

PARALLEL COMPUTING FOR DATA SCIENCE

INTRODUCTION

Jairo Cugliari

Master Informatique

Parcours Data Mining

UNIVERSITÉ
LUMIÈRE
LYON 2



Organisation

- Chaque séance : 1h CM + 2h Lab

Responsable

- Jairo Cugliari (Jairo.Cugliari@univ-lyon2.fr)

MCC

- Projet / Cours / Dossier en groupes (50%)
- Examen individuel (avec ordinateur 50%)

Books

- N. Matloff. Parallel Computing for Data Science: With Examples in R, C++ and CUDA (Chapman & Hall/CRC The R Series) Open textbook version
- M. Gorelick and I. Ozvald. High Performance Python. (O'Reilly)
- H. Wickley. Advanced R (Chapman & Hall/CRC The R Series) Open textbook version

E-learning

A tsunami of ressources, a glimpse:

- Blog by J. VanderPlas
- Blog by T. Oliphant
- Course by V. Miele (in French)

Why ?

```
int getRandomNumber()  
{  
    return 4; // chosen by fair dice roll.  
              // guaranteed to be random.  
}
```

Source: xkcd.com

What ?

HOW LONG CAN YOU WORK ON MAKING A ROUTINE TASK MORE EFFICIENT BEFORE YOU'RE SPENDING MORE TIME THAN YOU SAVE?
(ACROSS FIVE YEARS)

		HOW OFTEN YOU DO THE TASK					
		50/DAY	5/DAY	DAILY	WEEKLY	MONTHLY	YEARLY
HOW MUCH TIME YOU SHAVE OFF	1 SECOND	<div><div>1</div> DAY</div>	2 HOURS	30 MINUTES	4 MINUTES	1 MINUTE	5 SECONDS
	5 SECONDS	<div><div>5</div> DAYS</div>	12 HOURS	2 HOURS	21 MINUTES	5 MINUTES	25 SECONDS
	30 SECONDS	<div><div></div><div></div><div></div><div></div><div></div><div></div></div> 4 WEEKS	<div><div>3</div> DAYS</div>	12 HOURS	2 HOURS	30 MINUTES	2 MINUTES
	1 MINUTE	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div> 8 WEEKS	<div><div>6</div> DAYS</div>	<div><div>1</div> DAY</div>	4 HOURS	1 HOUR	5 MINUTES
	5 MINUTES	9 MONTHS	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div> 4 WEEKS	<div><div>6</div> DAYS</div>	21 HOURS	5 HOURS	25 MINUTES
	30 MINUTES		6 MONTHS	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div> 5 WEEKS	<div><div>5</div> DAYS</div>	<div><div>1</div> DAY</div>	2 HOURS
	1 HOUR		10 MONTHS	2 MONTHS	<div><div>10</div> DAYS</div>	<div><div>2</div> DAYS</div>	5 HOURS
	6 HOURS				2 MONTHS	<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div> 2 WEEKS	<div><div>1</div> DAY</div>
	<div><div>1</div> DAY</div>					<div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div> 8 WEEKS	<div><div>5</div> DAYS</div>

How ?



Main topics

1. Intro. Overview. Case study
2. Good practices for coding
3. Bottlenecks: seek & destroy
4. Some mathematics of parallel data analysis

Speakers

- Jairo Cugliari
- maybe an invited speaker also