

Engenharia de Dados com Hadoop e Spark



Data Science Academy



Data Science Academy

Bem-vindo



Data Science Academy

www.datascienceacademy.com.br



Data Science Academy

Usando MapReduce em Grandes Volumes de Dados

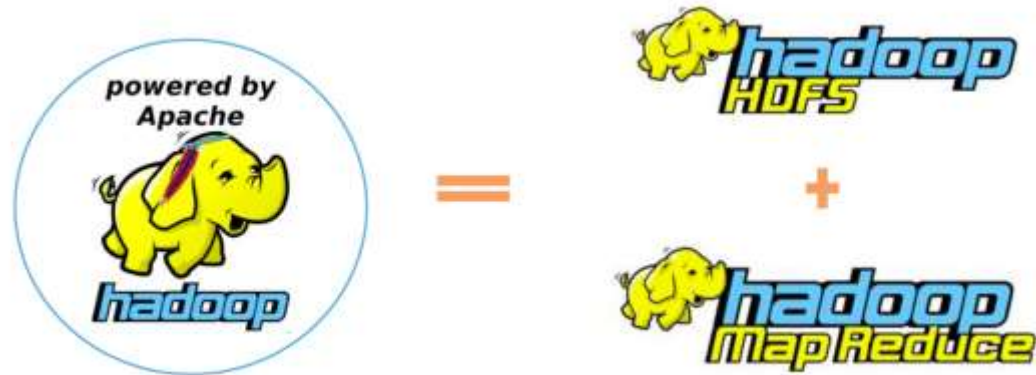


Data Science Academy

www.datascienceacademy.com.br

MapReduce

Hadoop



Data Science Academy

MapReduce

O que vamos estudar neste capítulo?

- Computação Distribuída
- Funcionamento do MapReduce
- Processamento de Dados Armazenados no HDFS
- Processamento de Big Data
- Criação e Monitoramento de Jobs MapReduce
- Processamento de Jobs MapReduce em Nuvem, com o Serviço AWS da Amazon



O que exatamente é processamento MapReduce?

Programação de Computadores

Python Fundamentos para Análise de Dados
Big Data Real-Time Analytics com Python e Spark





Computação Distribuída

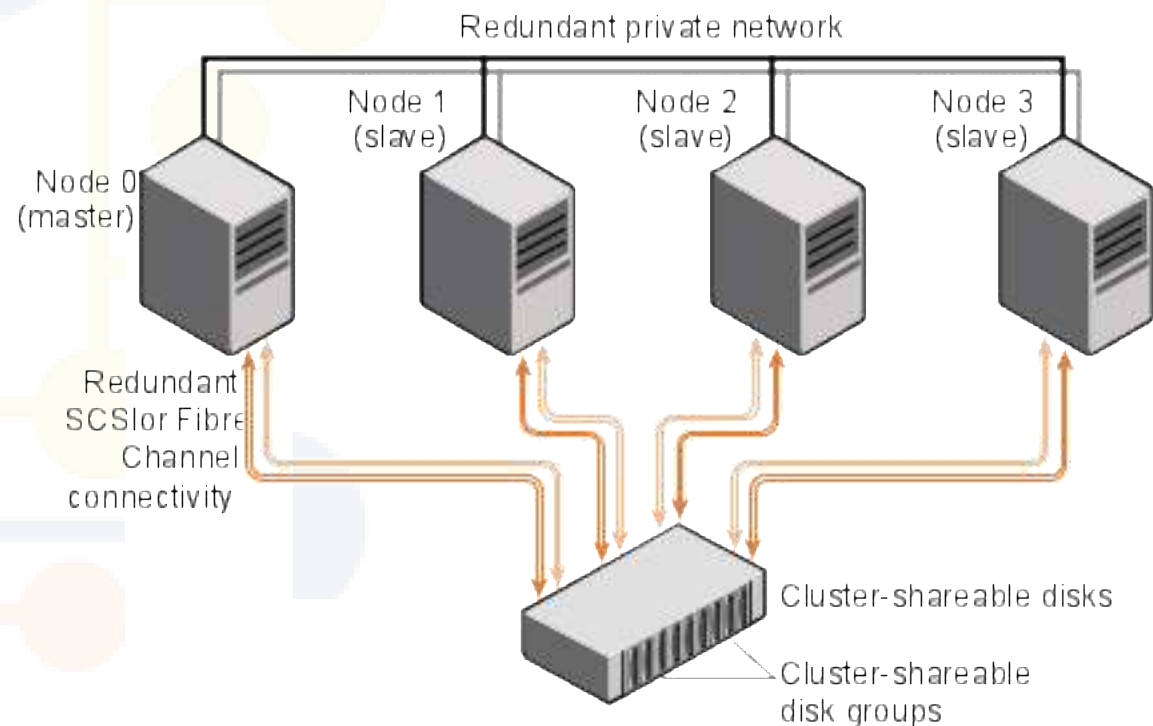


Data Science Academy

www.datascienceacademy.com.br

Computação Distribuída

Sistema de Processamento
Distribuído e Paralelo



Data Science Academy

Computação Distribuída

Uma tarefa qualquer, pode ser dividida em várias subtarefas e então executadas em paralelo.



Data Science Academy

Computação Distribuída

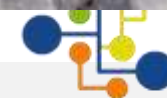


- Pesquisas científicas
- Previsões climáticas
- Descoberta de novas partículas
- Controle de epidemias
- Armazenamento e Processamento de Big Data



Data Science Academy

Computação Distribuída



Data Science Academy

Computação Distribuída

Sistemas Computacionais estão cada vez mais elaborados e complexos

Grande parte das máquinas interligadas por redes de computadores

Computação Distribuída

Sistemas Distribuídos

Maior poder de processamento
Maior carga, maior número de usuários
Melhor tempo de resposta
Maior confiabilidade



Data Science Academy

Computação Distribuída

A computação distribuída consiste na utilização de um conjunto de máquinas conectadas por uma rede de comunicação, atuando como um único sistema



Data Science Academy

Computação Distribuída

Computação Distribuída

- Executa aplicações através de máquinas diferentes, como se estas fossem uma só
- Tornou-se possível com a popularização das redes de computadores
- As máquinas podem estar interligadas por redes intranets, internet, redes públicas e privadas



Computação Distribuída

Computação
Distribuída

Vantagens

Utilizam melhor o poder de
processamento

Podem apresentar maior
confiabilidade

Apresentam melhor
desempenho

Permitem reutilizar serviços
já disponíveis

Permitem compartilhar
dados e recursos

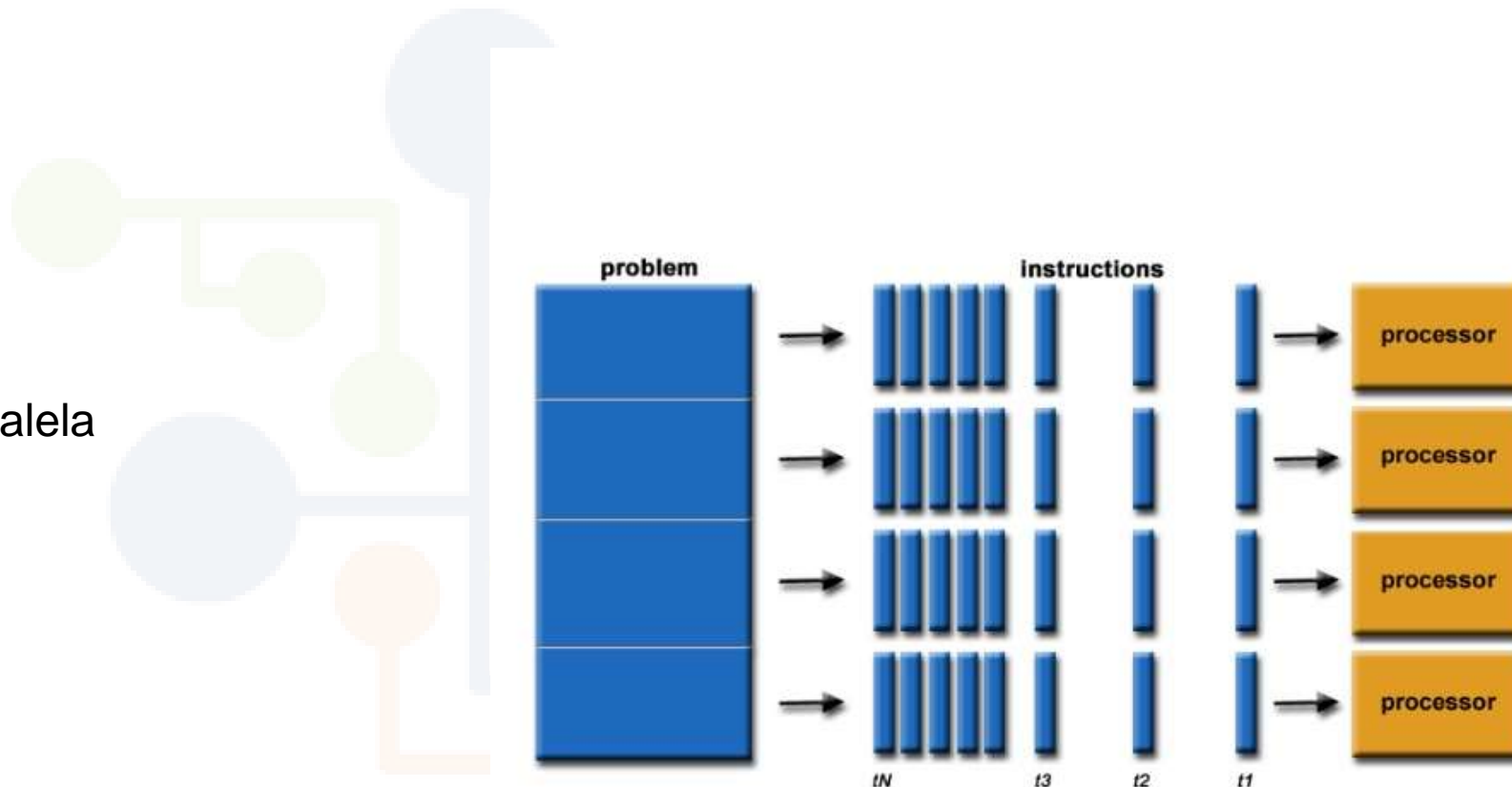
Permitem processar
grandes conjuntos de dados



Data Science Academy

Computação Distribuída

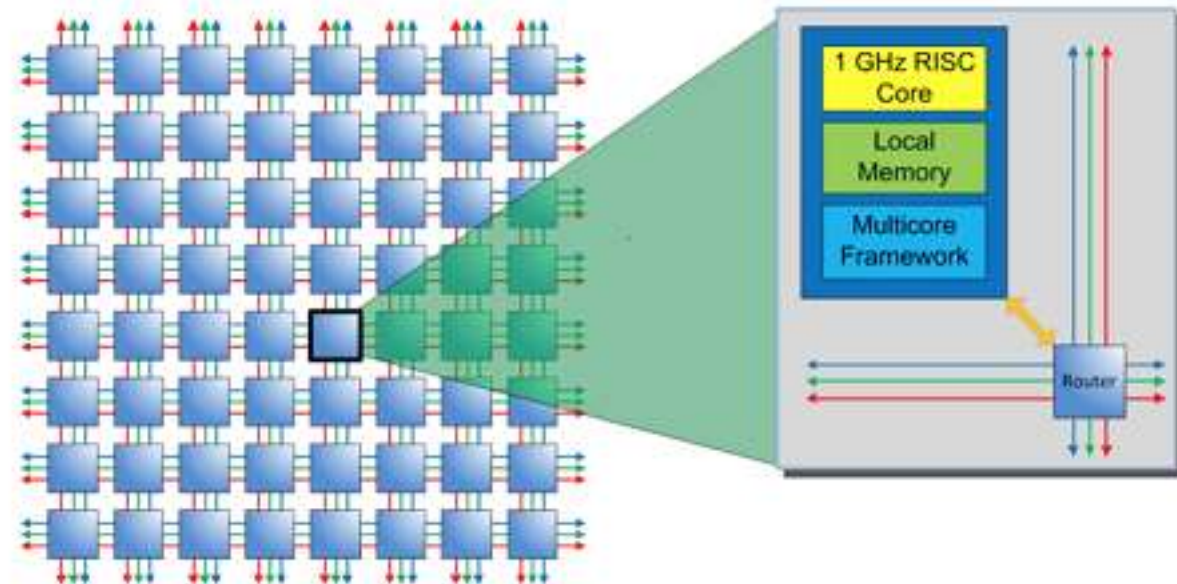
Computação Paralela



Data Science Academy

Computação Distribuída

Computação Paralela



Data Science Academy



Computação Distribuída

Cloud Computing



Data Science Academy

www.datascienceacademy.com.br

Computação Distribuída



Data Science Academy

Computação Distribuída



Data Science Academy

Computação Distribuída

Pay as you go



Data Science Academy

Computação Distribuída

Infraestrutura como um serviço
(IaaS)

Pay as you go



Data Science Academy

Computação Distribuída

Plataforma como um serviço
(PaaS)

Pay as you go



Data Science Academy

Computação Distribuída

Software como um serviço
(SaaS)

Pay as you go



Data Science Academy

Computação Distribuída

Mas o Big Data (sempre ele) trouxe mais um serviço para a nuvem



Data Science Academy

Computação Distribuída

Big Data como um serviço
(BDaaS)

Pay as you go



Data Science Academy

Computação Distribuída

Big Data como um serviço
(BDaaS)

Pay as you go



Data Science Academy

Computação Distribuída



Academy



Computação Distribuída



Data Science Academy

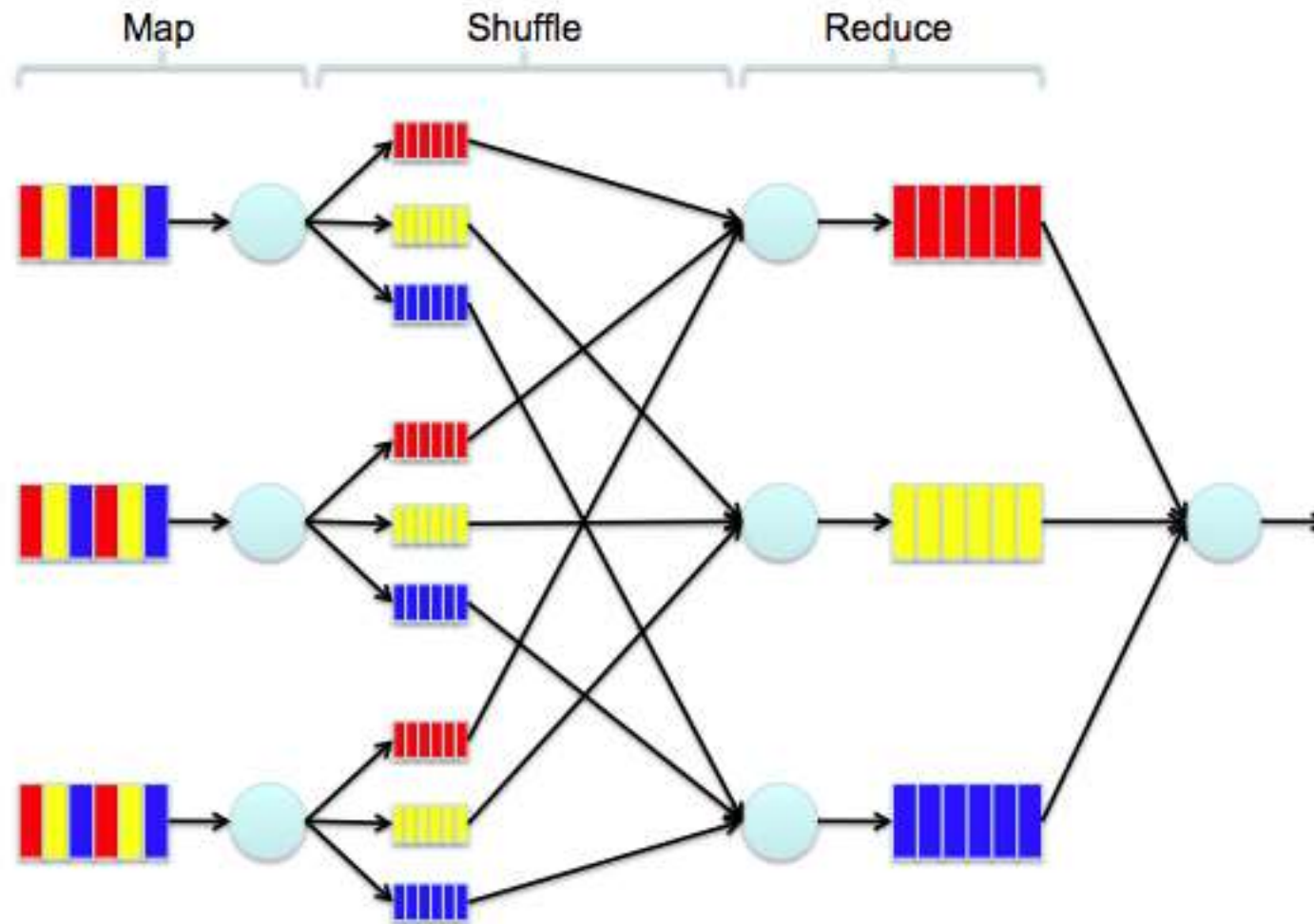


O Modelo de Programação MapReduce

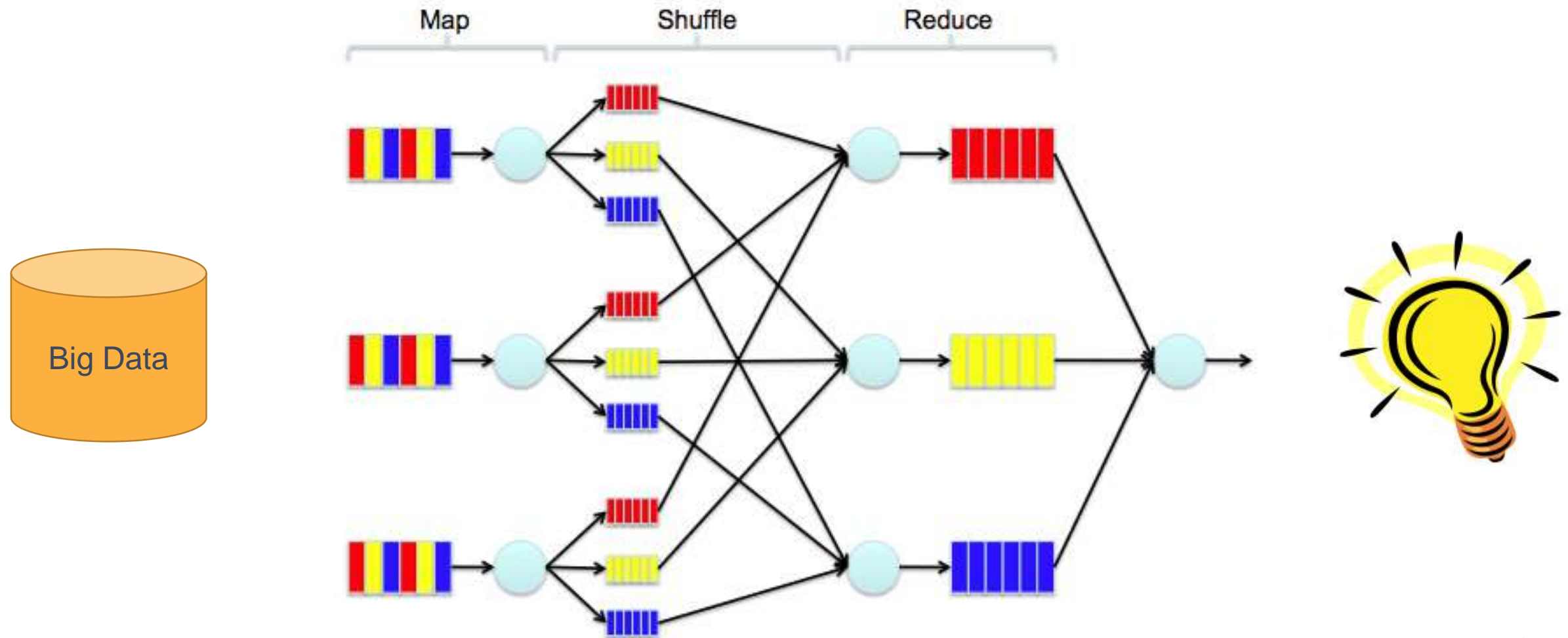


Data Science Academy

MapReduce



MapReduce



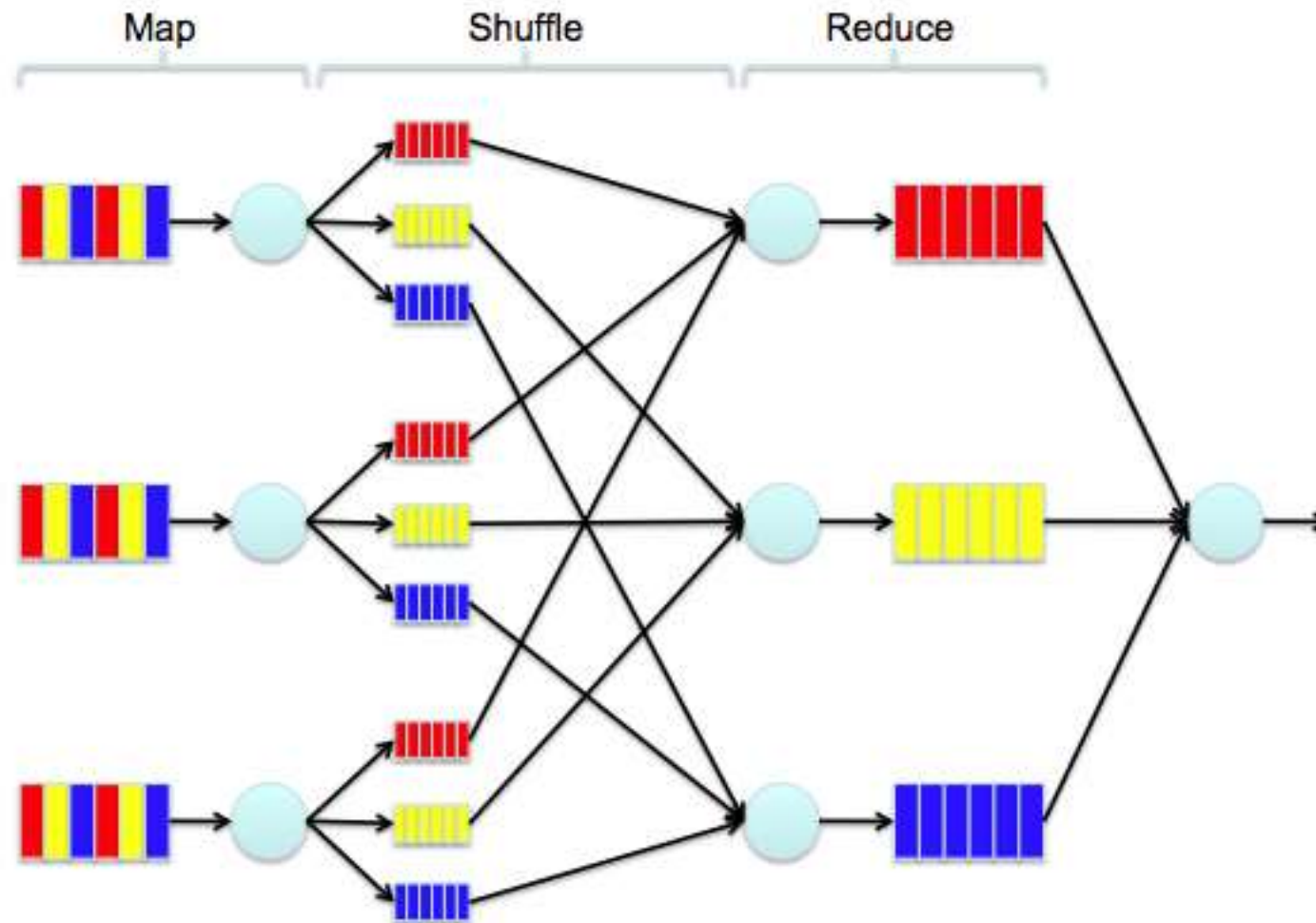
MapReduce

Como exatamente funciona o modelo MapReduce?

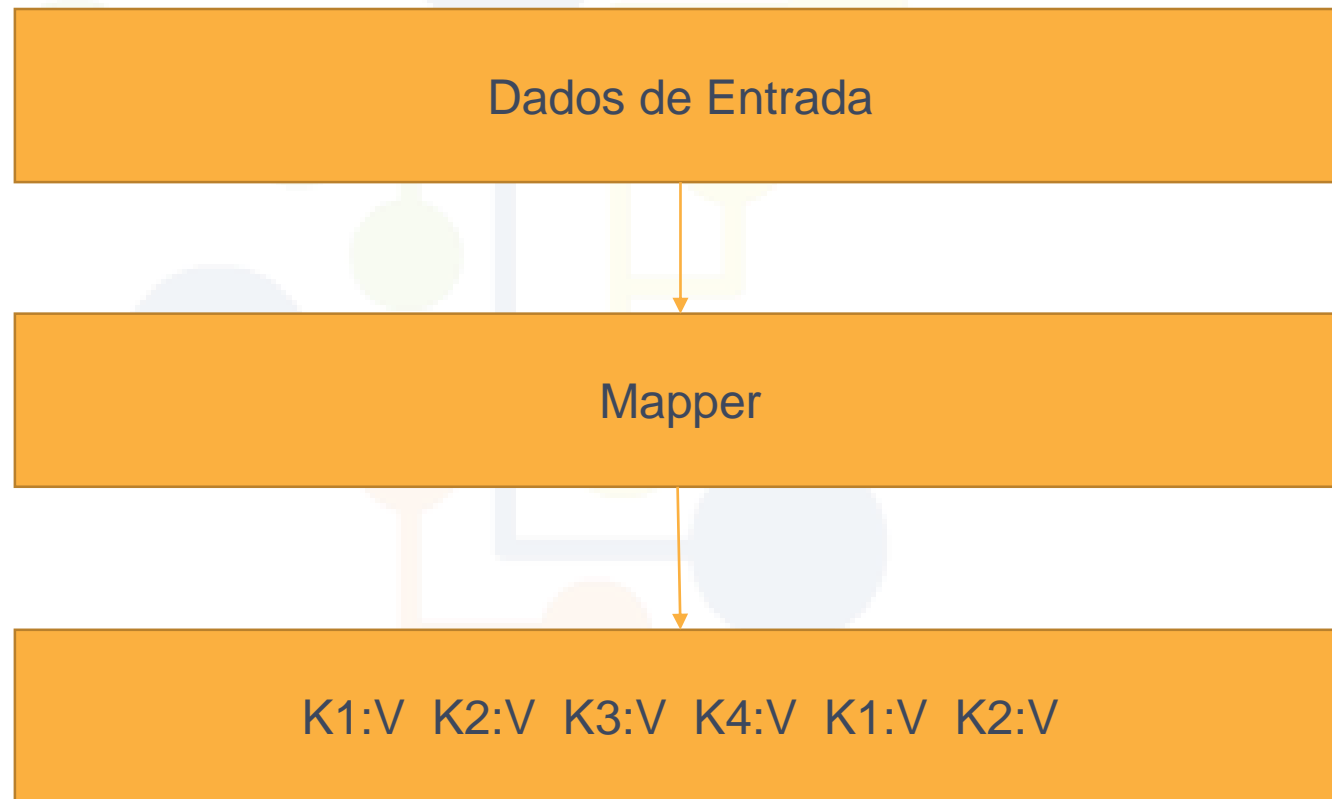


Data Science Academy

MapReduce



MapReduce



MapReduce

Quem define o que será a chave e o que será o valor?

Você, Cientista de Dados!



Data Science Academy

MapReduce

Dataset MovieLens (u.data)

userID	movieID	rating	timestamp
241	198	2	981769876
197	302	3	781769876
197	378	4	751769876
186	153	4	721769876
165	349	3	741769876
187	472	1	681769876
187	267	2	581769876

Quantos filmes cada pessoa assistiu?

Mapper

Key : Value

userID : movieID

241 : 198

197 : 302

197 : 378

186 : 153

165 : 349

187 : 472

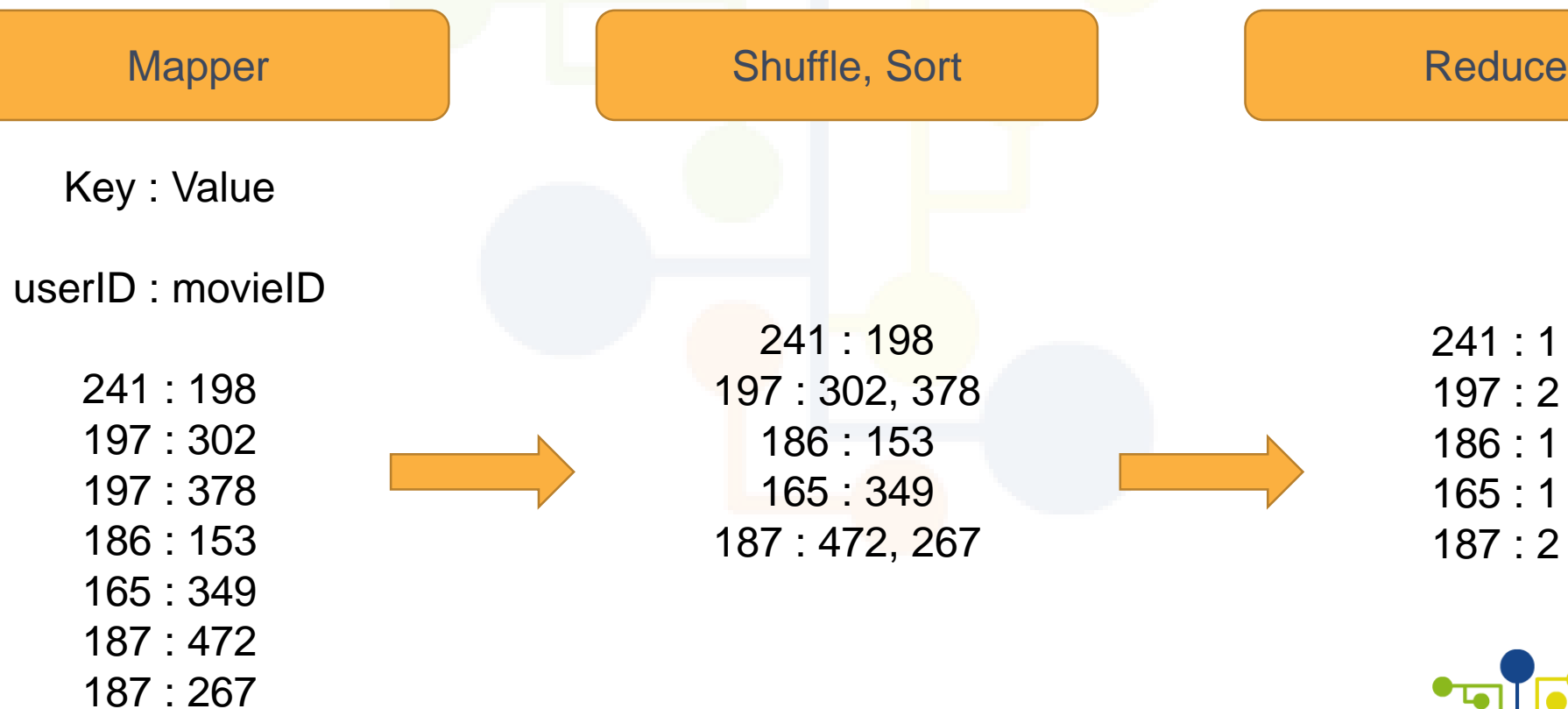
187 : 267



Data Science Academy

MapReduce

Quantos filmes cada pessoa assistiu?



Data Science Academy

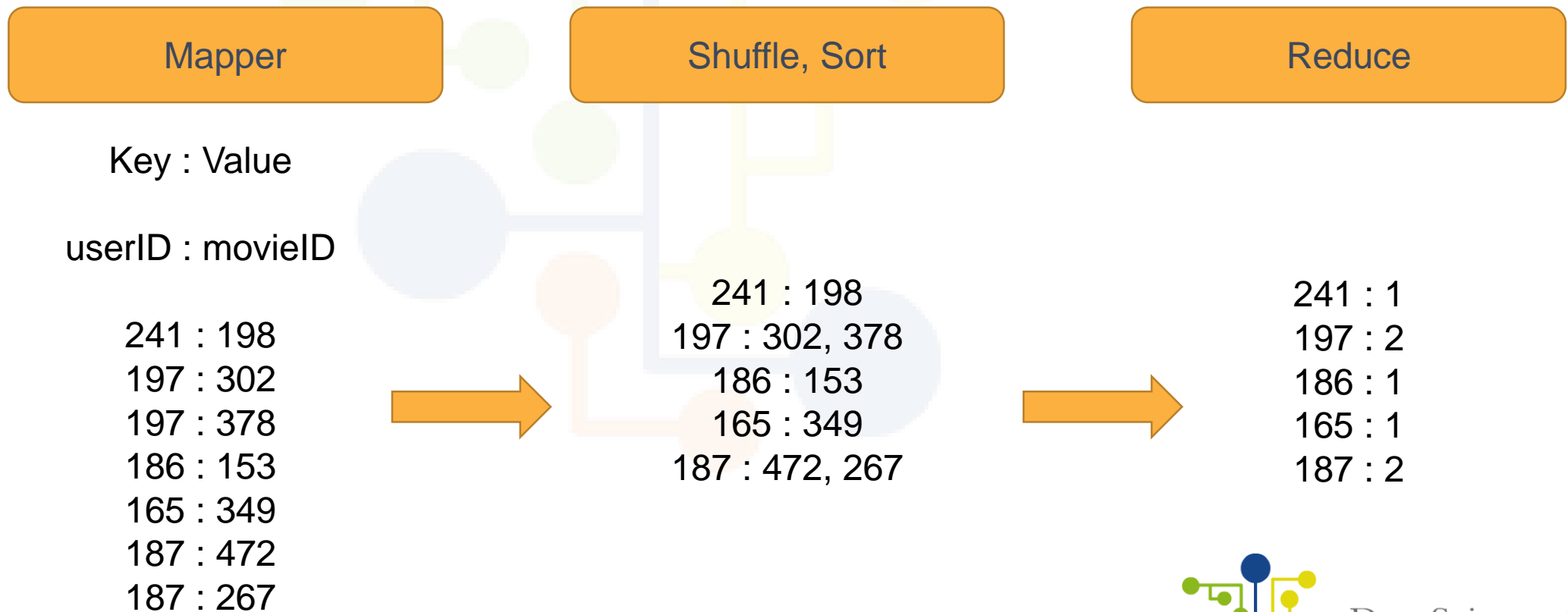
MapReduce



Data Science Academy

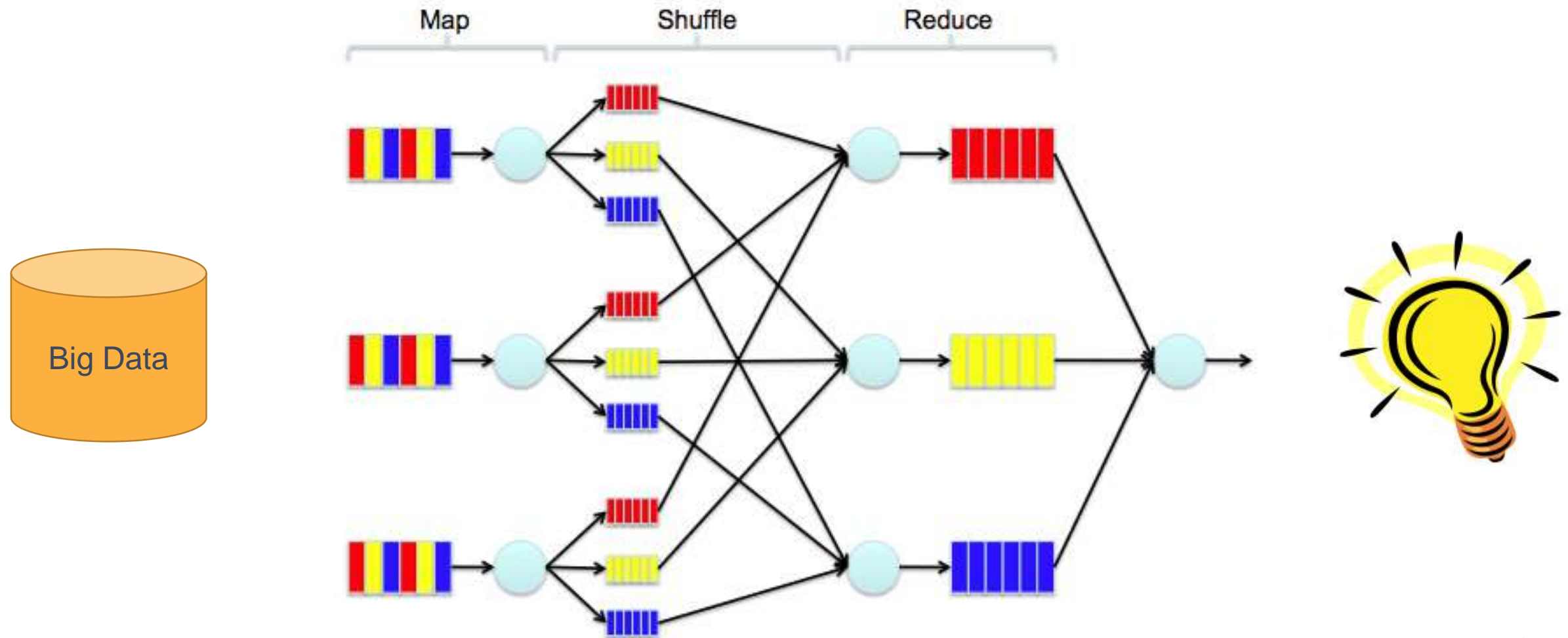
MapReduce

Quantos filmes cada pessoa assistiu?

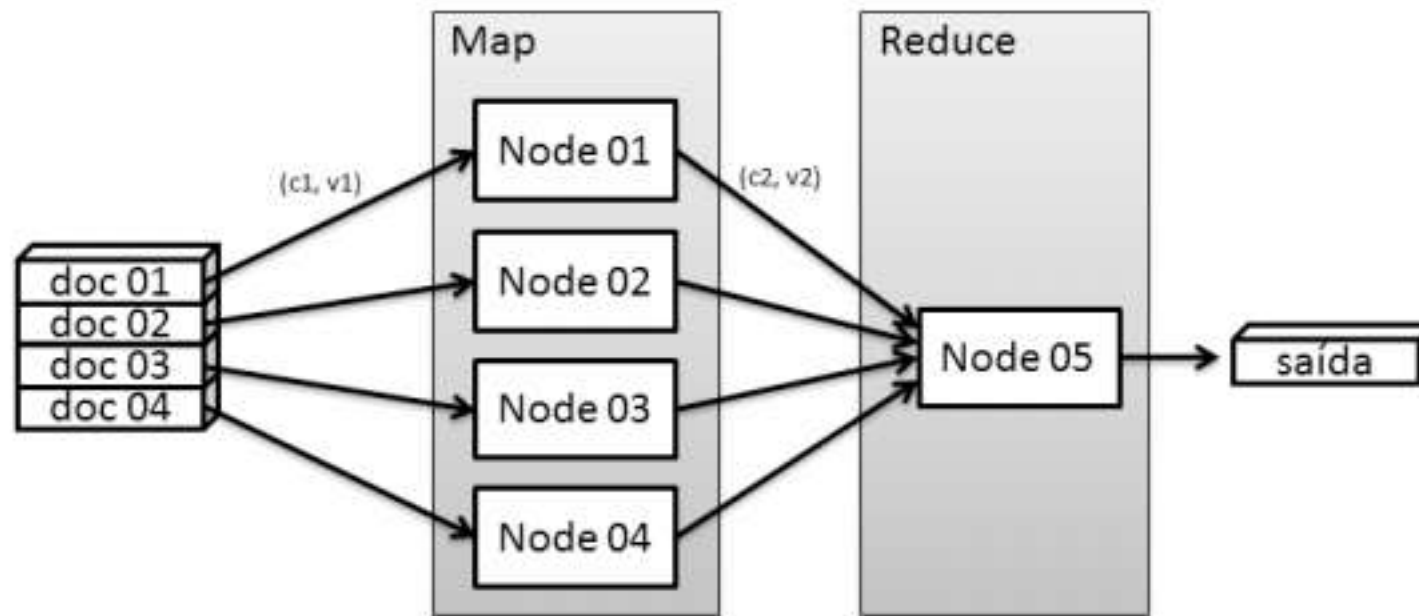


Data Science Academy

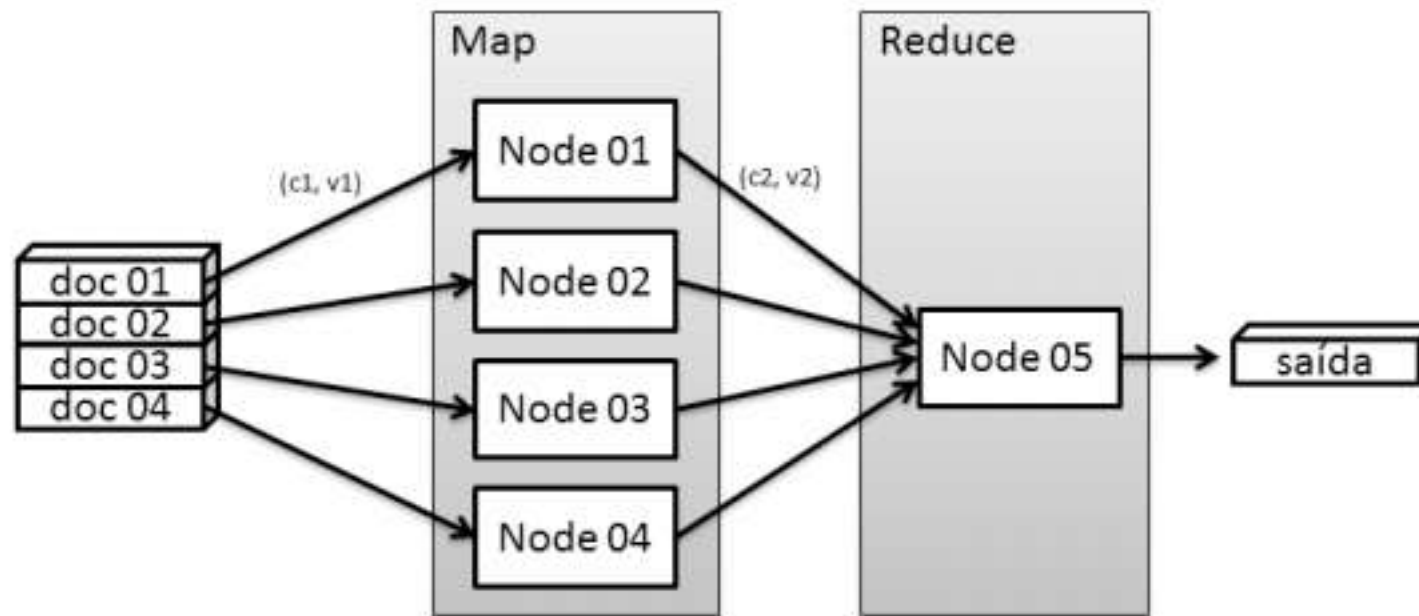
MapReduce



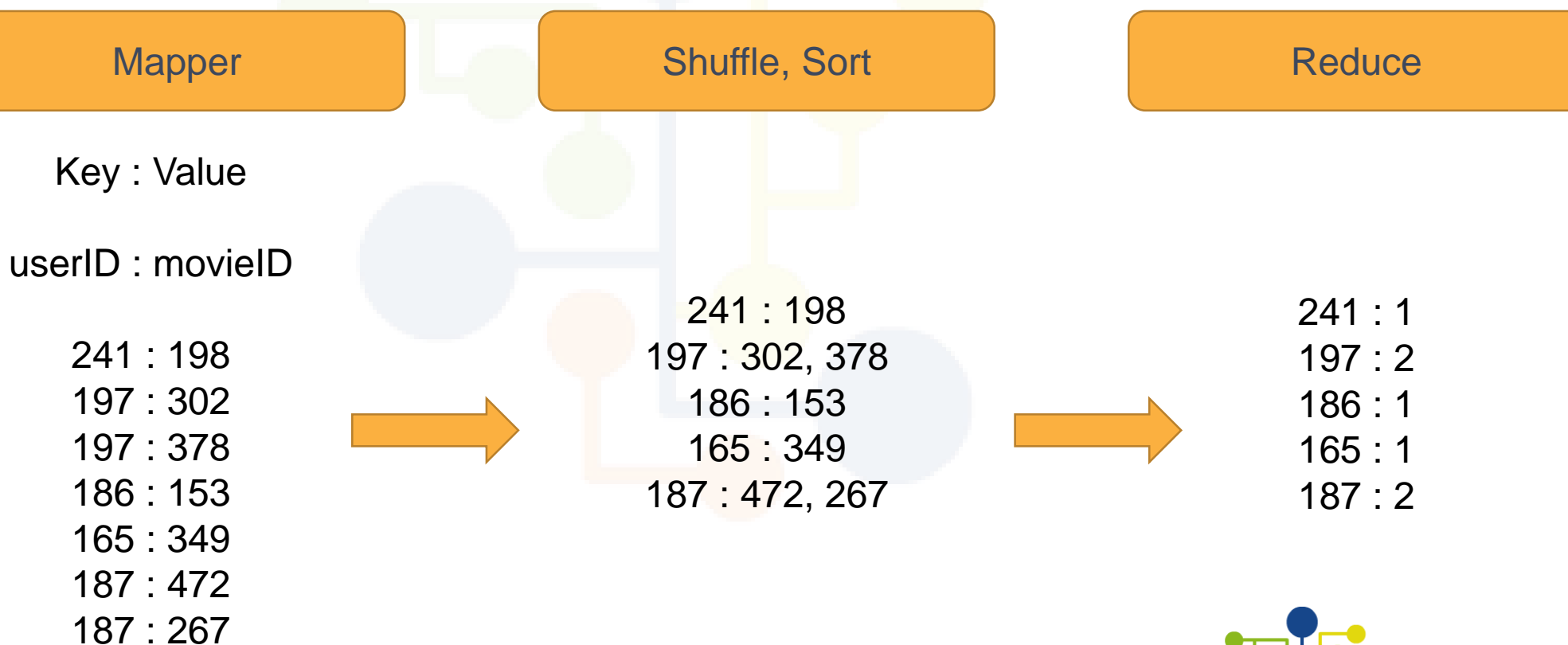
MapReduce



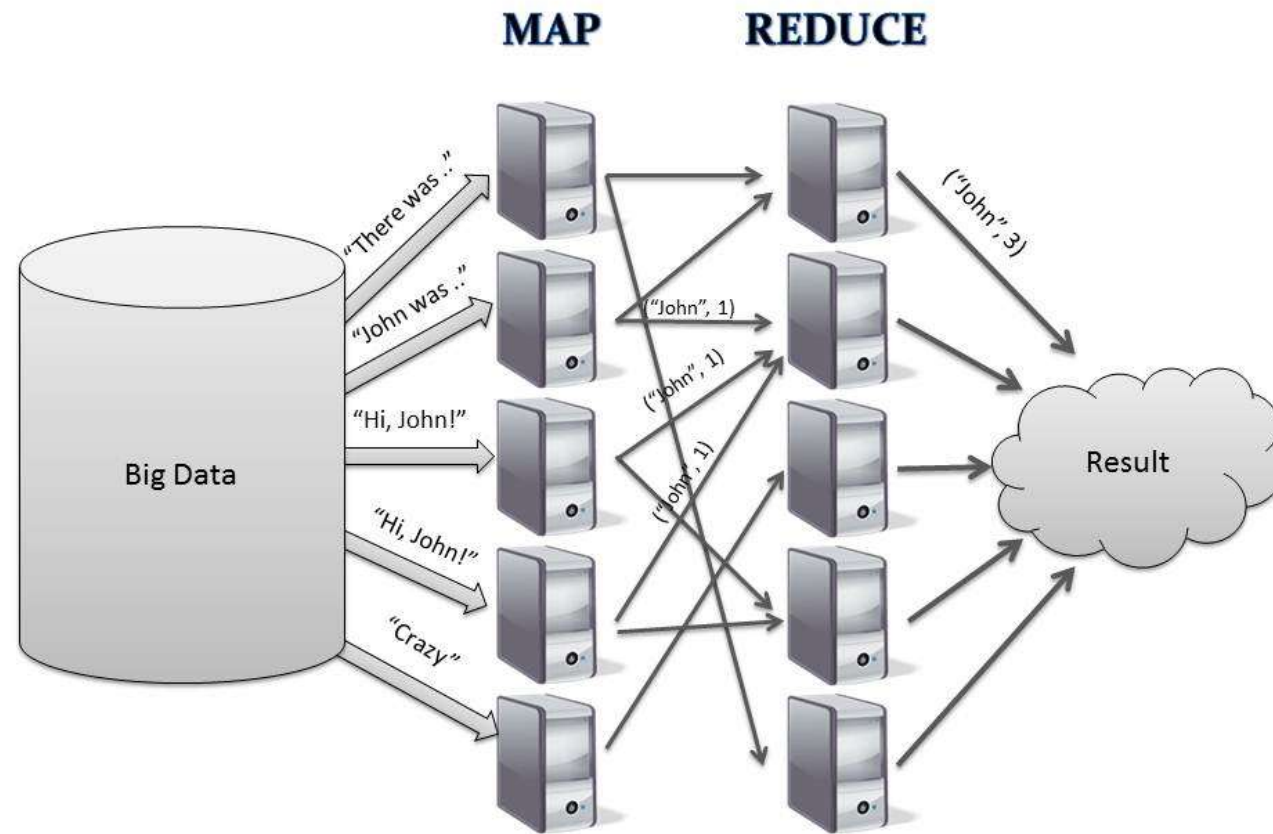
MapReduce



MapReduce



MapReduce



Como o MapReduce Utiliza a Computação Distribuída



Data Science Academy

MapReduce

A operação do MapReduce inclui



Comunicação eficiente
entre os nodes



Tratamento de erros e
problemas de performance



Computação paralela
através do cluster



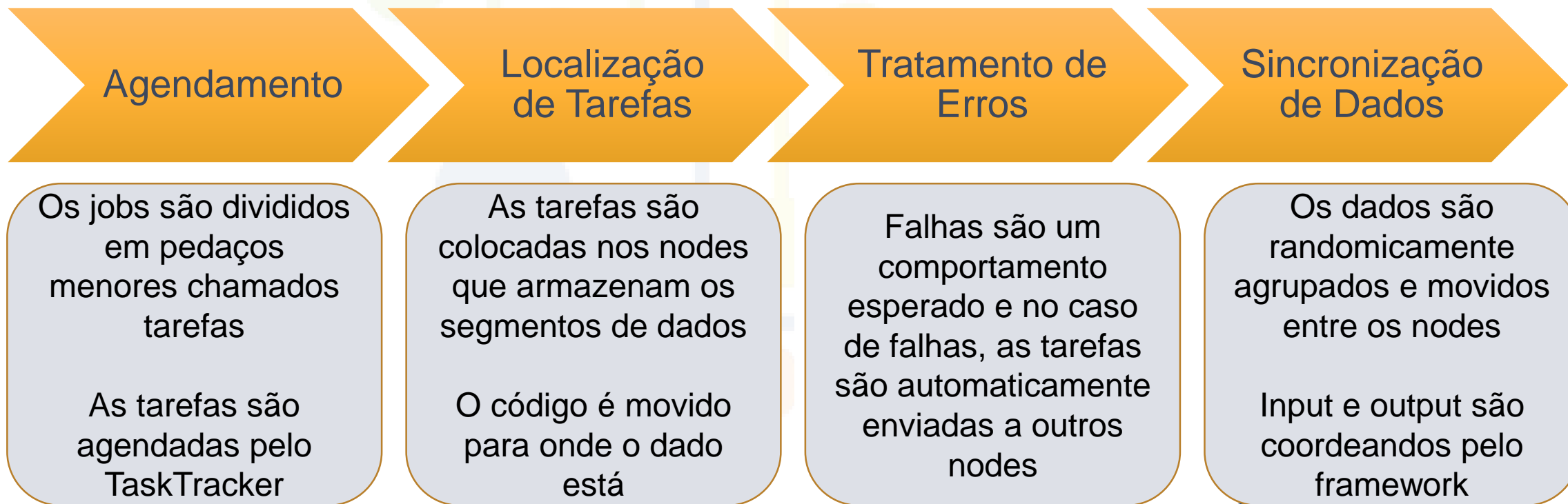
Funções de mapeamento e
redução



Data Science Academy

MapReduce

Workflow do MapReduce

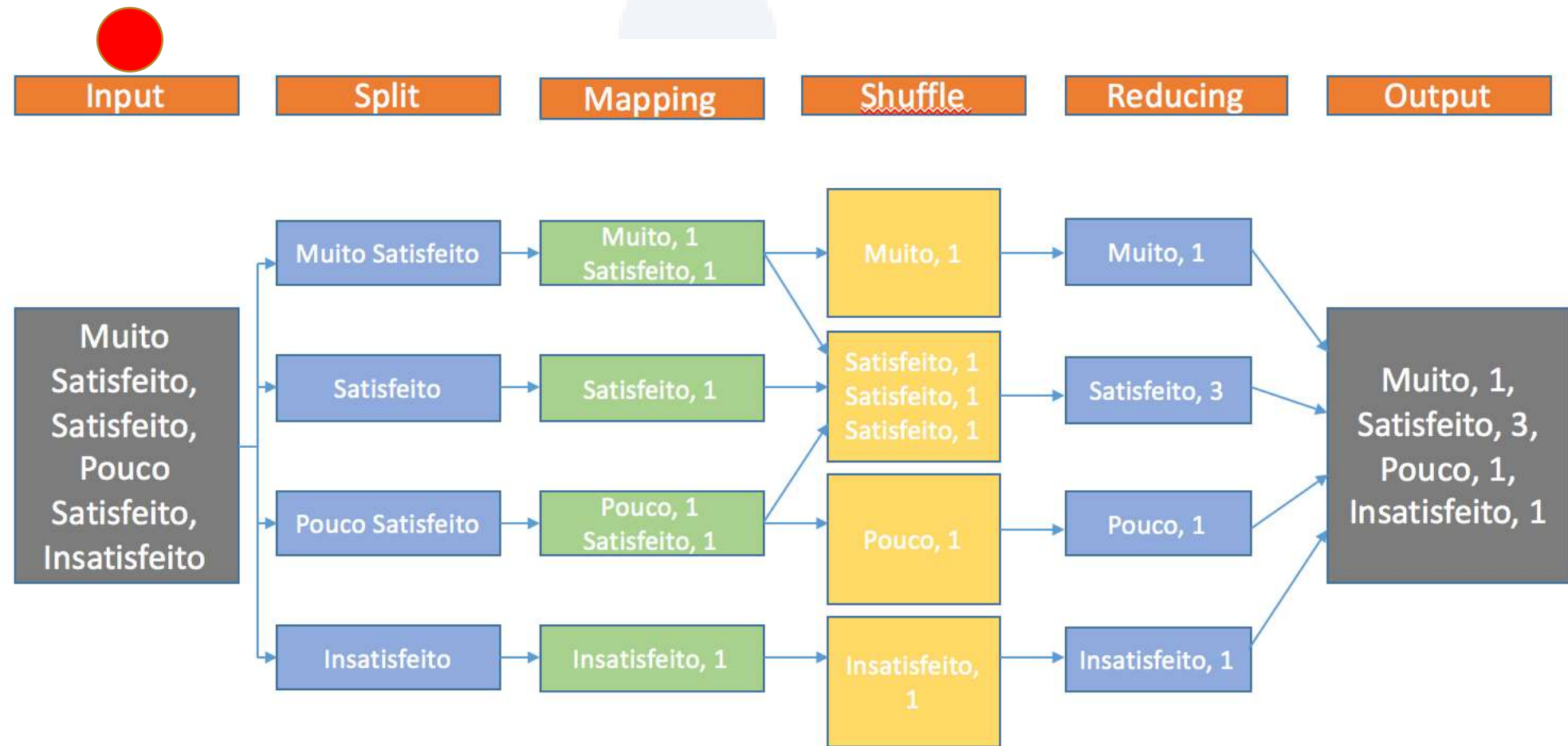


MapReduce

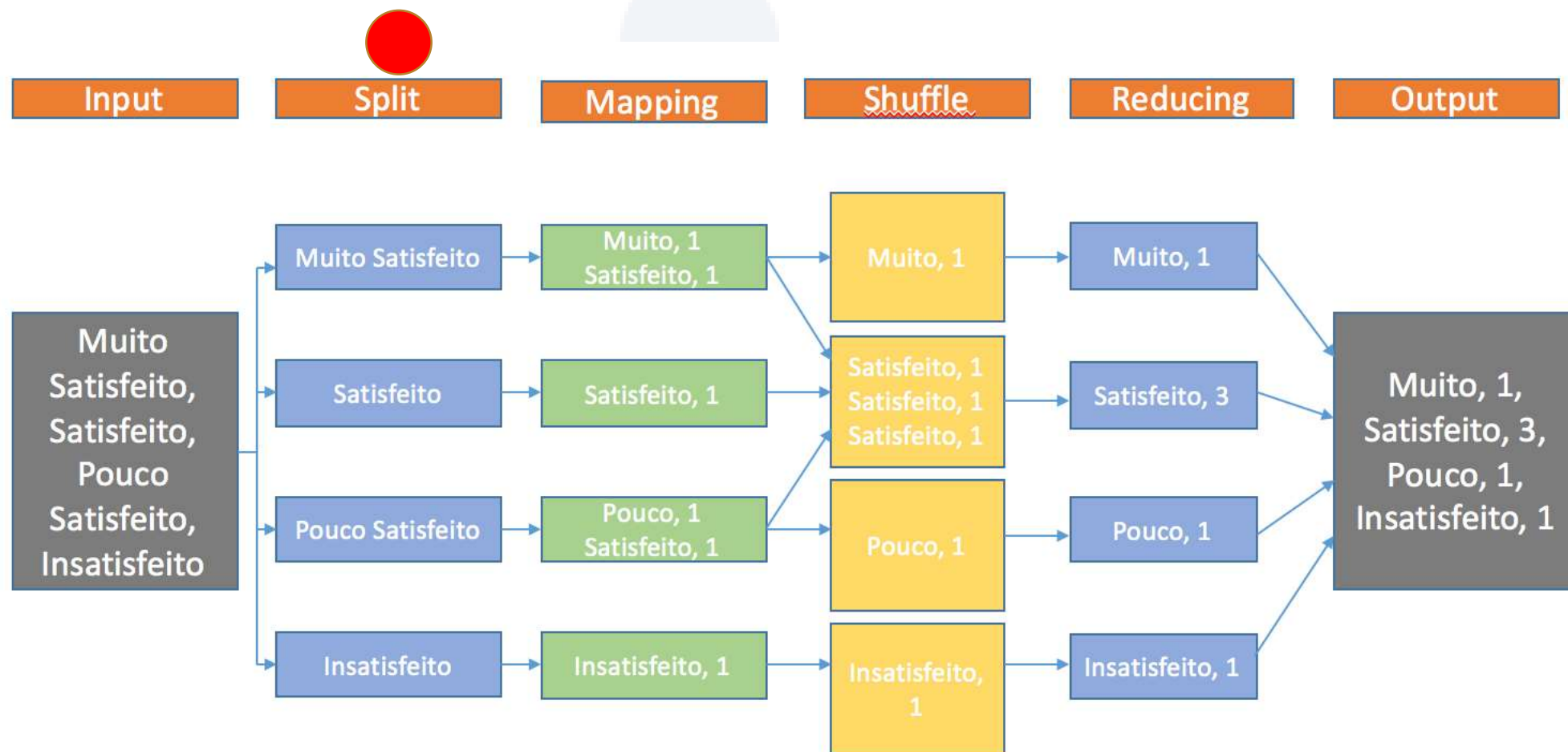
{Muito Satisfeito, Satisfeito, Pouco Satisfeito, Insatisfeito}



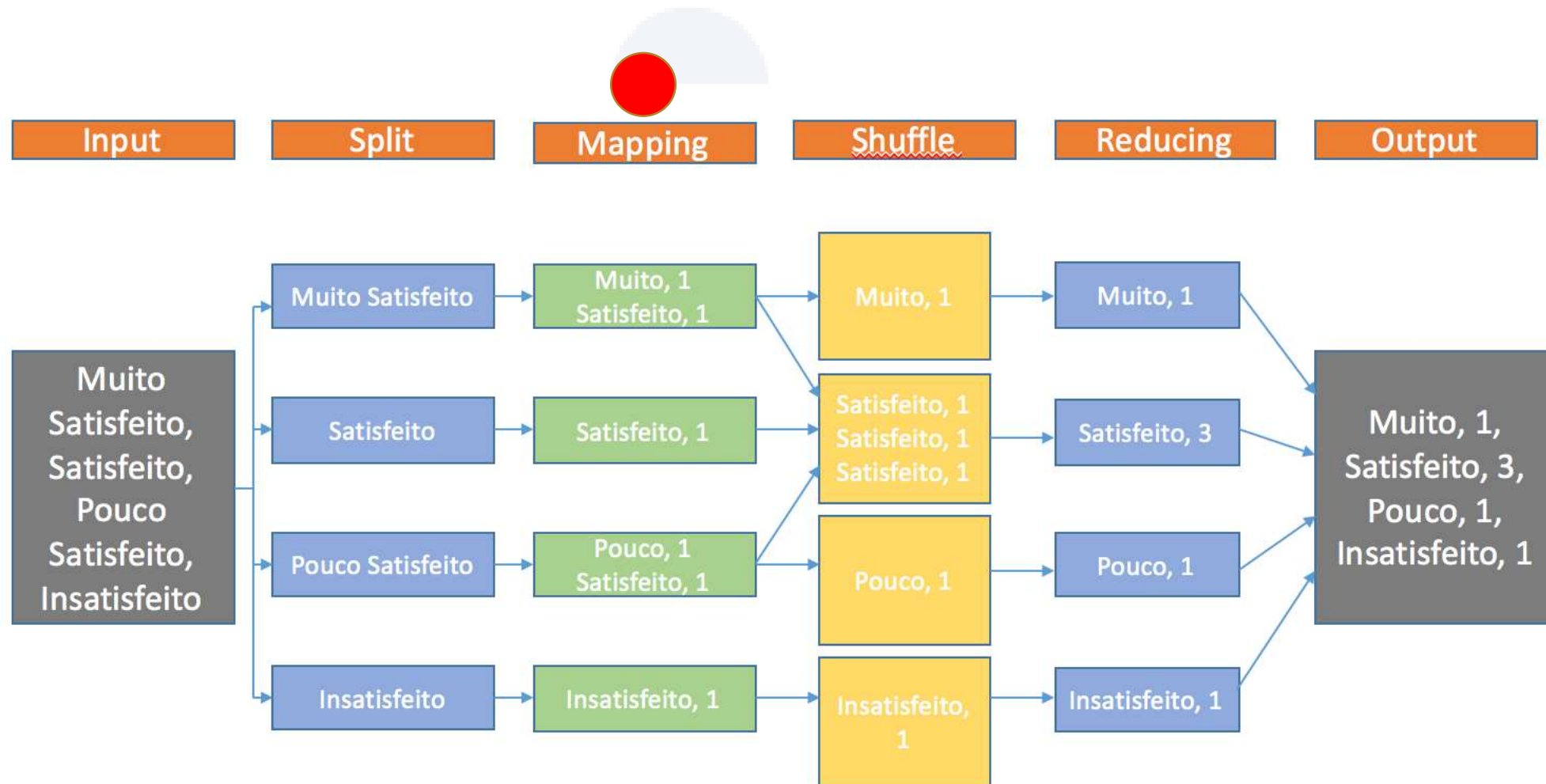
MapReduce



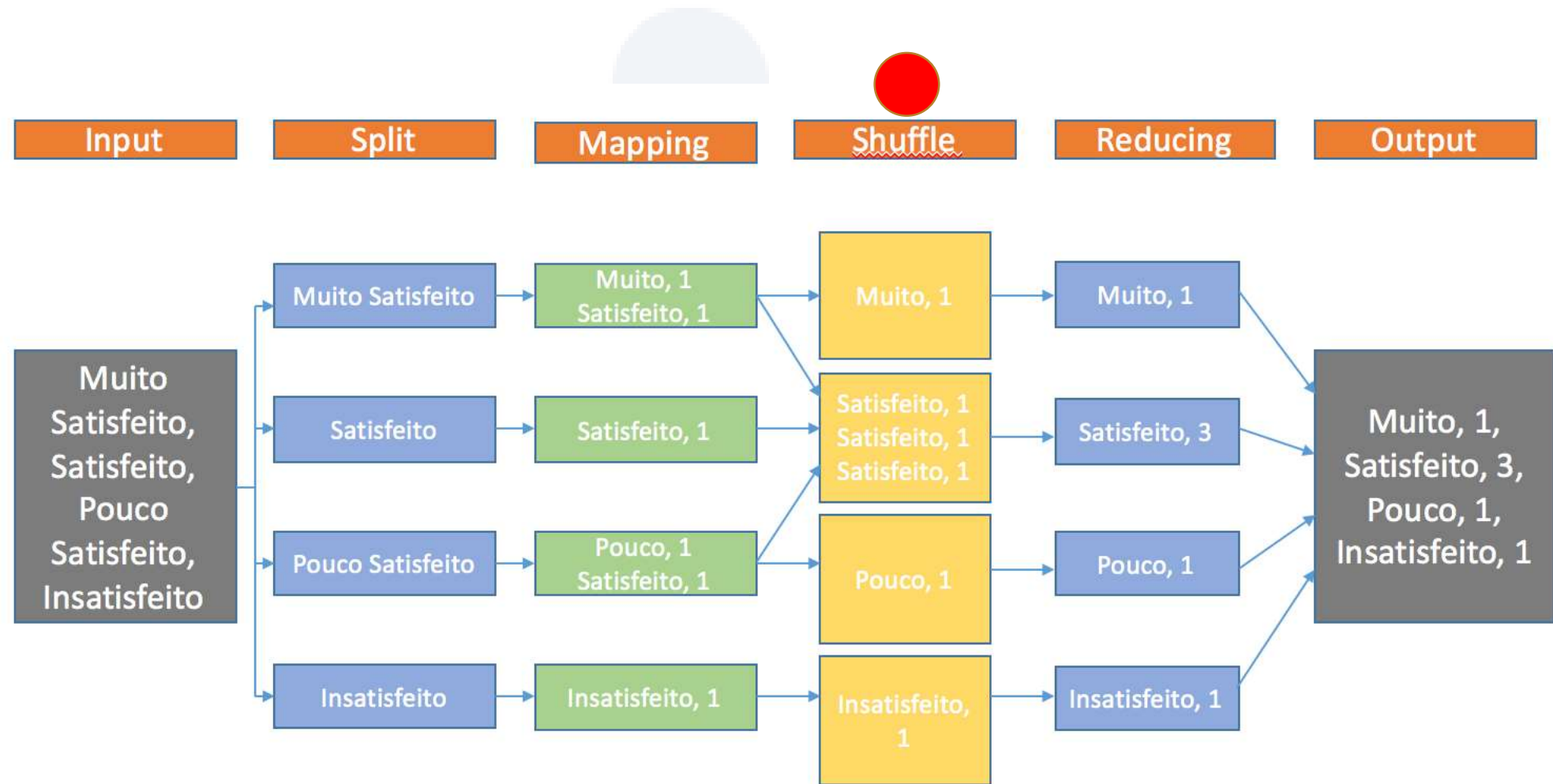
MapReduce



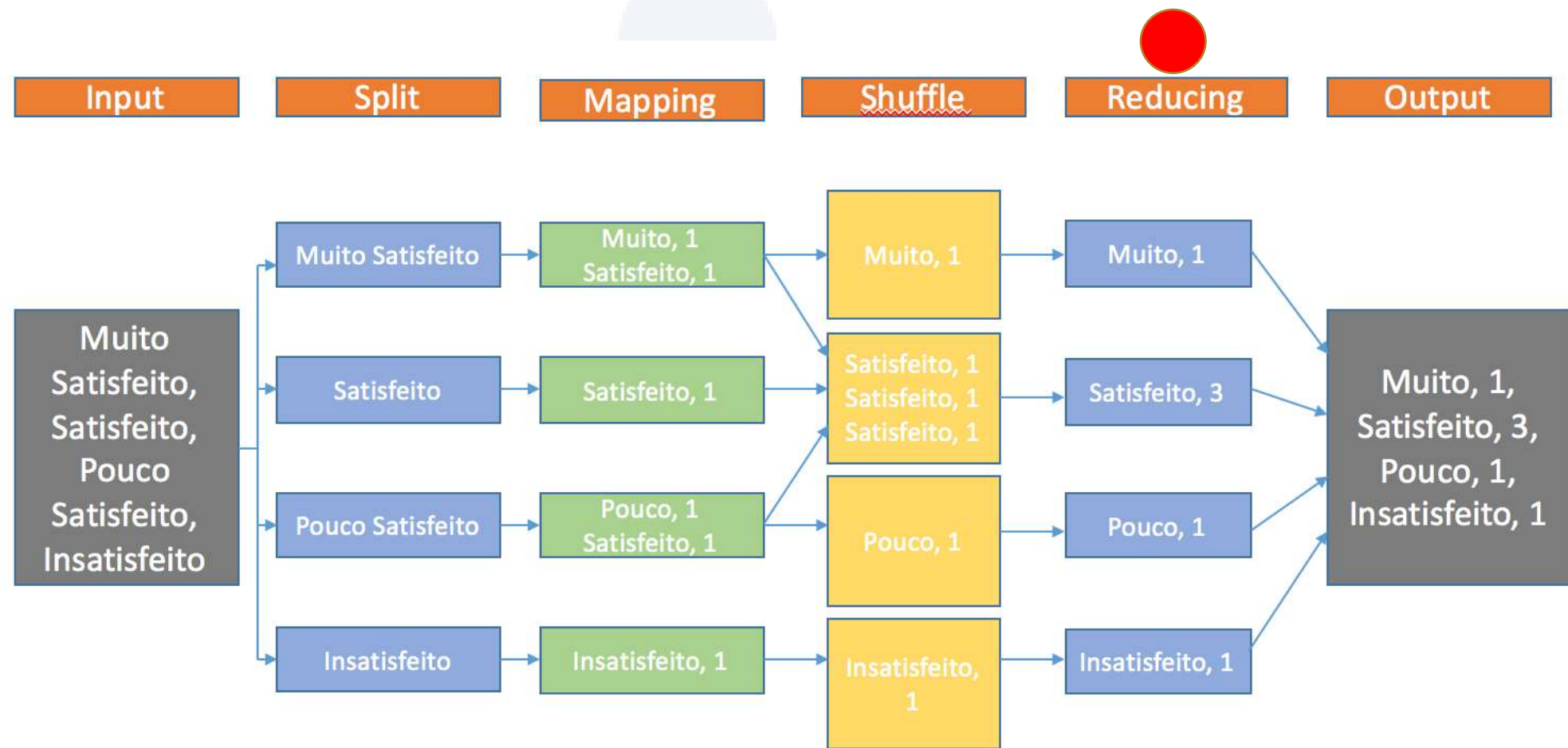
MapReduce



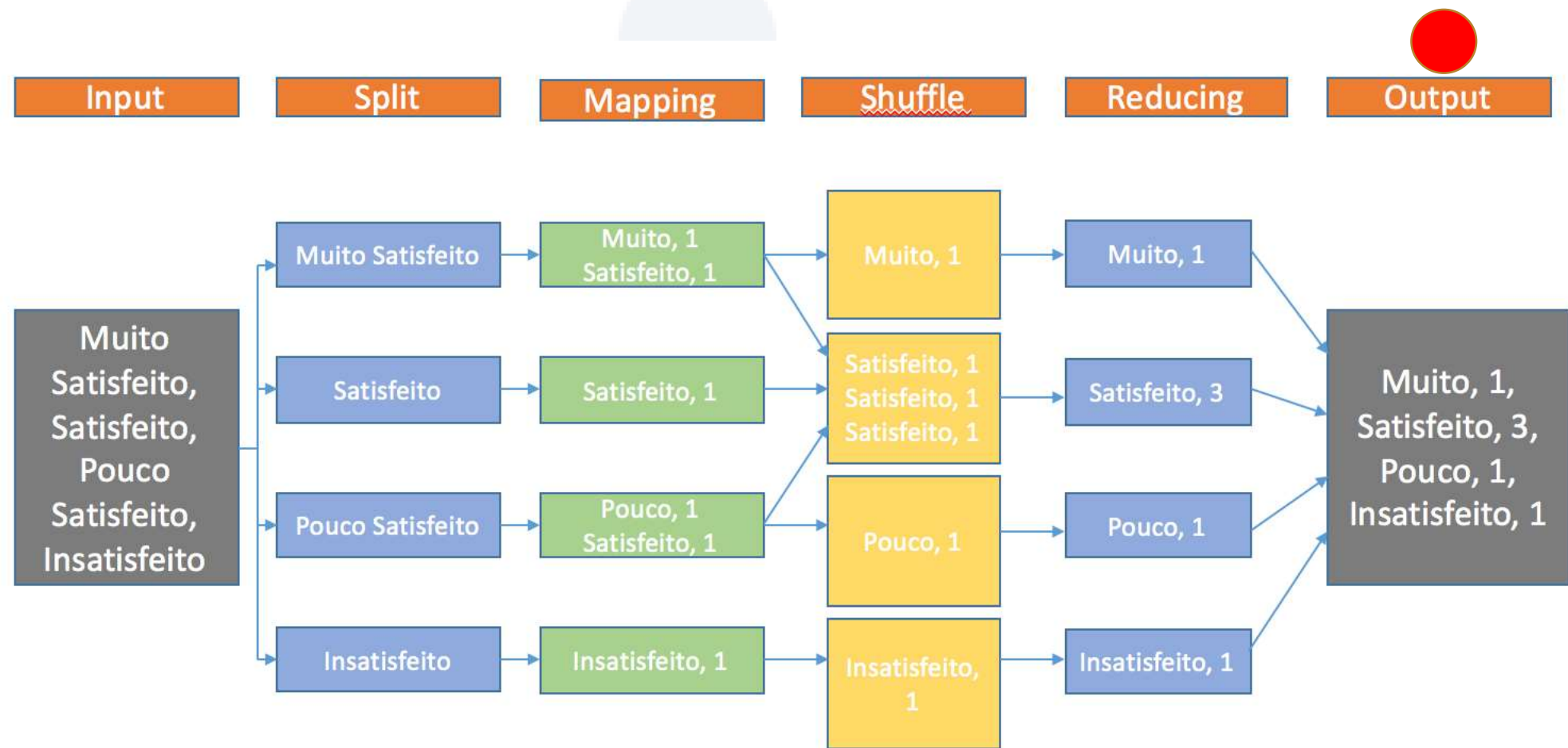
MapReduce



MapReduce



MapReduce



MapReduce

Características do MapReduce

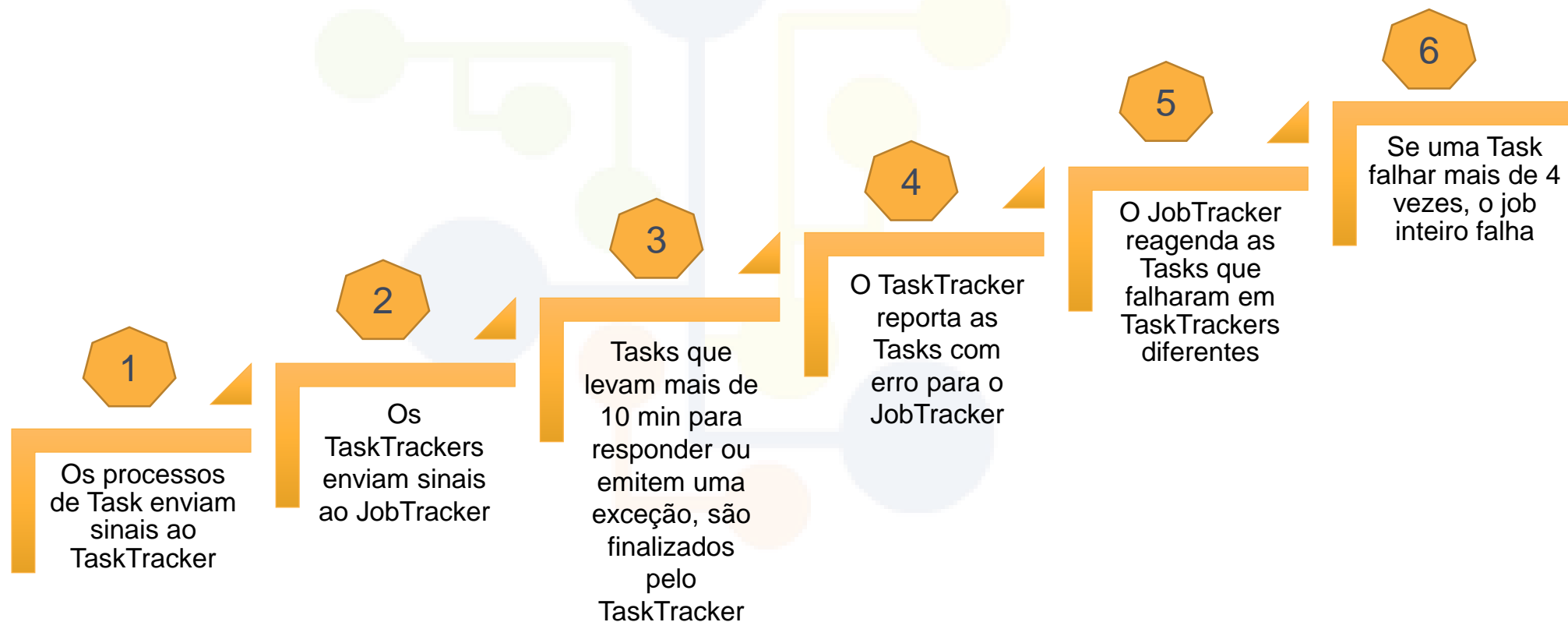
Algumas das principais características do MapReduce:

- Consegue trabalhar com grandes volumes de dados
- Funciona bem com o conceito WORM (Write Once and Read Many)
- Permite paralelismo
- As operações são realizadas próximas dos dados
- Hardware e storage de baixo custo podem ser usados
- O runtime fica responsável por dividir e mover os dados para as operações



MapReduce

Processo de Recuperação a Falhas do MapReduce



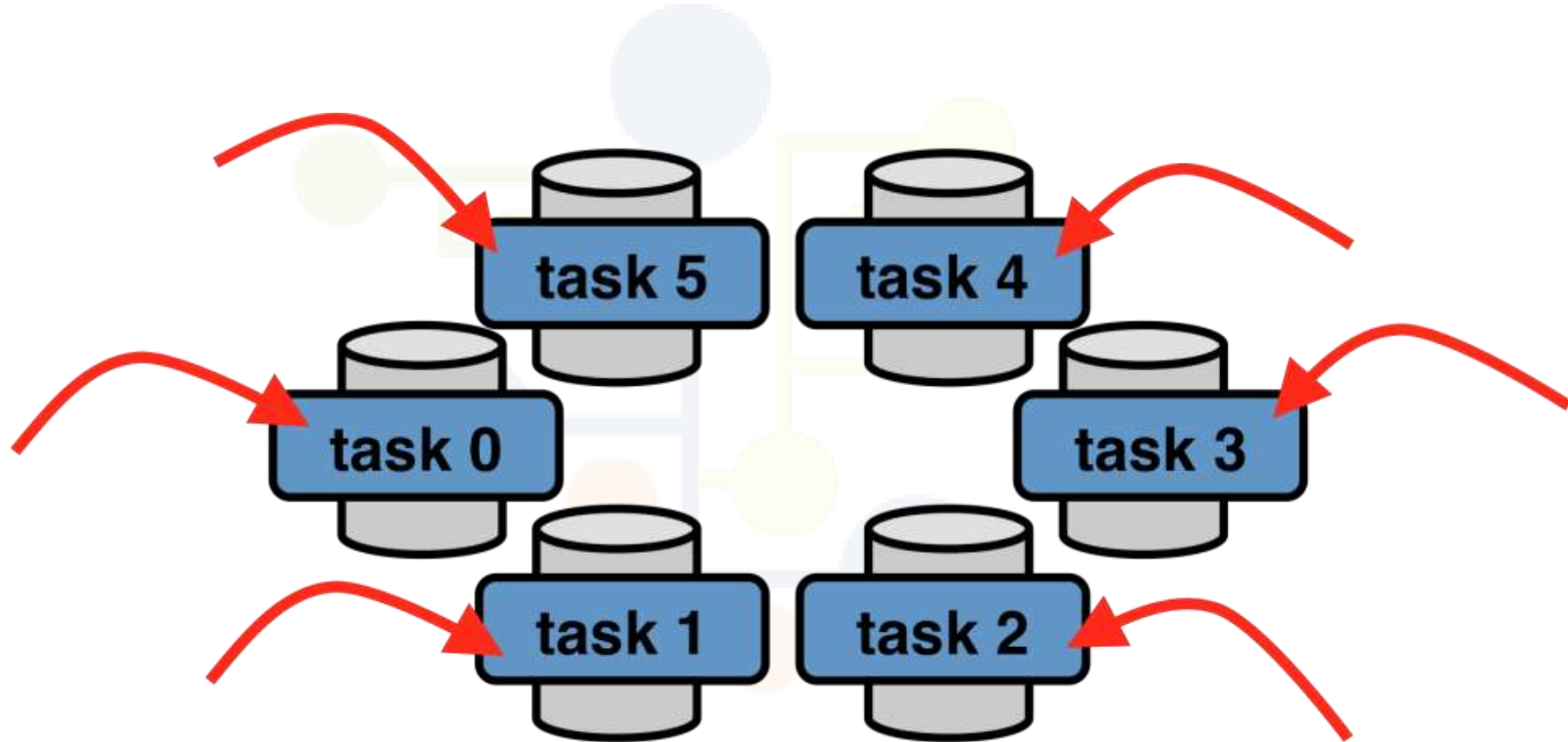


YARN

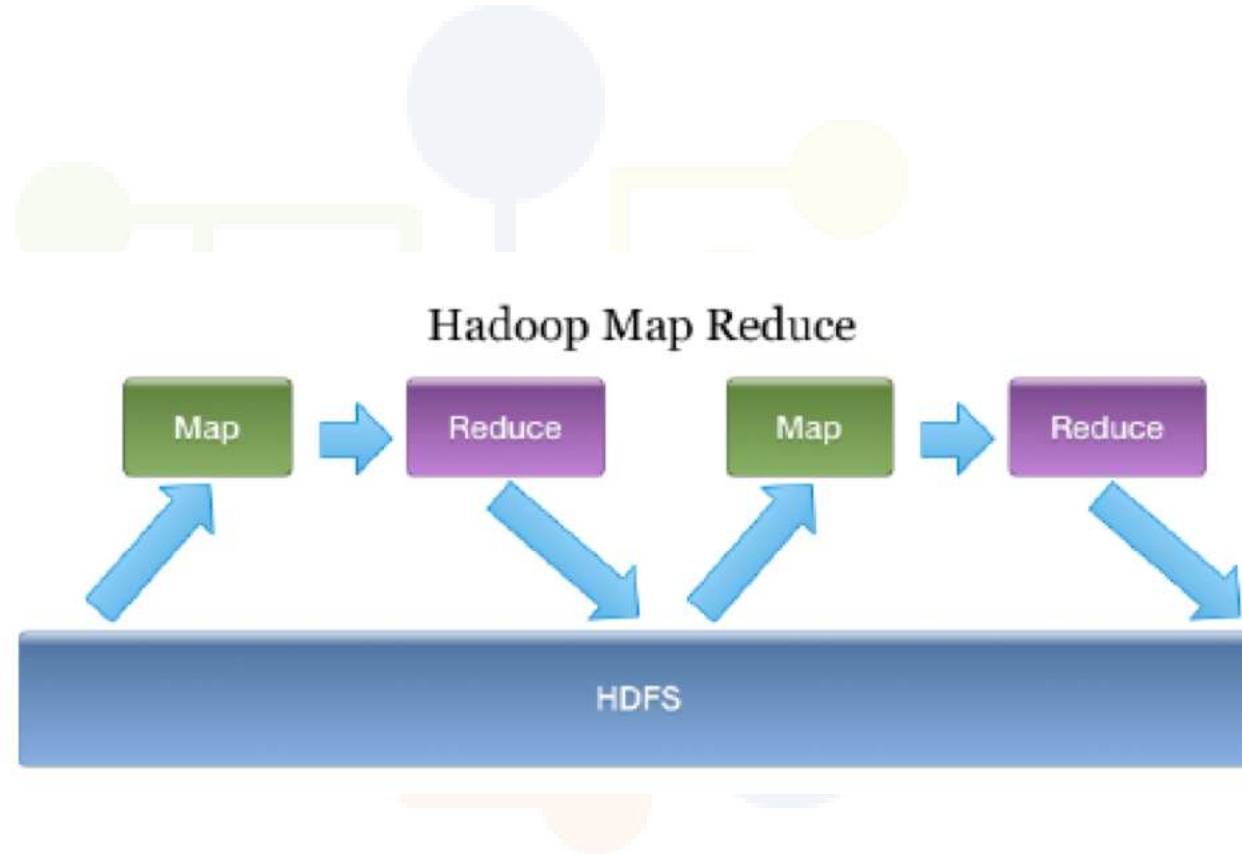


Data Science Academy

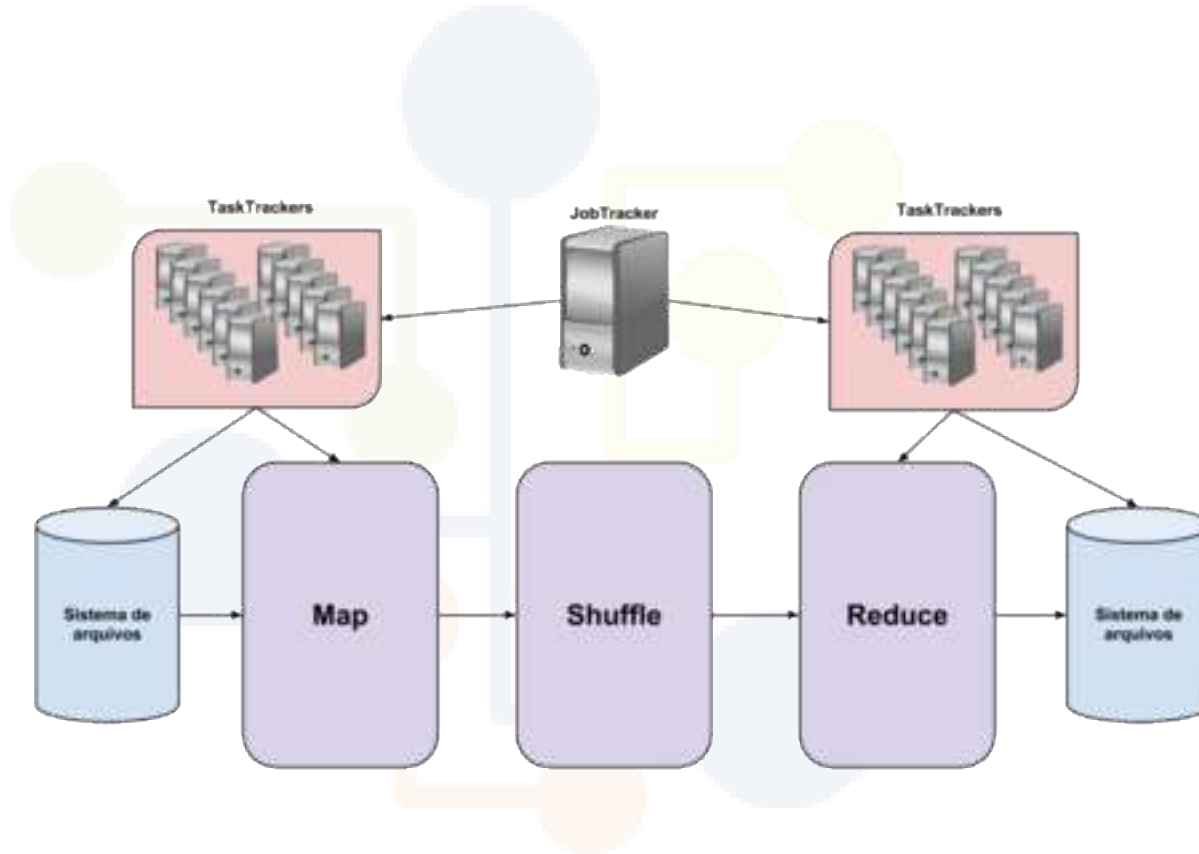
YARN



YARN



YARN



YARN

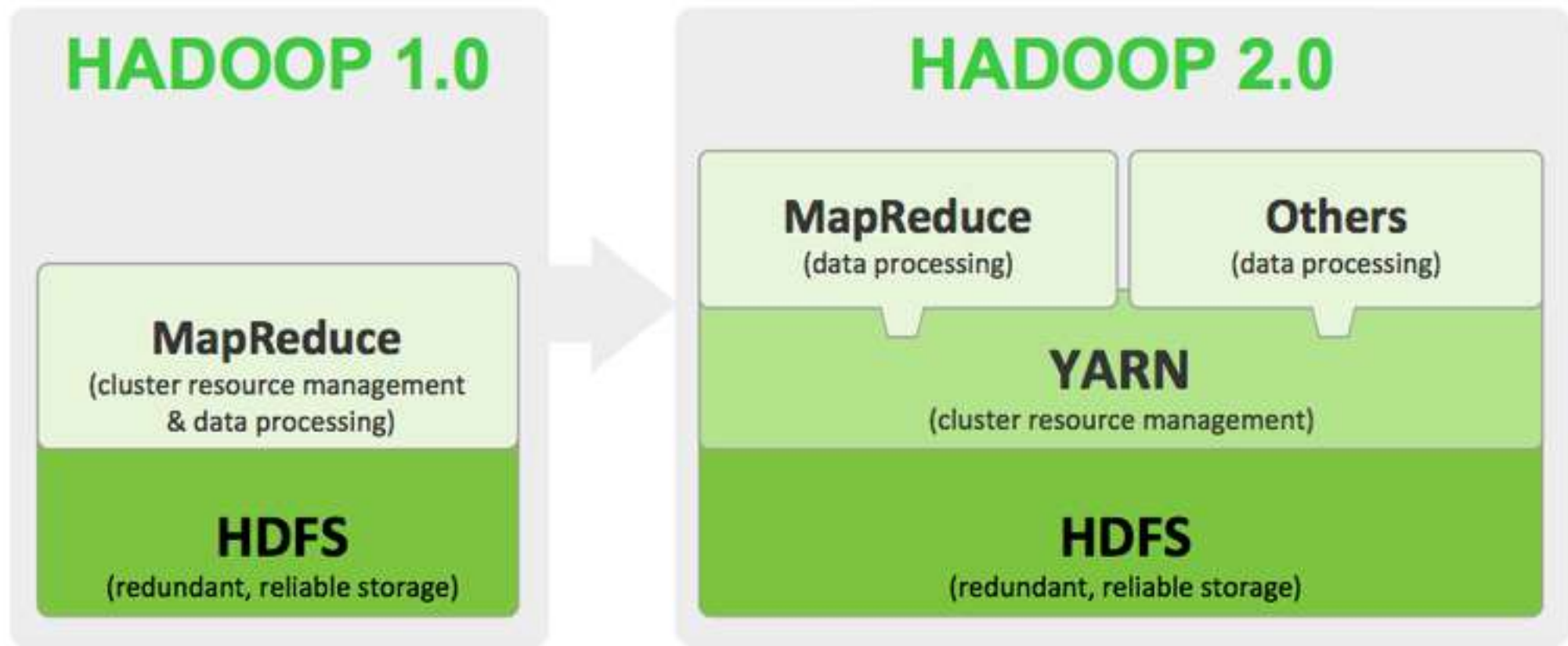
YARN

Mas uma nova ferramenta, chamada "**Yet Another Resource Negotiator**" (YARN) foi introduzida com o Hadoop 2.0 e a gestão de recursos foi transferida do MapReduce para o YARN

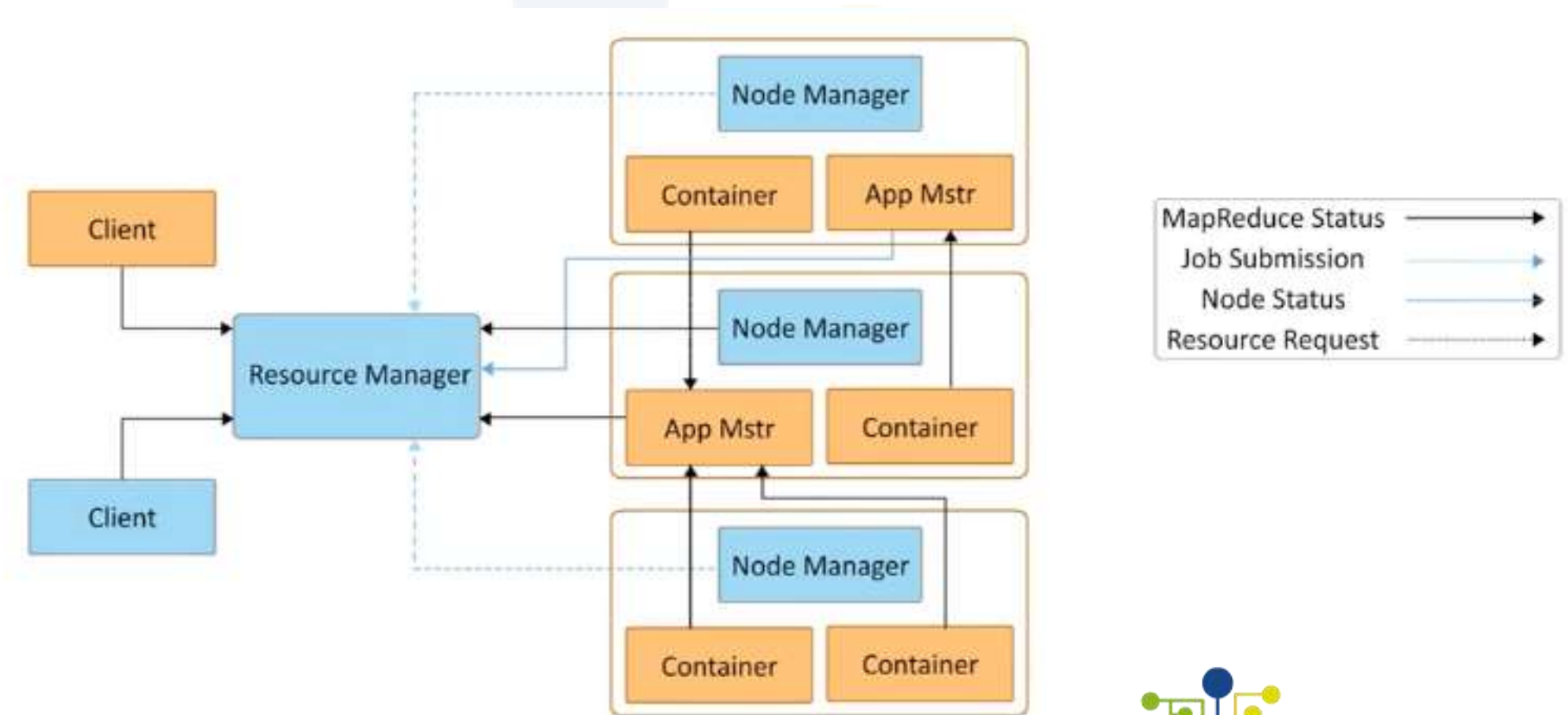


Data Science Academy

YARN



YARN



YARN

YARN



A necessidade do YARN surgiu a partir de problemas com a gestão de recursos em versões anteriores do MapReduce:

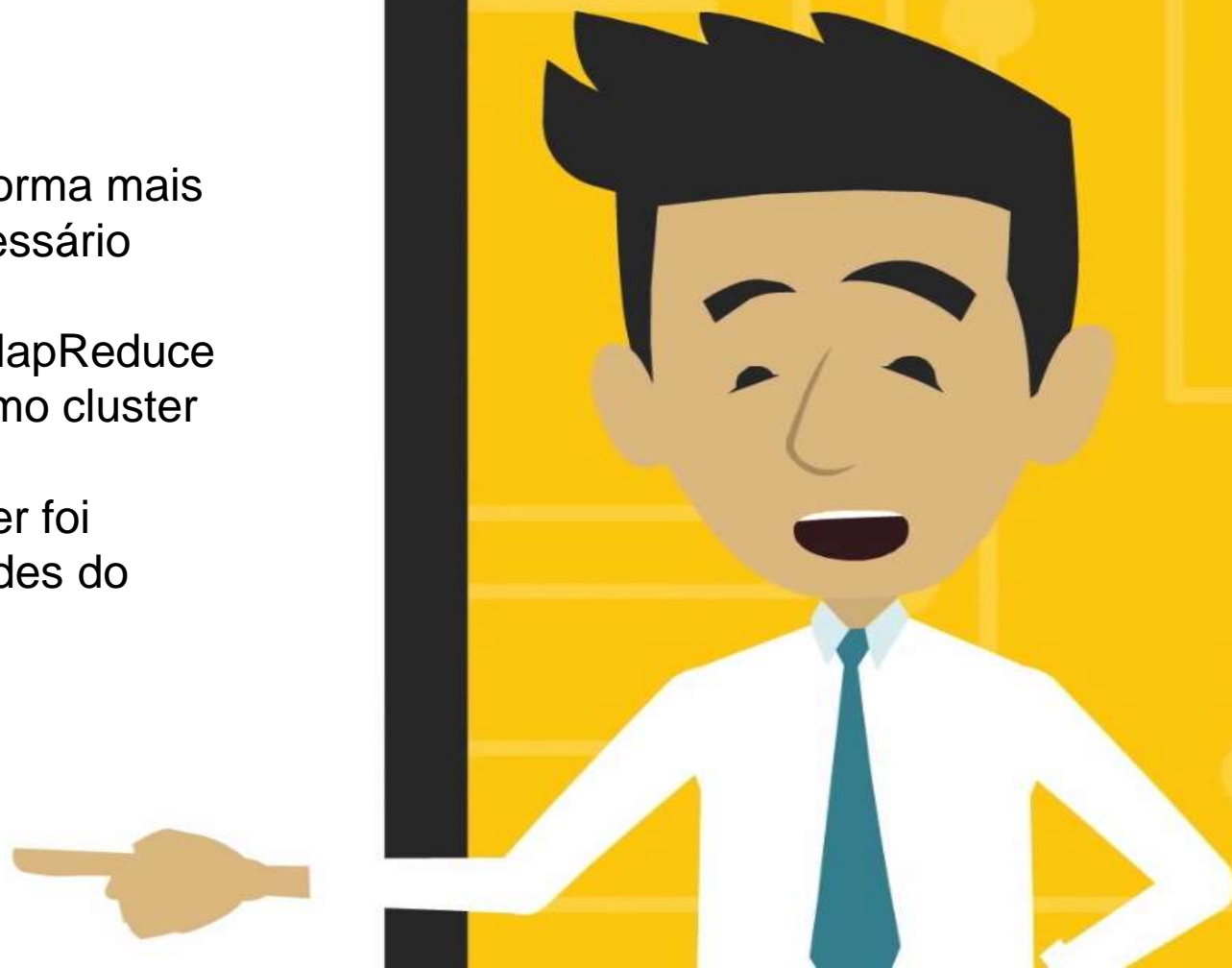
- ❖ A utilização do cluster era baixa durante as tarefas de MapReduce
- ❖ Recursos não eram compartilhados
- ❖ Havia problemas com a preferência de tarefas a serem executadas
- ❖ Havia apenas 1 JobTracker por cluster



Data Science Academy

Benefícios de utilização do YARN

- Os recursos são aplicados de forma mais eficiente e apenas quando necessário
- Aplicação MapReduce e Não-MapReduce podem ser executadas no mesmo cluster
- O conceito de Application Master foi introduzido e as responsabilidades do JobTracker redesenhadas



YARN

YARN Daemons

Resource Manager	Application Master
Executa no node Master	Executa nos nodes slaves
Agendador global de recursos	Comunica com o gerenciador de recursos



Data Science Academy

Obrigado



Data Science Academy

www.datascienceacademy.com.br