*Workshop 9 - Big Data Analytics*

Luca Morandini
Data Architect – AURIN Project
University of Melbourne
luca.morandini@unimelb.edu.au

# Starting Spark With Docker

- The instructions on how to build a mini-cluster (a Spark master and one Spark worker Docker containers), and an example of how to run a word-count MapReduce job is at the usual **https://github.com/AURIN/comp90024** repository, go under the **spark** directory and follow the README.

# Visualizing a Spark Cluster

● Spark *Web-UI* allows to graphically see nodes (port 8080)

# Visualizing Spark Jobs: Executors

- Spark *Web-UI* allows as well to see the division of jobs in stages and tasks, and the allocation of executors to nodes (port 4040)

# Visualizing Spark Jobs: Jobs

# Visualizing Spark Jobs: Single Job

# Visualizing Spark Jobs: Tasks of a Stage