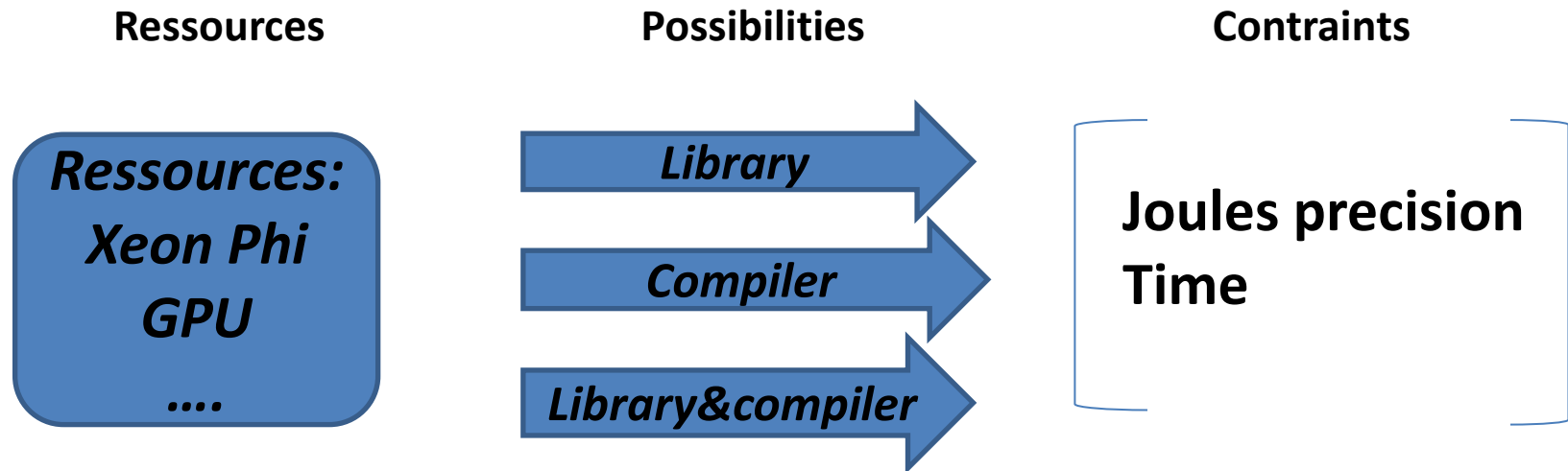
« I want to process and analyze Big Data on R without changing my programming habits ... »

- How to satisfy such requirements ?
- What is the best way to present my solution users ?
- Give them a lesson on their HPC or hide the complexity of using HPC tools ?

## Ways to do this...

- DSL: a programming language whose specifications are dedicated to a specific application domain
- Old fashion optimization : smart library
- New fashion : DSL -> smart compiler
- Smart library and smart compiler

## To sum up



➤ A graph problem



## Mach Project (Massive calculations on heterogeneous systems)

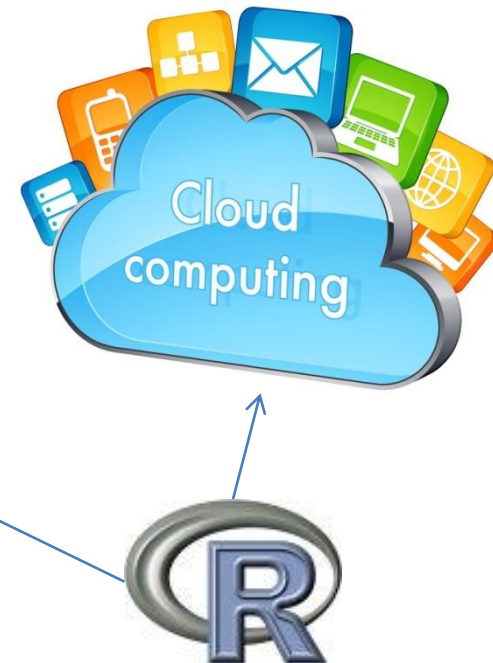
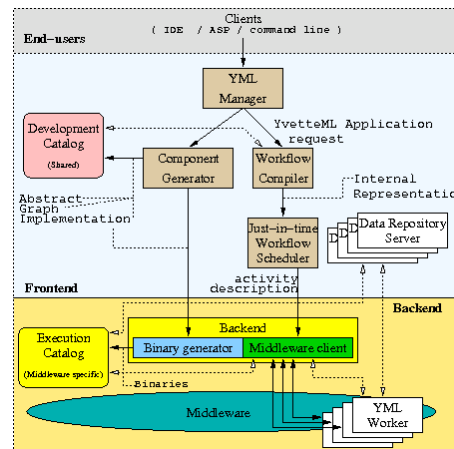
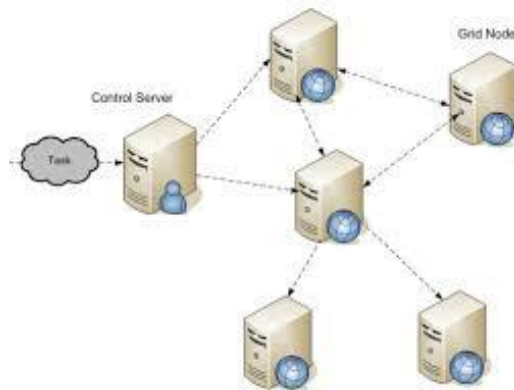
- European Project: private and public partners
- Duration: 3 years



# Architectures Targets

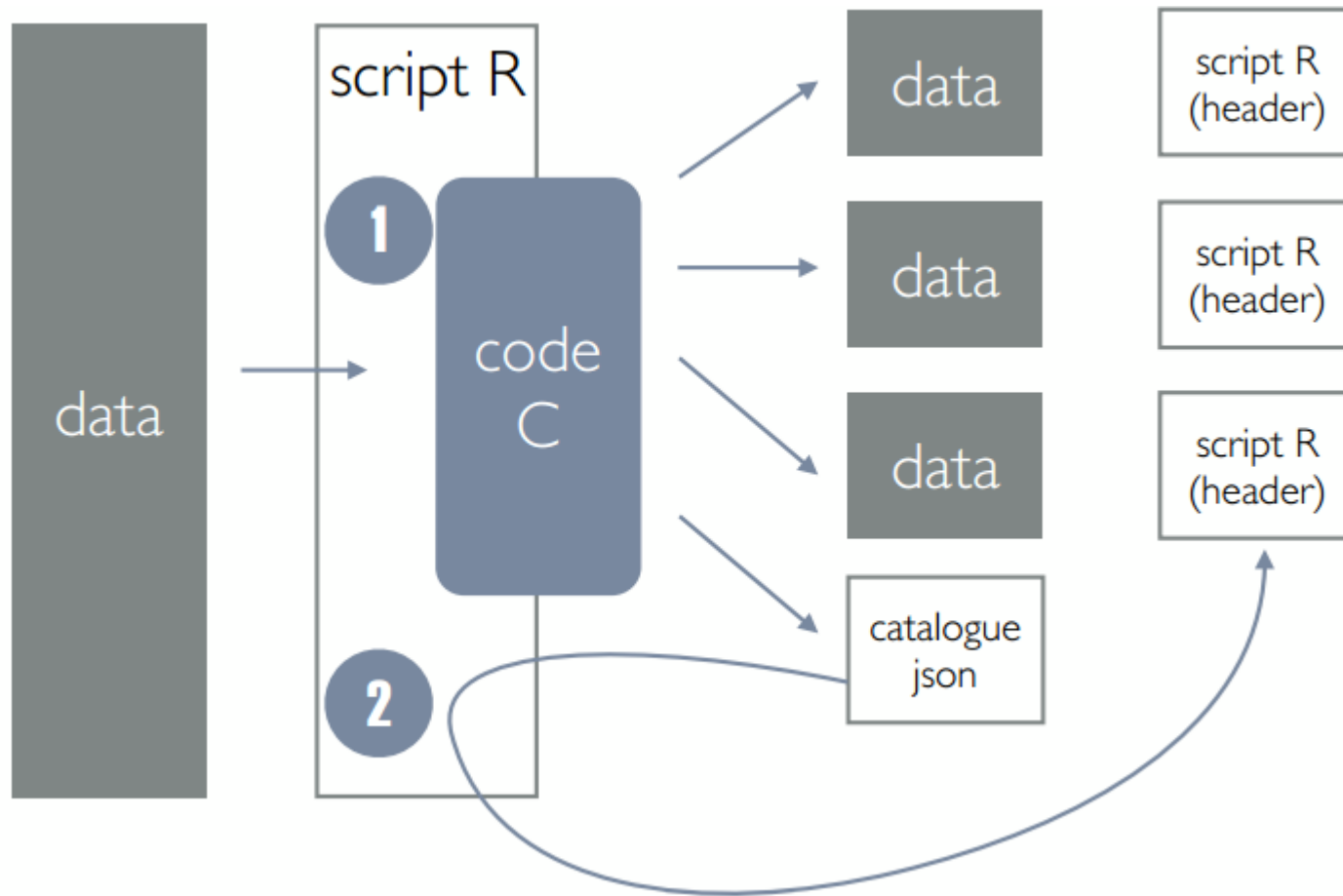


## Scheduler



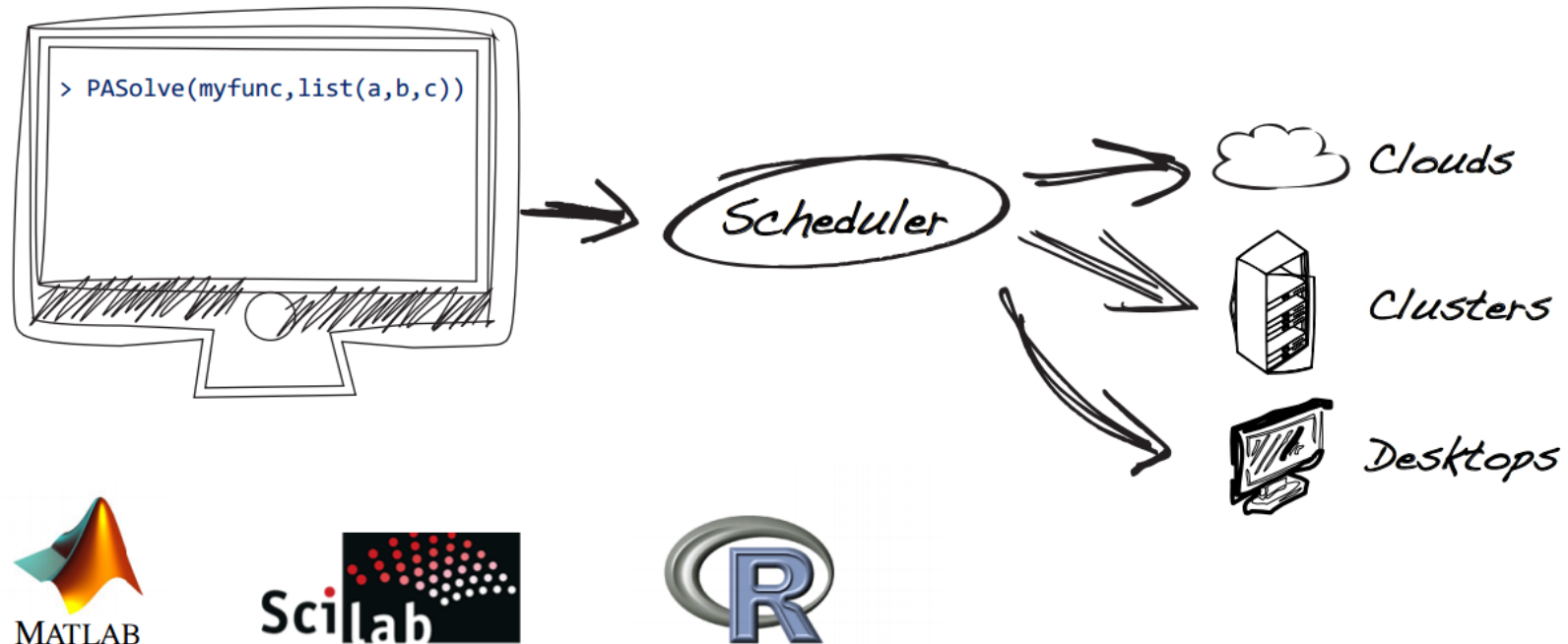
# At Current Day

## MapReduce/ 'R' Megapack



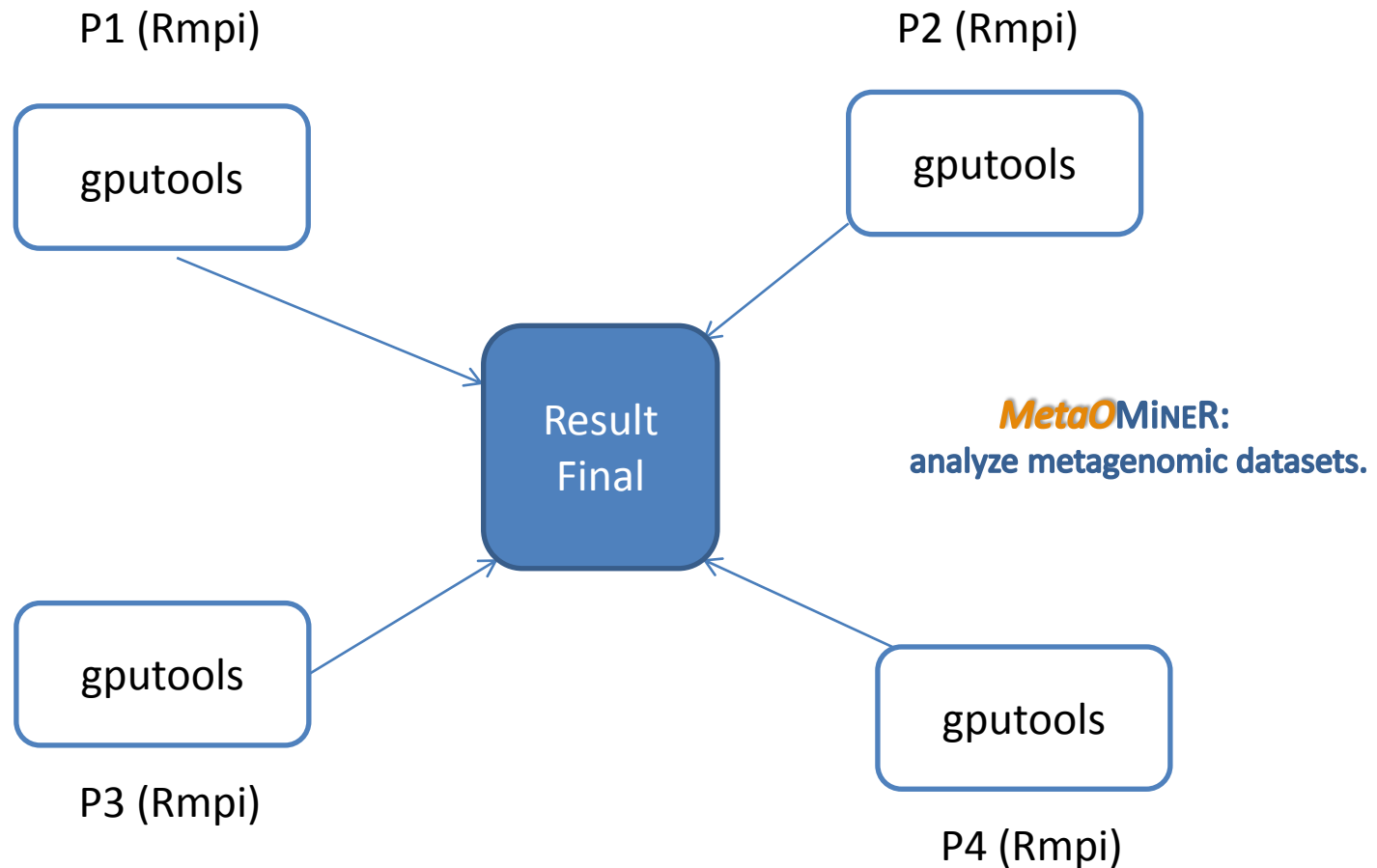


## Jobs soumission/ 'R' PARConnector

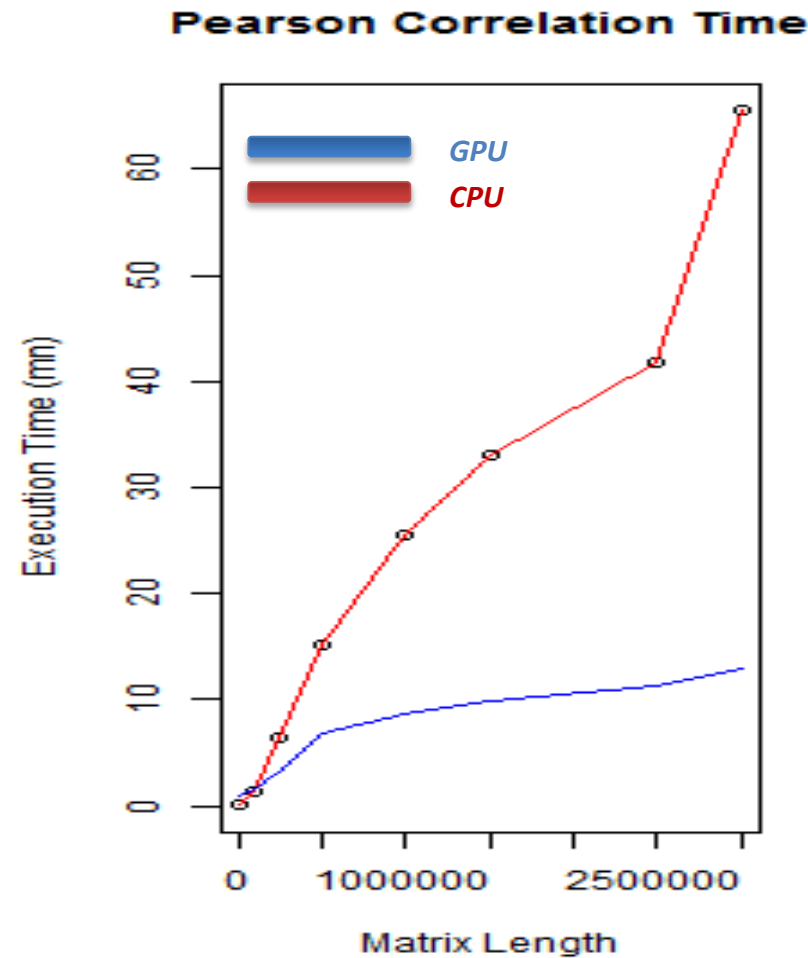


*Intégration transparente dans votre environnement scientifique*

## 'R'gpuStat



## Benchmark of Pearson correlation in GPU



## Conclusion

- Approach the problem from a graph is it the best solution?
- Big Data is -there becoming a barrier to innovation?

Efficient Maximum Flow Algorithms, By Andrew V. GOLBERG & Robert E. TARJAN,  
2014 ACM 0001-0782/14/08

## Thanks



***Jad Abou Ghantous***  
***Anne-Sophie Alvarez***

***Jean-Michel Batto***

***Magali Berland***

***Nahid Emad***

***Amine Ghozlane***

***Vincent Heuschling***

***Emmanuelle Le Chatelier***

***Pierre Léonard***

***Nicolas Pons***

***Florian Plaza***

***Edi Prifti***

