

Introduction to Bigdata

Distributed Systems

Hadoop and its role in Bigdata

Spark Introduction.

Benefits of using PySpark

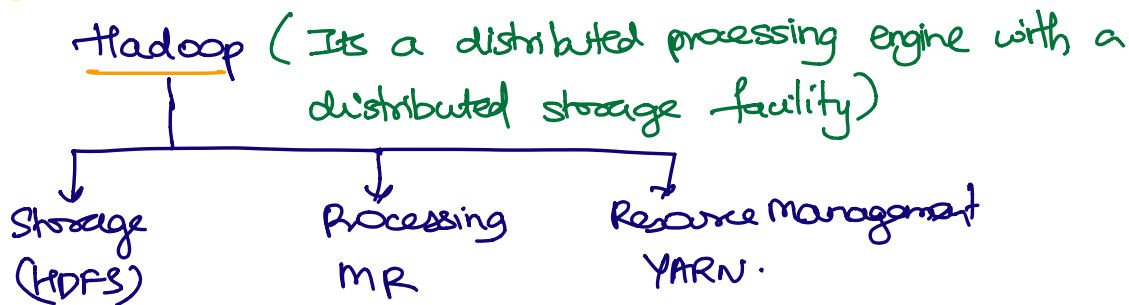
What is Bigdata? → data with huge size.

201201hardy.txt → 507 MB

Notepad → Failed → **Bigdata**

Waldpad → Had some latency but → **NORMAL DATA**
loading was successful

Distributed Architectures.



- Benefit (Parallel Processing)

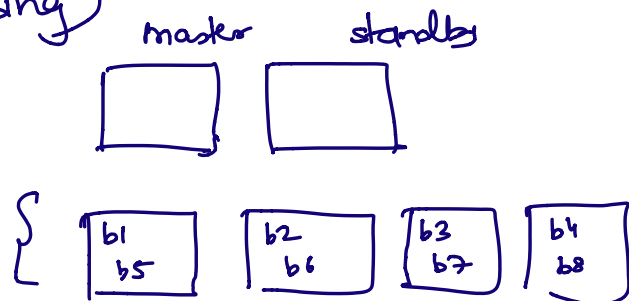
file.txt → 512 MB (8 blocks)

Hadoop splits your data in the

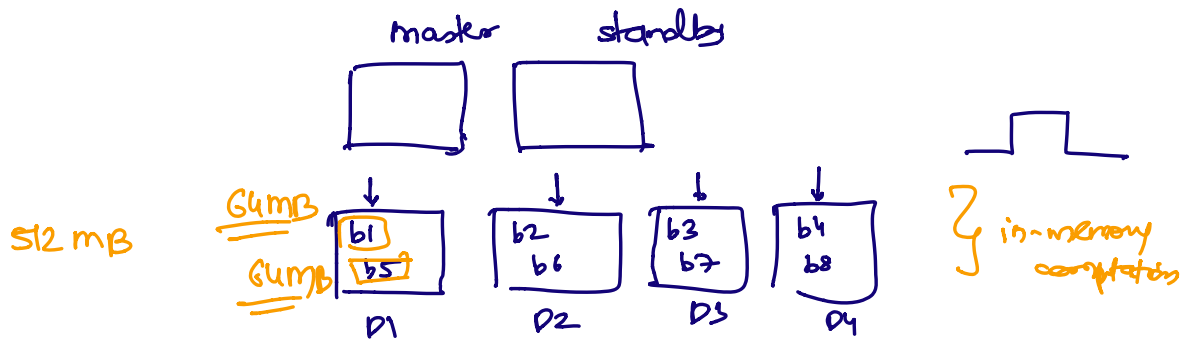
form of
blocks

maintain
data and
perform processing

B.S = 64 MB



when it comes to Processing, all systems will participate



Thus I advise parallel processing

