



PARTNERSHIP FOR
ADVANCED COMPUTING IN EUROPE

MaRS - Matrix of RNA-Seq



12èmes Journées

Fabien PIERRAT – Bertrand CIROU

November, the 25th 2016

ACOBIOM – CINES – SHAPE





The MaRS Project

MaRS - Matrix of RNA-Seq

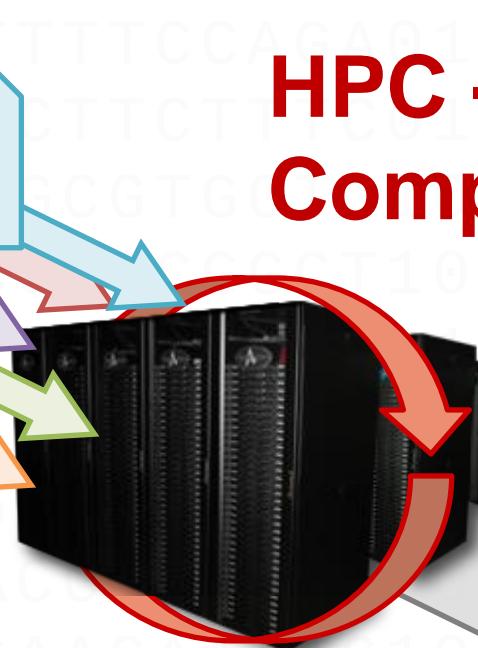
- A large set of **gene expression data (RNA-Seq)** is published & available
- Represents a big **source of information for:**
 - the discovery of new **Biomarkers**
 - the validation of **Biomarkers**
- But, due to
 - **Huge amount of data, huge size of files**
 - **various methods of analysis used**
- These data are:
 - **challenging to compare**
 - **unexploited all together**

The MaRS Project

HPC - High Performance Computing (only 1 method)



Gene Expression profiles
(plenty of files)



MaRS
(only 1 matrix)

	Lib01	Lib02	Lib03	Lib04	Lib05
Gene001	0	5	76	10	0
Gene002	23	0	2	6	9
Gene003	2	95	2	0	2
Gene004	11	4	0	33	1
Gene005	0	5	5	0	53



RNA-Seq collection

- RNA-Seq libraries are collected from public databases (NCBI, EBI)

- We selected libraries:

- with comparable data :

- same type of data

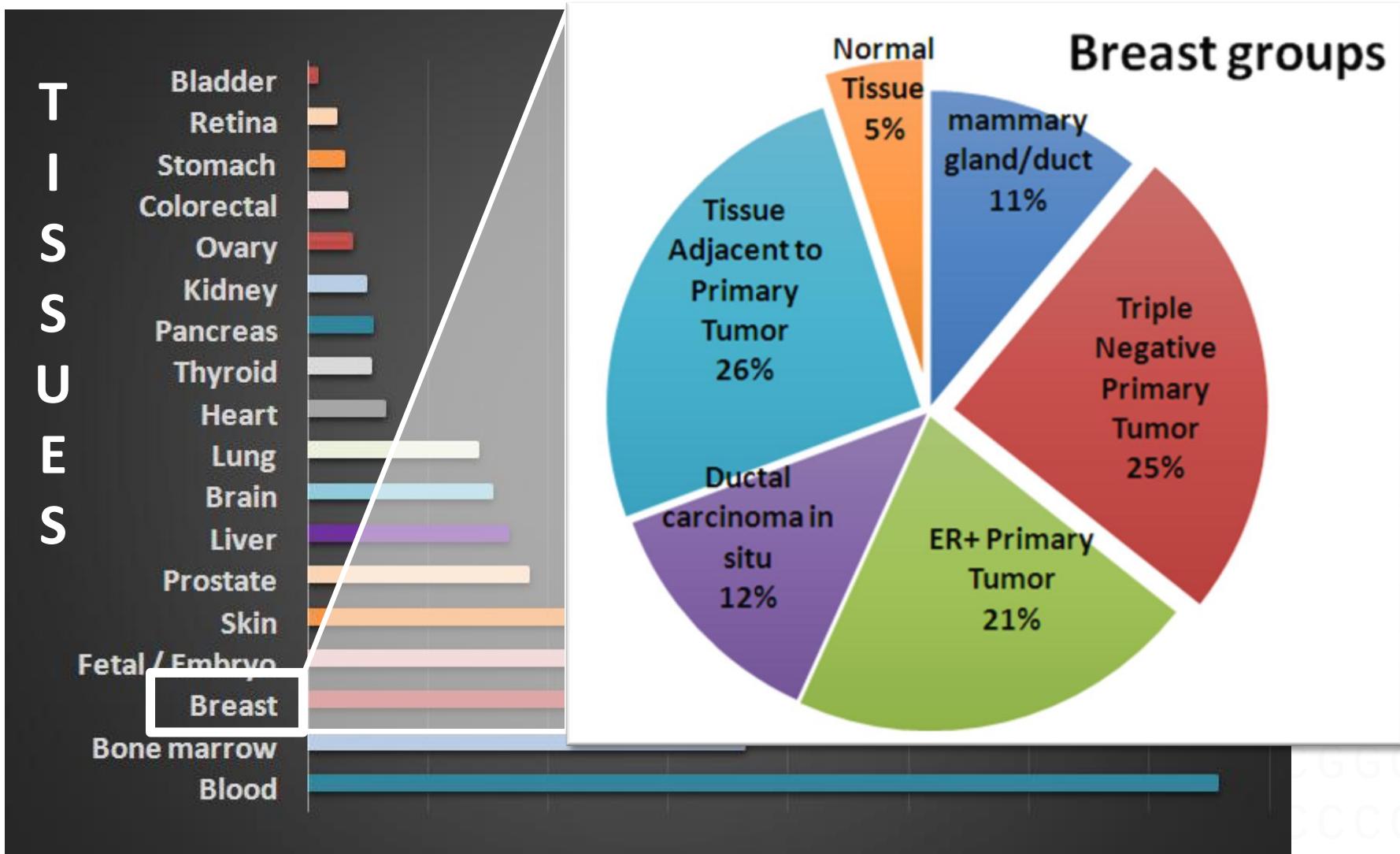
- same technology of sequencing

- from multiple sources (laboratories, countries)

→ ~27000 Human RNA-Seq selected

(from multiple organs / pathologies / ethnic groups...)

RNA-Seq diversity





HPC Toolchain

compressed FASTQ



1

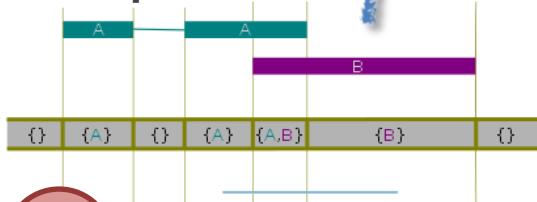
Downloading

	Lib01	Lib02	Lib03	Lib04	Lib05
Gene001	0	5	76	10	0
Gene002	23	0	2	6	9
Gene003	2	95	2	0	2
Gene004	11	4	0	33	1
Gene005	0	5	5	0	53

5

Matrix building MaRS

HTSeq



4

Counting

3

Mapping

TopHat

A spliced read mapper for RNA-Seq

Bow

TIE

Organisation & Alignment

Trimmomatic
Picard

build passing

Quality & Filters

Cleaning

Why HPC?

① Download

6 Go &
0.5 hour

- ② Cleaning
- ③ Mapping
- ④ Counting

4 hours
on 12 CPU =
48 hours/core

+ ⑤ Matrix Building

MaRS

120 To &
416 Days

X
27000

1.2 million
hours/core



Many Challenges

- **Very Large amount of data:**
 - Downloading
 - Processing

→ HPC on cluster (**Occigen**)
- **Software's Installation, Settings on the cluster**
- **Optimization**
 - compiler choice
 - packages choice
- **Jobs duration & Limitations management**
 - Adaptation of software's parameters for all libraries

Tests & Benchmarks

Performance
GCC / Intel ratio



Compiler GCC vs compiler Intel

**Best with
Intel**

**Best with
GCC**

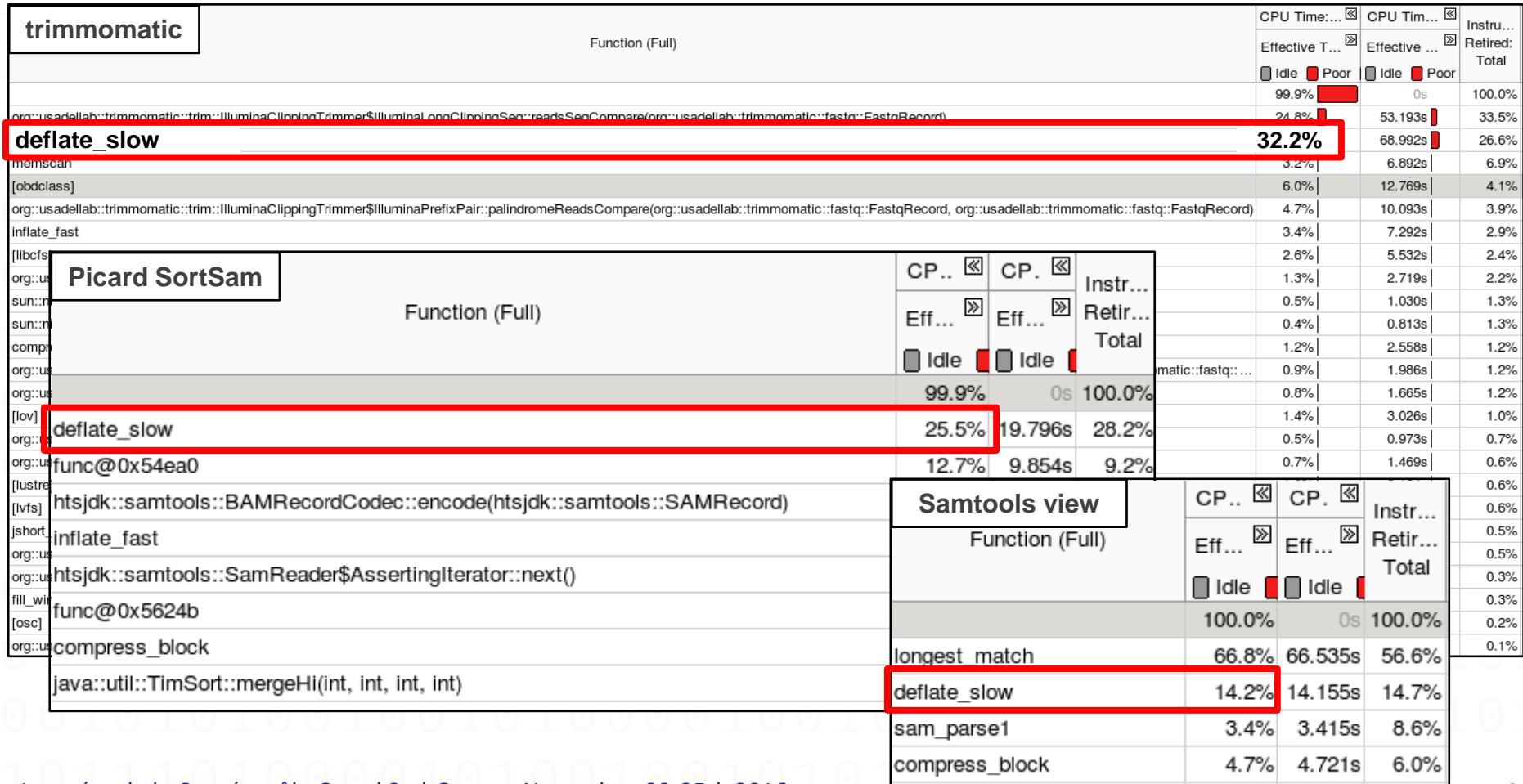
Processing Steps

Libraries

- SRR1644117
- SRR1955853
- ERR649045
- SRR1313202
- SRR1644118
- SRR1313135
- SRR1523680
- SRR2014238

Performance & Optimization

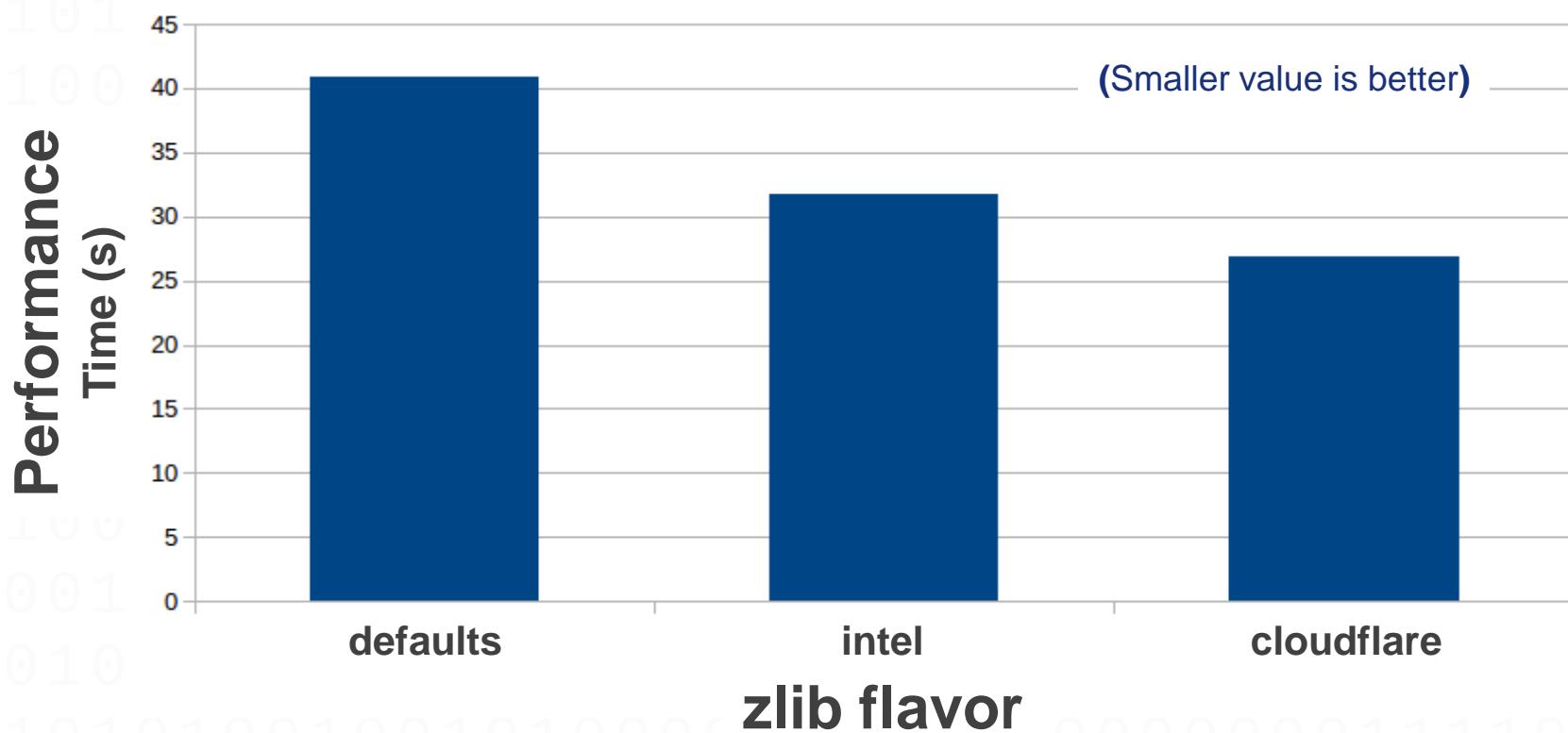
Vtune analysis highlights the high consumption of compression & decompression steps



Performance & Optimization

zlib decompression performance

~750 MB compressed lib





Conclusion & Perspectives

- **NOW : → 95% of the Human MaRS matrix is built**
- **NEXT :**
 - **MaRS Exploration**
Matrix hosting & Queries allowing:
 - Web interface / Database
 - integrated analysis tools
 - **RNA-Seq Data production increases continuously**
 - Integration of new profiles in MaRS
 - **Similar work on other species: Mouse**
- **BIOMARKERS discovery**



Partners



ACOBIOM

- French SME
- Discovery & validation of new biomarkers



SHAPE
SME HPC Adoption Programme in Europe



- Multiple missions:
HPC, digital archiving, hosting
- *Occigen* cluster: 2.1 Pflops
(Top 100 of the most powerful computers in the world)

PRACE - SHAPE

SME HPC Adoption Programme in Europe



Acknowledgments

**THANK YOU for
your attention**



THANKS to all collaborators:

- Fabien PIERRAT
- Laurent MANCHON
- Roman BRUNO
- David PIQUEMAL
- Bertrand CIROU
- Victor CAMEO PONZ
- Francis DAUMAS*

