



HDFS/MapReduce

Eng. Software / PSPD

Prof. Fernando W Cruz



Sumário

- HDFS
- MapReduce/Hadoop



Contextualizando...

- Uso de **ambientes computacionais** com grande capacidade de armazenamento e processamento

- *Clusters* de computadores
- Computação em nuvem (*cloud computing*)



modelo que permite acesso ubíquo, conveniente e sob demanda via rede a um conjunto compartilhado e configurável de recursos computacionais (como redes, servidores, armazenamento, aplicação e serviços) que pode ser rapidamente fornecido e liberado com esforços mínimos de gerenciamento ou interação com o provedor de serviços

modelo composto por uma coleção de computadores dispostos de forma paralela e distribuída e interconectados por redes de alta velocidade



Contextualizando...

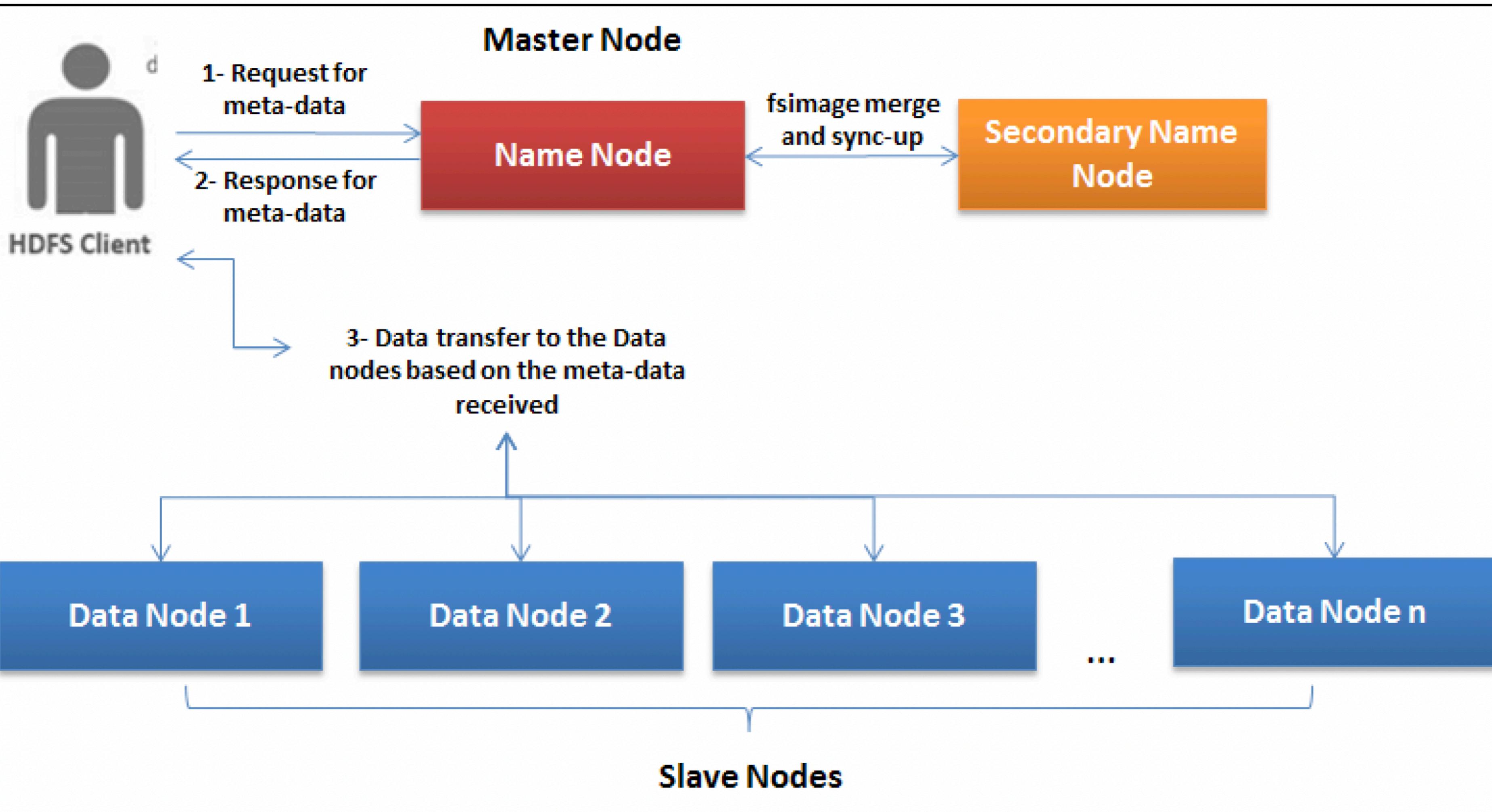
- Uso de sistemas de arquivos distribuídos
 - HDFS (*Hadoop Distributed File System*)



sistema de arquivos distribuído para armazenar grandes quantidades de dados, com alta tolerância a falhas, e capaz de ser empregado em *hardware* de baixo custo
(baseado no Google File System)

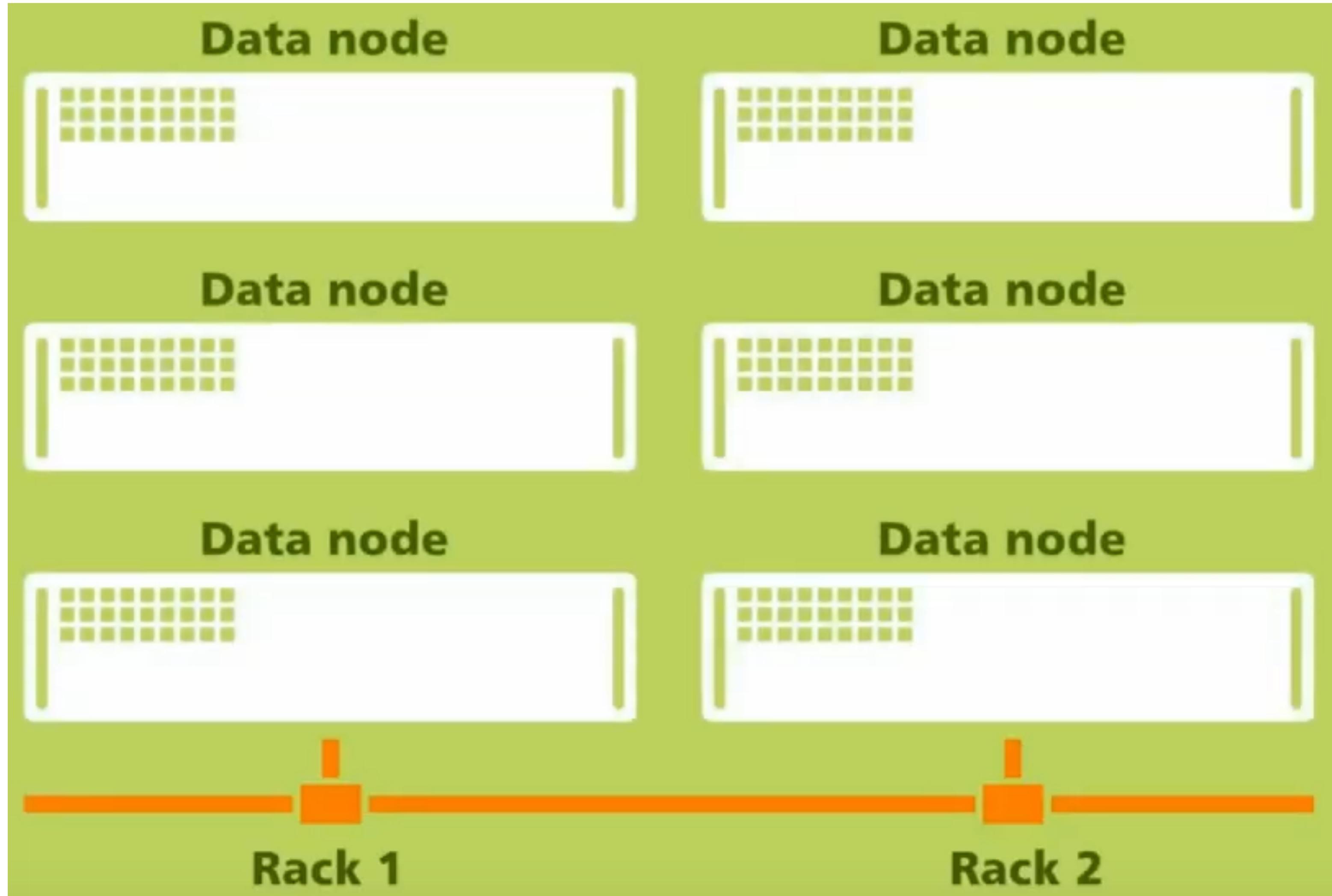


HDFS - visão dos nós de um cluster



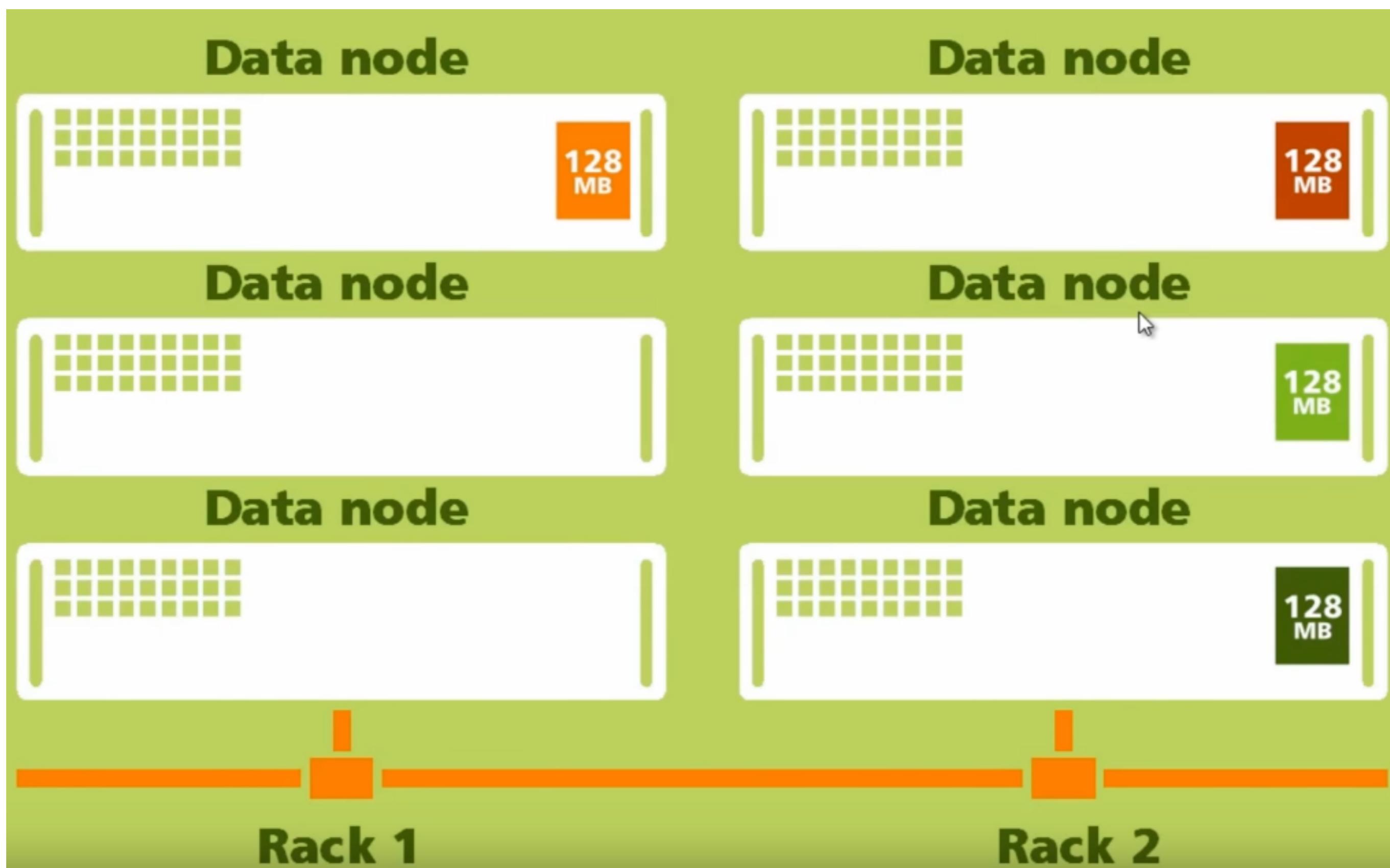
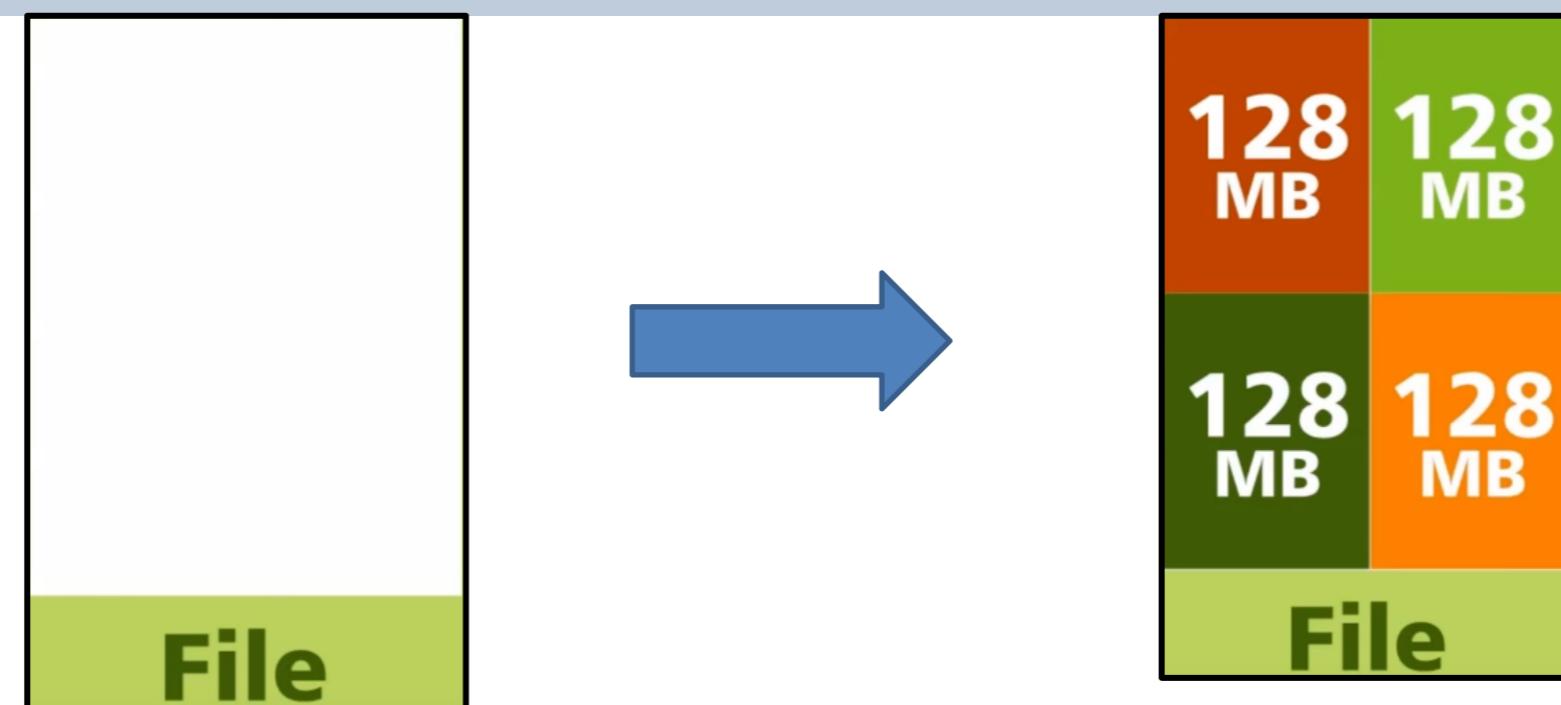


HDFS - visão dos nós de um cluster





HDFS





HDFS

Data node



Data node



Data node



Data node



Data node



Data node



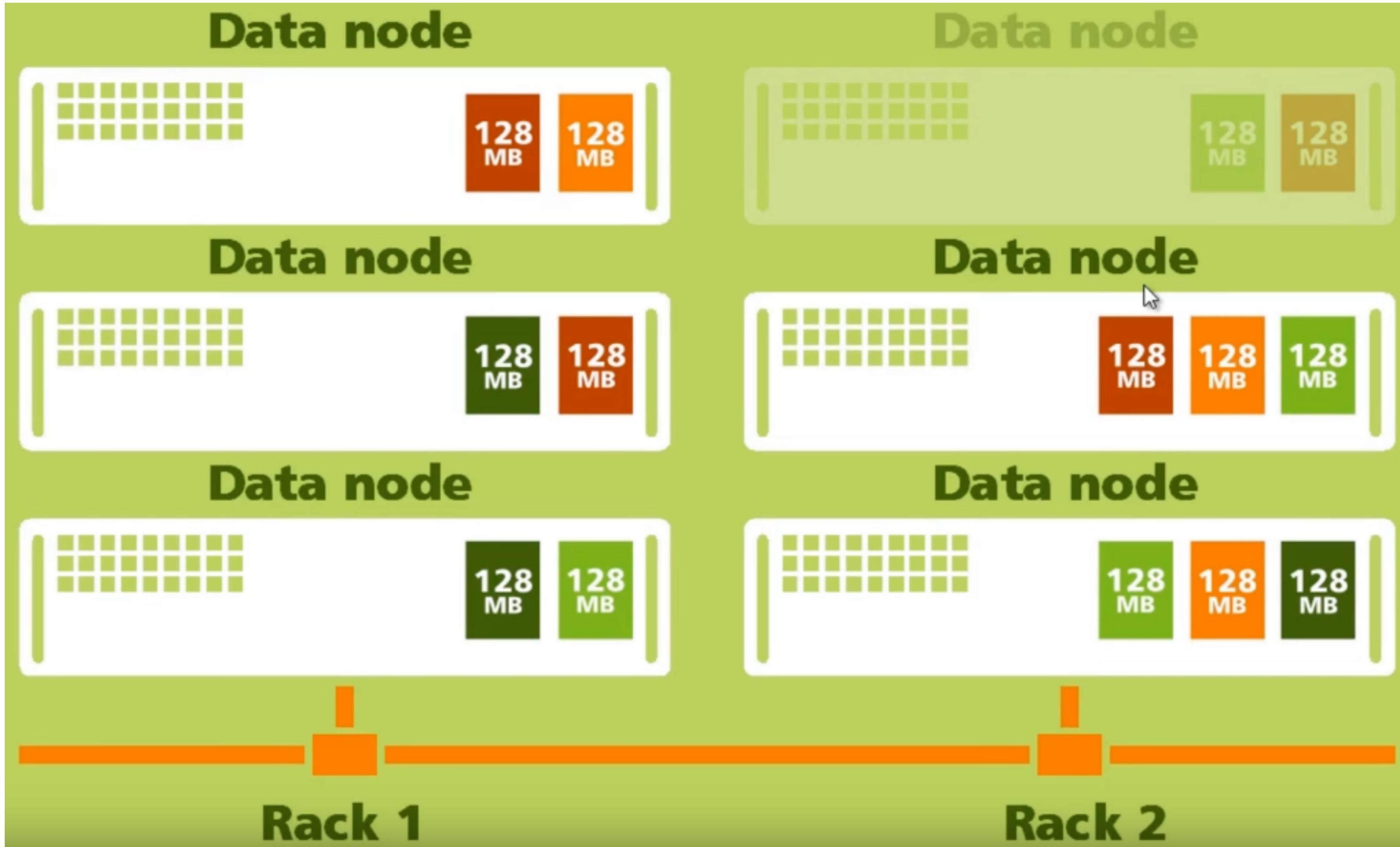
Rack 1



Rack 2



HDFS



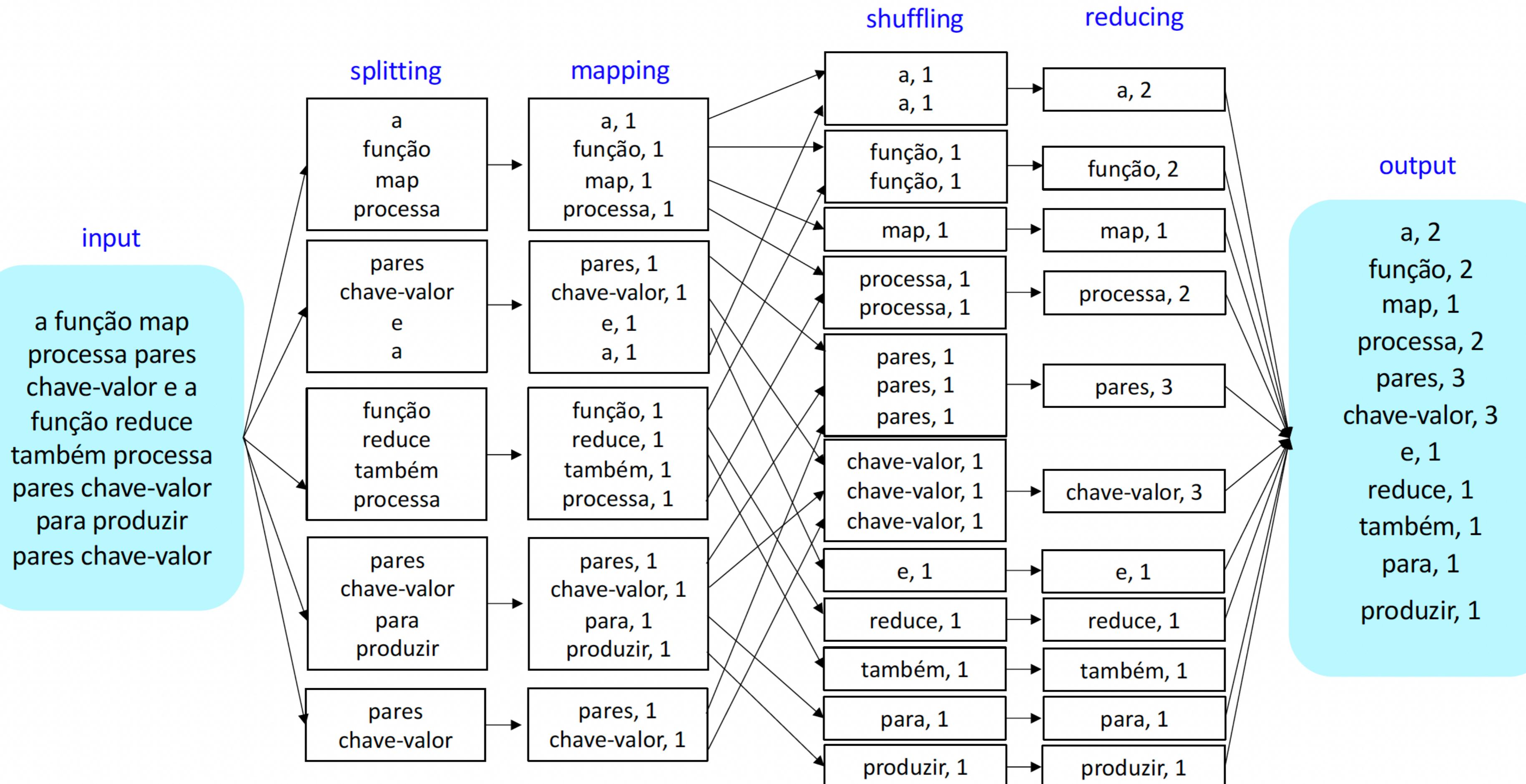


MapReduce

- Algoritmo clássico proposto pelo Google
- Modelo de programação funcional
- A programação envolve dois estágios: Map e Reduce
- No Map:
 - A entrada são os dados brutos
 - A saída é um conjunto de pares <chave, valor>
- No Reduce:
 - A entrada são os diversos pares <chave, valores[]>
 - A saída são novos pares reduzidos com <chave, valores_contabilizados>

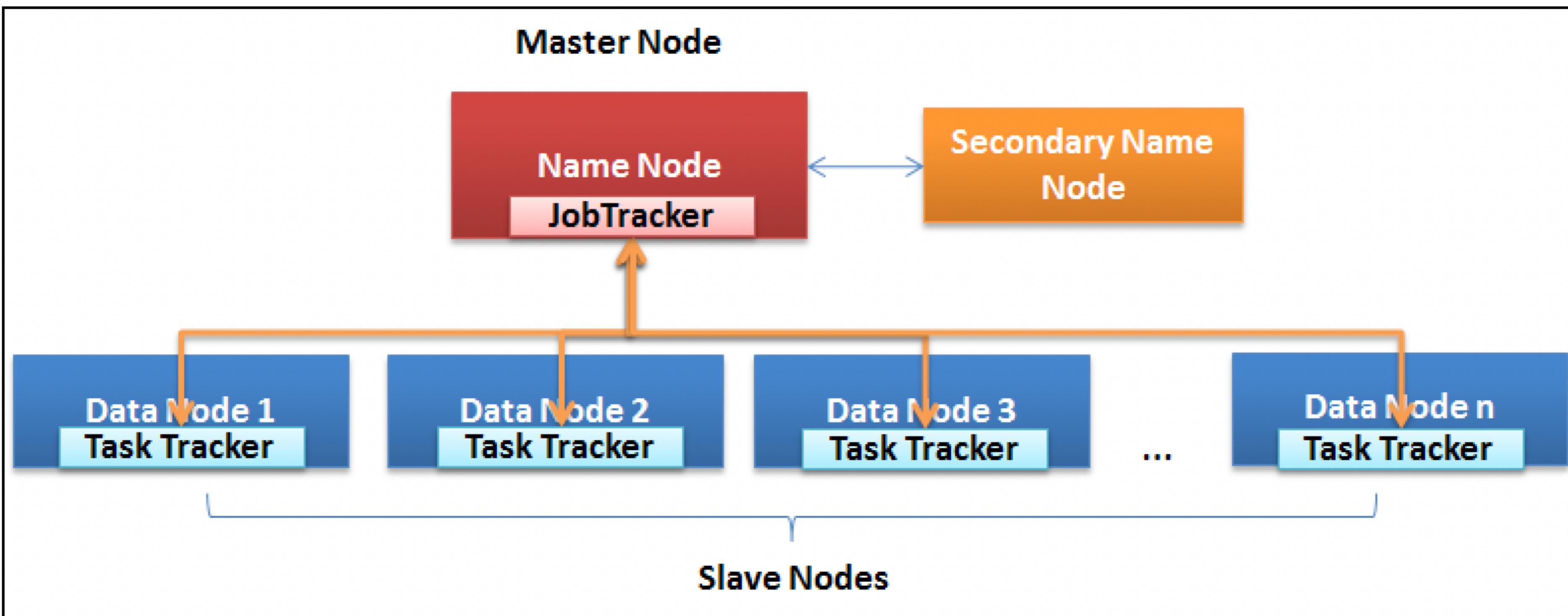


MapReduce - Exemplo





Arquitetura MapReduce



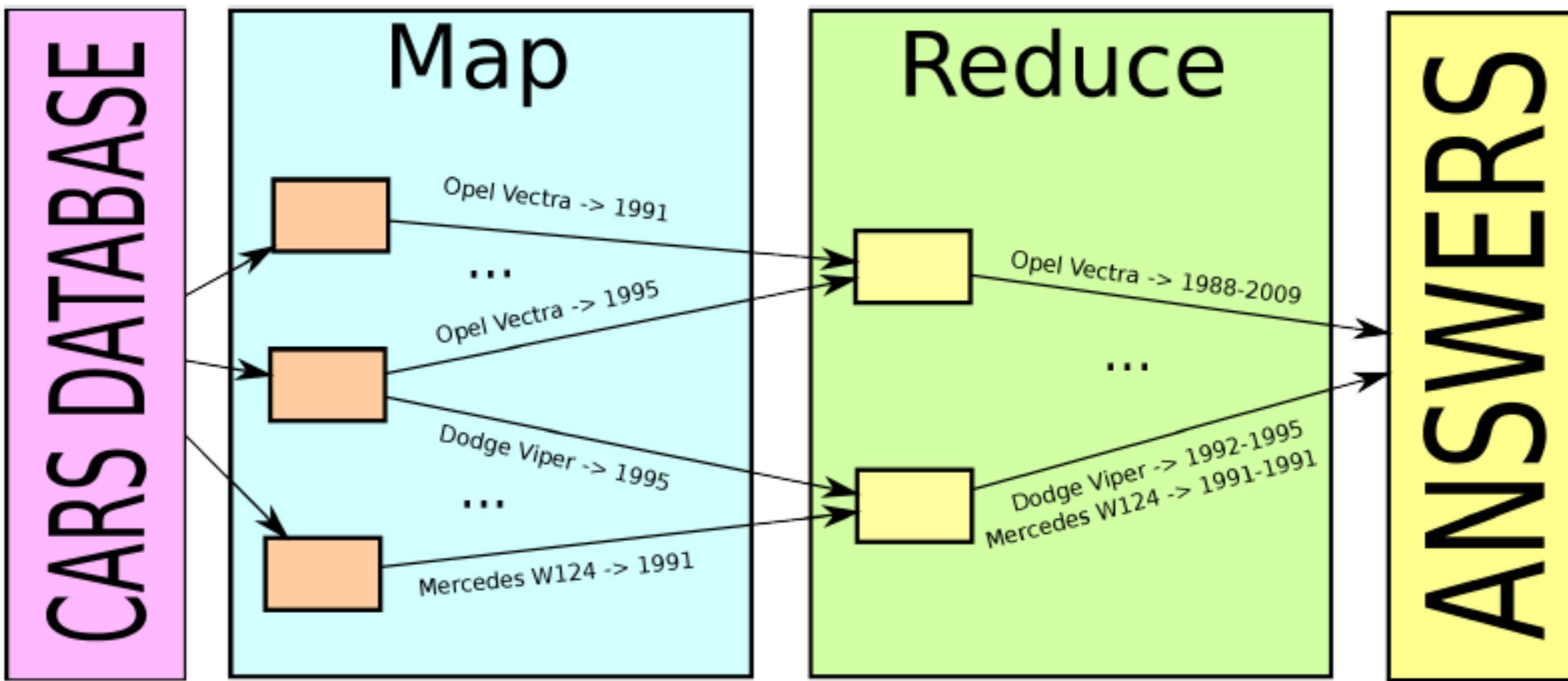


MapReduce - exemplos

- Dada uma base de carros, registrados por ano.
- Input:
 - 1991:
 - Opel Vectra, ABS
 - Mercedes W124, ABS, Airbag
 - 1995:
 - Opel Vectra, Ac, TC
 - Dodge Viper, ABS, airbag
 - ...
- Descobrir quais carros foram produzidos em qual ano
 - Opel Vectra: 1988-2009
 - Dodge Viper: 1992-1995
 - ...



MapReduce - exemplos



- No mapeamento (MAP)
 - A geração de resultados é independente
 - Tarefa paralela perfeita
- Na redução (REDUCE)
 - Pode ser executada em paralelo
 - Exige boa distribuição de entradas necessárias
- O algoritmo MapReduce é
 - Altamente paralelo
 - Largamente adotado
 - Boa solução para tarefas em batch



MapReduce – aplicações práticas



https://upload.wikimedia.org/wikipedia/commons/a/aa/Logo_Google_2013.Official.svg

https://upload.wikimedia.org/wikipedia/commons/8/8e/Hadoop_logo.svg

<https://upload.wikimedia.org/wikipedia/en/f/f8/CouchDB.svg>

https://upload.wikimedia.org/wikipedia/en/e/eb/MongoDB_Logo.png

https://upload.wikimedia.org/wikipedia/en/8/8e/Riak_distributed_NoSQL_key-value_data_store_logo.png

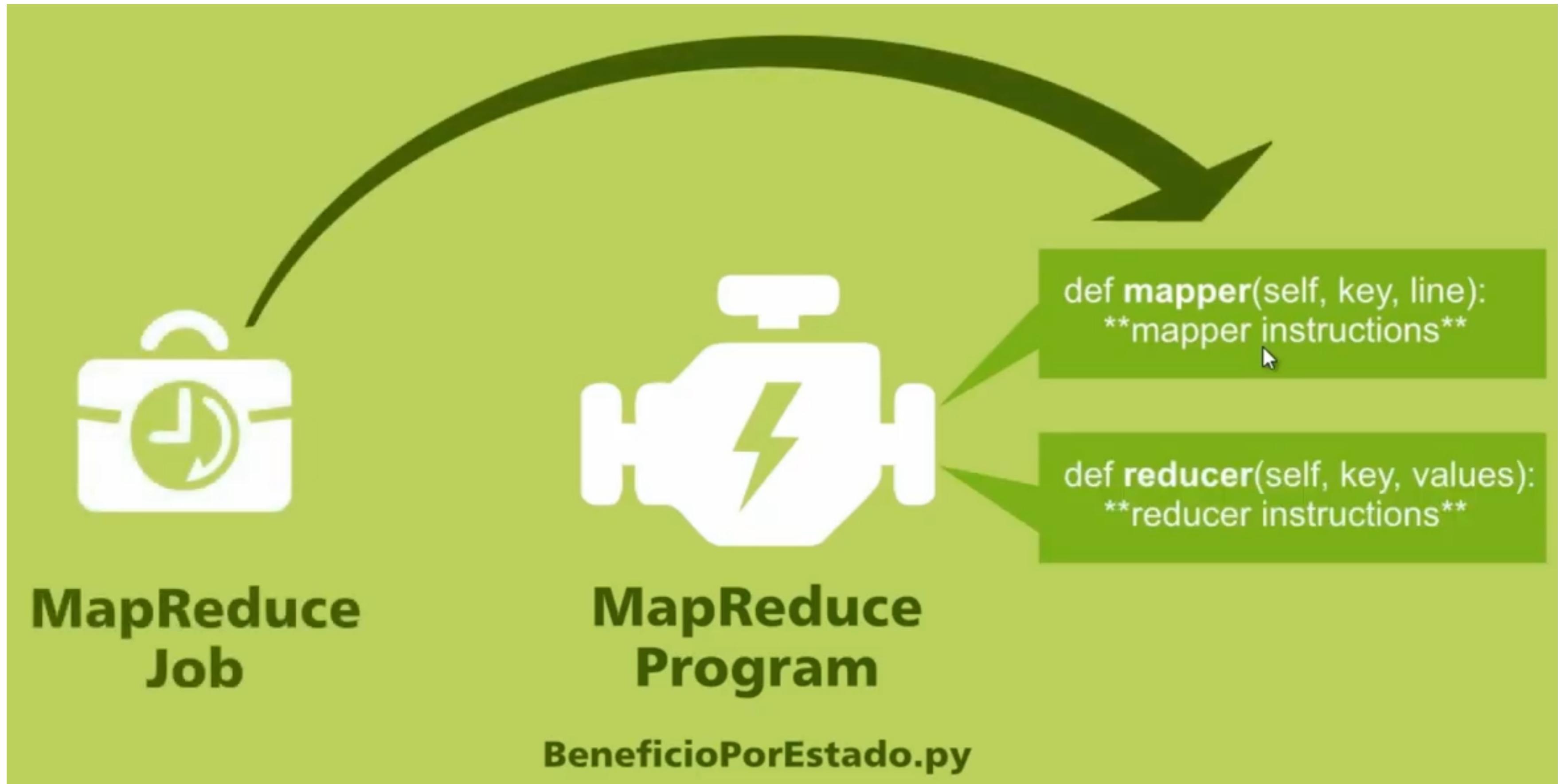


MapReduce





MapReduce



MapReduce Job implementa 3 fases:

- (i) Map
- (ii) Sort
- (iii) Reduce



MapReduce

Exemplo (leitura de dados do bolsa família (arquivo Input 1.7GB)

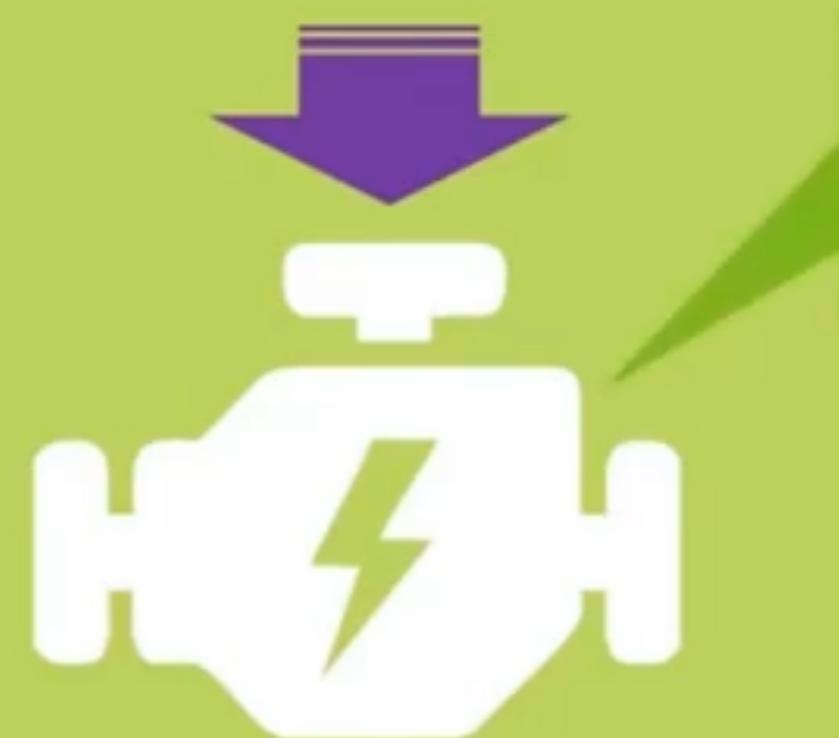
Fase 1: Map

Dado bruto (CSV)

UF	Nome Município	Nome Favorecido	Valor	Competência
SP	CAMPINAS	ANGELA SERAFIM	180.00	01/2016



MapReduce Job
Fase: Mapper



MapReduce Program
BeneficioPorEstado.py

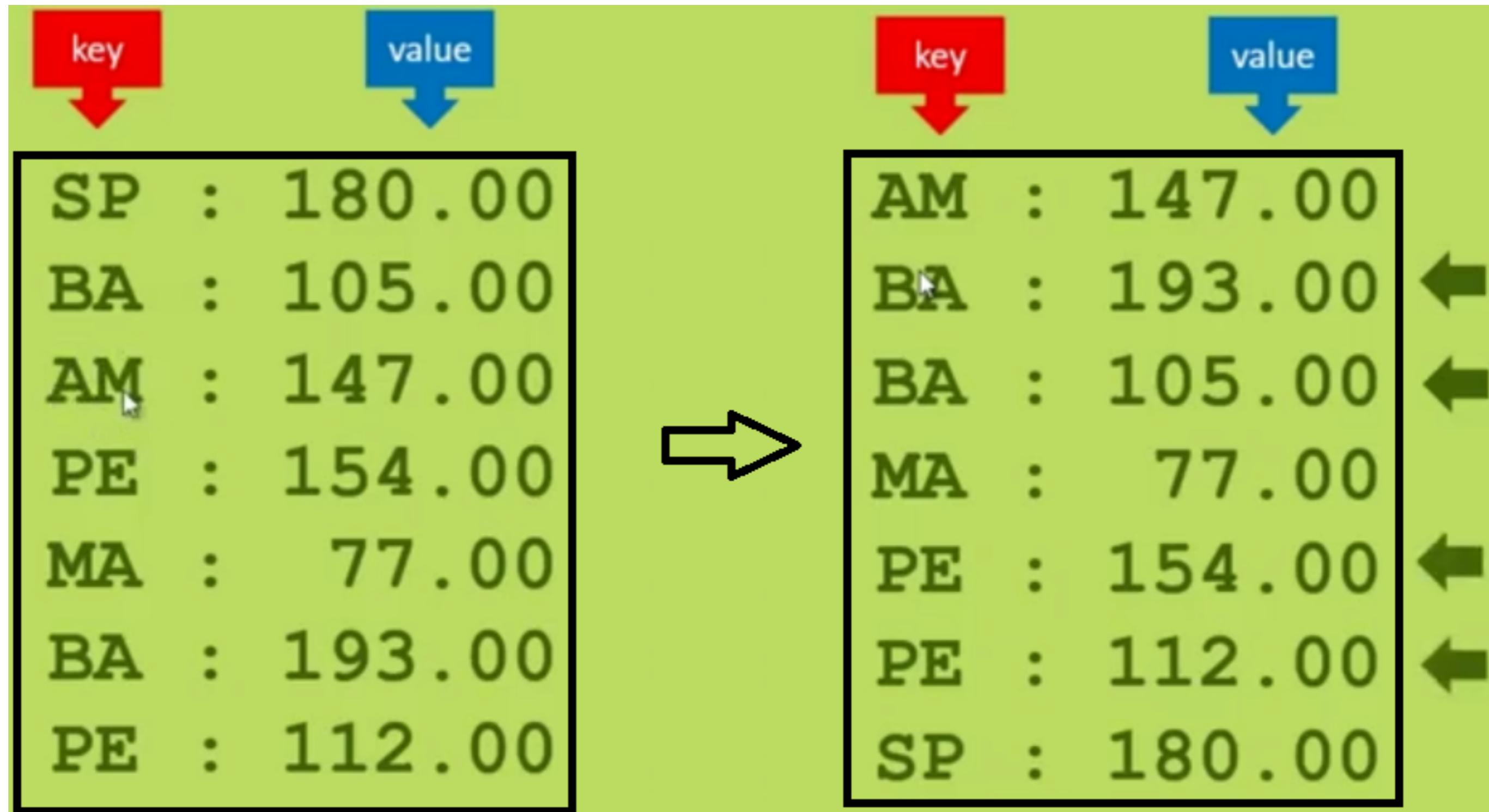
```
def mapper(self, key, line):  
    **mapper instructions**
```

key
value
SP : 180.00



MapReduce

Fase 2: Sort



AM : 147.00
BA : 193.00, 105.00
MA : 77.00
PE : 154.00, 112.00
SP : 180.00



MapReduce

Exemplo (leitura de dados do bolsa família (arquivo Input 1.7GB)

Fase 3: Reduce

```
AM : 147.00
BA : 193.00, 105.00
MA : 77.00
PE : 154.00, 112.00
SP : 180.00
```



MapReduce Program
BeneficioPorEstado.py

```
AM : 147.00
BA : 298.00
MA : 77.00
PE : 266.00
SP : 180.00
```



PSPD

Prof. Fernando W Cruz