



Big Data e Computação em Nuvem

Aula 10

Spark DataFrames - Cache

Prof. Michel Fornaciali, PhD.

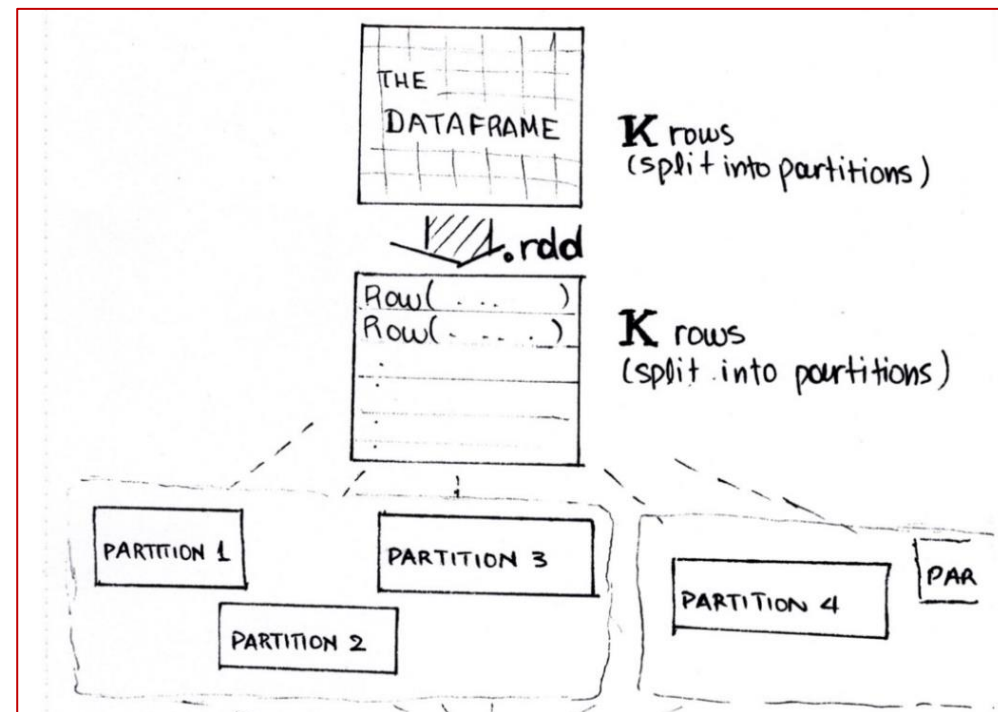
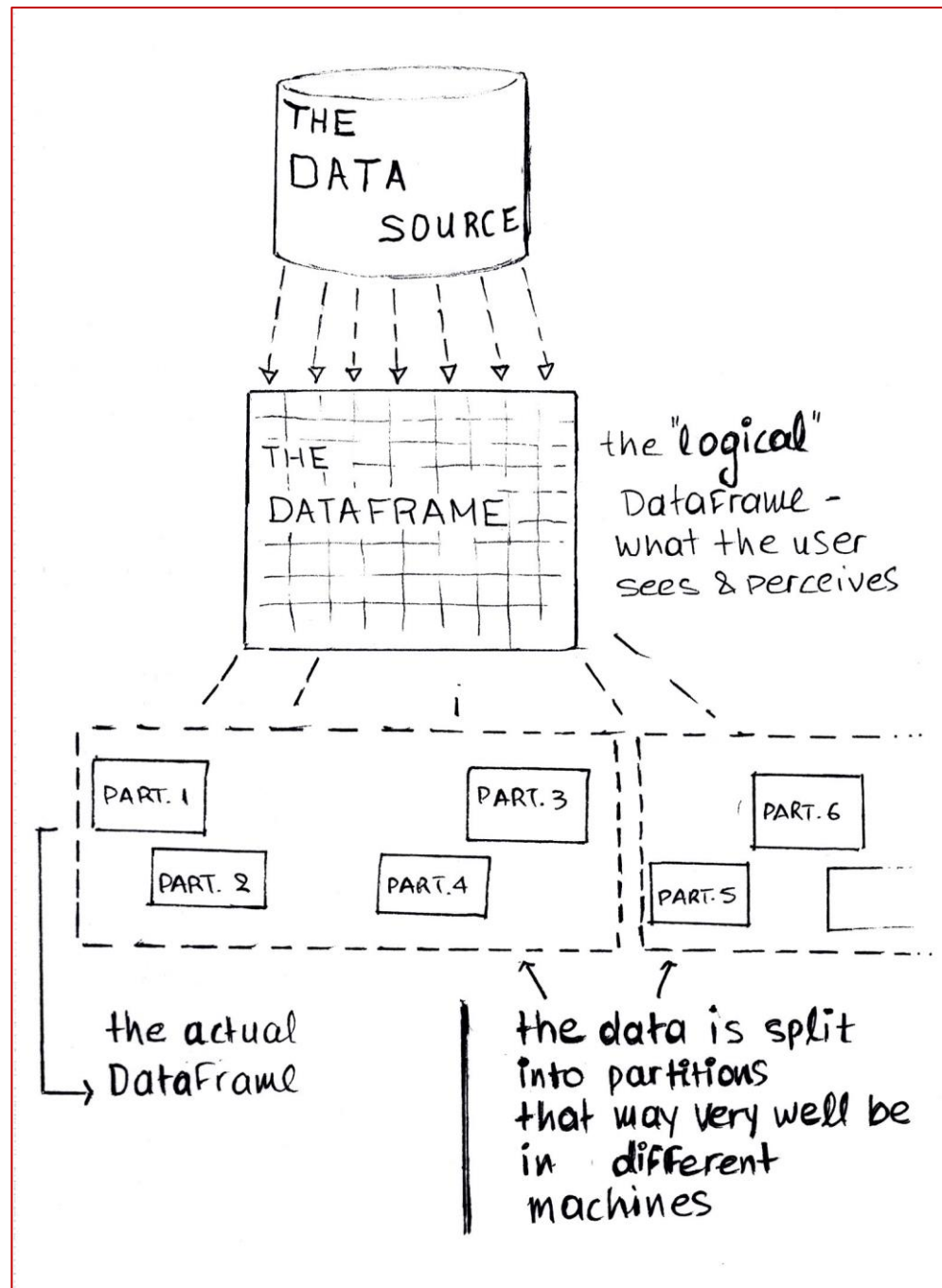
Prof. Thanuci Silva, PhD.

Contatos:

MichelSF@insper.edu.br

thanucis@insper.edu.br

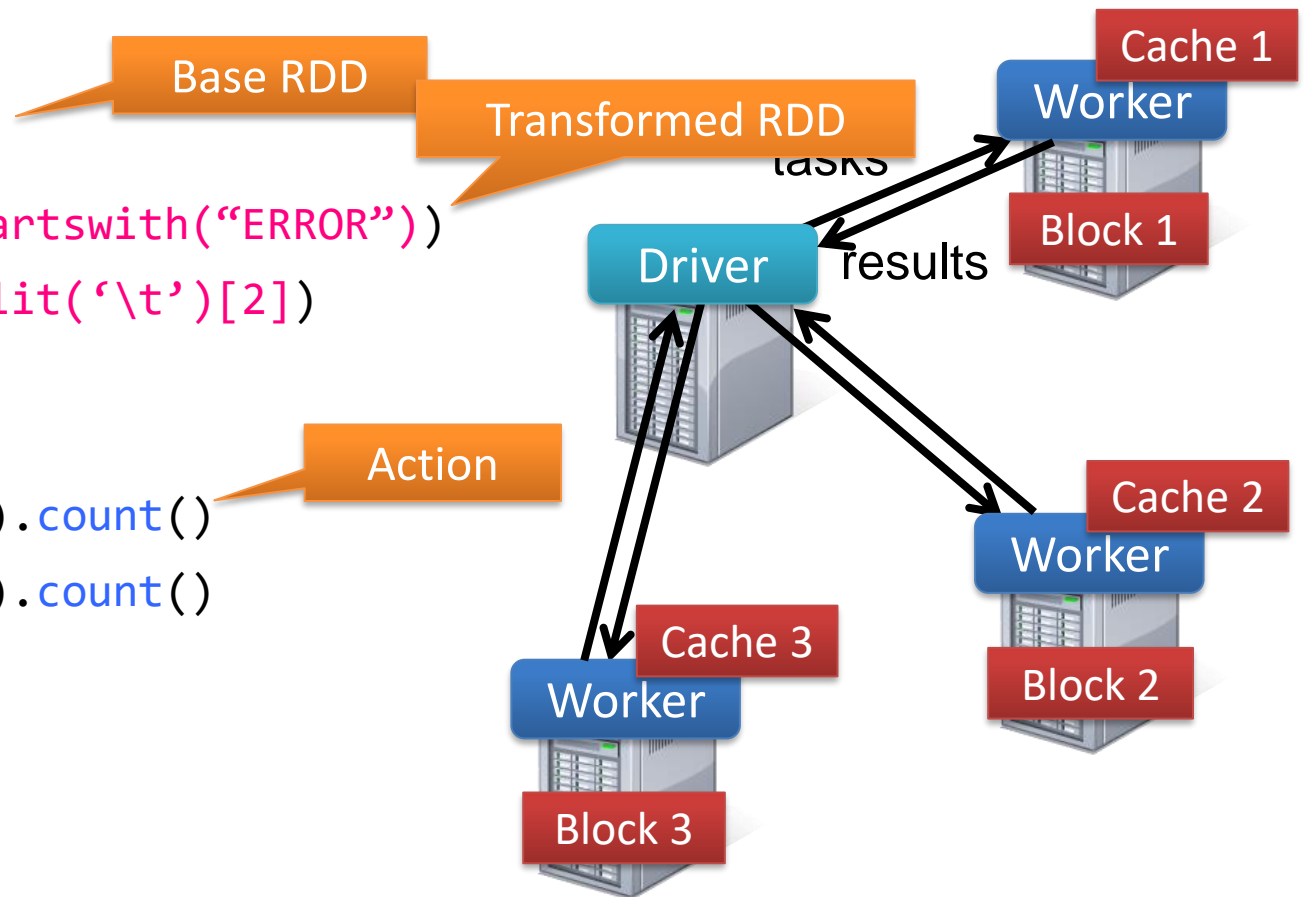
DataFrame



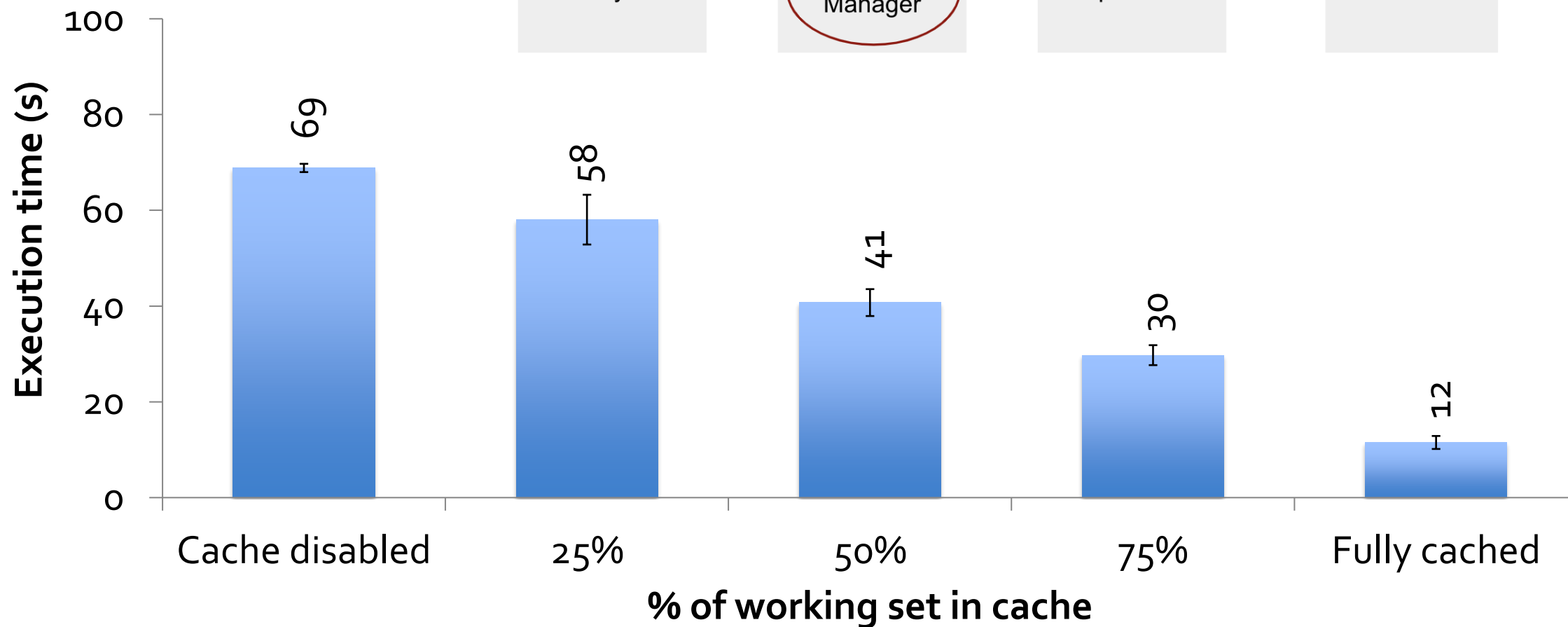
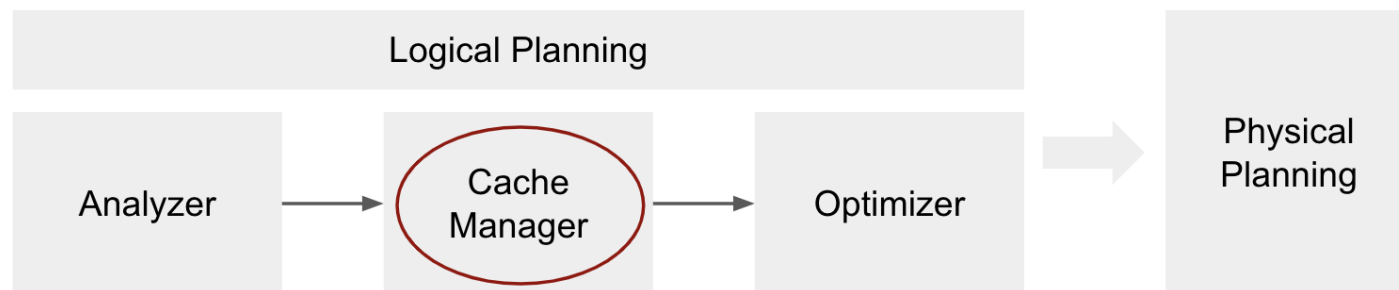
Spark – Mecanismo de Cache

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split('\t')[2])
messages.cache()

messages.filter(lambda s: "foo" in s).count()
messages.filter(lambda s: "bar" in s).count()
. . .
```



Scaling Down Use of Cache



Uso do cache

```
inputRDD = sc.textFile("words.txt").cache()

print("Number of words: ",inputRDD.count())
print("Number of distinct words: ", inputRDD.distinct().count())
```

Inspire