

# Sequenciamento de DNA utilizando Hadoop (MR) (Java Vs. C++)

RAFFAEL BOTTOLI SCHEMMER

UFRGS - PPGC - POD 2014/I

# Sumário

- ▶ Introdução.
- ▶ Motivação e Objetivos.
- ▶ Sequenciamento de DNA.
- ▶ Modelo Map Reduce.
- ▶ Implementação.
- ▶ Funcionamento da aplicação.
- ▶ Resultados.
- ▶ Conclusões e trabalhos futuros.

# Introdução.

- ▶ Forte interesse do mercado por computação intensiva em dados.
- ▶ Cluster favorece o modelo de persistência distribuído:
  - ▶ Suporte a tolerância a falhas.
  - ▶ Suporte a leitura e escrita paralela de informação.
  - ▶ Alta disponibilidade de persistência.
- ▶ Sequenciamento de DNA:
  - ▶ Alto volume de geração de dados.
  - ▶ Demanda de computação intensiva em processamento de dados.

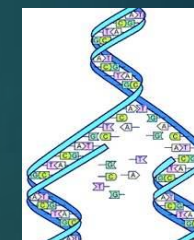
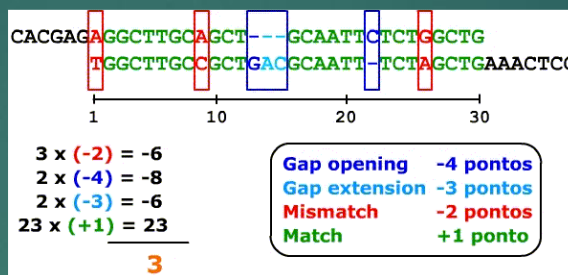
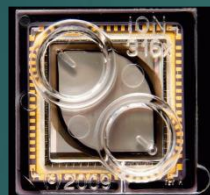
# Motivação e Objetivos.

- ▶ Hadoop:
  - ▶ Ferramenta de big data voltada para:
    - ▶ Persistência distribuída de informações.
    - ▶ Suporte a tolerância a falhas de tempo real.
    - ▶ Suporte a processamento distribuído (MR).
  - ▶ Implementado em Java:
    - ▶ RMI implementa o mecanismo de troca de mensagem do FW.
    - ▶ Rede pode se tornar um gargalo dado o modelo MR utilizado.
    - ▶ Sockets é suportado através de Hadoop Pipes (C++).
- ▶ Autor deste trabalho acredita que:
  - ▶ Uso de Hadoop Pipes pode melhorar o tempo de execução e o consumo de recursos da aplicação.

# Motivação e Objetivos.

- ▶ GPPD:
  - ▶ Desenvolve cooperação com HCPA.
  - ▶ Desenvolveu um TG que propôs uma aplicação de sequenciamento de DNA usando de MR utilizando o Hadoop.
- ▶ Este trabalho têm como proposta:
  - ▶ Entender o trabalho implementado.
  - ▶ Rescrever o código Java em C++ com suporte a Hadoop Pipes.
  - ▶ Avaliar ambas as implementações.
  - ▶ Uso do cluster GradeP.

# Sequenciamento de DNA.



mapt.NA12156.alter.bam x

00000000	1f 8b 08 04 00 00 00 00 ff 06 00 42 43 02 00	.....ÿ..BC..
00000010	11 3c ad 7d 7b d4 25 57 55 e7 97 ee 3c 3a dd 4d	.< }{Ô\$WUc.î<:ŸM
00000020	ba ab 72 3a a9 74 27 9d 9c c3 45 ce 2c e5 5b f5	°xr:®t'!..ÆI,â[ð
00000030	ae 5b 3a 38 54 7d c5 7c 05 63 94 c7 47 14 99 11	ø[:®T)â].c.ÇG...
00000040	0c 20 2c 5e 42 08 24 60 20 a9 ca 01 0b a3 3c 04	,,^B.ê' øË...f<.
00000050	1a e5 29 82 38 3a a2 20 33 12 11 11 61 8d 0b 94	.â).8:ø 3...a...



ABCB11 p2\_444\_2\_gtc\_gyc  
this mutation was not found in the current DataBase.

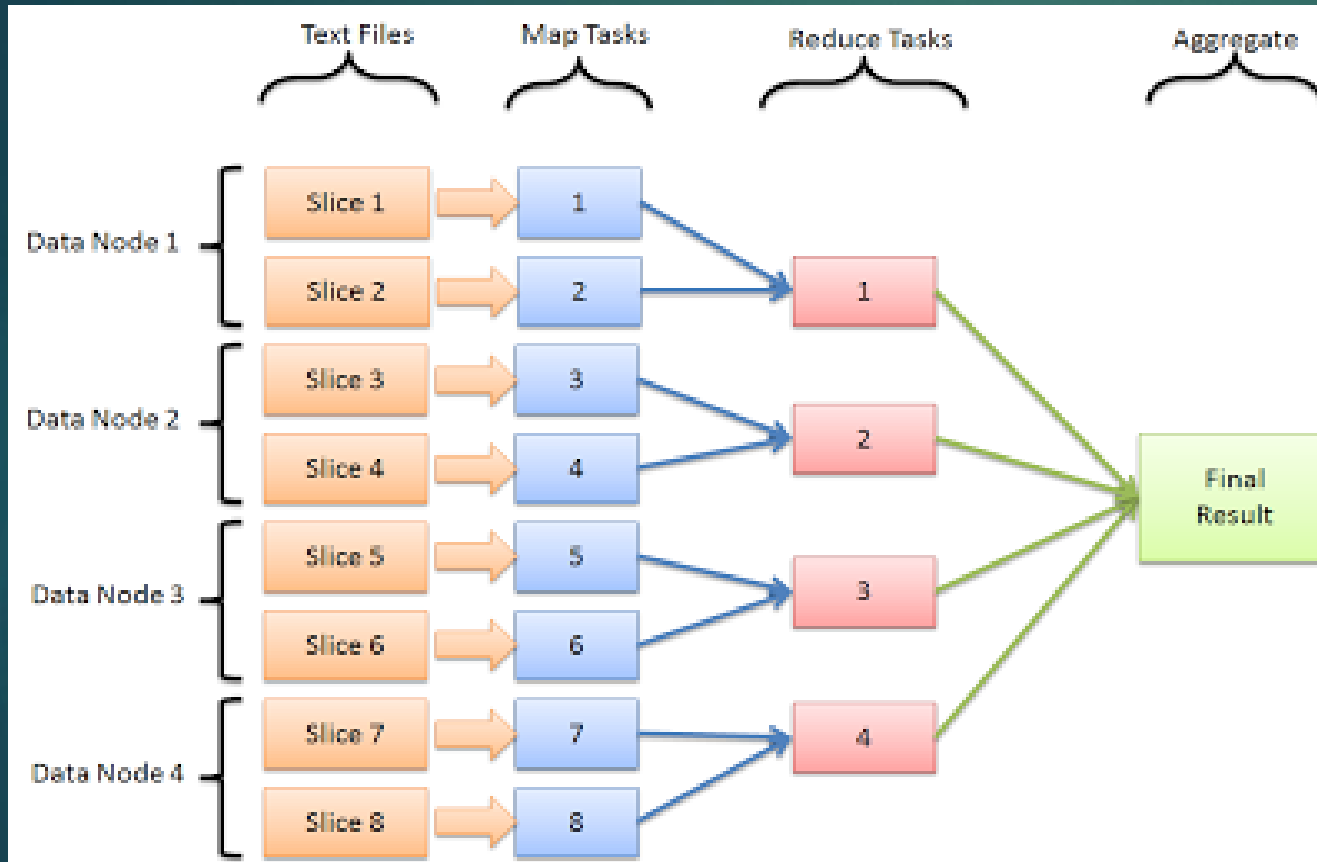
ABCB11 p52\_297\_2\_gag\_grg  
DB data: 297 2 A>G Glu-Gly  
11568372 pathogenic germline  
<http://cmim.org/entry/603201#0002>

1

2

3

# Modelo Map Reduce.





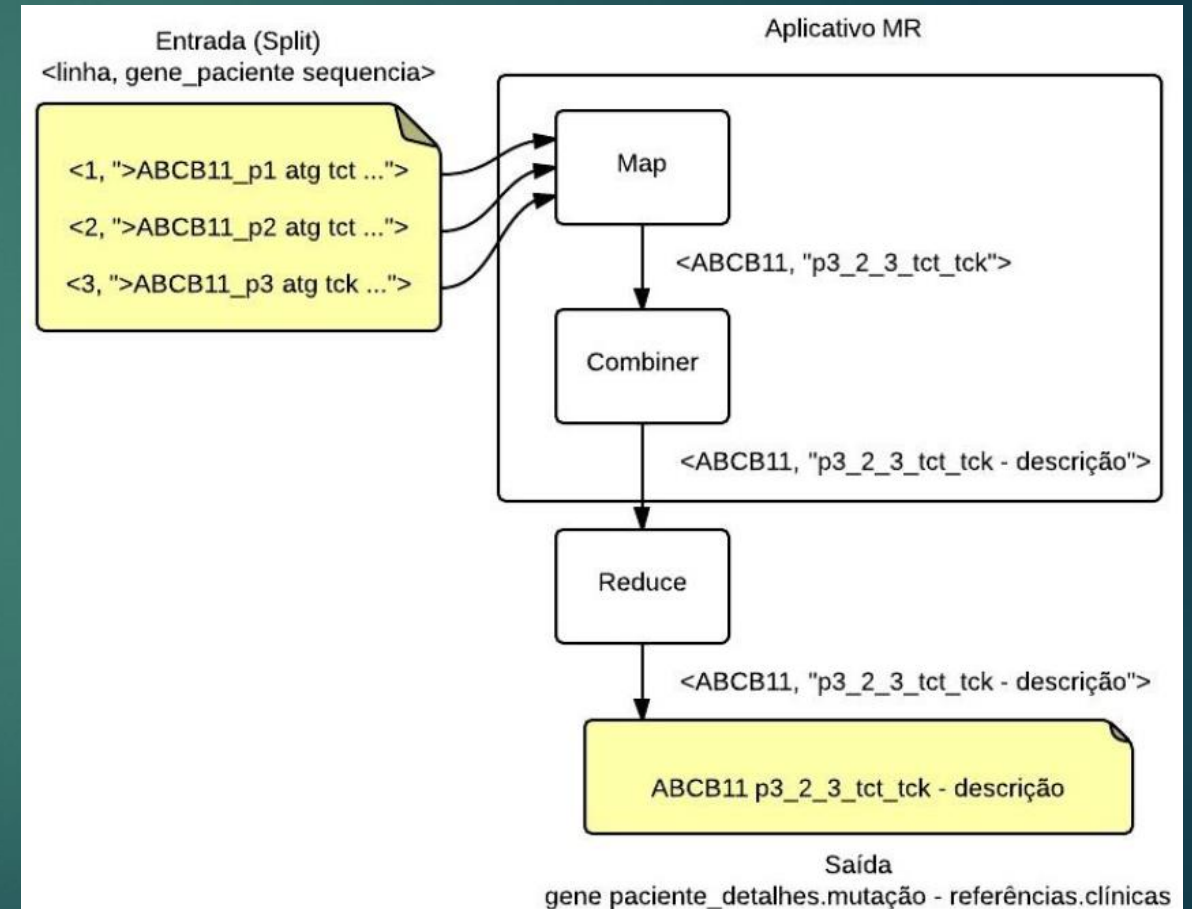
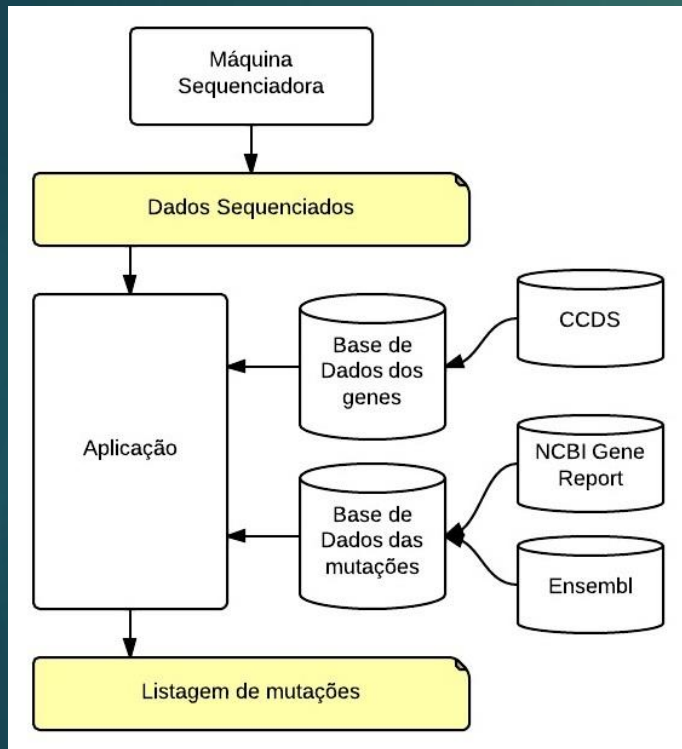
# Implementação.

- ▶ Dificuldades encontradas:

- ▶ Número de linhas de código: 253 (Java) x 476 (C++) : +88%
- ▶ Tamanho do código executável 6.6KB (Jar) 1.3MB (Bin).
- ▶ Pouca ou nenhuma documentação existente para C++.
- ▶ Poucas as estruturas suportadas em C++ para manipulação de Strings.
- ▶ Código C++ declarou em código as bases de referência:
  - ▶ Java utilizou a cache distribuída do Hadoop.
- ▶ Várias bibliotecas auxiliares foram necessárias ao C++.
- ▶ Processo de compilação e execução extremamente complexo.
- ▶ Foi necessário correção e recompilação de bibliotecas (util, pipes).
- ▶ Processo de instalação e configuração do Hadoop é o mesmo que Java.



# Funcionamento da aplicação



# Resultados.

- ▶ Cluster utilizado:

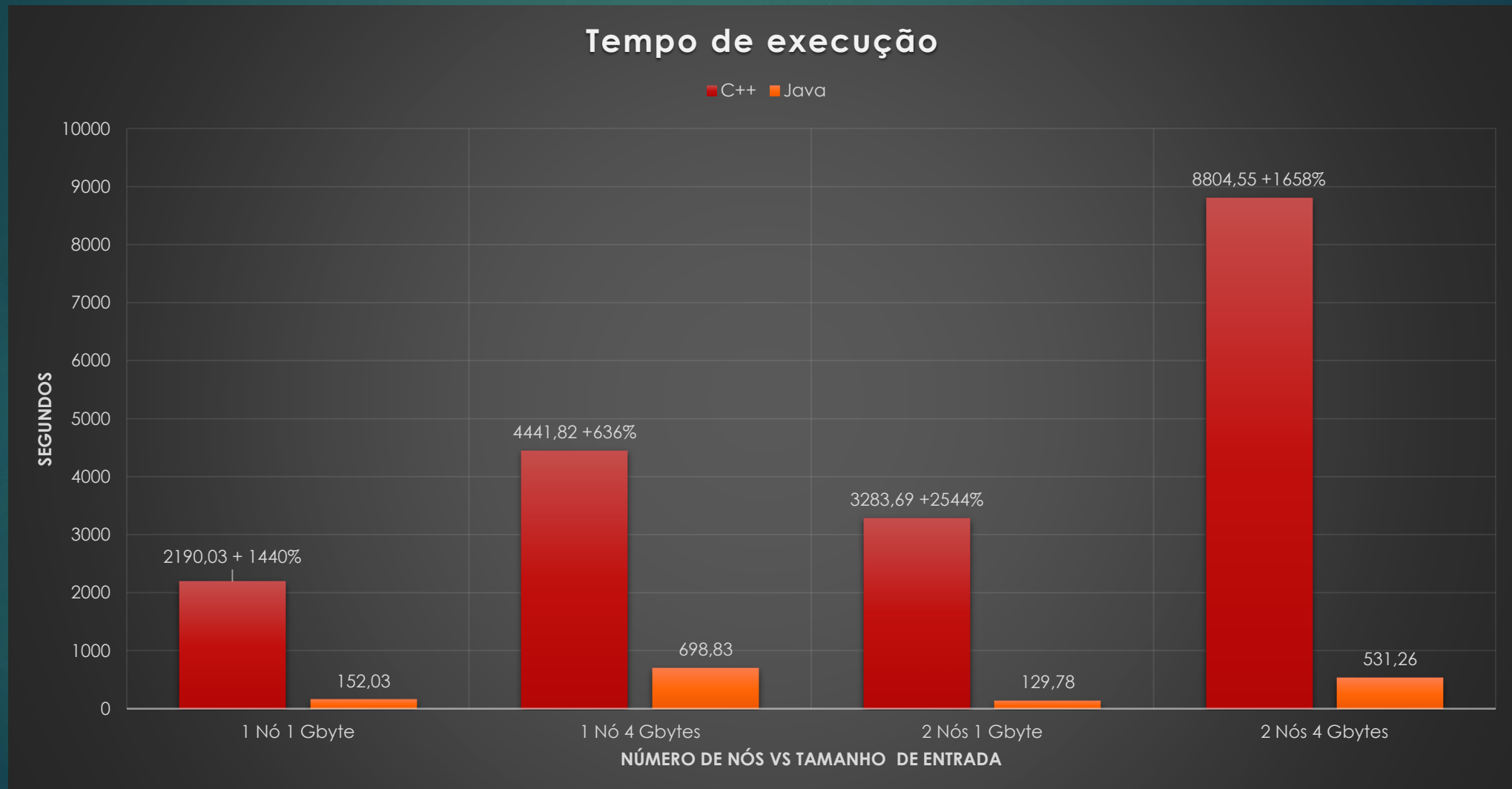
- ▶ 4 Máquinas GPPD:

- ▶ 4 Nós SMP (16 Cores):

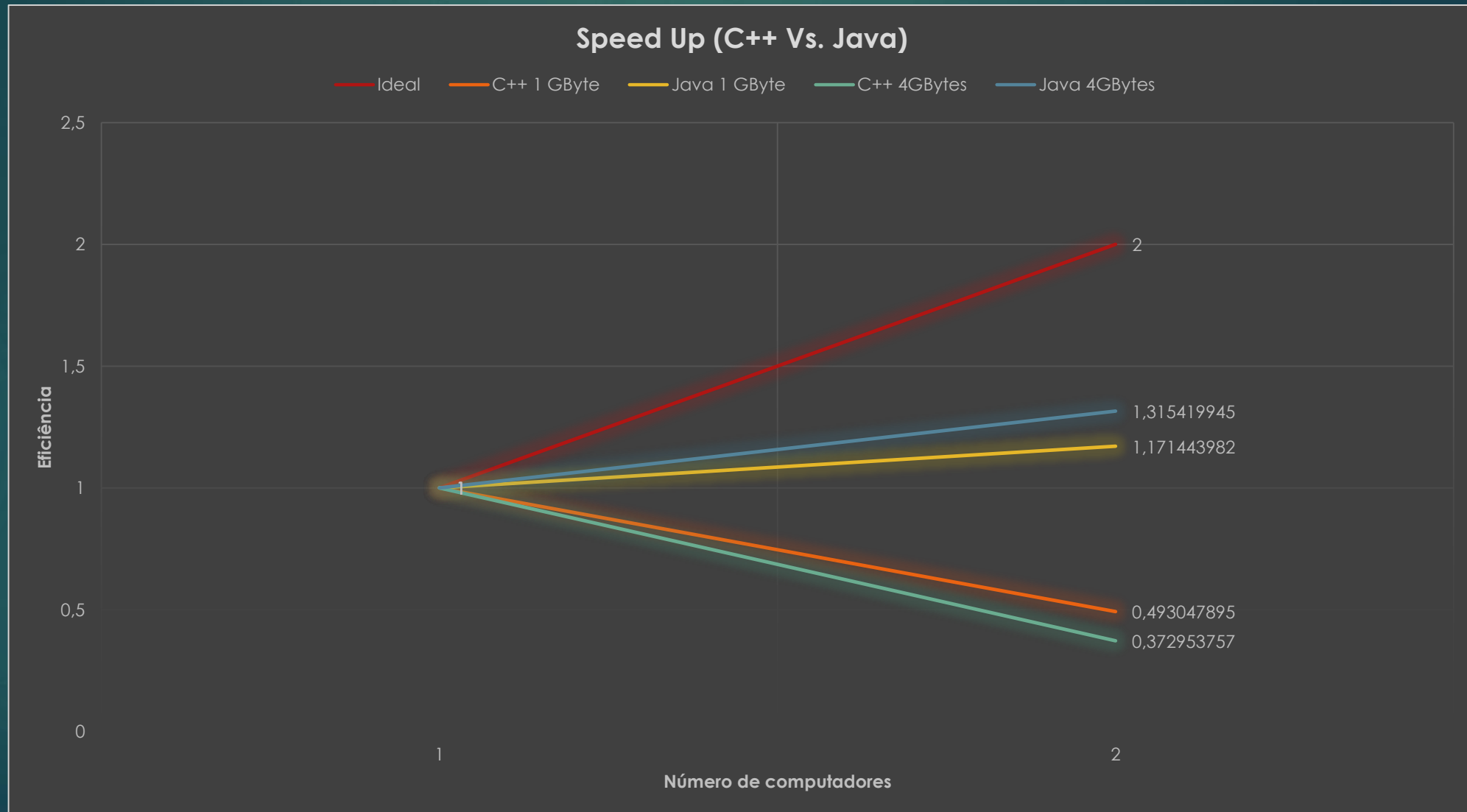
- ▶ 3 C2Q Q8200 (2.33GHz) (4 CPUs) com 4Gbytes cada.
      - ▶ i5 3570 (3.2GHz) (4 CPUs) com 6Gbytes cada.
      - ▶ Rede Ethernet 100Mbps (Compartilhada).
      - ▶ Discos Sata-I ~40Mbytes (Read/Write).

- ▶ Hadoop 1.0.4.

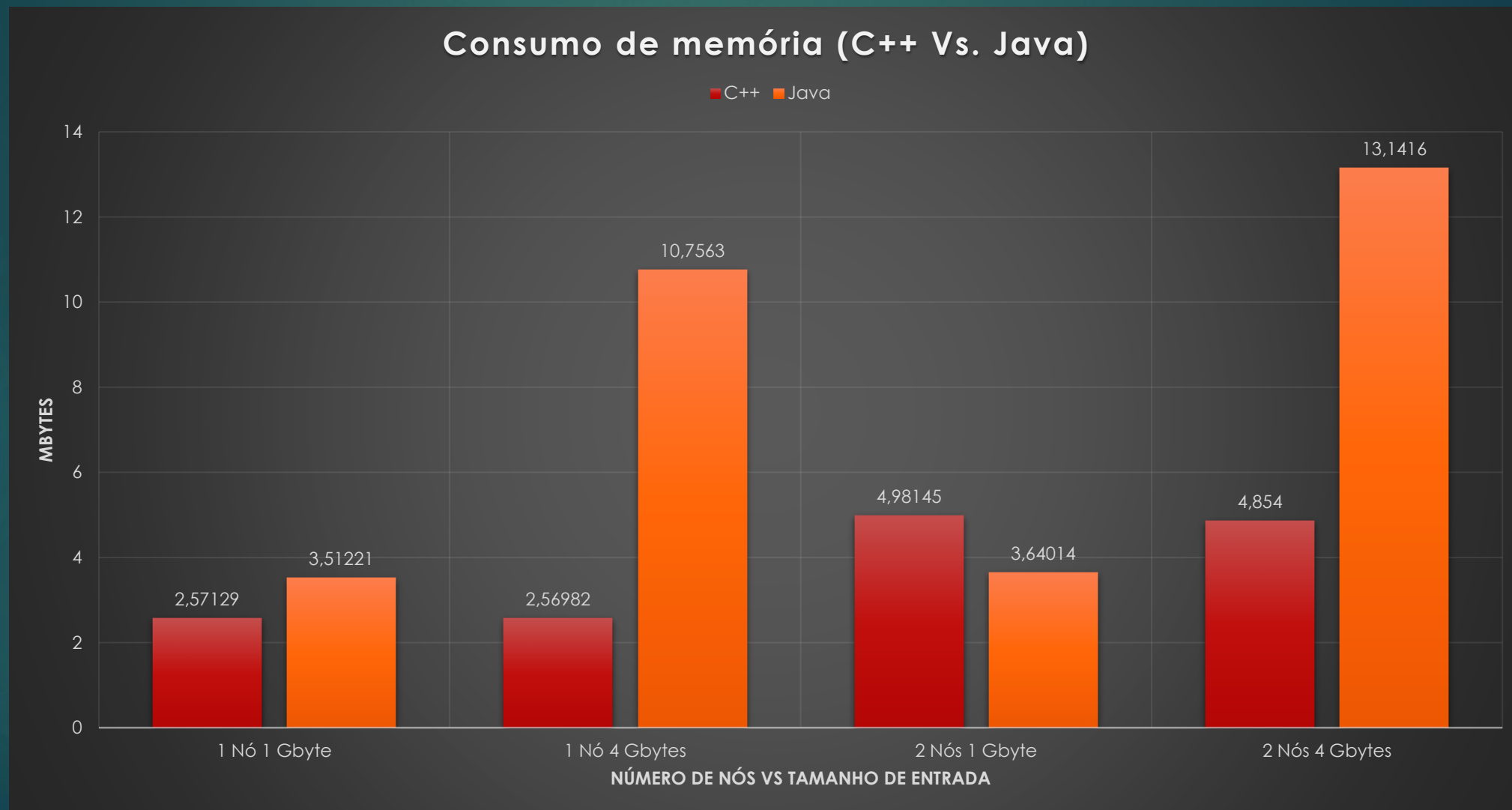
# Resultados. (Tempo de execução)



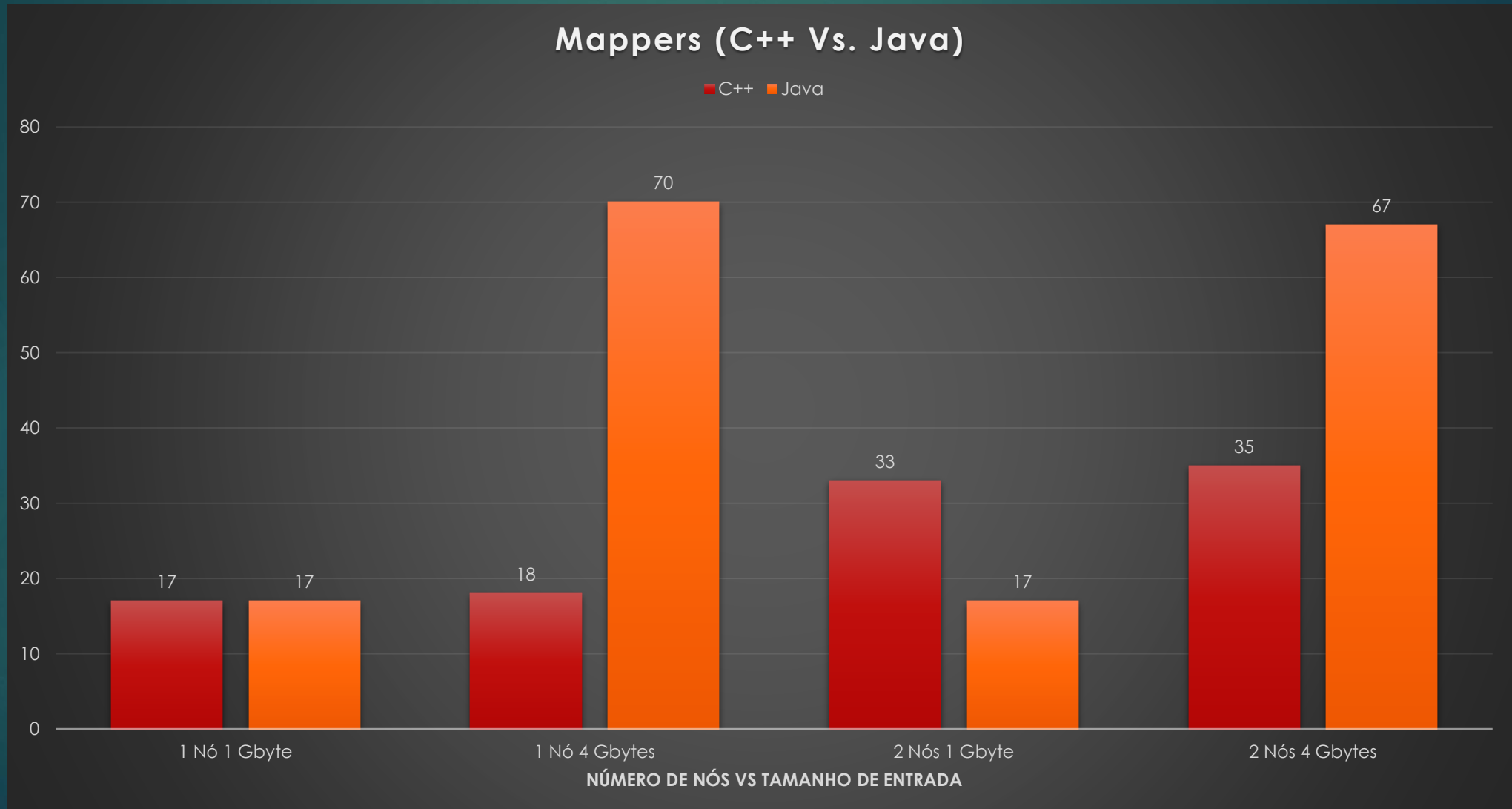
# Resultados. (Speed Up)



# Resultados. (Consumo de Memória)

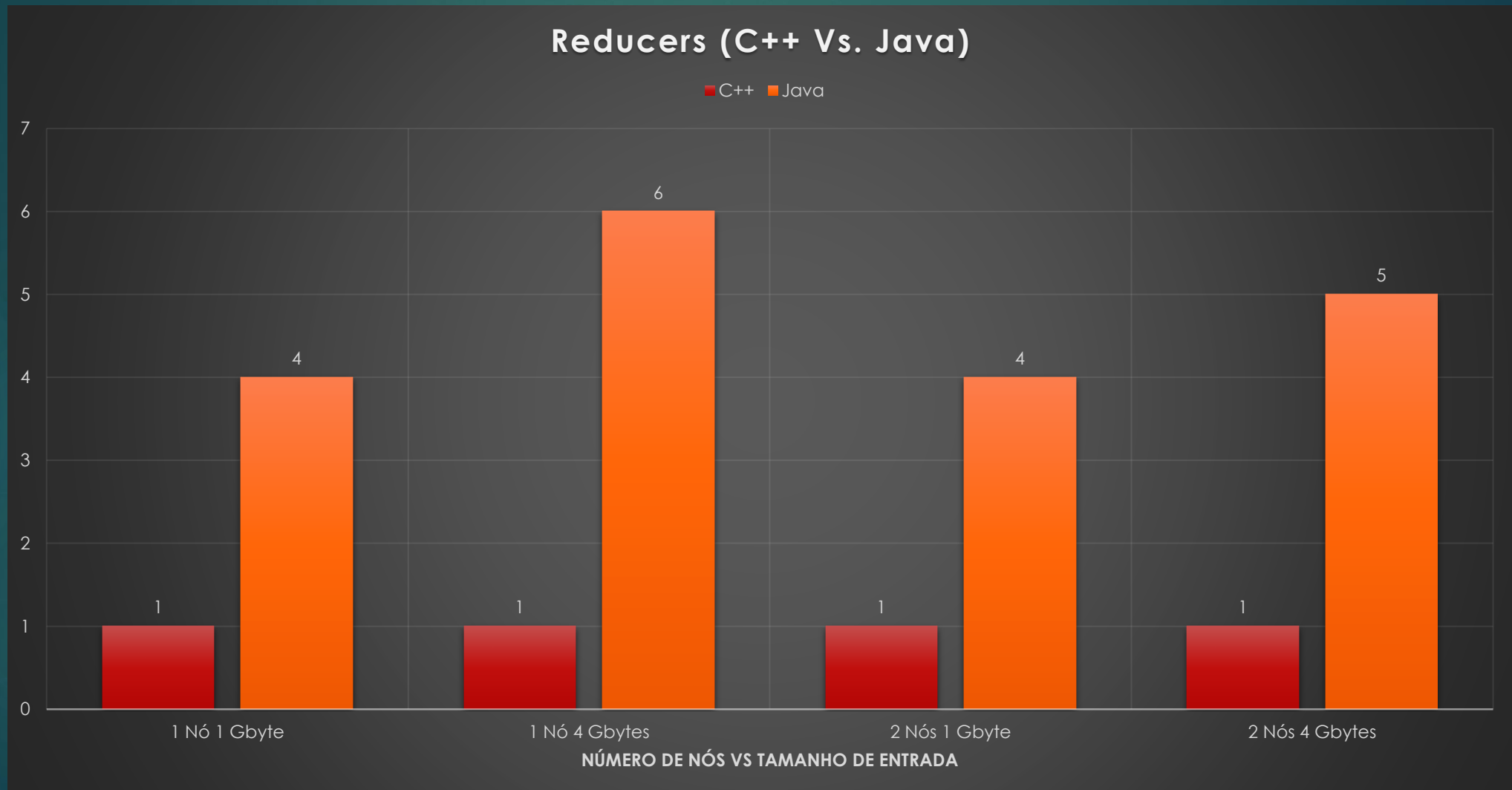


# Resultados. (Mappers)

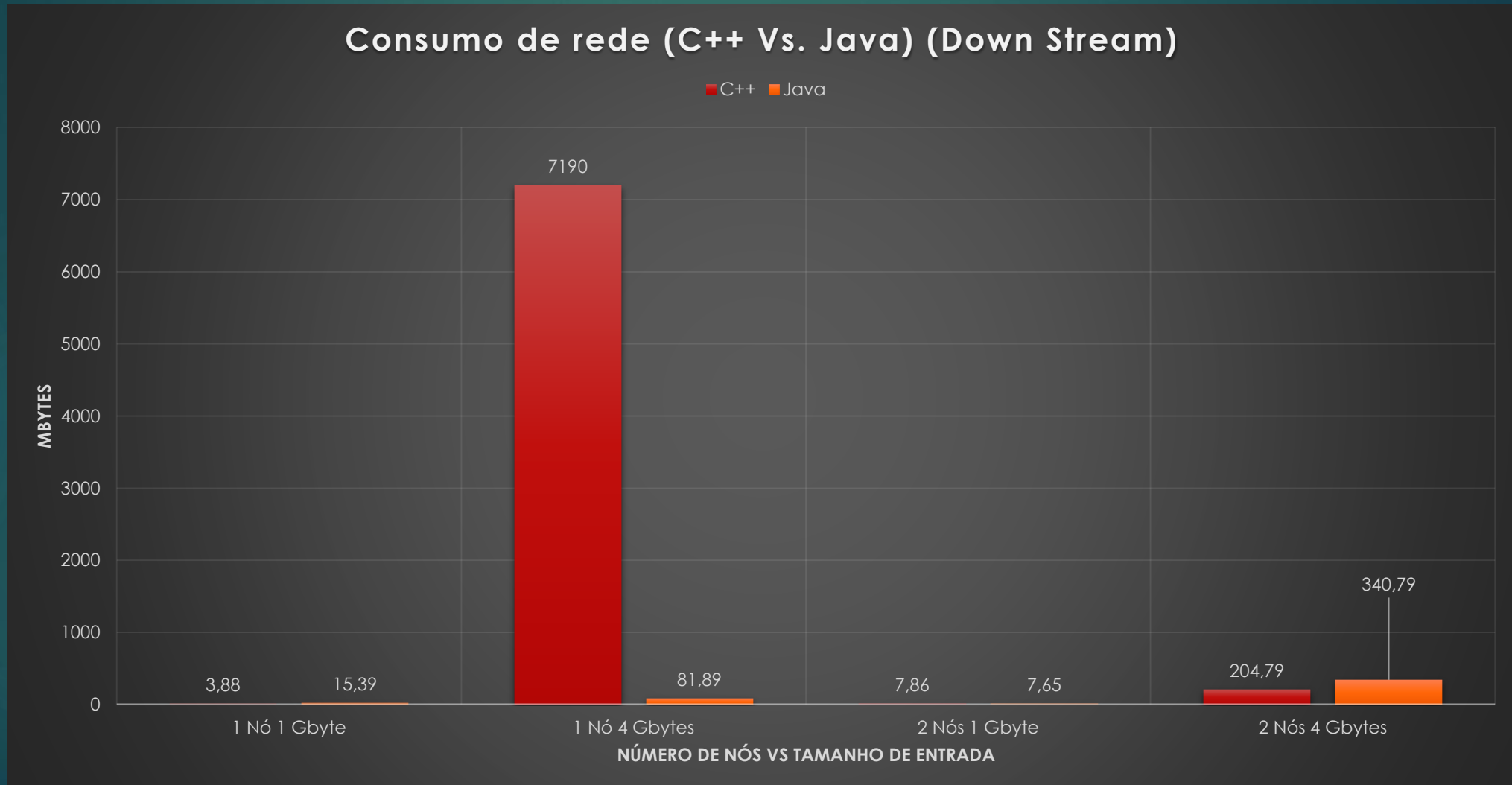




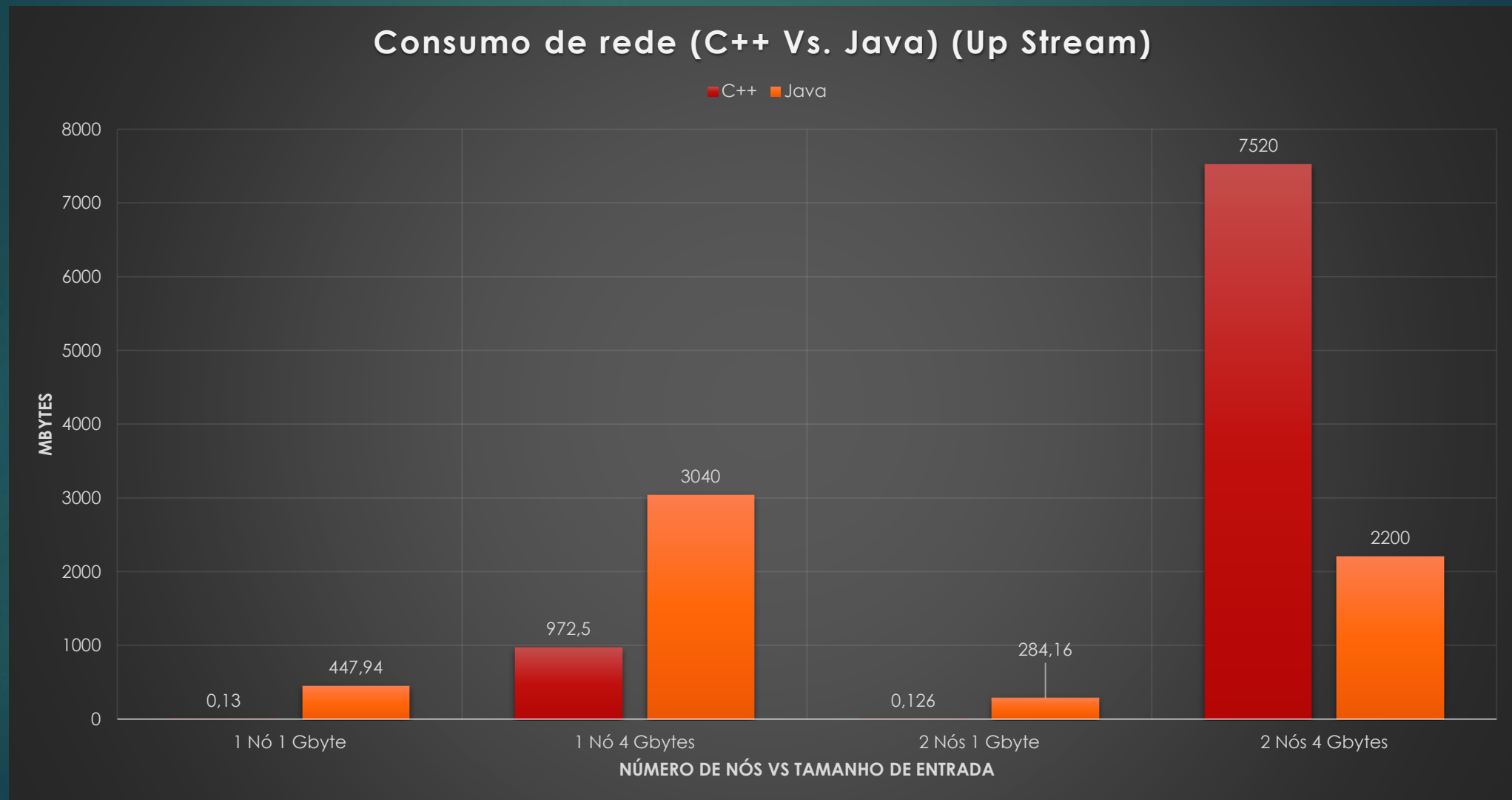
# Resultados. (Reducers)



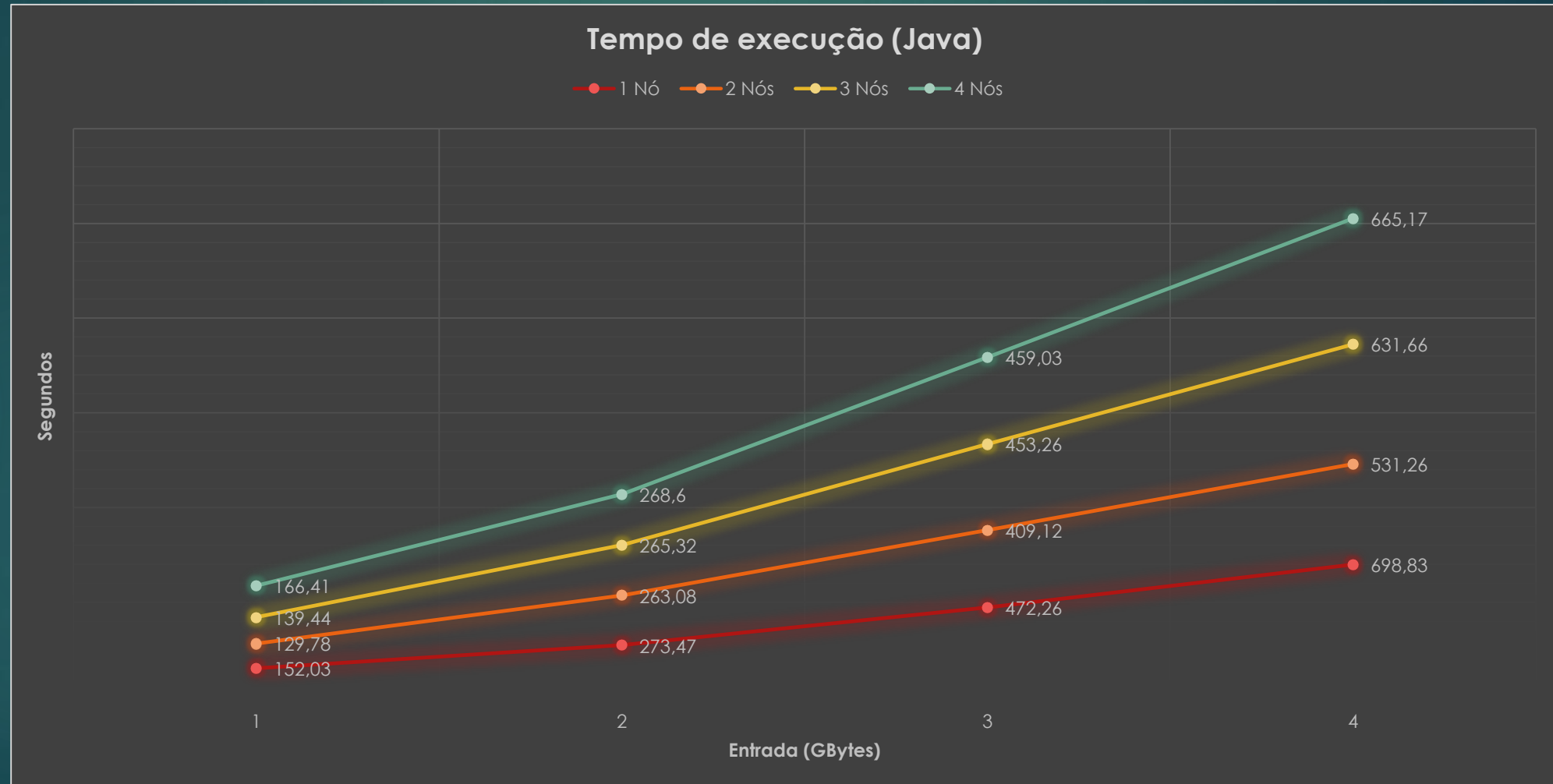
# Resultados. (Network Down)



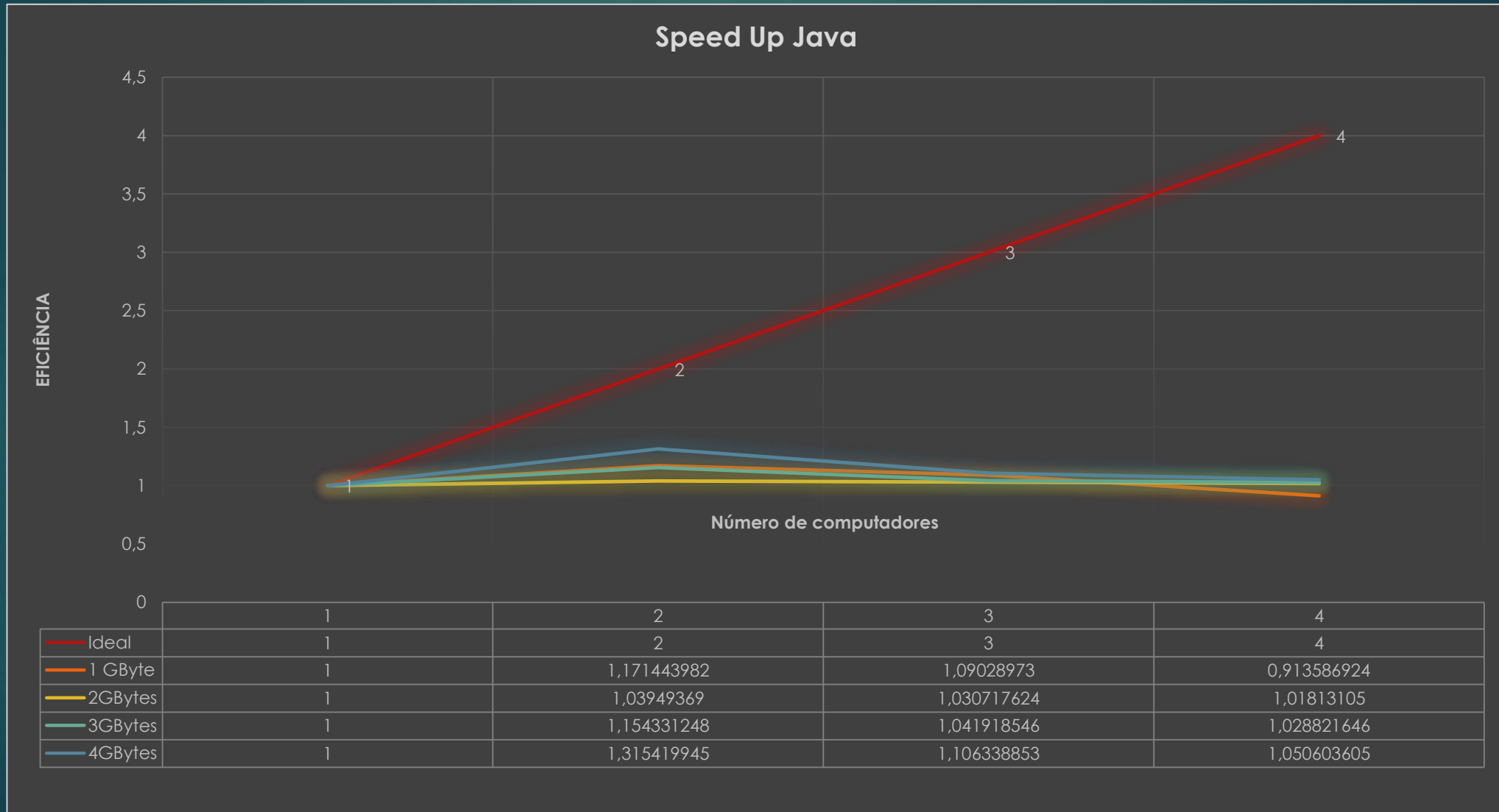
# Resultados. (Network Up)



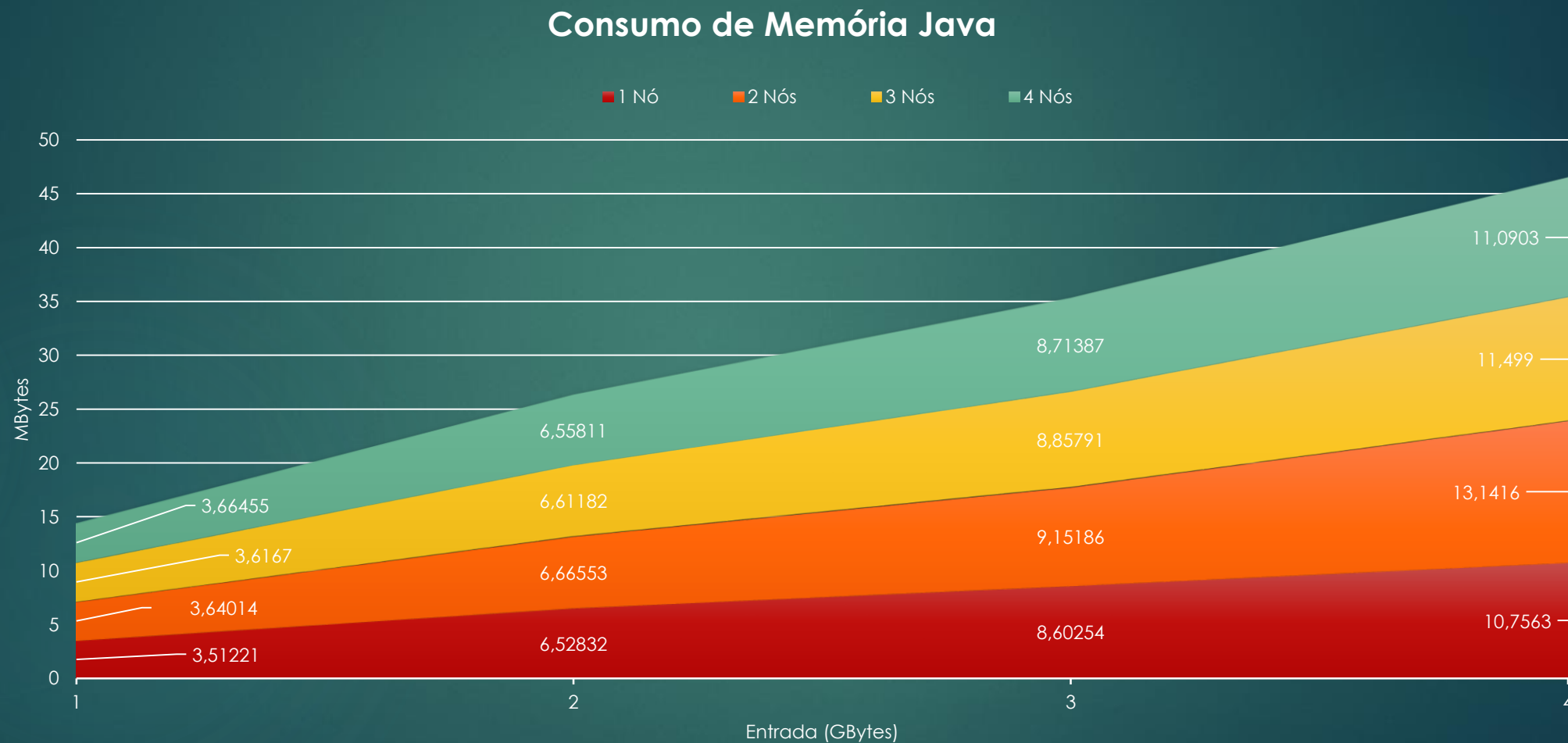
# Resultados. (Tempo de execução)



# Resultados. (Speed Up)

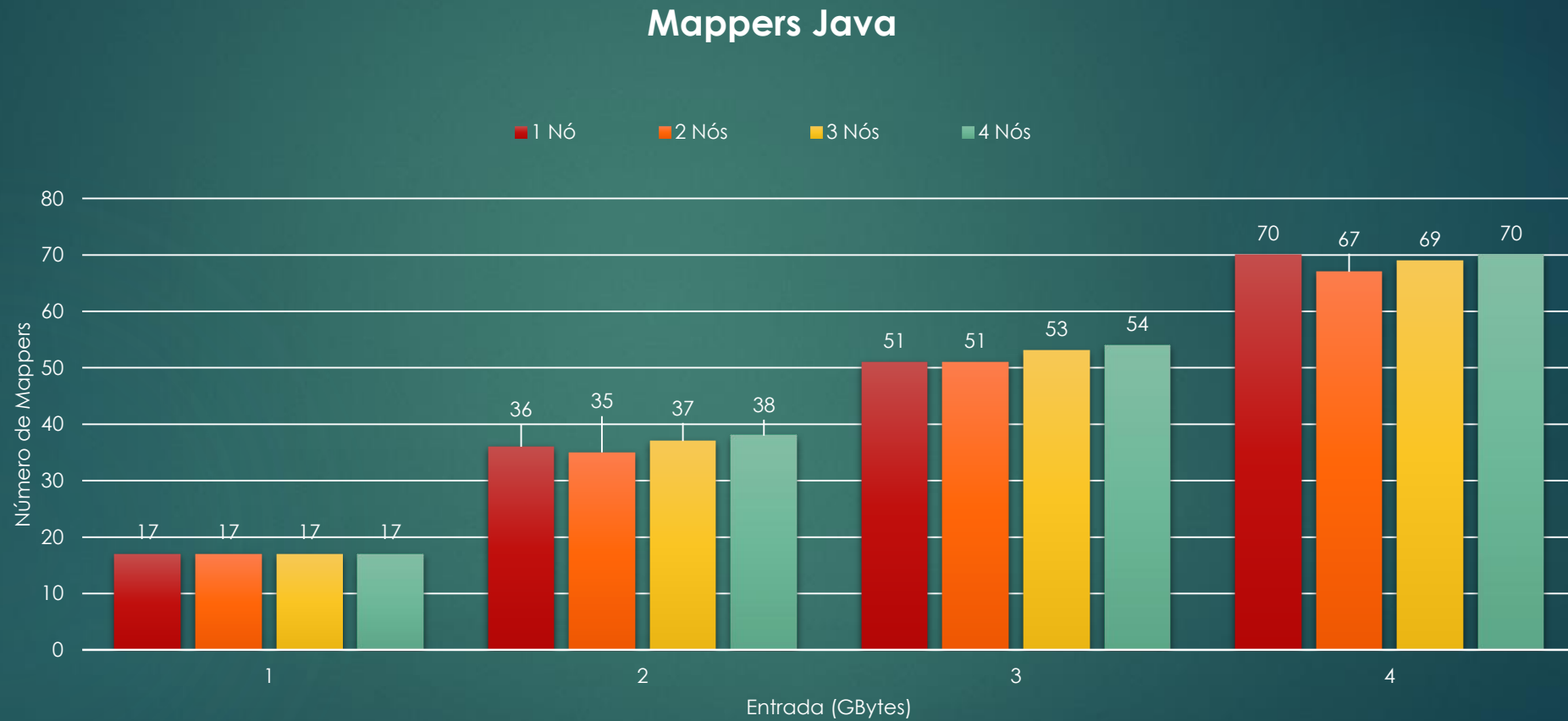


# Resultados. (Consumo de Memória)

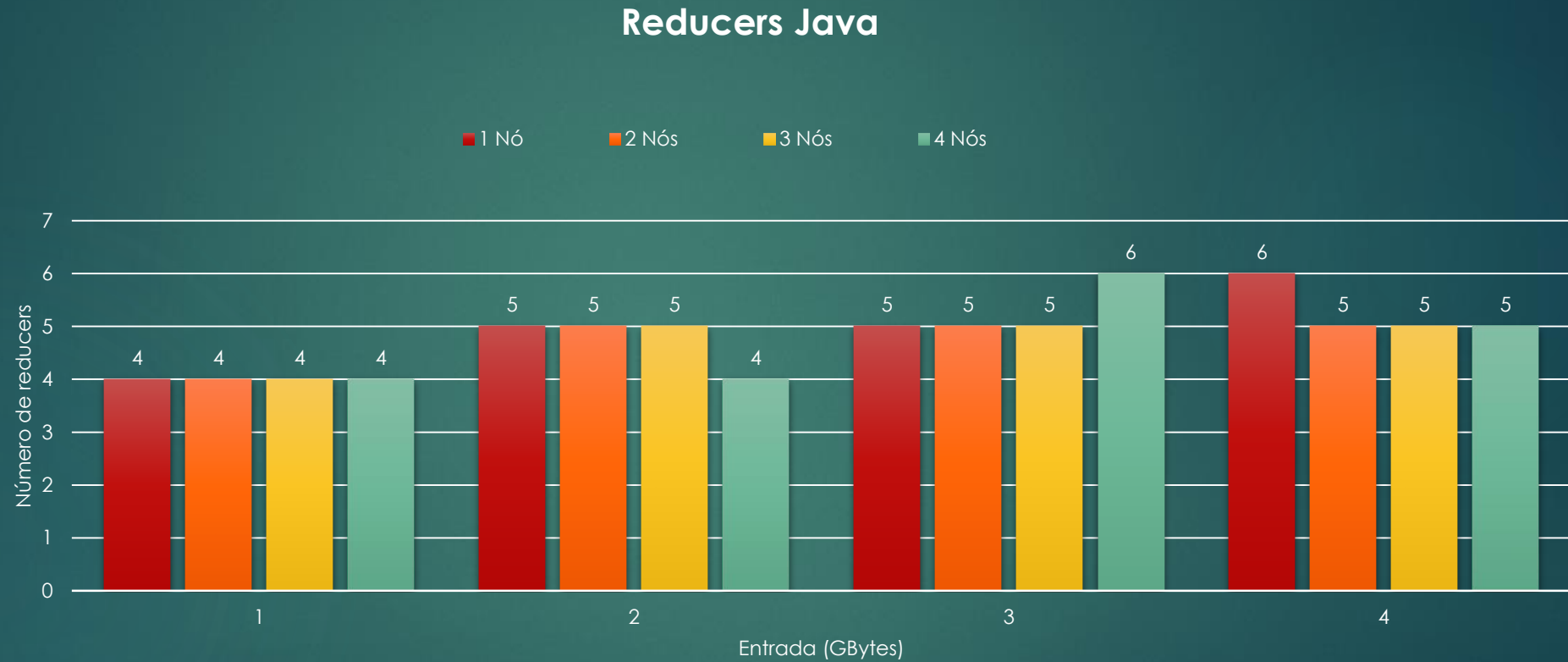




# Resultados. (Mappers)

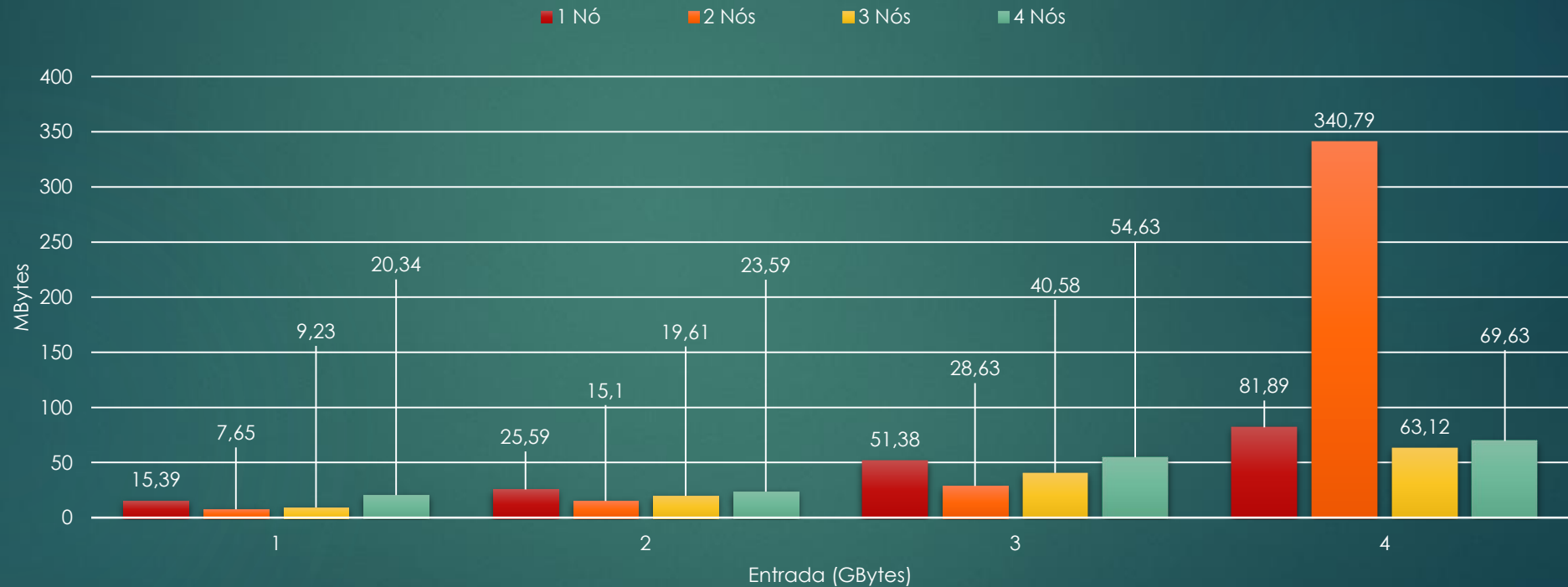


# Resultados. (Reducers)



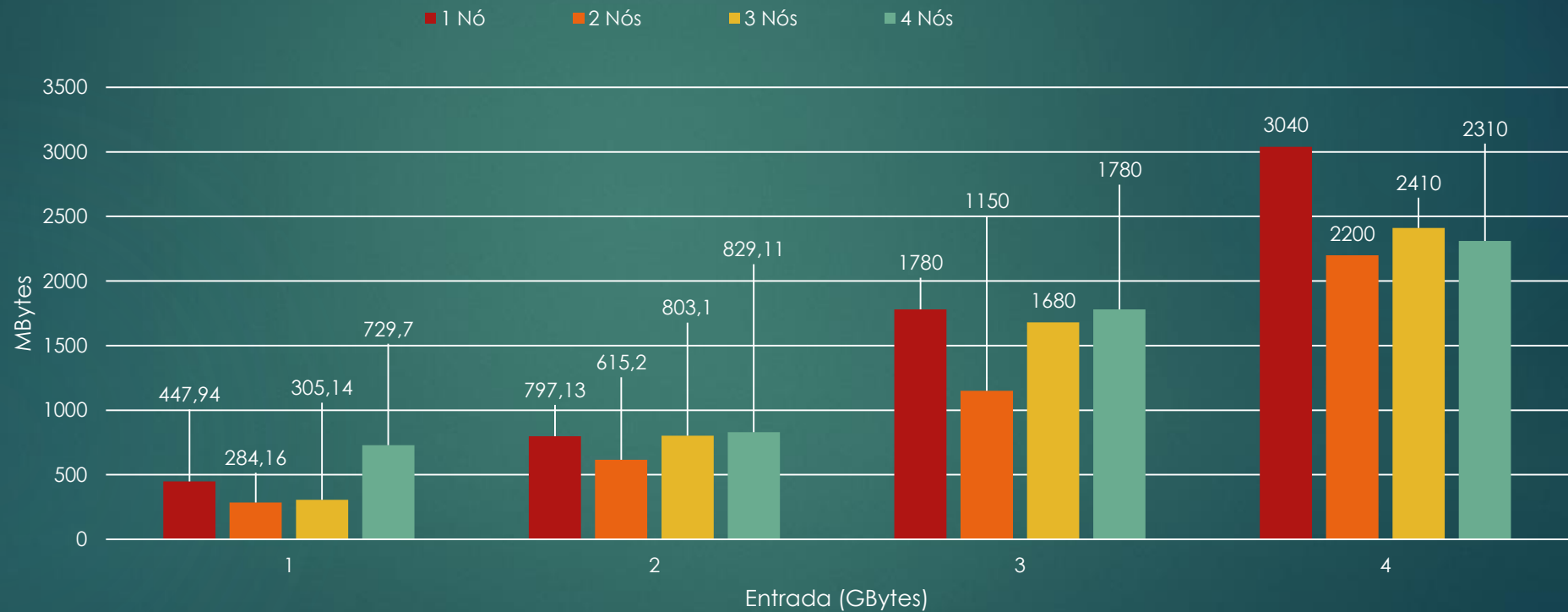
# Resultados. (Network Down)

## Consumo de rede Java (Down Stream)



# Resultados. (Network Up)

## Consumo de rede Java (Up Stream)



# Conclusões e trabalhos futuros.

- ▶ Com relação a arquitetura distribuída:
  - ▶ Cluster pode não dar suporte a aplicação.
  - ▶ Volumes de entrada são custosos para executar/avaliar.
    - ▶ Dado o tempo de execução.
  - ▶ Configuração de um ambiente distribuído é custoso.
- ▶ Com relação ao uso de C++:
  - ▶ Resultados demonstraram que C++ não apresenta competitividade comparado a Java, tanto no consumo de recursos como no tempo de execução.
  - ▶ C++ não apresenta a mesma facilidade de programação e uso do Hadoop quando comparado a Java.
  - ▶ C++ limita a portabilidade da aplicação para ambientes heterogêneos.

# Conclusões e trabalhos futuros.

- ▶ Com relação a aplicação de sequenciamento de DNA:
  - ▶ É limitada as bases de referencia.
  - ▶ Não dá suporte a busca de outros tipos de genes.
  - ▶ Pelas avaliações pode não ser escalável.
  - ▶ Amostras sequenciadas podem vir sobre outros formatos (BAM?) e não estarem alinhadas.
- ▶ Com relação a trabalhos futuros:
  - ▶ Detectar sobre quais aspectos o Map Reduce pode ser interessante para a aplicação proposta.
  - ▶ Execução deste trabalho no cluster gradeP.
  - ▶ Identificar outros aspectos da aplicação a serem utilizados pelo modelo como por exemplo:
    - ▶ Processo de alinhamento.
    - ▶ Busca nos arquivos de referencia pelos genes a serem sequenciados.
    - ▶ Busca nas bases de mutações por genes utilizados durante o sequenciamento.





Perguntas ?