

# Programmation Big Data

## Présentation du module

Sylvain Tenier

Département TIC - Esigelec

Février 2016

# Objectifs

- ① Programmation concurrente
- ② Programmation distribuée évolutive
- ③ SQL sur systèmes distribués
- ④ Systèmes “In memory”

# Prérequis

- ① Tronc commun : JAVA, SQL, gestion de projet
- ② TIC : Git
- ③ Dominante : Threads JAVA

# Déroulé type d'une séance

- ❶ Présentation et restructuration (30 min)
  - Présentation de la séance
  - Restructuration des séances précédentes
    - Envoi de vos questions par mail avant la séance
- ❷ Réalisation (3h)
  - Conception (groupe de 4)
    - Groupes différents des autres modules !
  - Réalisation
  - Création des slides de restitution
- ❸ Restitution (30min)
  - Projet Git à jour sur Github
  - Présentation orale de 5 min maximum (bibliographie, description de l'approche retenue, répartition des tâches, résultats obtenus, perspectives)
  - Tableau comparatif des performances

# Dataset 1 : MovieLens

- Le groupe de recherche GroupLens a collecté plusieurs millions de notes attribuées à des films par des utilisateurs du site MovieLens
- Vous devez exploiter ces données pour répondre aux questions suivantes :
  - ① Quelle est la proportion de films ayant 1, 2, 3, 4, 5 étoiles ?
  - ② Quel utilisateur a noté le plus de films ?
  - ③ Quel est le film le plus populaire ?
- Vous répondrez à ces questions en utilisant une technologie différente à chaque séance

## Dataset 2 : Yelp

- Yelp est un service permettant de trouver et noter des services de proximité. La société a mis à disposition un “dataset académique” comprenant les évaluations des 250 sociétés les plus proches de 30 universités :
  - [https://www.yelp.com/academic\\_dataset](https://www.yelp.com/academic_dataset)
- Vous devez exploiter ces données pour répondre aux questions suivantes :
  - 1 Quelle est la proportion d'établissements par nombre d'étoiles ?
  - 2 Pour chaque catégorie (useful, funny, cool), quel utilisateur a apporté le plus de votes ?
  - 3 Quel est le restaurant le mieux noté pour chaque université ?
- Vous répondrez à ces questions en utilisant une technologie différente à chaque séance