

# Comentário do artigo “Apache Spark”

Bruno Susana, Gonalo Mascarenhas, Guilherme Correia, Henrique Cruz, Ricardo Rosa,

FCT UALG

**Abstract**— O presente documento visa expor a opinio crtica do grupo GC sobre o trabalho desenvolvido pelo grupo G3, que teve como objetivo de estudo o “Apache Spark”. O trabalho ser avaliado em diferentes aspetos e ser atribuída uma classificao entre 0 e 10 em cada aspecto. No final ser realizado um pequeno comentrio sobre o porqu de cada nota atribuída.



## 1 INTRODUO

O artigo em questo tem como tema o “Apache Spark”, este *Framework*  utilizado sobretudo em *Big Data*, termo este aplicado ao tratamento e anlise de grandes volumes de informao e de dados que seriam demasiado demorosos e complexos de analisar pelos mtodos mais tradicionais. O artigo em questo comea por abordar os conceitos de sistema paralelo e sistema distribuído na sua generalidade, apresentando vrios conceitos relacionados, dos quais as diferenas entre programaco concorrente e sequencial, regio crtica e *cloud and grid computing*.

De seguida  apresentado ao leitor o ponto fulcral do artigo que consiste no *framework* que est sendo anlise, aqui  apresentada alguma informao relevante tal como a origem, a arquitetura, as bibliotecas implementadas e o seu funcionamento.

Aps esta introduo -nos apresentado alguns casos de uso em que este *framework*  utilizado e exemplos de grandes empresas que empregam esta tecnologia nos dias atuais, mostrando desta forma que esta tecnologia continua bem presente atualmente.

No final  feita uma comparao de *performance* em que se compara o *framework Spark* com o seu maior concorrente no mercado, *Hadoop-Map Reduce*, onde so avaliados alguns pontos fulcrais quando se escolhe qual a ferramenta mais adequada a cada caso.

Aps uma primeira anlise no realizada ainda com rigor,  possvel constatar que o trabalho produzido pelo grupo apresenta bastante qualidade em termos de estrutura e em

termos de escrita onde so apresentados vrios contedos que foram lecionados em aula.

## 2 AVALIAO

Onde so abordados os sistemas paralelos e distribuídos os autores falam sobre os vrios envolventes de forma clara e explicita, com o auxlio de imagens e tabelas para ajudar o leitor a compreender melhor os temas base do artigo.

Onde  abordado o tema *cloud computing*, os autores referem o impacto que este teve e  feita uma comparao com *grid computing* de forma a entender as semelhanas entre os conceitos. Desta forma,  possvel ao leitor entender de forma suprflua a diferena entre os conceitos, no obstante o facto de serem feitas boas referncias e introduo a temas que so falados mais  frente no artigo.

No ponto 4.2 no que se refere  histria do Apache,  mencionado que ficou um Top-Level Apache Project em 2013, no entanto isso aconteceu em Fevereiro de 2014 [2].

De seguida -nos introduzido o *framework* em estudo com uma esclarecedora e sucinta introduo onde so expostas as bases do *Spark* e a sua arquitetura. Aqui  mostrado como funcionam os *RDD*’s associados aos *framework*.

Na parte do *higher-level libraries*, os autores optaram por explicar as utilidades das bibliotecas no *Apache spark*. Comeado pelo *Spark sql*, os autores explicam de forma geral em que consiste e ao longo do texto so indicadas as linguagens de programaco utilizadas no desenvolvimento desta biblioteca, a figura 7  bem escolhida demonstrando complementando graficamente a explicao dada.

Na parte sobre *ghrapx*, o tema poderia ser mais completo e mais explicativo, apesar do que j  apresentado ser o suficiente como forma de introduo do tema.

Relativamente ao *Mlib*, o tema encontra-se bem explicado, os autores optaram por no colocar cdigo, desta forma  mais simples para o leitor.

No que toca ao tema da comparao do Apache Spark

- 
- Bruno Susana, aluno de Licenciatura em Engenharia Informtica, na Universidade do Algarve; email: [a61024@ualg.pt](mailto:a61024@ualg.pt)
  - Gonalo Mascarenhas, aluno de Licenciatura em Engenharia Informtica, na Universidade do Algarve; email: [a64533@ualg.pt](mailto:a64533@ualg.pt)
  - Guilherme Correia, aluno de Licenciatura em Engenharia Informtica, na Universidade do Algarve; email: [a61098@ualg.pt](mailto:a61098@ualg.pt)
  - Henrique Cruz, aluno de Licenciatura em Engenharia Informtica, na Universidade do Algarve; email: [a61099@ualg.pt](mailto:a61099@ualg.pt)
  - Ricardo Rosa, aluno de Licenciatura em Engenharia Informtica, na Universidade do Algarve; email: [a62461@ualg.pt](mailto:a62461@ualg.pt)

com o Apache Hadoop no geral foram feitas boas comparações, com bons fundamentos entre cada característica de cada framework de forma clara e compreensível explicando como funciona cada aspeto de cada framework.

Finalmente quando são referidas as aplicações, os autores escolheram temas de grande relevo no mercado e bastante apelativos à leitura, acompanhado de exemplos que captam o leitor.

### 3 CONFORMIDADE COM AS DIRETRIZES DO AUTOR

Em termos de escrita, é utilizada um tipo de escrita agradável visualmente e de fácil leitura, o tamanho da mesma também não gera qualquer dificuldade visual. Não se detetou nenhum tipo de abreviação ou anotação que gerasse dúvida na ideia que se pretende transmitir.

Existe um pequeno reparo em que uma das figuras é repetida e a numeração das mesmas encontra-se errada, de tal forma que o leitor primeiro encontra a figura 3 e de seguida a figura 2 sendo que ambas as figuras constam também com a mesma descrição.

A figura 4 encontra-se num local inadequado, esta inserida no tema *Spark*, mas deveria estar em *Cloud computing and Grid computing*.

A classificação das figuras não se encontra uniforme ao longo do documento, temos como exemplo a figura 1 e a figura 6 que não apresentam um mesmo tamanho na identificação da respetiva figura.

Na figura 8 é apresentada uma fase denominada *Spark engine* que não é explicada, a contrário das outras fases que tiveram uma explicação clara e sucinta.

A figura 9 encontra-se bastante reduzida, dificultando a sua leitura em formato físico.

### 4 CLASSIFICAÇÃO

De seguida são atribuídas classificações de 0 a 10 em diferentes aspetos qualitativos:

Gráficos: 7

Referências: 7

Organização (Ortografia e gramática): 8

Descrição do caso de estudo e enquadramento: 9

Clareza de organização: 9

Completude: 9

### 5 DICAS ÚTEIS

Na parte do Graphx os autores poderiam ter colocado uma figura a mostrar como funciona a biblioteca de forma a transmitir melhor para o leitor. [1]

Deveriam ter indicado de onde vieram algumas referências sobre as aplicações.

Na parte da comparação entre o Spark e o Hadoop poderiam ter usufruído de mais figuras para ter uma comparação mais clara e mais consolidada em alguns dos aspetos e não apenas na

performance, acrescentando assim uma maior legibilidade e consolidação da informação descrita.

### 6 CONCLUSÃO

Com este comentário, podemos concluir que o artigo analisado é um artigo bem estruturado, que apresenta os conteúdos de forma linear, começando com algo geral e ao longo do artigo torna-se mais concreto. Os autores também expõem o conteúdo de forma que o leitor a entenda sem ter qualquer tipo de dúvida.

Por outro lado, têm alguns erros que são necessários para compreensão do leitor e também possuem erros a nível de conteúdo.

Nota final: 9

### REFERÊNCIAS

- [1] spark, a. (2020). *GraphX Programming Guide*. Obtido de spark apache: <https://spark.apache.org/docs/latest/graphx-programming-guide.html>
- [2] wikipedia "Apache Spark". (23 de Abril de 2021). Obtido de [https://en.wikipedia.org/wiki/Apache\\_Spark](https://en.wikipedia.org/wiki/Apache_Spark)