

# 3IF4260 - Calcul à Haute Performance pour l'analyse de données

Responsables : **Stephane VIALLE**

Langues d'enseignement : **FRANCAIS**

Campus où le cours est proposé : **CAMPUS DE PARIS - SACLAY**

Nombre d'heures d'études élèves (HEE) : **40**

Nombre d'heures présentielles d'enseignement (HPE) : **24**

Année académique : **2024-2025**

Niveau avancé : **non**

## Présentation, objectifs généraux du cours :

Ce cours propose d'étudier les environnements de programmation distribuée "Spark" et "MPI", avec suffisamment d'approfondissement pour une mise en oeuvre sur des problèmes d'analyse de données, avec des exécutions sur clusters de PC et avec une approche permettant un "passage à l'échelle". Des TP avec mesures et analyses de performances jalonnent le déroulement du cours.

## Période(s) du cours (n° de séquence ou hors séquence) :

SM11

## Prérequis :

- Cours commun de 1A : Systèmes d'Information et Programmation (1CC1000)
- Cours commun de 1A : Algorithmique & Complexité (1CC2000)
- Cours de 3A : Systèmes concurrents et répartis (3IF1040)
- Cours de 3A : Architectures matérielles (3IF4010)

## Plan détaillé du cours (contenu) :

Ce cours comprend 2 parties avec TP et un examen final : CM 12h00, TP 12h00 (total : 24 HPE)

- **Algorithmique et programmation distribuée selon un schéma "map-reduce" en Spark** : CM 6h00, TP 6h
  - Rappel de programmation Spark RDD
  - Concept et programmation SQL en Spark (Spark-SQL)
  - Concept et mécanismes de traitements de flux de données en Spark (Spark Structured Streaming)
  - TPs de Spark distribué sur cluster de PC (TP1.1 : Spark-SQL, TP1.2 : TP Spark structured streaming)
- **Algorithmique et programmation distribuée par envois de messages en MPI** : CM 6h00, TP 6h
  - Rappels de déploiement sur cluster de PC multi-coeurs de code MPI+OpenMP
  - Entrées/Sorties disques depuis un programme MPI
  - Distribution d'un algorithme de clustering par k-means
  - Communications asynchrones, recouvrement calculs-communications
  - TPs de MPI sur clusters de PCs (TP2.1: k-means distribué en mpi4py, TP2.2: produit de matrices avec recouvrement calculs-communications en C+MPI+OpenMP)

## Déroulement, organisation du cours :

Les concepts vu en cours seront mis en oeuvre et approfondis dans 4 TP. Chaque solution développée sera exécutée sur un cluster de PC et ses performances seront mesurées et analysées. Des optimisations des codes et des algorithmes seront effectuées si nécessaire pour obtenir des solutions performantes et aptes au passage à l'échelle.

- Cours : 12h00
- TP sur cluster de PC, avec rapports et évaluations : 12h00

## Organisation de l'évaluation :

Evaluation à partir des TP des parties MPI et Spark:

- 50% : Comptes rendus des TP de la partie 1
- 50% : Comptes rendus des TP de la partie 2

Rmq : Le contenu et le nombre de pages des comptes rendus sont contraints, afin de forcer les étudiants à un effort de synthèse et de clarté

En cas d'absence non justifié à une séance de TP la note de 0 sera appliquée, en cas d'absence justifiée à une séance de TP la note finale sera calculée à partir des autres séances de TP.

L'examen de rattrapage sera un examen écrit ou oral, qui constituera 100% de la note de rattrapage.

## Moyens :

**Equipe enseignante :**

**S. Vialle** et **G. Quercini** (CentraleSupélec)

**Ressources de calcul :**

les TP se dérouleront sur les clusters de PC du Data Center d'Enseignement du campus de Metz, accédées à travers Internet.

## Acquis d'apprentissage visés dans le cours :

A l'issue de ce cours, les étudiants sauront :

- **AA1** : concevoir et implanter un algorithme MPI avec communications par circulation de données et par communications collectives,
- **AA2** : concevoir et implanter un algorithme MPI avec recouvrement des calculs et des communications,
- **AA3** : concevoir et implanter un algorithme Spark traitant un flux de données,
- **AA4** : mesurer les performances et optimiser des codes Spark distribués sur clusters de PC,
- **AA5** : tester l'aptitude d'un code à "passer à l'échelle".

## Description des compétences acquises à l'issue du cours :

- **C2**: Develop in-depth skills in an engineering field and a family of professions (in connection with learning outcomes AA1, AA2 and AA3).
- **C6**: Be operational, responsible, and innovative in the digital world (in connection with learning outcomes AA4 and AA5).
- **C7**: Know how to convince (in connection with the TP reports).

