

UFPR - Universidade Federal do Paraná  
SEPT - Setor de Educação Profissional e Tecnológica

## IAA - Especialização em Inteligência Artificial Aplicada

### IAA013 - Big Data

Prof. João Eugenio Marynowski – [jeugenio@ufpr.br](mailto:jeugenio@ufpr.br)

# Roteiro

- Apresentação do Professor
- Introdução ao Módulo
- Apresentação dos Alunos
- Big Data
- Utilização da VM Hadoop
- Ambiente Distribuído
- Programação MapReduce

# Quem Sou Eu?

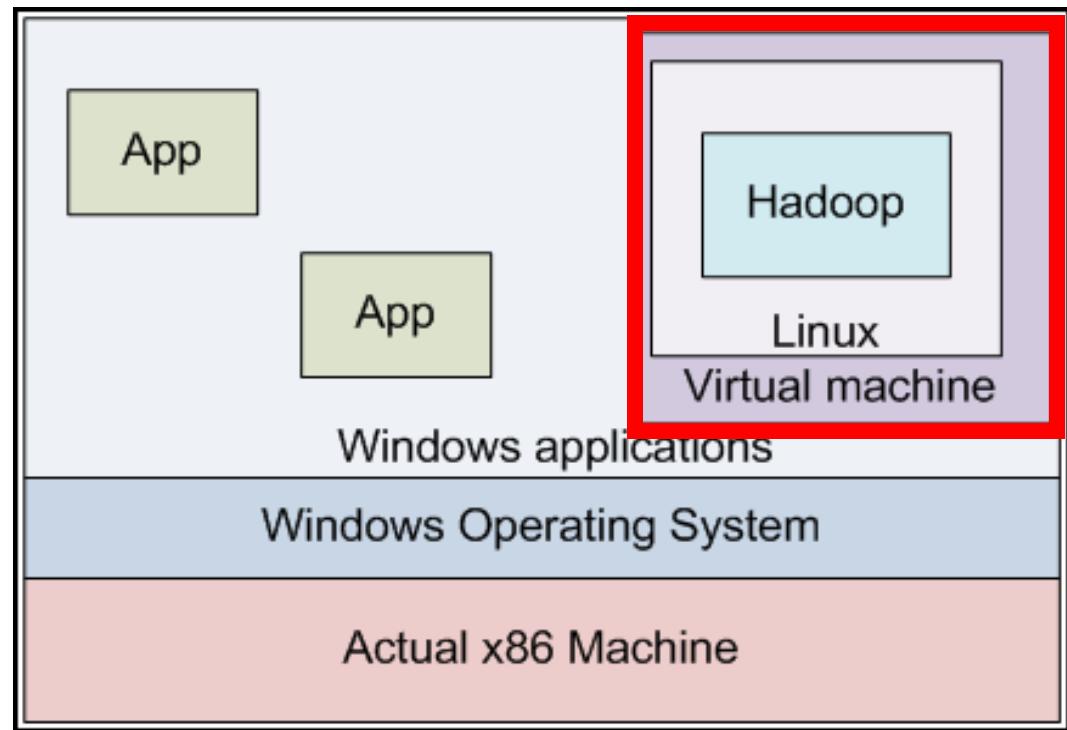
- João Eugenio Marynowski – jeugenio@ufpr.br
- Histórico/Experiência
  - '01 Bacharel em Ciência da Computação, UNIOESTE-Foz
  - '04 Mestre em Ciência da Computação, UFPR
  - Programador, Professor, Analista de Seg., DBA, Gerente TI, Empresário
  - '13 Doutor em Ciência da Computação, UFPR
  - '14 e '15 Pós-doutorado PUCPR/PPGIa, Segurança
- '16 Professor UFPR/SEPT/TADS
- Professor e pesquisador em BD e SD (Big Data)
  - SOLID e CORE
  - Extensão

# Especialização em Inteligência Artificial Aplicada (IAA)

- Introdução à Inteligência Artificial
- Estatística Aplicada I
- Linguagem de Programação Aplicada
- Linguagem R
- Estatística Aplicada II
- Arquitetura de Dados
- Aprendizado de Máquina
- Visão Computacional
- **Big Data**
- Ciência de Dados Apl. Saúde
- Frameworks de IA
- Empreendedorismo e Inovação
- Metodologia Científica
- Tópicos em IA
- Seminários
- TCC

# Sistemas para as Aulas

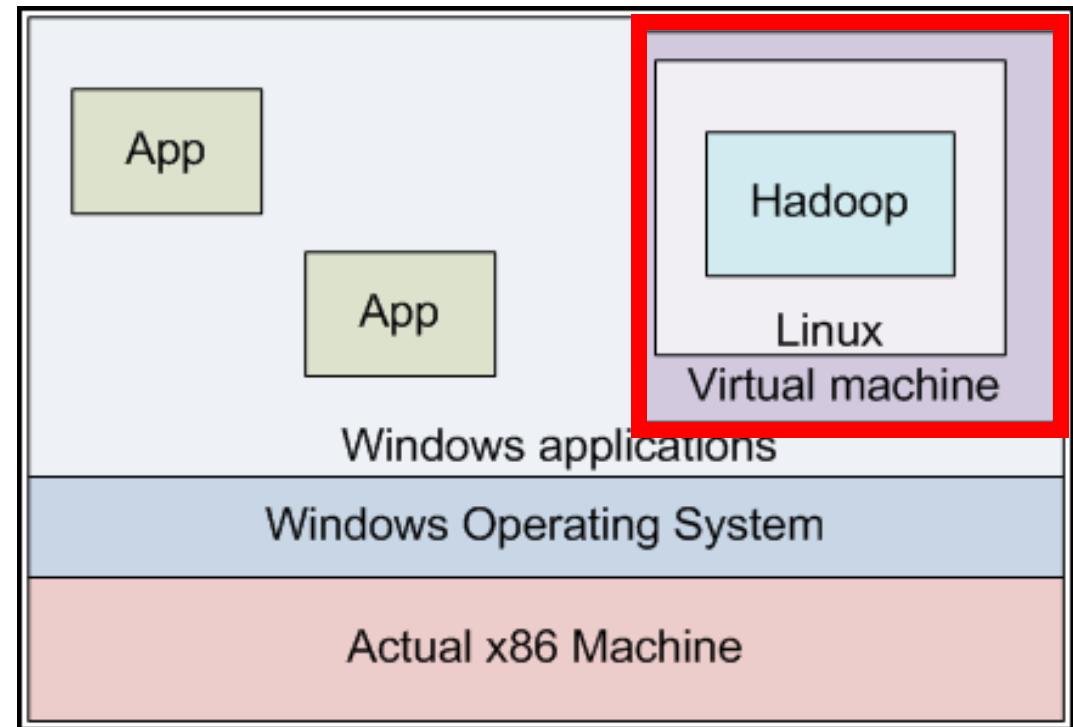
- Sandbox Hortonworks HDP 2.1



- <https://www.cloudera.com/downloads/hortonworks-sandbox/hdp.html>
- Bloqueada sem autenticação
- Permite obter Sandbox mais recentes, mas exige muito recurso (20GB)
- Então disponibilizei: Hortonworks\_Sandbox\_2.1.ova

# Sistemas para as Aulas

- Sandbox Hortonworks HDP 2.1
- VirtualBox
- Putty e WinScp
- IDE Java (Netbeans)
- VirtualBox
  - File, Import Appliance
    - Hortonworks\_Sandbox\_2.1.ova
    - CPU [2], RAM [4096 ou 2048]
    - Import



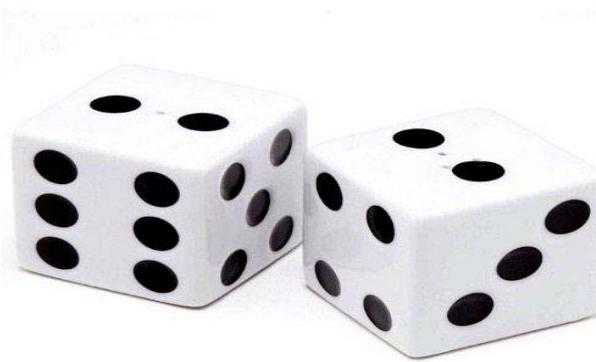
# *Big Data*



# ***Big Data***

- Muitos dados vem em mente...

- Muitos dados vem em mente...







# Big Data na Mídia



- Ah... Big Data é dados!!!
- Mas não estamos na “Era da Informação”?

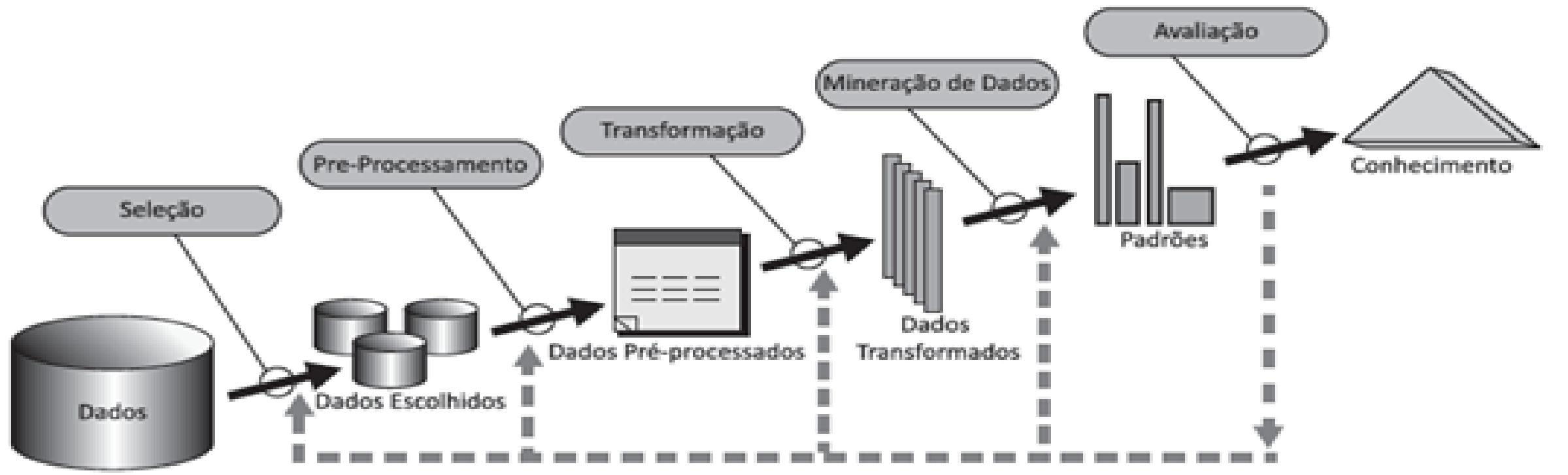
- <http://vimeo.com/31298658>

# Dado e a Tomada de Decisão



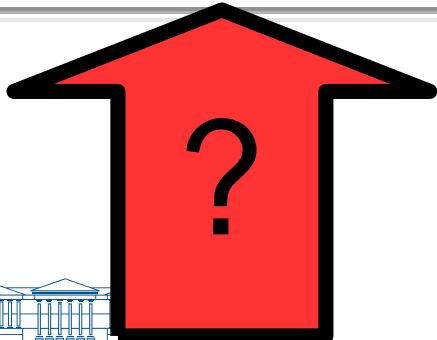
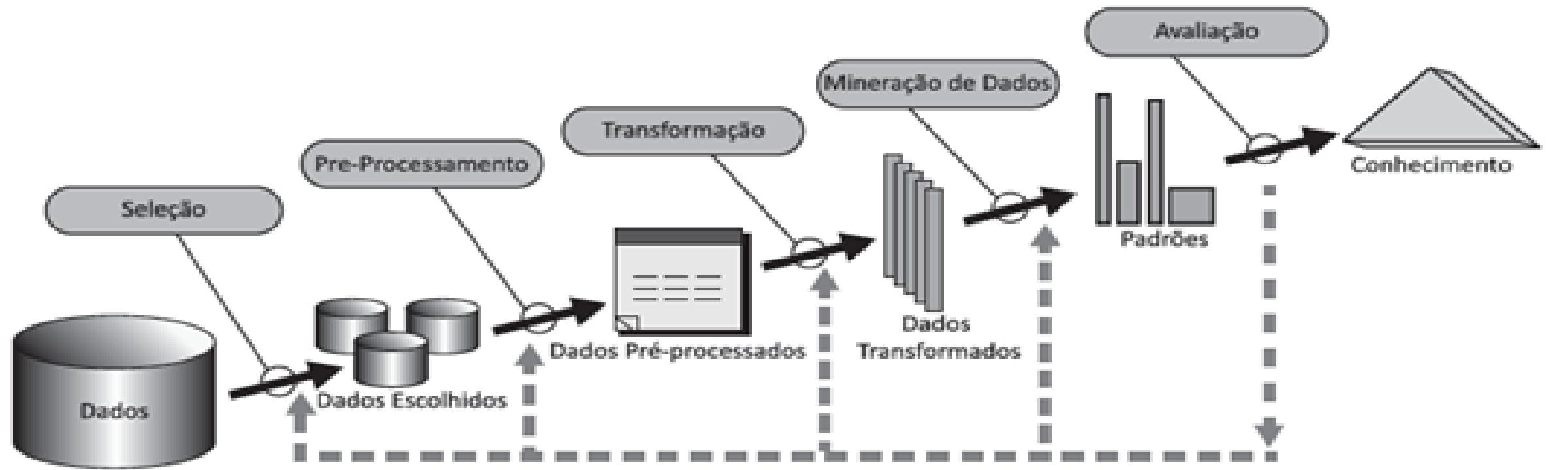
# Processo de Descoberta de Conhecimento

## KDD - Knowledge Discovery Databases



# Processo de Descoberta de Conhecimento

## KDD - Knowledge Discovery Databases

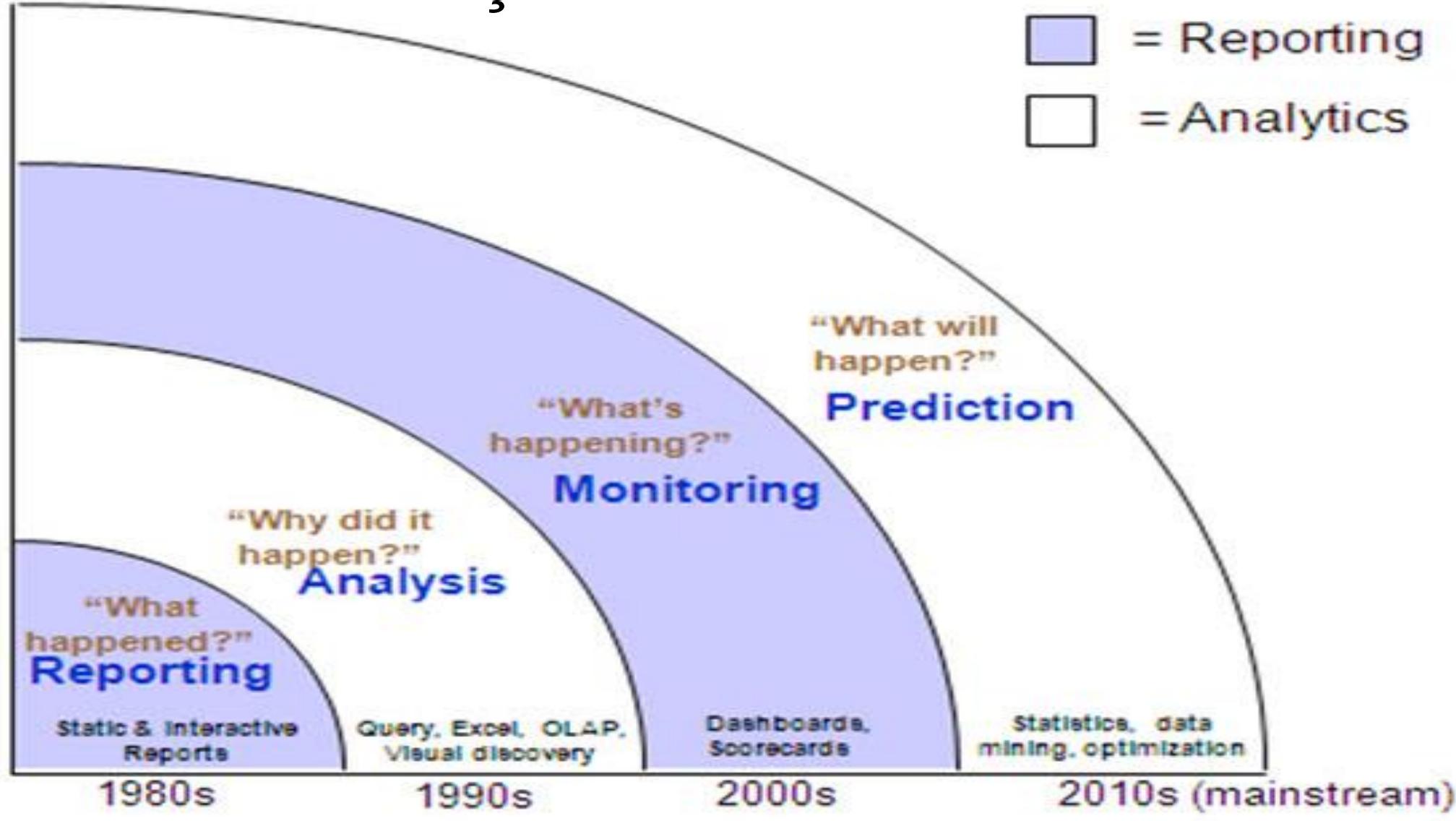


# Evolução Histórica de BI

High

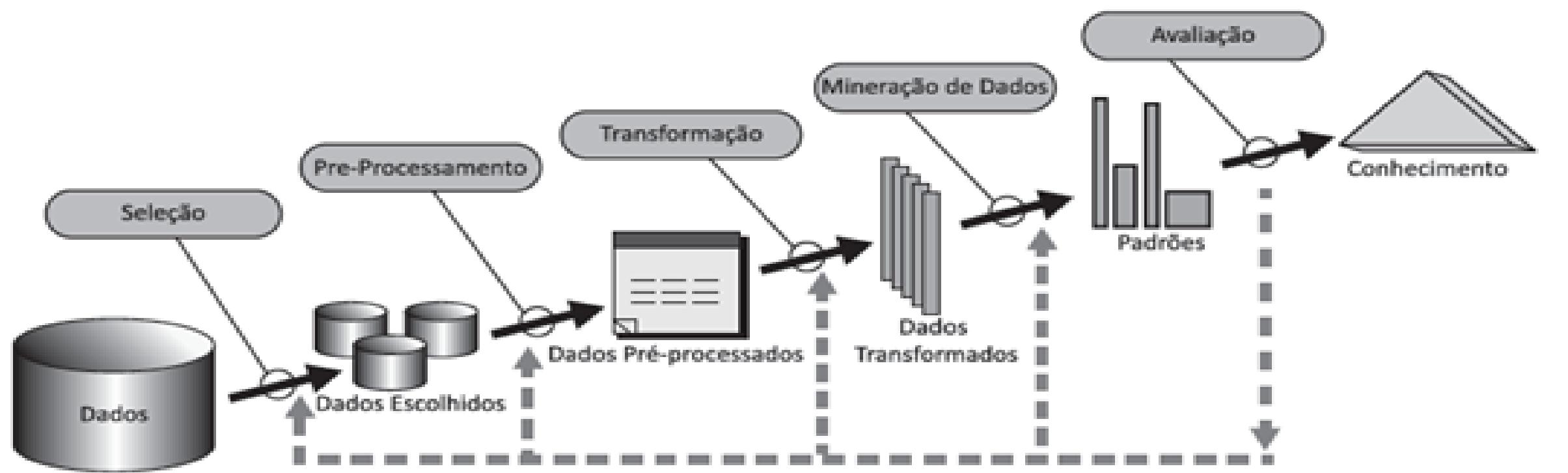
Business Value

Low



# Processo de Descoberta de Conhecimento

## KDD - Knowledge Discovery Databases



- bi mensagens/semana – 7mil/s!
  - mi transações/s
  - bi de fotos
- TB, PB, EB, ZB,...

# Problema

**Como processar um grande volume de dados?**

**1 PB – 1 computador (60MB/s) → 192 dias!**

É necessário utilizar **VÁRIOS** computadores...

**- FACILMENTE!**

# Vários Computadores...



facebook

twitter

Google  
YAHOO!

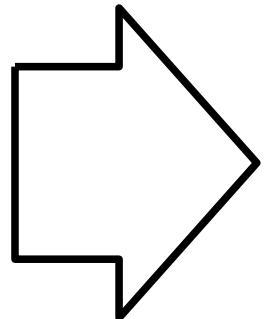
Linked in



100, 1.000, 10.000, ...

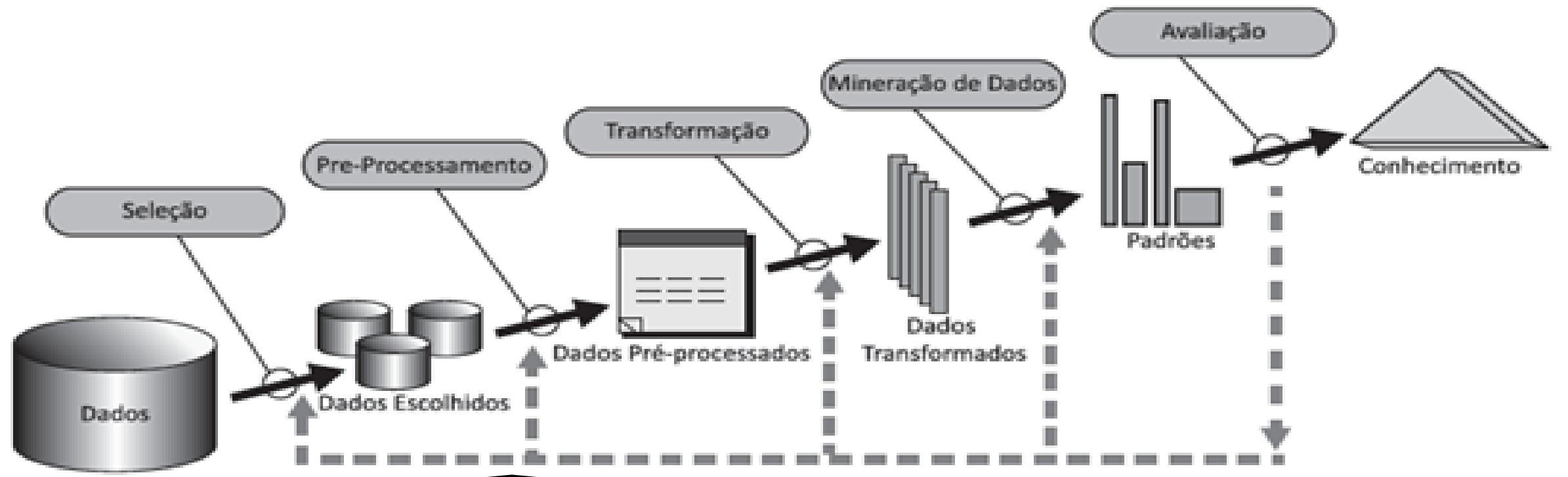
# Sistemas Big Data

***Sistemas***



# Processo de Descoberta de Conhecimento

## KDD - Knowledge Discovery Databases



ETL?

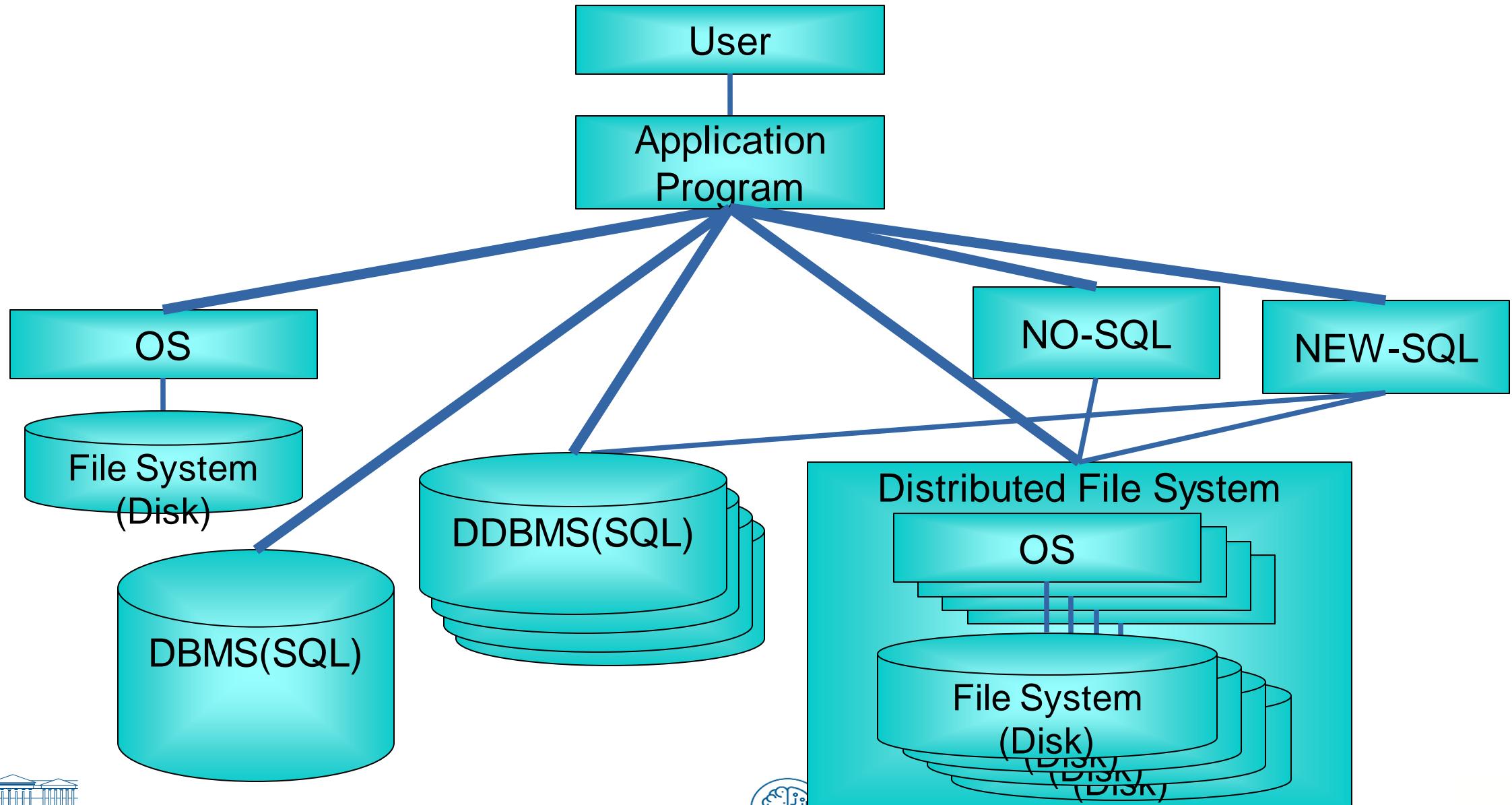
# Modelos de Dados para Big Data

- Relacional, Objeto e O-R
- Dimensional/Multidimensional (DW)
- Dados Geográficos (GIS)
- Colunar
- Chave-valor
- Documento
- Grafo (Hierárquico e Rede)
- Não estruturados
  - Excel/CSV, Áudio, Vídeo, Logs, ...

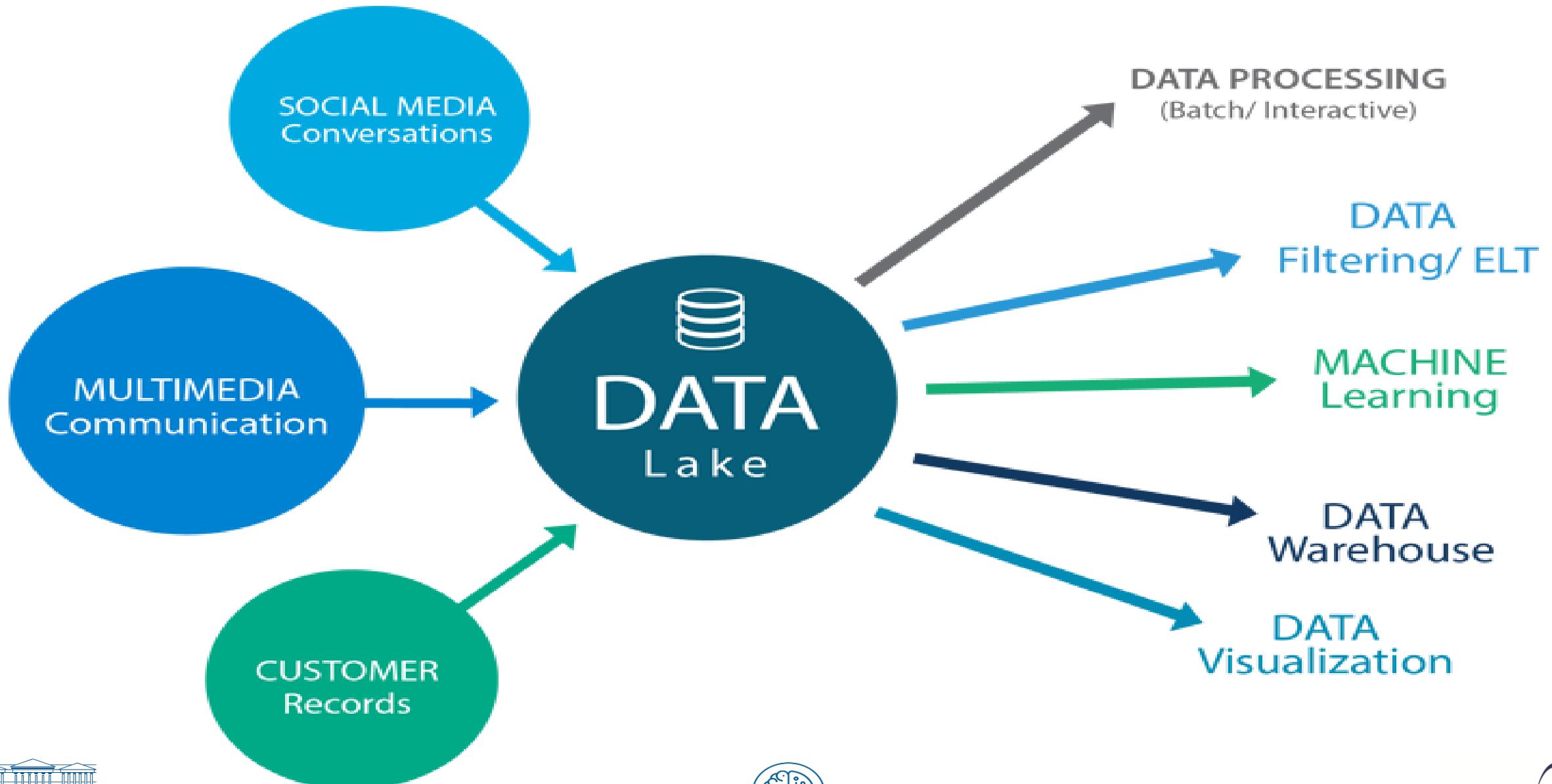
NoSQL

NewSQL

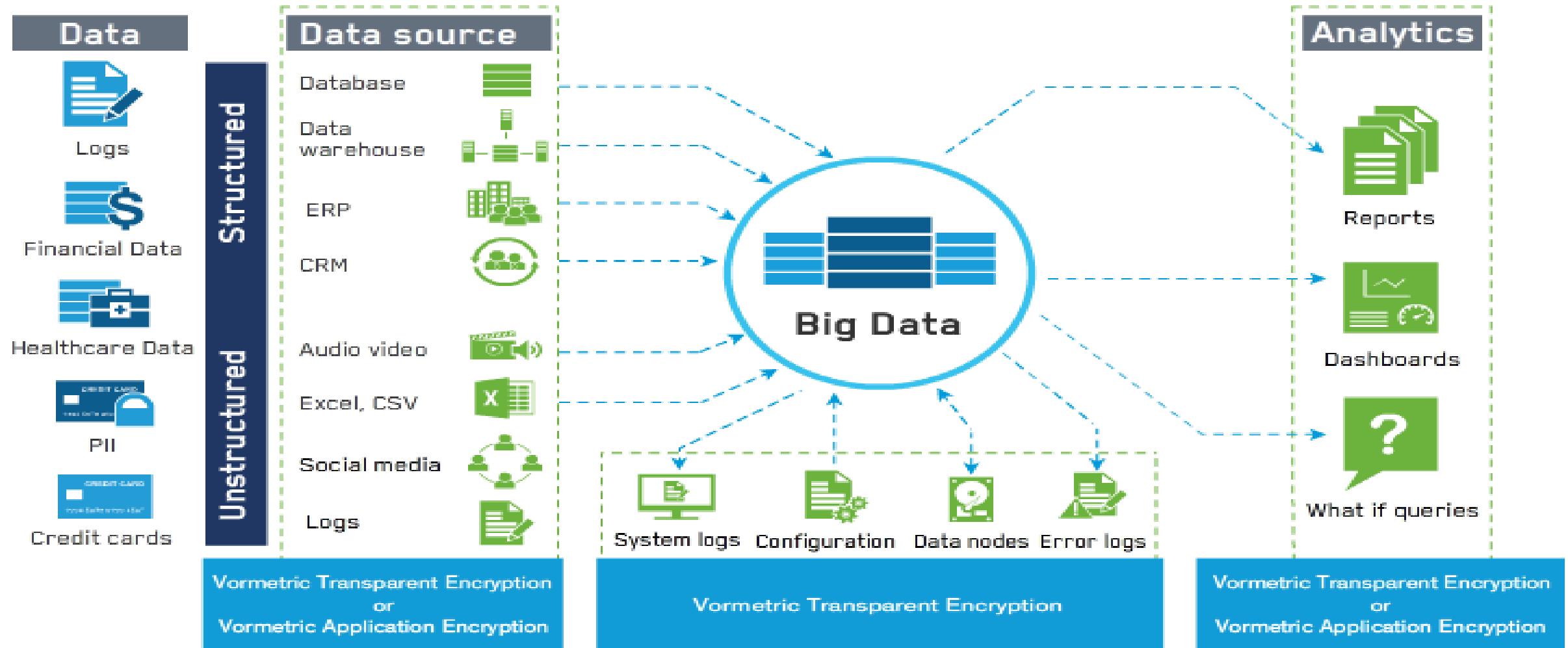
# Sistemas Big Data



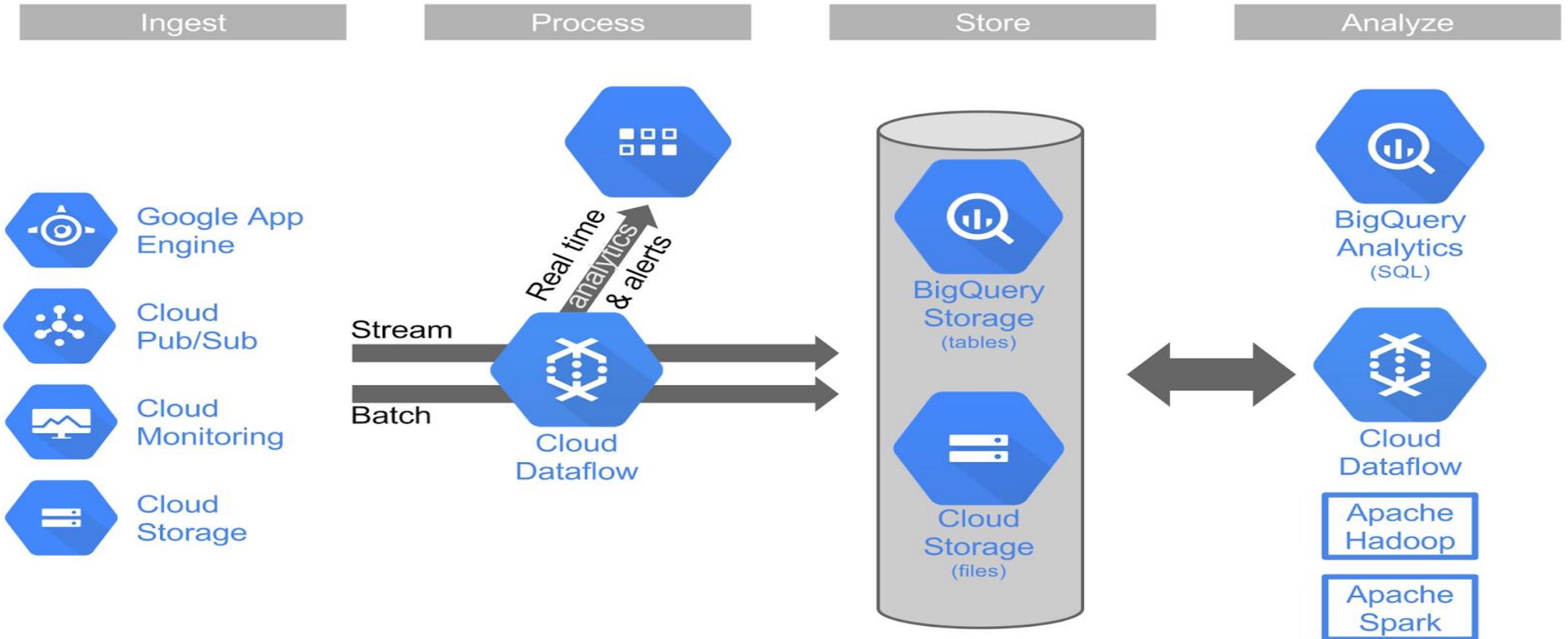
# Data Lake



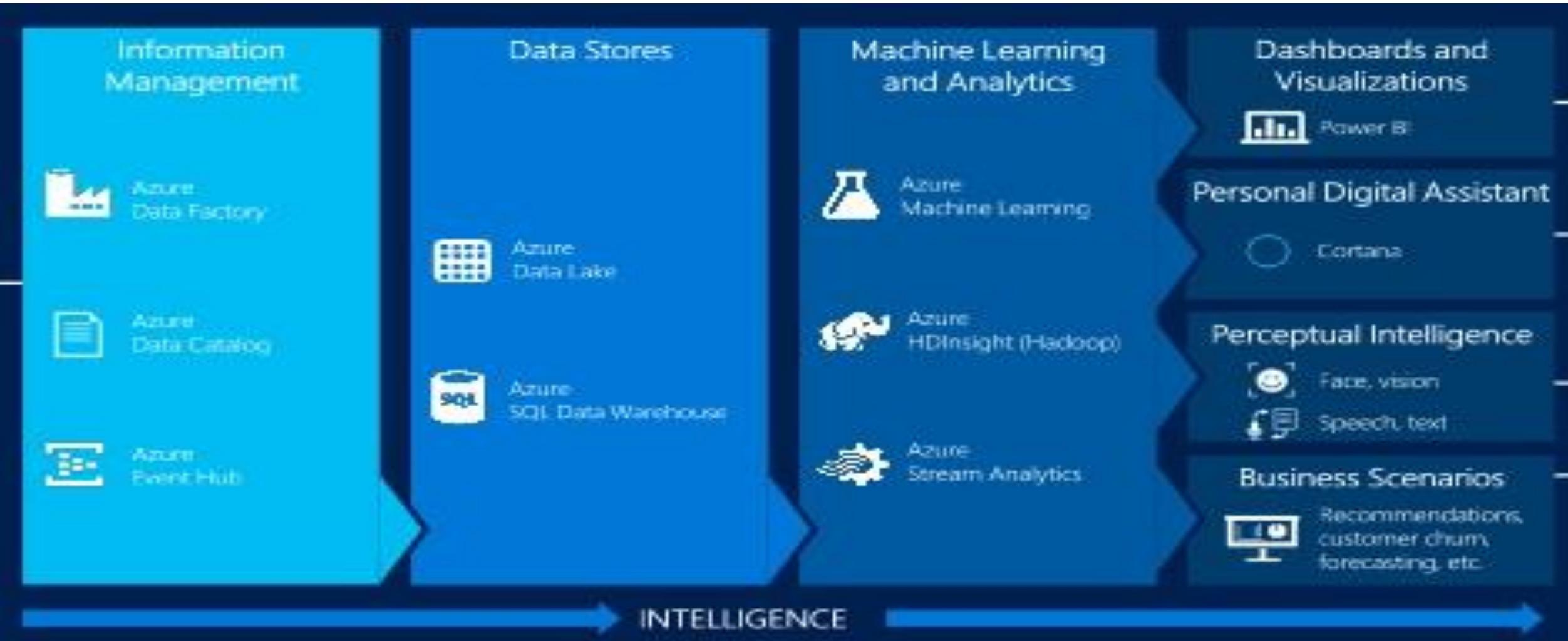
# Big Data e Data Lake



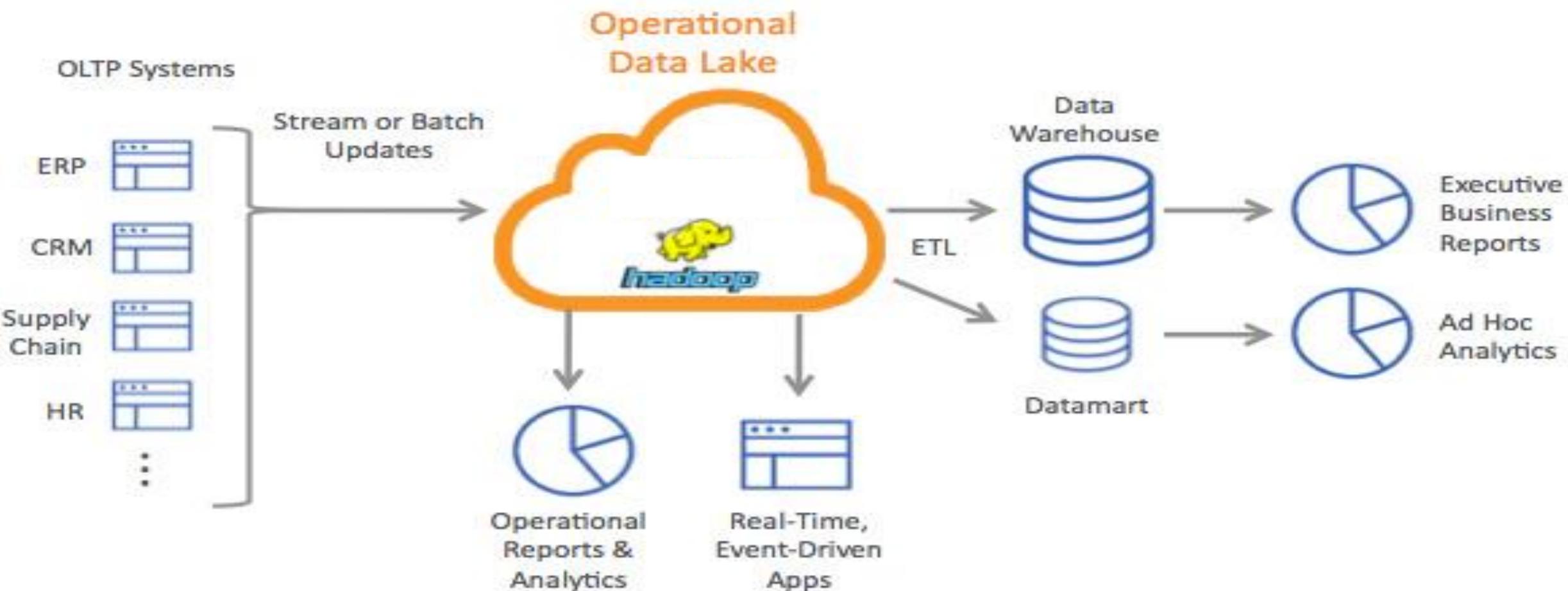
# Google Big Data Cloud



# Microsoft Data Analytics



# Data Lake



# Programa

- Introdução à Big Data
- Data Engineering (Gerenciamento/Ferramentas para Big Data)
- NoSQL e NewSQL
- Dados em movimento – Processamento de Streaming
  - Utilização da VM hortonworks 2.1
  - Ambiente distribuído: DFS, MapReduce/Hadoop e Ambari
  - Personalização de aplicações MR
  - Ecosistema Hadoop HDP 2.1

# Referências

- Hurwitz, J.; Nugent, A.; Dr. Halper, F.; Kaufman, M. **Big Data Para Leigos**. Alta Books Editora, 2016.
- White, T. (2012). **Hadoop: The Definitive Guide**, 3rd Edition (3rd ed., p. 688). O'Reilly Media. <http://it-ebooks.info/book/635/>
- Sakr, S., Liu, A., & Fayoumi, A. G. (2013). **The Family of MapReduce and Large-Scale Data Processing Systems**. ACM Computing Surveys, 46(1), 1–44.
- SILBERSCHATZ, A. KORTH, H; SUDARSHAN, S.. **Sistema de Banco de Dados**. 6.ed. Campus, 2012.
- Diversos Artigos
- Dean, J., & Ghemawat, S. (2004). **MapReduce: Simplified Data Processing on Large Clusters**. In Proc. of the OSDI - Symp. on Operating Systems Design and Implementation (pp. 137–149). USENIX.

# Avaliação

- Presenças
  - 70% (aulas e/ou atividades)
- Atividades
  - Entrega via UFPR Virtual
  - Nota  $\geq 70$

# Quando?

9 às 12h e das 13:30 às 17h

Sábados:

19/11

26/11

03/12

# Sandbox no Virtualbox

- File, Import Appliance
  - Hortonworks\_Sandbox\_2.1.ova
  - CPU [2], RAM [4096 ou 2048]
  - ... Import

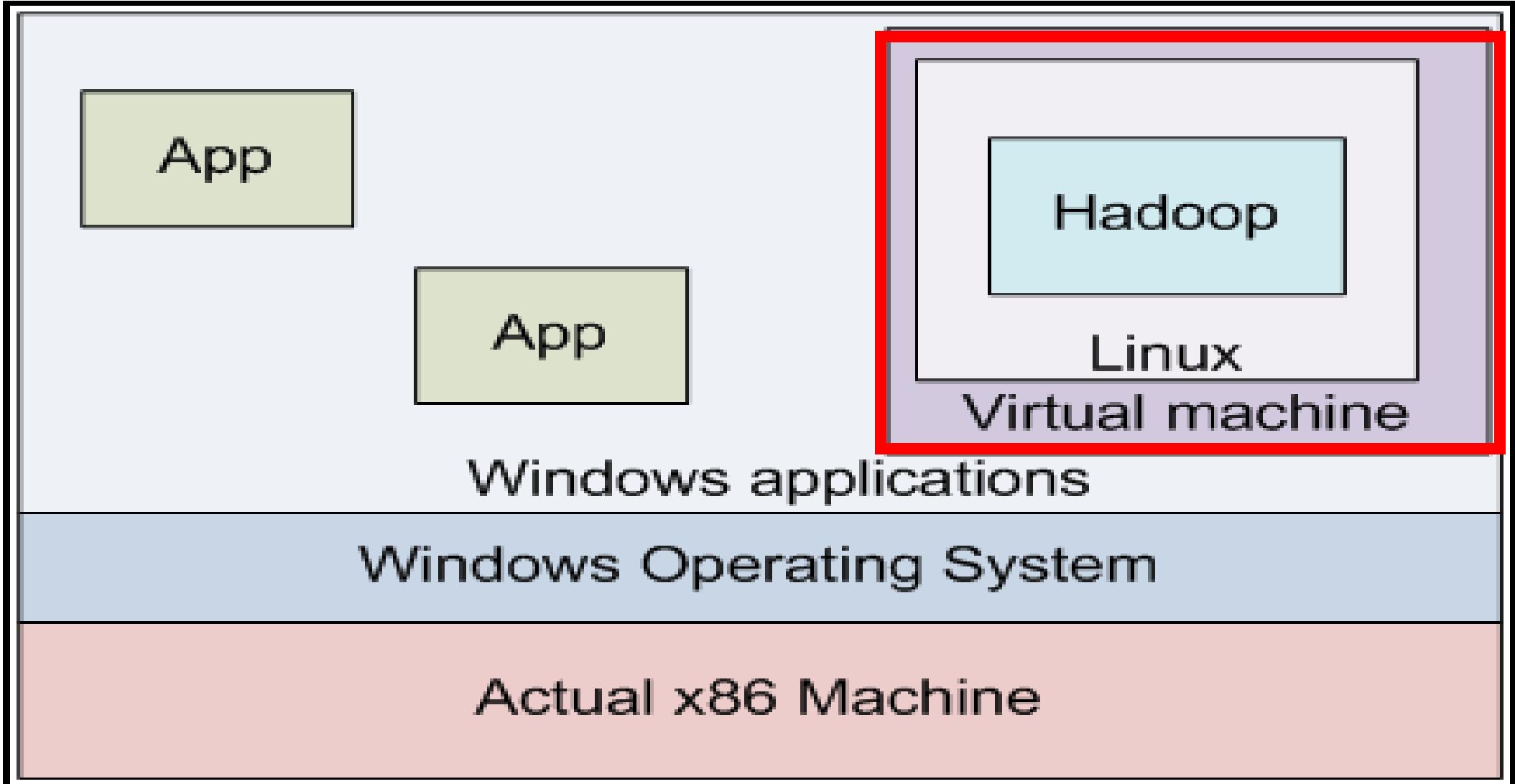
# Quem São VocêS?

- Nome
- Experiência relacionada com Big Data
- Objetivos futuros

# Sandbox no Virtualbox

- File, Import Appliance
  - Hortonworks\_Sandbox\_2.1.ova
  - CPU [2], RAM [4096 ou 2048]
  - ... Import
- ...
- **Start**

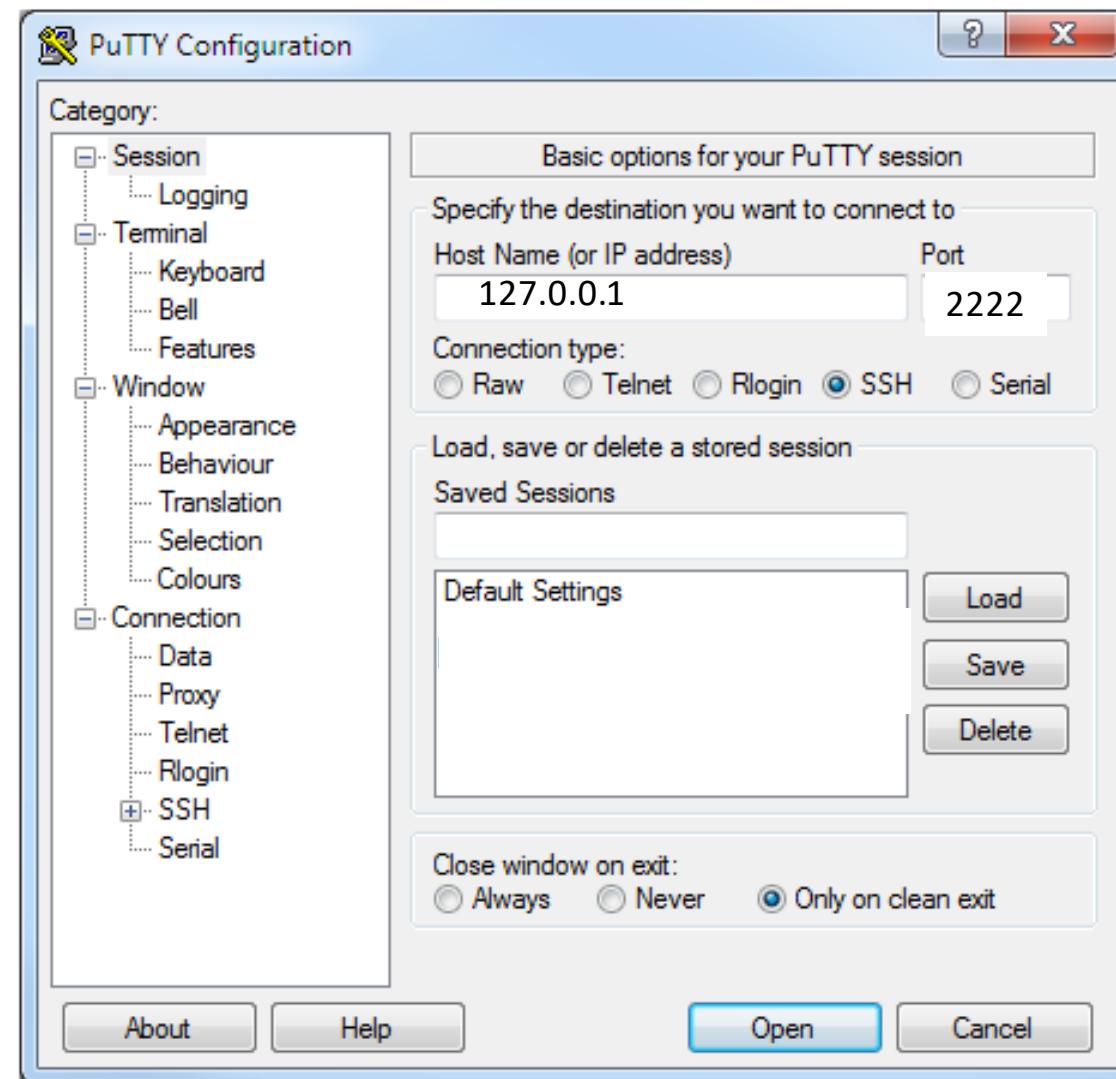
# Sandbox HDP 2.1



# Acesso à Sandbox HDP

- Linux: CentOS, usuário **root** senha **hadoop**
- Acesso Local
  - Terminal → Alt + F5
- Acesso “Remoto” via ip 192.\*.\*.\* ou 127.0.0.1
  - ssh root@<IP> [-p 2222]
  - **Putty (Windows)**

# Putty



## Atividade 1 (Salvem prints da tela para envio)

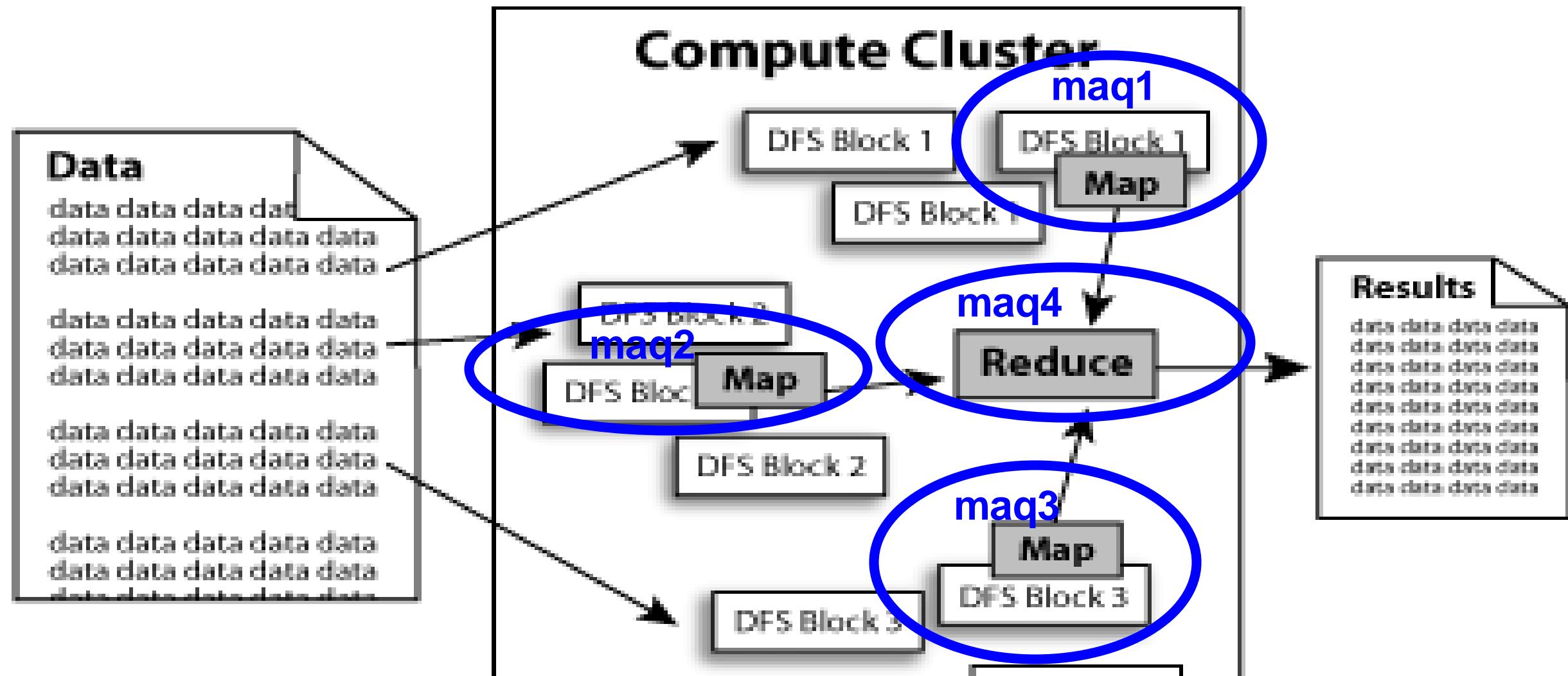
- 1) \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 1 100
  - 2) \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 1 500000
  - 3) \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 5 100000
- Comentem sobre a diferença de execução dos exercícios 2 e 3  
(Quais diferenças? Tempo? Por que?)

# Atividade 1 (Salvem prints da tela para envio) FEITO!

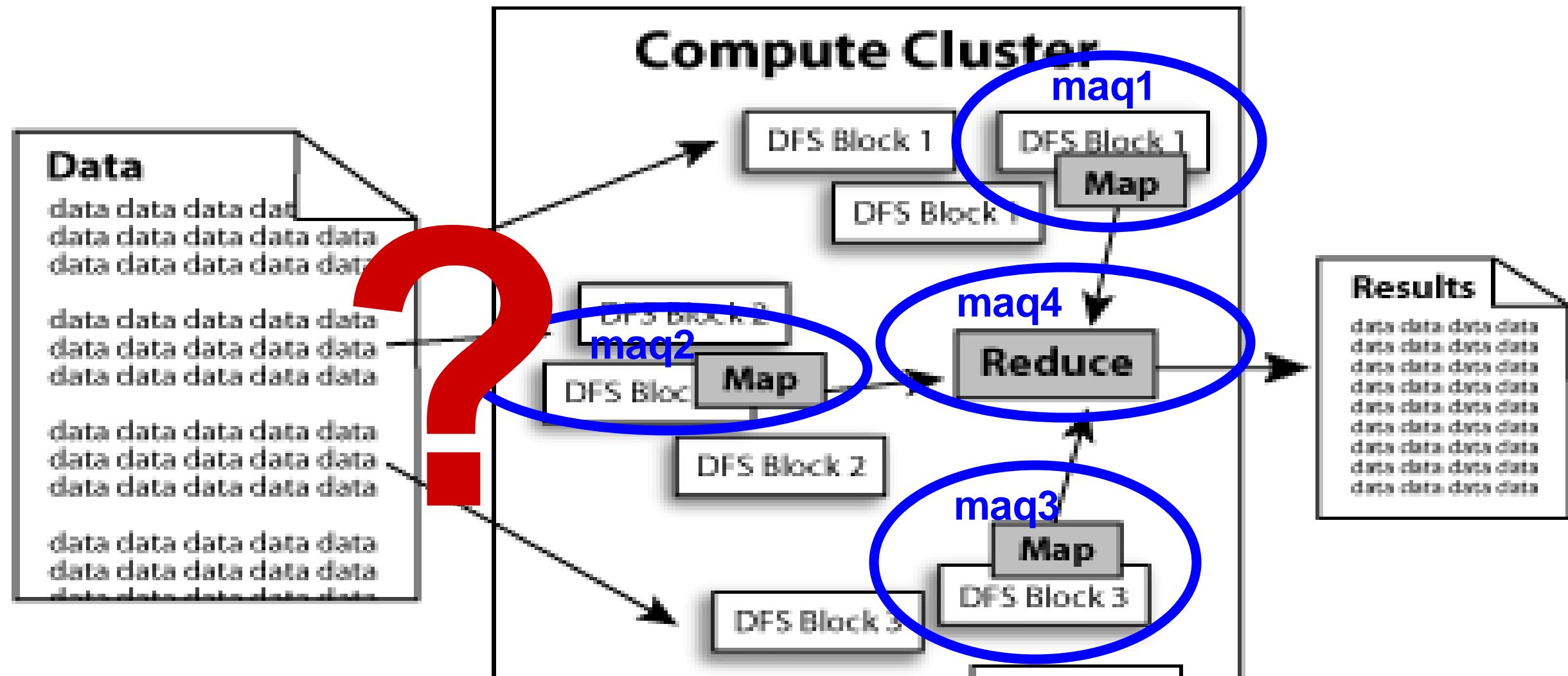
- 1) \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 1 100
  - 2) \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 1 500000
  - 3) \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar pi 5 100000
- Comentem sobre a diferença de execução dos exercícios 2 e 3  
(Quais diferenças? Tempo? Por que?)

# Como processo um arquivo?

# Distributed File System (DFS) + MapReduce



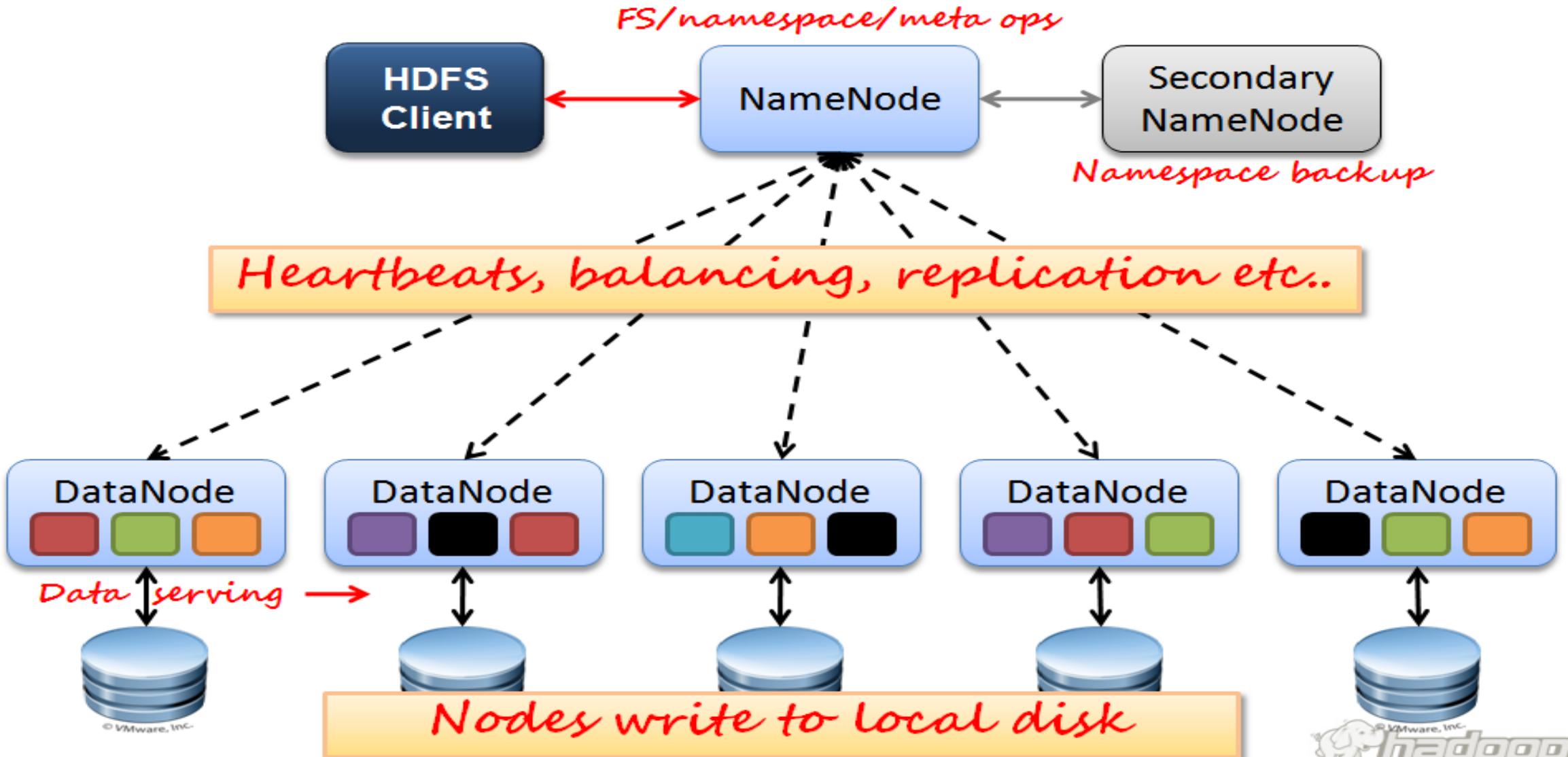
# Distributed File System (DFS) + MapReduce

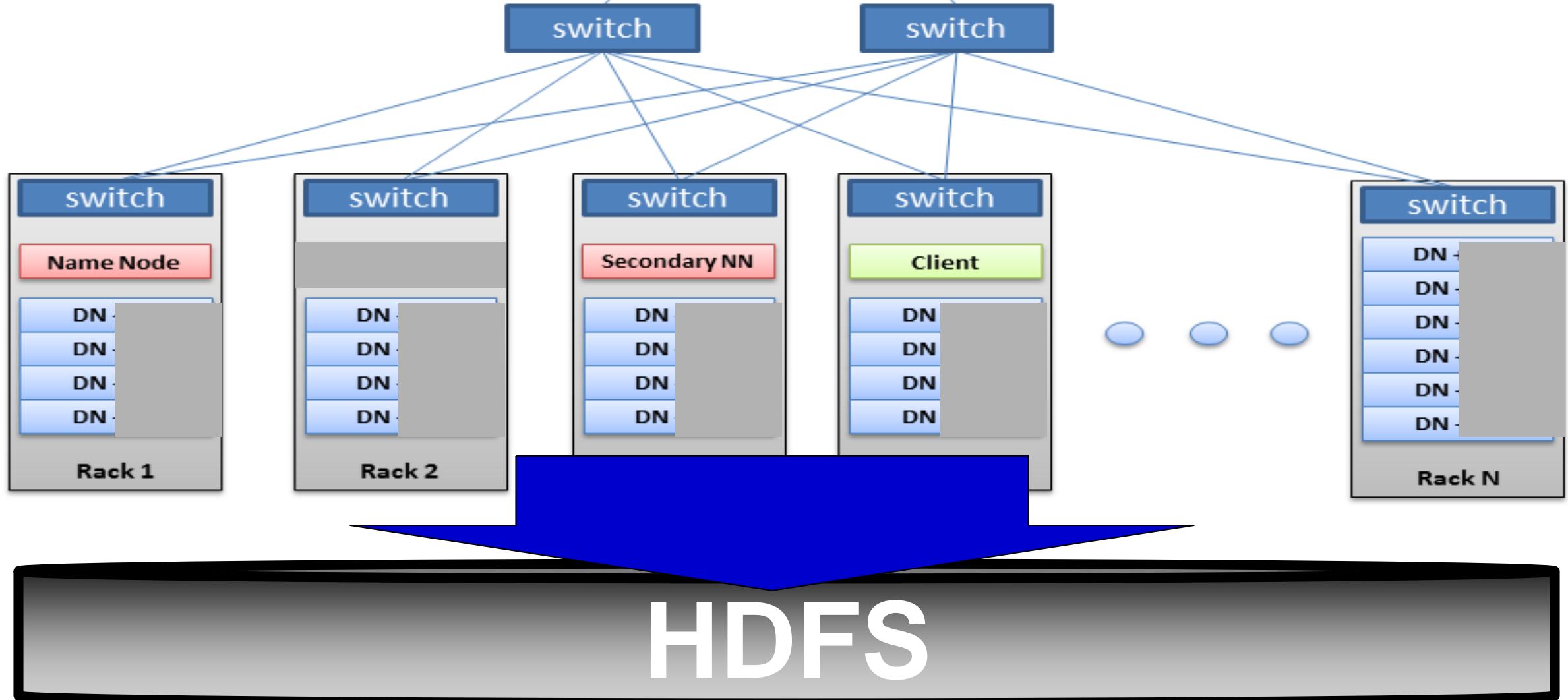


# GFS e HDFS

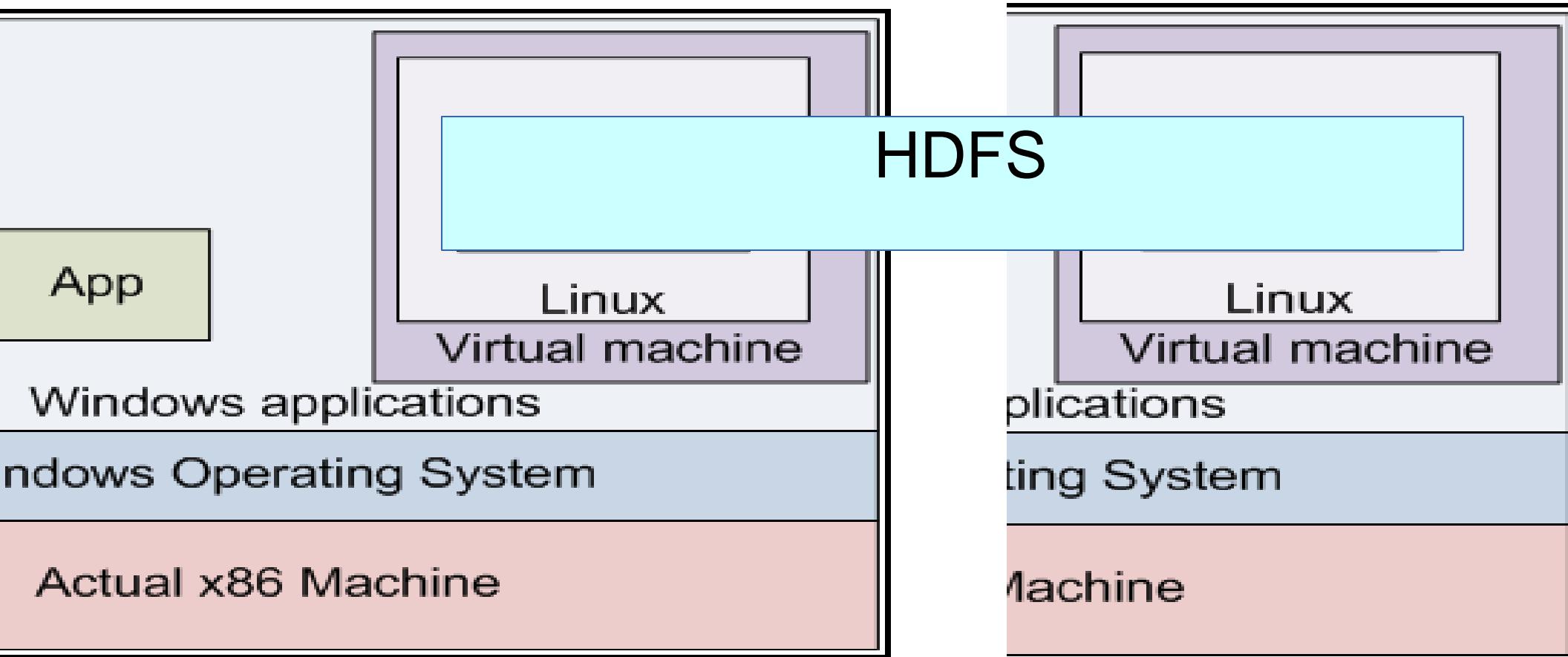
- **The Hadoop Distributed File System (HDFS)**
  - Permite acesso como um disco único
  - Executa sobre sistema de arquivo nativo (EXT3, EXT4, XFS,...)
  - Tolerante a falhas (Máquina, Rede, Disco, ...)
- Usa máquinas “comuns”
  - 2003: 669 dual 1.4 GHz, 2GB, 80GB
  - 2010: +25mil/3500 quad xeon 2.5GHz, 16GB, 4x1TB, 25 PB

# Arquitetura do HDFS





# Sandbox Hadoop x2, 3, 4,...



# Características do HDFS

- Armazenamento de grandes arquivos
  - GB, TB, PB, EB, ...
  - Milhões ao invés de bilhões de arquivos
  - (arquivo > 128 MB – bloco padrão)
- Escreve um e lê vários (write once and read-many)
- Leitura Sequencial (streaming)
  - Ruim para leitura aleatória
  - Operação de adição ao final (Append)
- Alta vazão (throughput)
  - Alta latência para pequenos pedaços (chunks)
  - (HBase trata bem muitos chunks)

# Comandos do HDFS

- HDP 2.1 → \$ hdfs dfs

- Shell commands format

```
$ hdfs dfs -<command> -<option> <path>
```

```
$ hdfs dfs -ls /
```

```
$ hdfs dfs -ls hdfs://localhost/to/path/dir
```

```
$ hdfs dfs -ls file:///to/path/file3
```

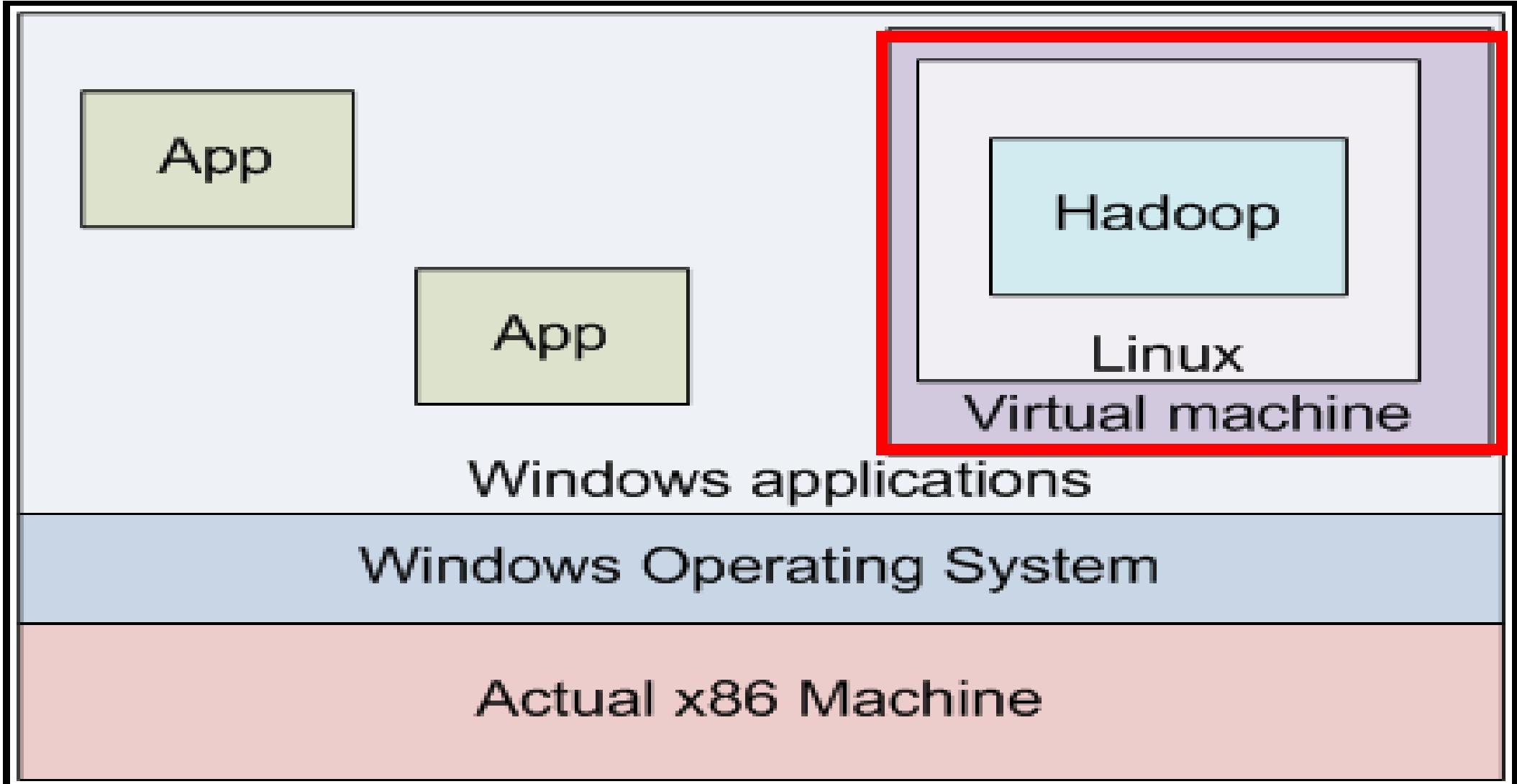
- Unix like (cat, du, rm, cp, mkdir, mv, put, get, tail, chmod,...)

```
$ hdfs dfs -help
```

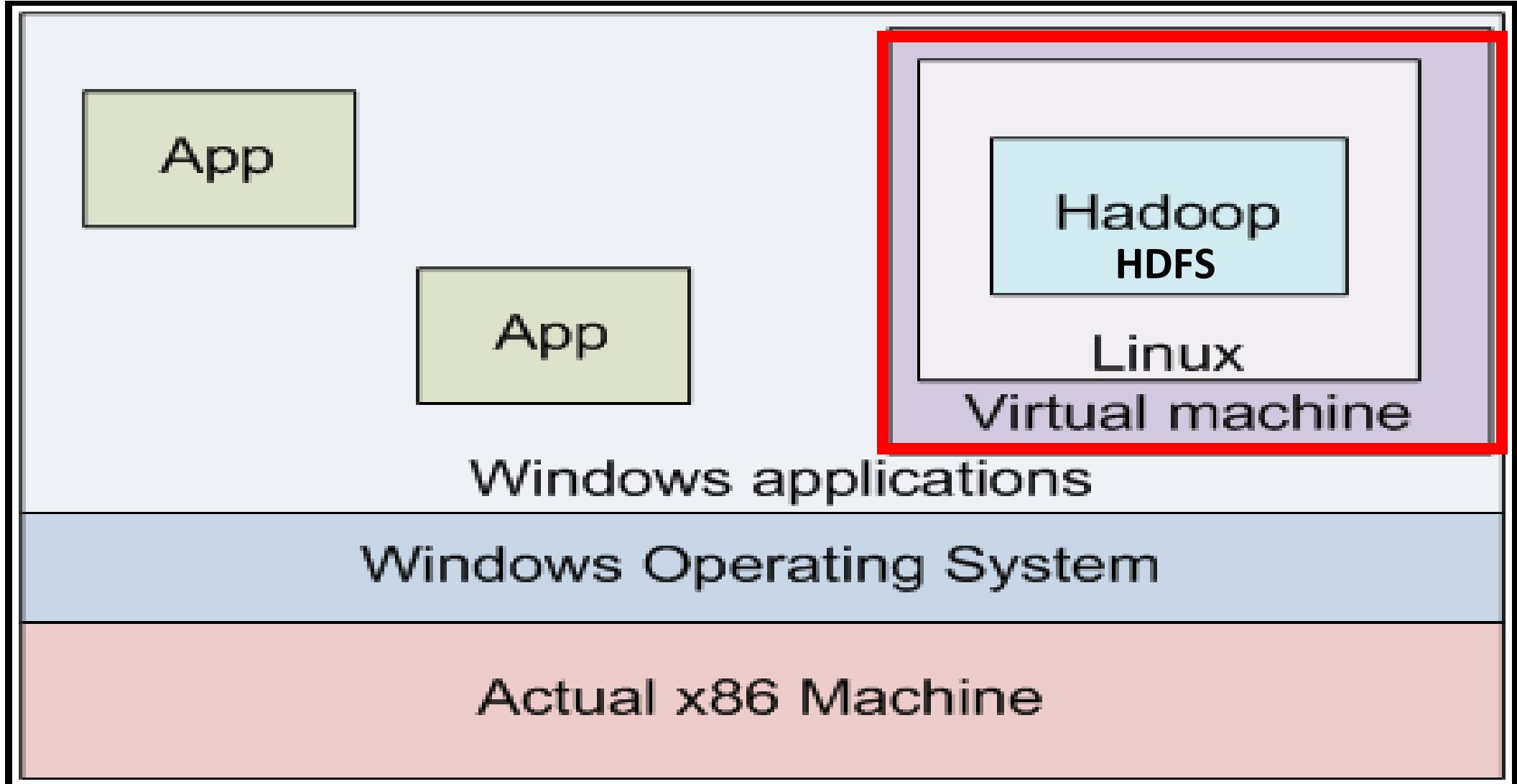
```
$ hdfs dfs -help <command_name>
```

# Sandbox Hadoop

HDFS ???



# Sandbox Hadoop



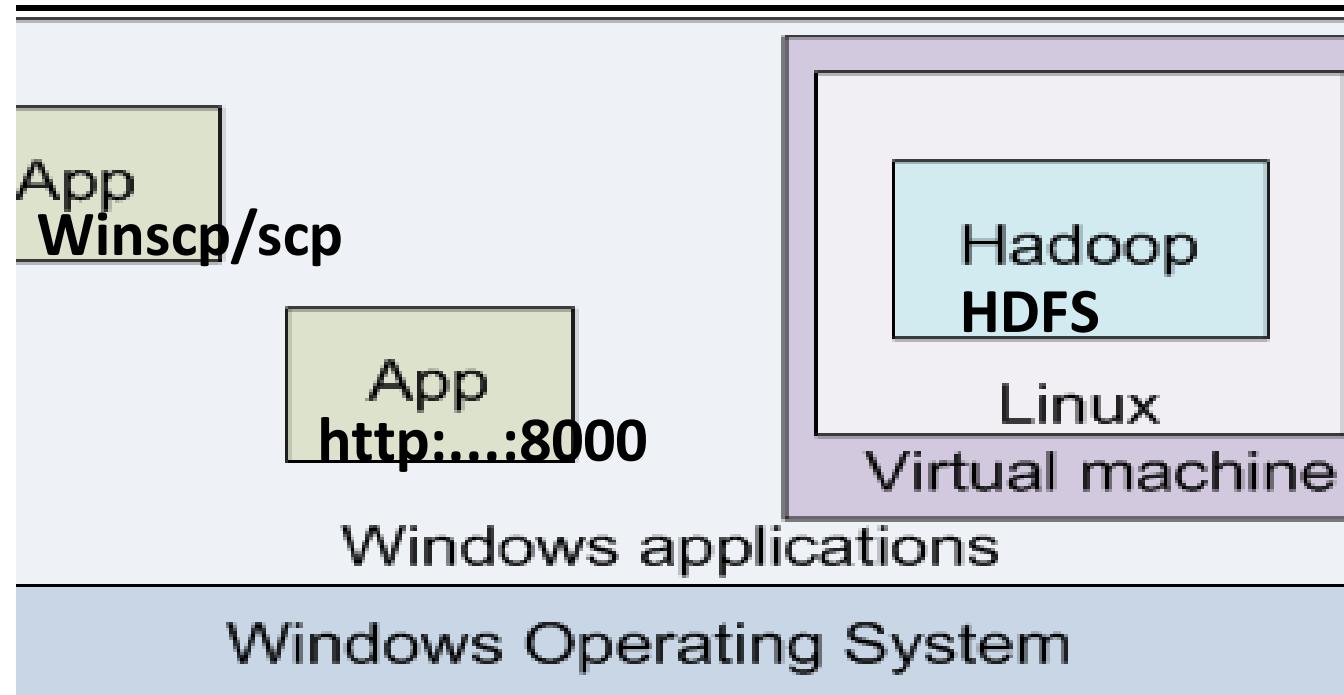
## Atividade 2 – Executar o Contador de Palavras (wordcount)

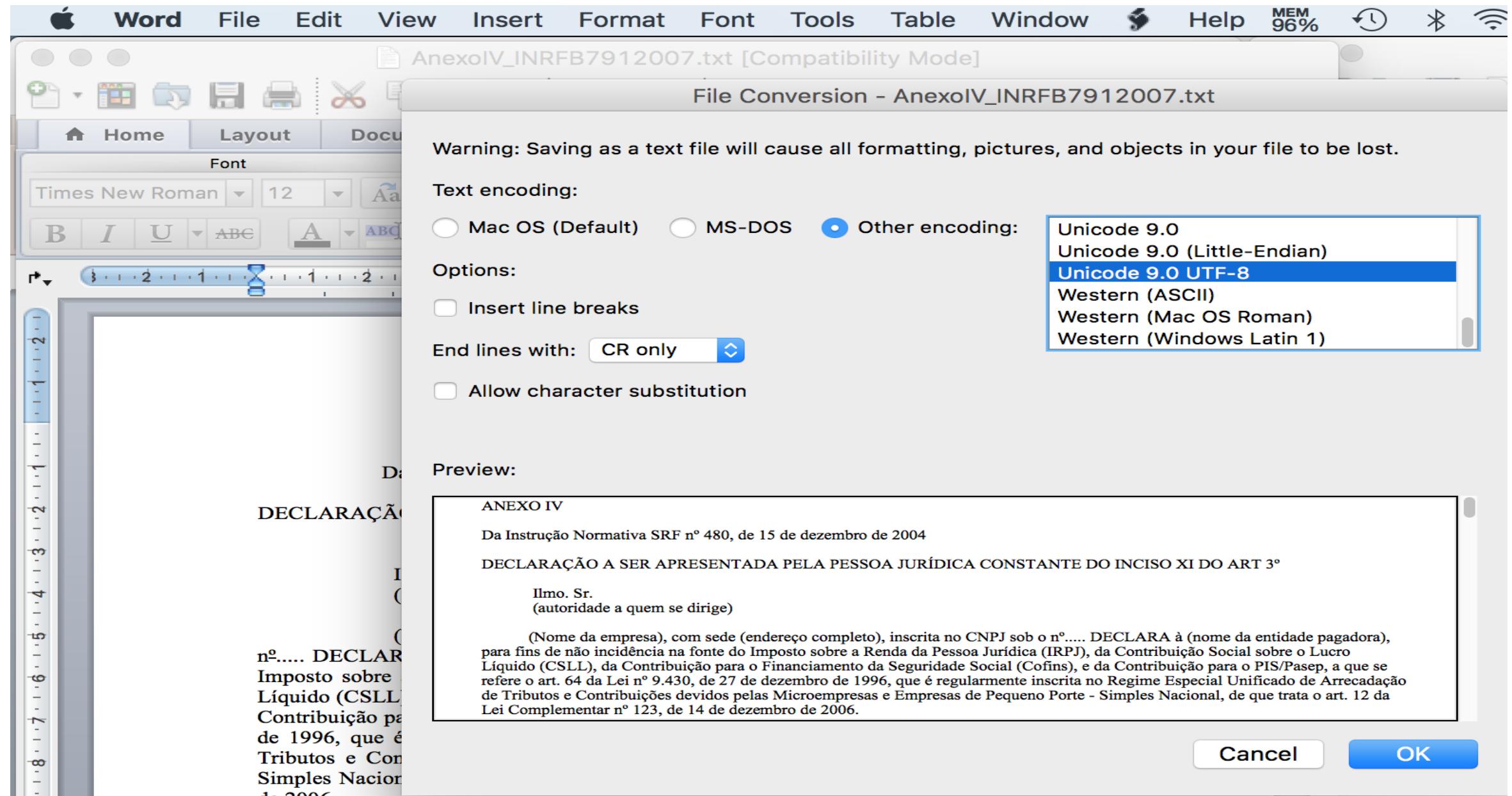
```
$ hdfs dfs -ls /
$ hdfs dfs -put /var/log/mysql.log /
$ hdfs dfs -ls /
$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-
mapreduce-examples.jar wordcount /mysql.log /saída
$ hdfs dfs -cat /saída/part*
$ hdfs dfs -get /saída /root/
$ ls -l /root/saída
$ cat /root/saída/part*
```

# Como processo um arquivo da minha máquina?

# Desktop, VM, HDFS

- Desktop → VM
  - Windows: Winscp (portable)
  - Linux: scp
- VM → HDFS
  - \$ hdfs dfs -[put|get] \*
- Desktop → HDFS
  - HUE, filebrowser
  - <http://<ip>:8000/filebrowser>





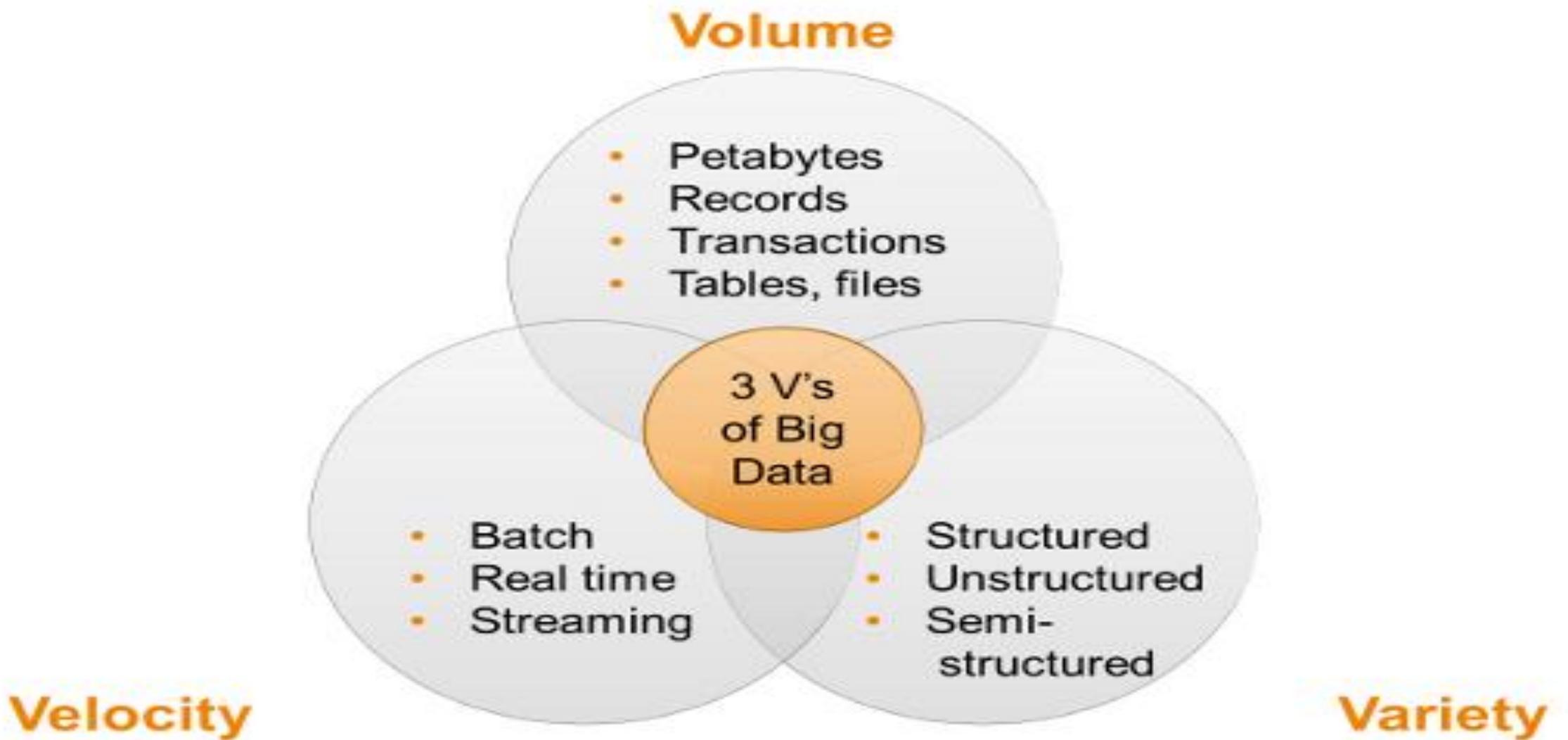
## Atividade 3

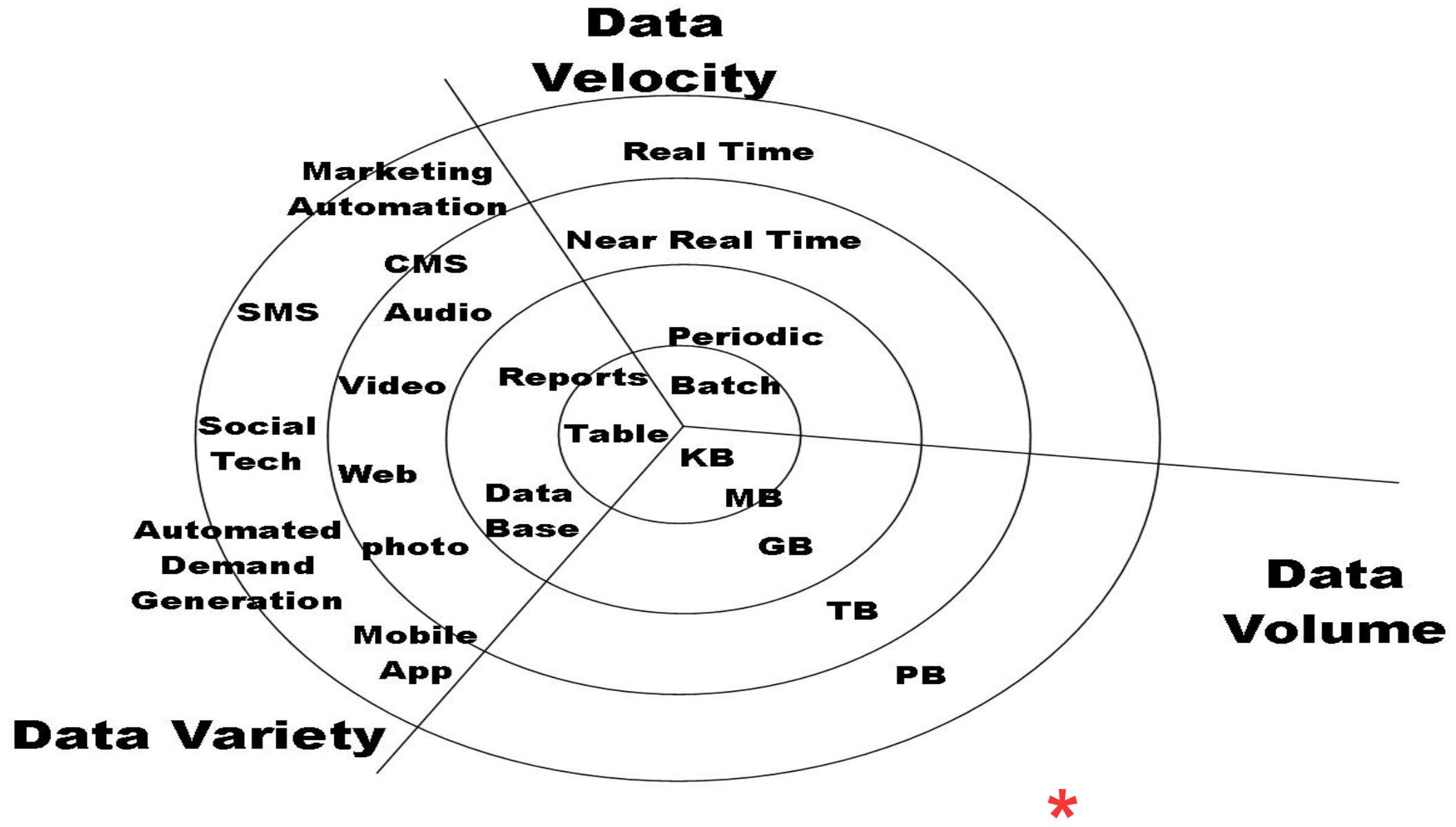
Executar o *wordcount* de um arquivo do seu computador.

# BIG DATA



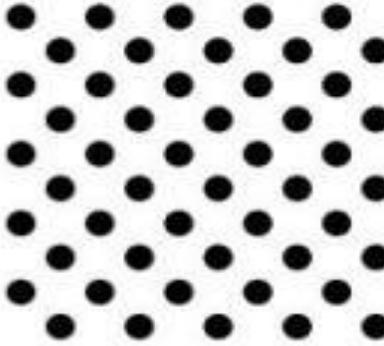
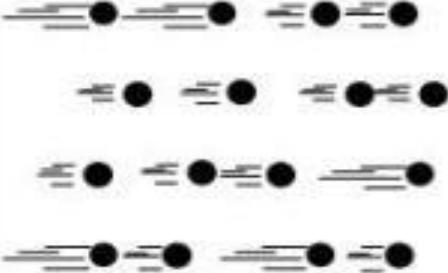
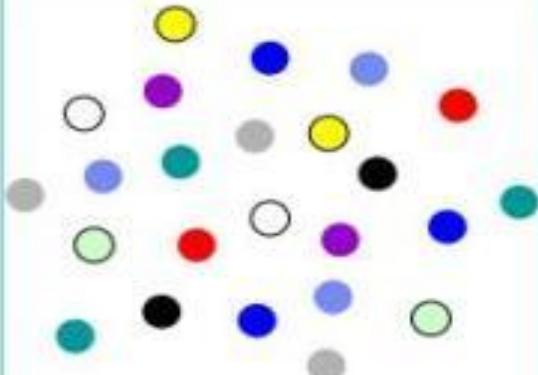
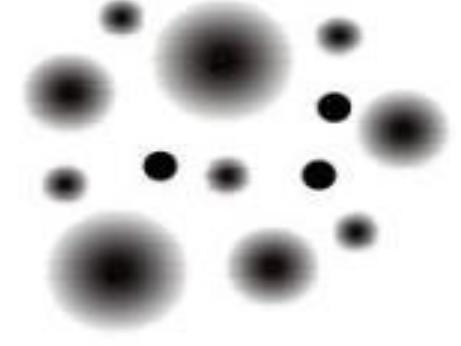
# Big Data '00



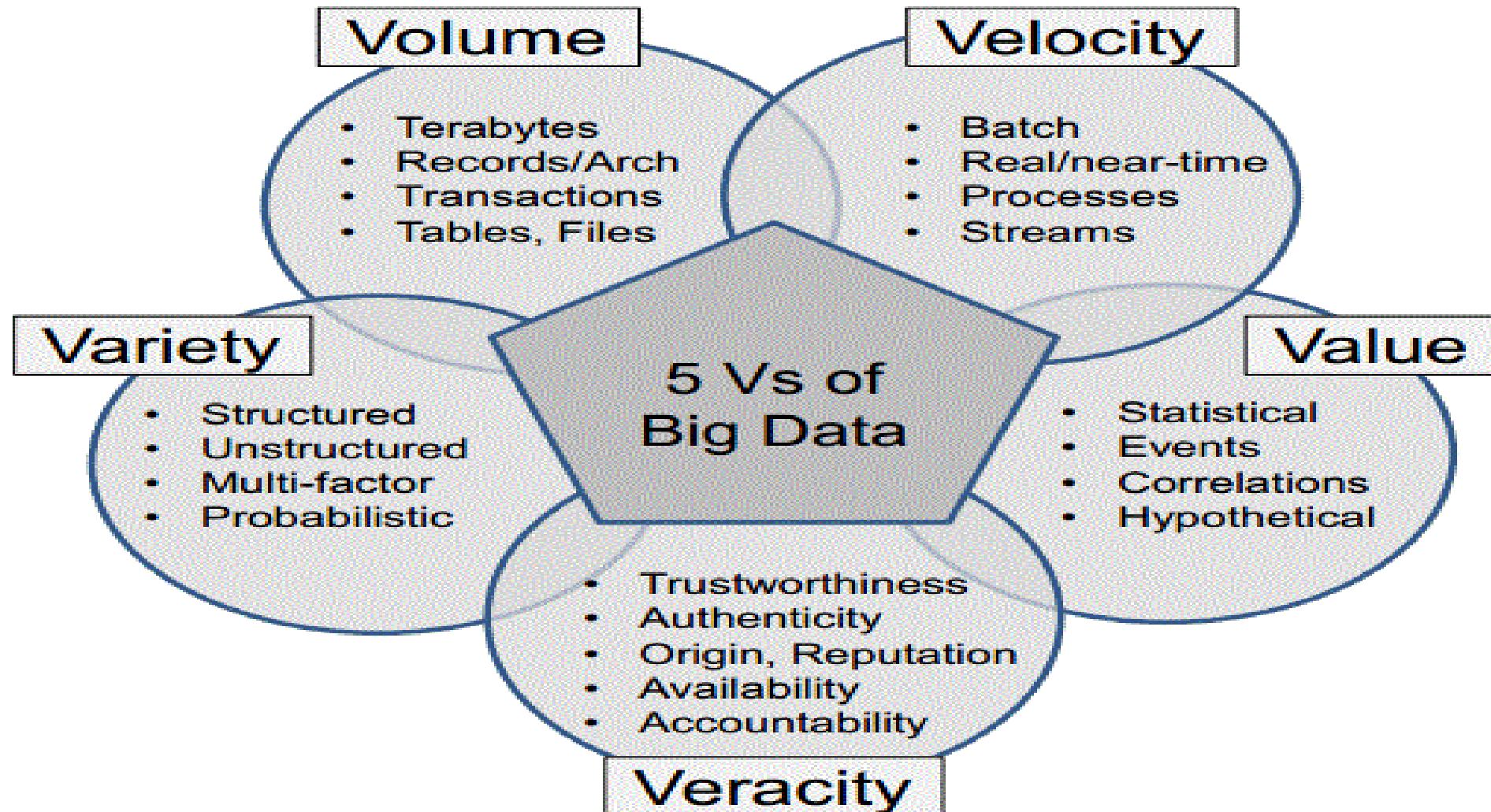


\*

# 4 Vs

Volume	Velocity	Variety	Veracity*
			
<b>Data at Rest</b>  Terabytes to exabytes of existing data to process	<b>Data in Motion</b>  Streaming data, milliseconds to seconds to respond	<b>Data in Many Forms</b>  Structured, unstructured, text, multimedia	<b>Data in Doubt</b>  Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations

# 5Vs



# 6 Vs ou Mais...

- Volume
- Variety
- Velocity
- Veracity
- Value
- [ Variability | Viability | Visibility | ... ]

# Buzzword

## 3Vs – Volume, Variedade e Velocidade



Processamento de grande volumes de dados  
não-estruturados para tomada de decisão em “tempo real”

# Caracterização de Big Data

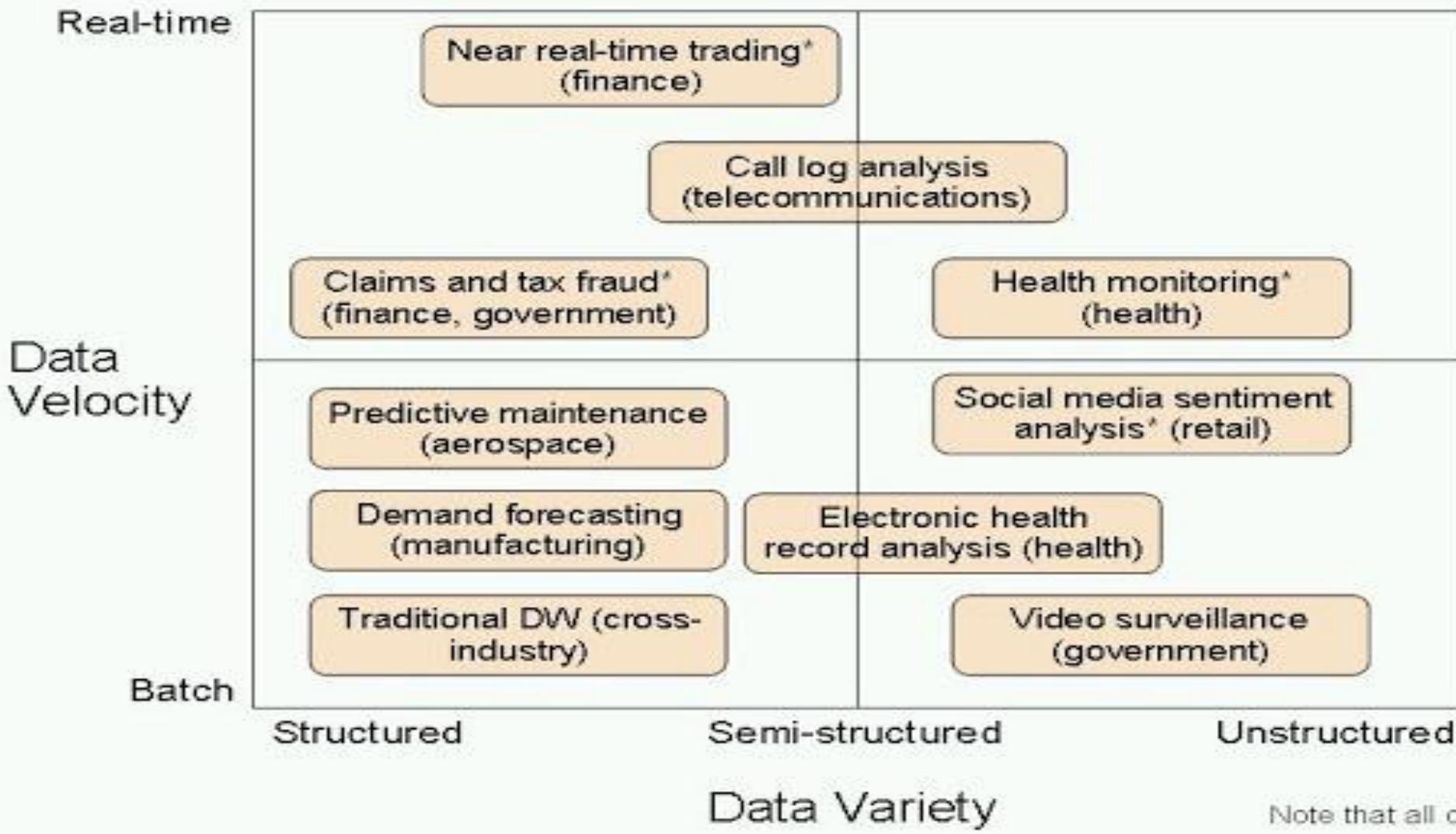
- Big Social Data
- Big Urban Data
- Big Medical Data
- Big ... Data

# Big Human Data

Atleta Extremo – Diabético  
(Iron Man triatleta e montanhista)  
<http://bcove.me/w0nks9qq>

- Milhares de sensores biométricos
  - Localização (GPS), altitude, temperatura
  - Temperatura corporal, níveis de insulina, caloria, sódio, batimento, ...
- Saúde
- Habilidades
- Performance

# Velocidade e Variedade



# Volume

- Twitter: 1 bi/semana – 6939/s!
- WallMart: 1 mi transações/s
- Facebook: 40 bi de fotos
- IDC: universo digital terá 40 ZB em 2020 (1ZB = 1Bi TB)
- IBM: 2,5 quintilhões dados/dia – 90% criado em 2 anos
- **NYC Taxi – 50GB/ano sem trajeto ou 2TB/ano com trajeto**

# Revisando Big Data

- Big Ideas: How big is Big Data (portuguese)
  - <https://youtu.be/xWx9THKXd6M>

# Então...

- Como armazenar e processar Big Data?
  - Volume
  - Variedade
  - Velocidade

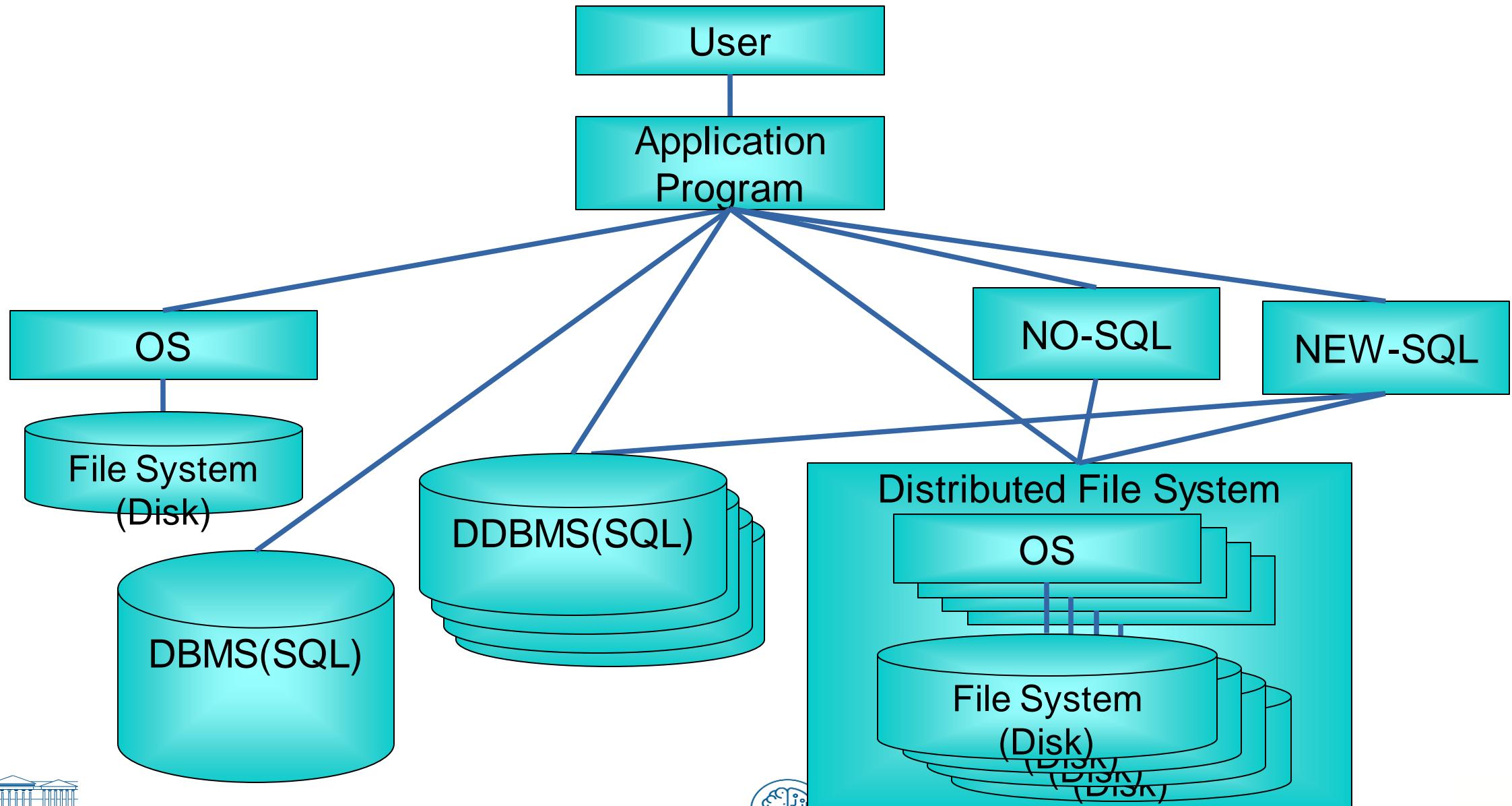
- Normalização → Desnormalização
  - Exclusivo → Redundante
  - Combinação (Join) → Acesso direto
  - “Lento” → “Rápido”
- Essencial → “Tudo”
  - Histórico (log)
  - Tendência

“Grava “tudo” pois poderá ser útil”

# Modelos de Dados para Big Data

- Relacional, Objeto e O-R
- Dimensional/Multidimensional (DW)
- Dados Geográficos (GIS)
- Colunar
- Chave-valor
- Documento
- Grafo (Hierárquico e Rede)
- Não estruturados
  - Excel/CSV, Áudio, Vídeo, Logs, ...

# Sistemas Big Data



# Componentes de Sistemas Big Data

## Data Analysis & Platforms



## Databases / Data warehousing



## Operational



## Multivalue database



## Business Intelligence



## Data Mining



## Social



## Big Data search



## Graphs



## KeyValue



## Document Store



## Object databases



## Multimodel



## XML Databases



Created by: www.bigdata-startups.com

DATA & ALLIANCESCAPE 2020

## INFRASTRUCTURE



ANALYTICS & MACHINE INTELLIGENCE



#### APPLICATIONS – ENTERPRISE



The image is a horizontal collage of logos for various database management systems, organized into six main categories: NoSQL Databases, NewSQL Databases, Graph DBs, MPP DBs, Serverless, and Cluster SVCS. Each category contains multiple logos of different database brands.



The image is a horizontal collage of logos for various fintech companies, arranged in five distinct sections. From left to right, the sections are: 'LEGAL' (with logos for Ravel, Qseed, DISCO, Cognacore, Onyo, and RSSS); 'REGTECH & COMPLIANCE' (with logos for Byd, Next ID, TRADEFISH, SCALEFACTORY, builddoer, builddoer, and exopps); 'FINANCE' (with logos for Anaplan, ZUUD, CAPTRANA, and TRADEFISH); 'AUTOMATION & RPA' (with logos for UPI Patch, Blue Prism, DIBER, Maitreya, VIVADO, DataPump, Neuronyx, ALLiVANT, eXtreme, and ZEALONIS); and 'SECURITY' (with logos for TANDEM, BlackBerry, eSigner, BlackCAT, FORTINET, Redbeam, DATAVERSUS, Proton, AegisLab, ANDROMALI, ReversingLabs, INVECTRA, ROSETTE, Lockwise, Qualys, Webroot, Fortinet, eSigner, RADAR, and McAfee). Each section has a thin black border around its respective group of logos.

A horizontal collage of logos for various data management companies, including Talend, pentaho, Alteryx, DataStax, Paxata, UNIFI, dotformz, Data Integration, MuleSoft, Trifacta, Stitch, Fivetran, Import.io, Matillion, Pendulum, Informatica, Atalant, Collibra, dremio, Fishtech, Okera, Dataworld, and PrivaCERA.



ADVERTISING	EDUCATION	REAL ESTATE	GOVT & INTELLIGENCE	COMMERCE	FINANCE - LENDING	INSURANCE
AppNexus  MediaMath	Unacademy	Redfin	VTS	Palantir	STITCH FIX	Root
criteo  IAS	Babbel	Broker	OpenGov	DODGINN	NowGood	Monrovia
ORACLE MARK	BENEFITON	Orchard	MARK43	ANDURIL	STANDARD	ZEST
albert  gumgum	declaro	iRealogy	FISCAL NOTE	FiscalNote	Upgrade	SAC
Opfer  theTradeDesk	KORTBI	SPACE MAKER	Geoprym	AYASDI	KENSIC	CAPE
TAPAB			Quid  PRIMER	AEGOSPAN	NUMERO	EvolutionIQ
				100Credit	AQUARIUS	EQUITY
					ZEST	jewell



The banner displays a grid of company logos. The first row contains logos for flatiron, AYRUS, METABOTA, banyon, 3DMed, ReBal, TEMPUS, adjuvant, Inc., AlQura, LUMINA, Z, PAUSE, Olive, Imp, Health, Resonacal, LUMINA, eSpringHealth, imitro, Caption Health, entic, Zebra, and zebra. The second row contains logos for color, Verily, Genentech, DNAfit, verily, GUESS, Neuro, Aptiv, Zogenix, instacart, Remedy, Amgen, Amantia, Neo, Optum, G7, Ford, Veeva, Xodisk, Imprivé, Cognex, TESSX, VisusHealth, Bluebeam, and semios. The third row contains logos for FLONERS, Genmural, AutoX, Nauto, Blue River, Kudu, TARANIS, Prospera, ALICE, BHARTI, and Lightrock.

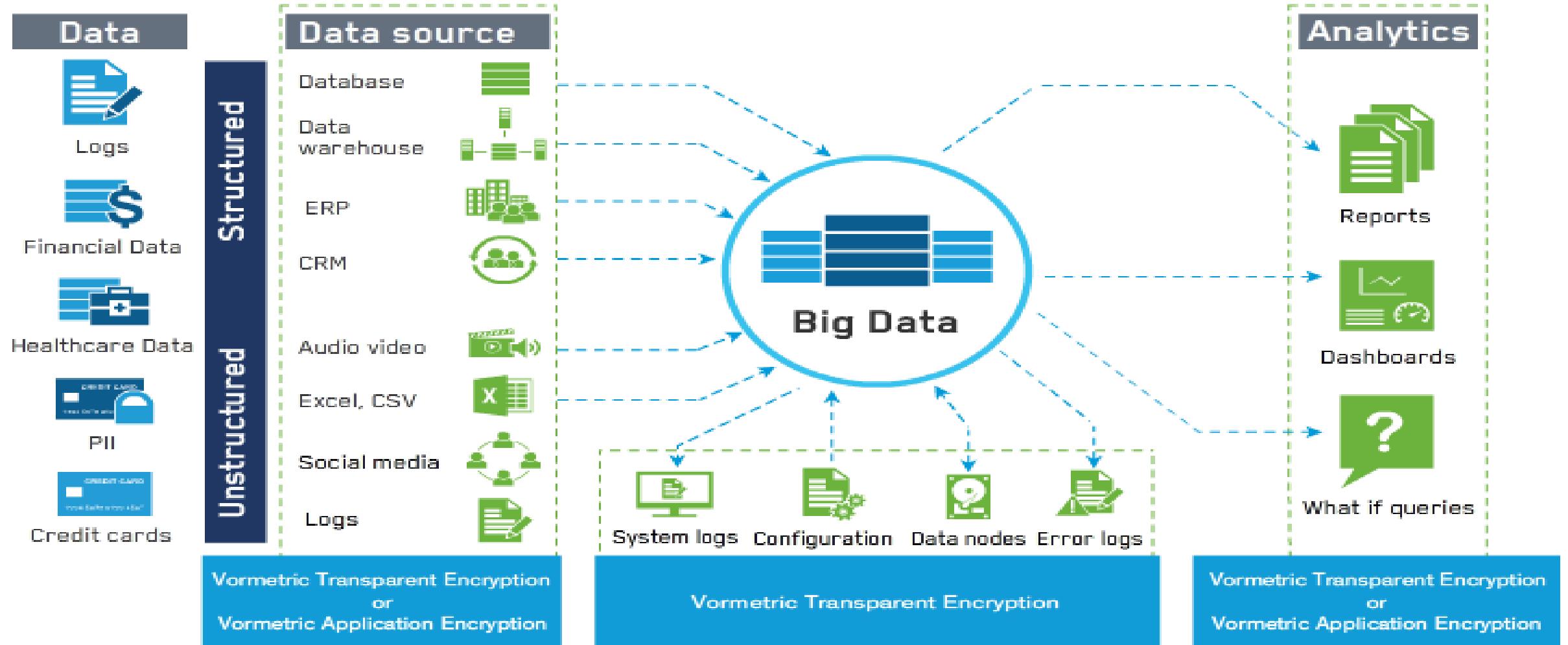


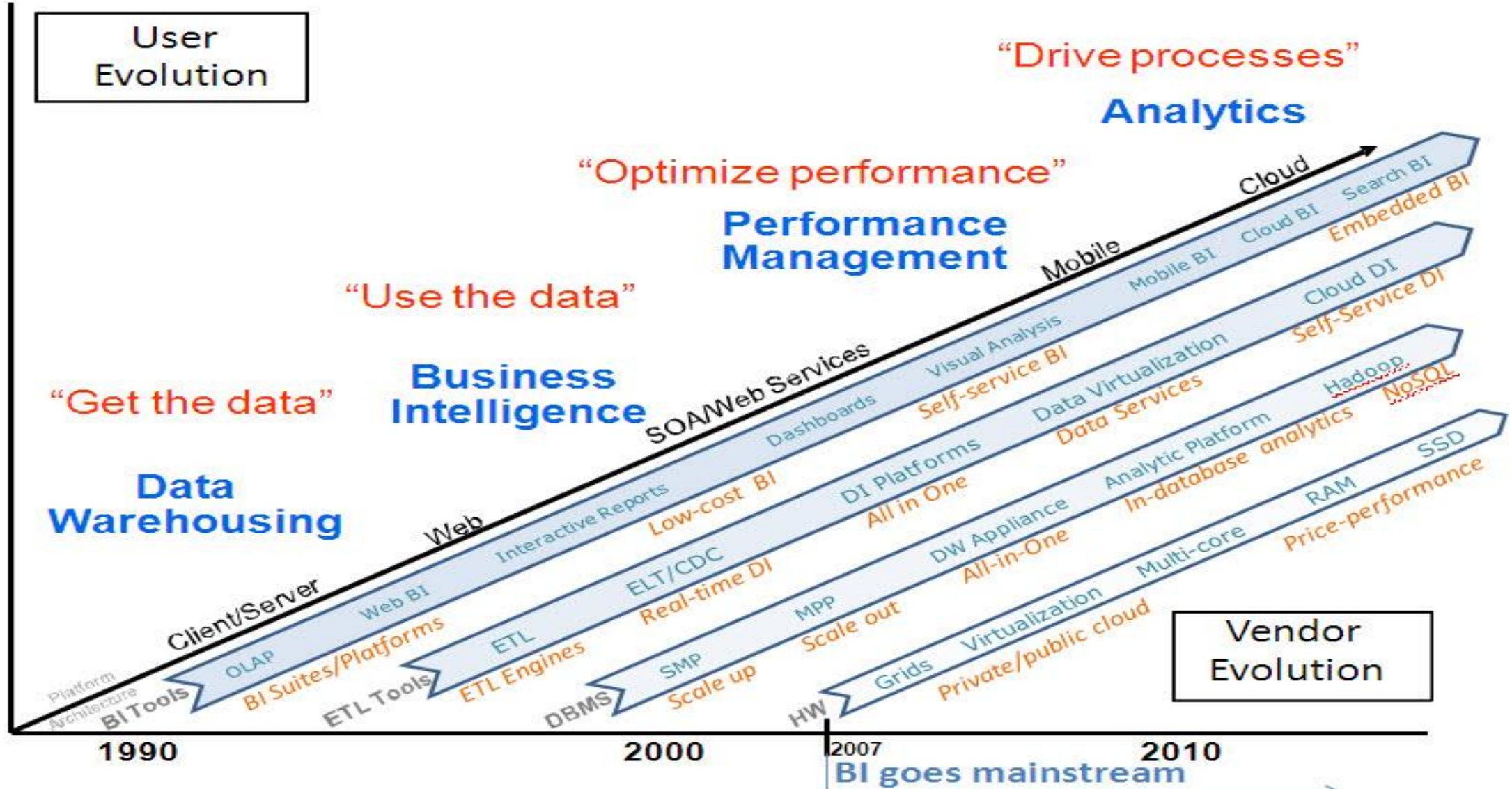
The banner displays a grid of logos for fintech and data companies. The first column under 'DATA MARKETPLACES & DISCOVERY' includes AWS Data Exchange, DAWEX, and data.world. The second column under 'FINANCIAL & ECONOMIC DATA' includes Bloomberg, Thomson Reuters, Dow Jones, S&P Capital IQ, ICB Insights, Plaid, Qualtracs, and several smaller logos for fintech companies like Shockwave, xignite, Bluebeam, and earnest. The third column under 'AIR / SPACE' includes Orbital Insight, DataRobot, and Zipline.



The banner features several sections: 'OTHER' with links to Data.gov and Booz Allen Hamilton; 'DATA SERVICES' with QuantumBlack, Kaggle, ElectraAI, Fractale, and DataCamp; 'INCUBATORS & SCHOOLS' with PredictiveHive, General Assembly, DataFlair, DataCamp, Galvanize, and BetaIS; and 'RESEARCH' with OpenAI, Facebook Research, NASA, MIRI, Vector Institute, and Alluvial Institute.

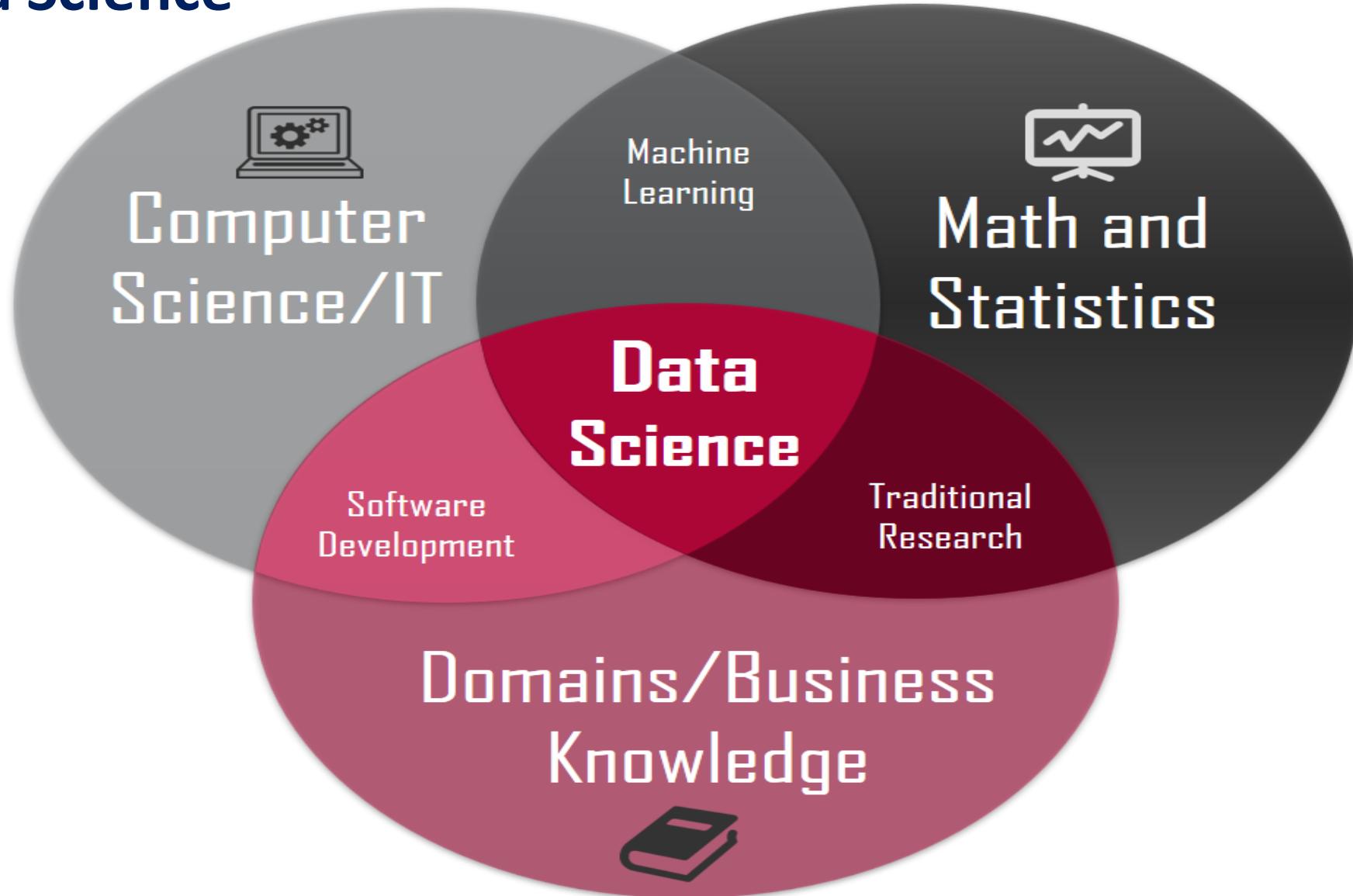
# Big Data e Data Lake







# Data Science



Especialização em  
Inteligência Artificial

# MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants

# DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

# PROGRAMMING & DATABASE

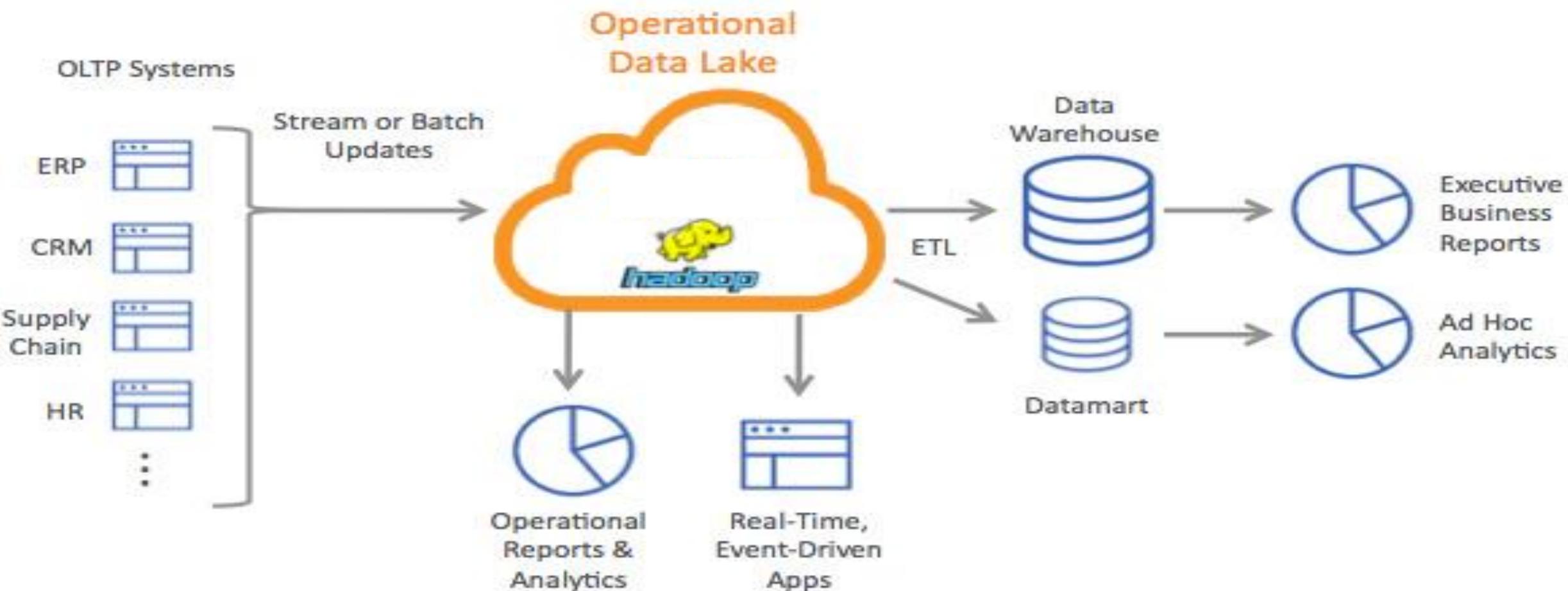
- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS



# COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

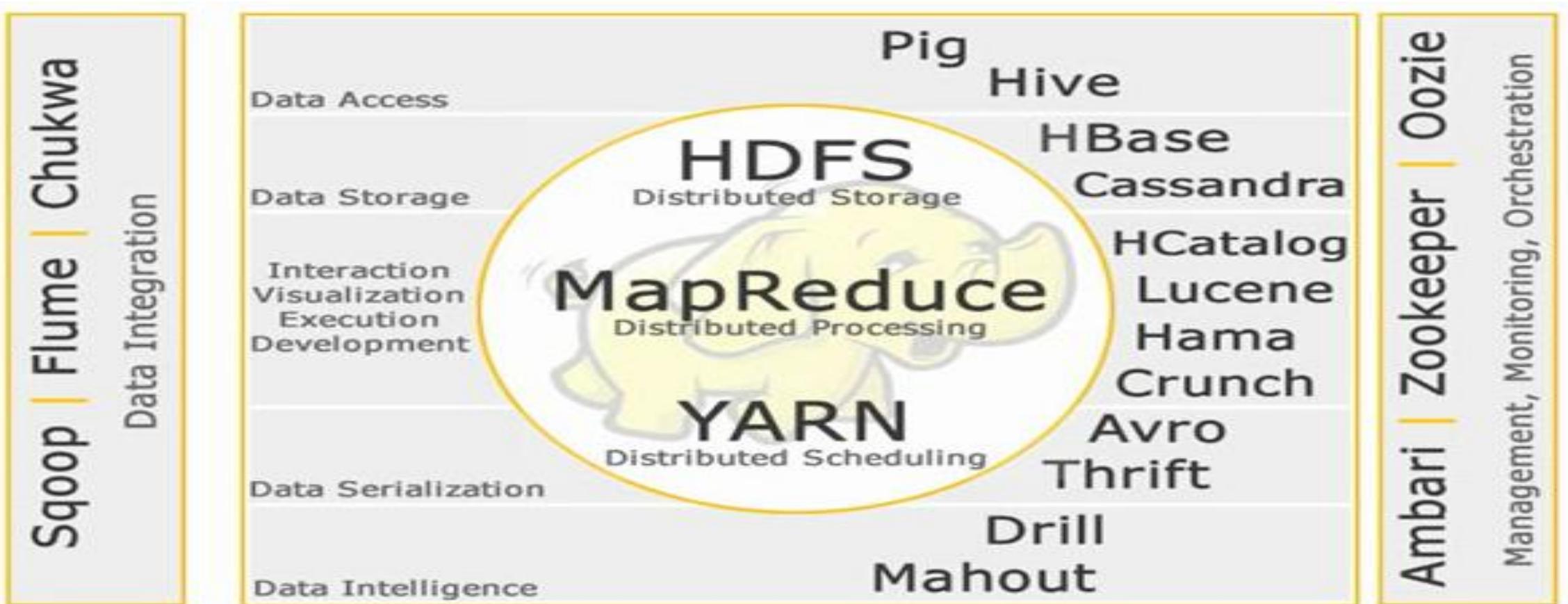
# Data Lake



# Hadoop e o Ecosistema Hadoop



<http://hadoop.apache.org>



# Referências

- ...
- Dean, J., & Ghemawat, S. (2004). **MapReduce: Simplified Data Processing on Large Clusters**. In Proc. of the OSDI - Symp. on Operating Systems Design and Implementation (pp. 137–149). USENIX.

# MapReduce (MR) [by Google'04]

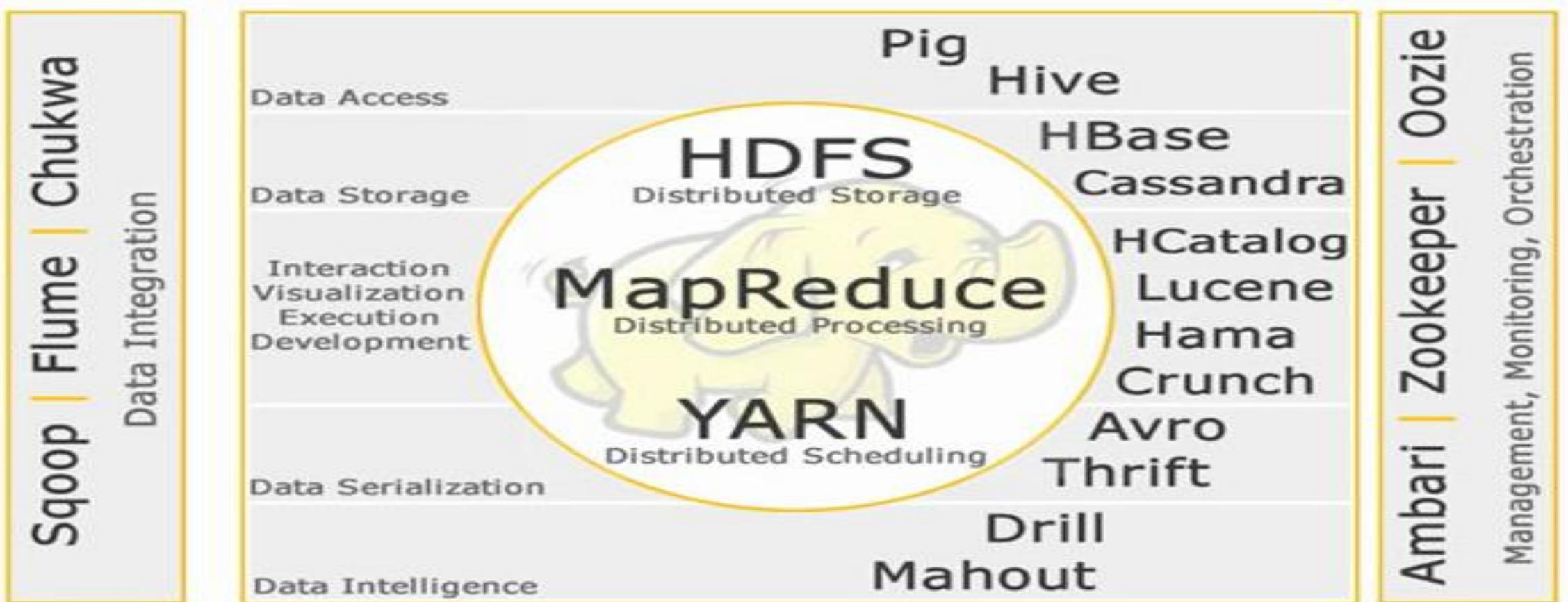
- Modelo de programação e uma implementação/framework
  - Programação personalizada
  - Distribuição e Paralelização automática
  - Milhares de máquinas “comuns”
  - Balanceamento de carga
  - Otimização de rede e transferência de arquivos
  - Tolerância à falhas
  - Ambiente de desenvolvimento comum
    - Melhorias beneficiam todos usuários

# Hadoop e o Ecosistema Hadoop



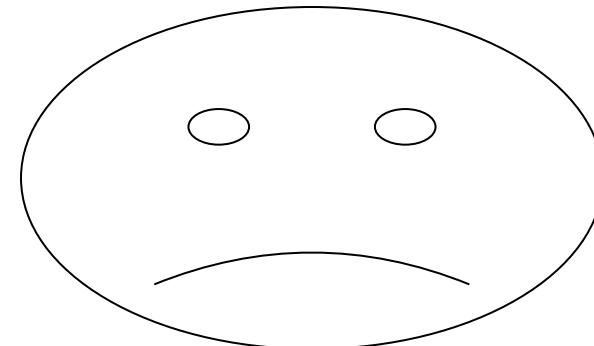
[Yahoo et al., 2006]

<http://hadoop.apache.org>



# Mão na Massa...

- Instalar Hadoop (na mão)
  - <http://hadoop.apache.org/>
  - Local, pseudo-distributed, distributed
- E os outros sistemas?
  - Pig, Hive, Cassandra, Hbase,...
- Atualizações?



# Distribuições Hadoop

- Problemas de incompatibilidade entre versões dos sistemas Hadoop (eco-sistema)
  - Pacotes Linux, VMs, tarballs
  - Scripts adicionais para manipular sistemas Hadoop
  - Suporte/Backport de características e atualizações feitas pela Apache
  - Distribuidores empregarão sistemas Hadoop e ajudarão a atualizar/manter o Apache Hadoop

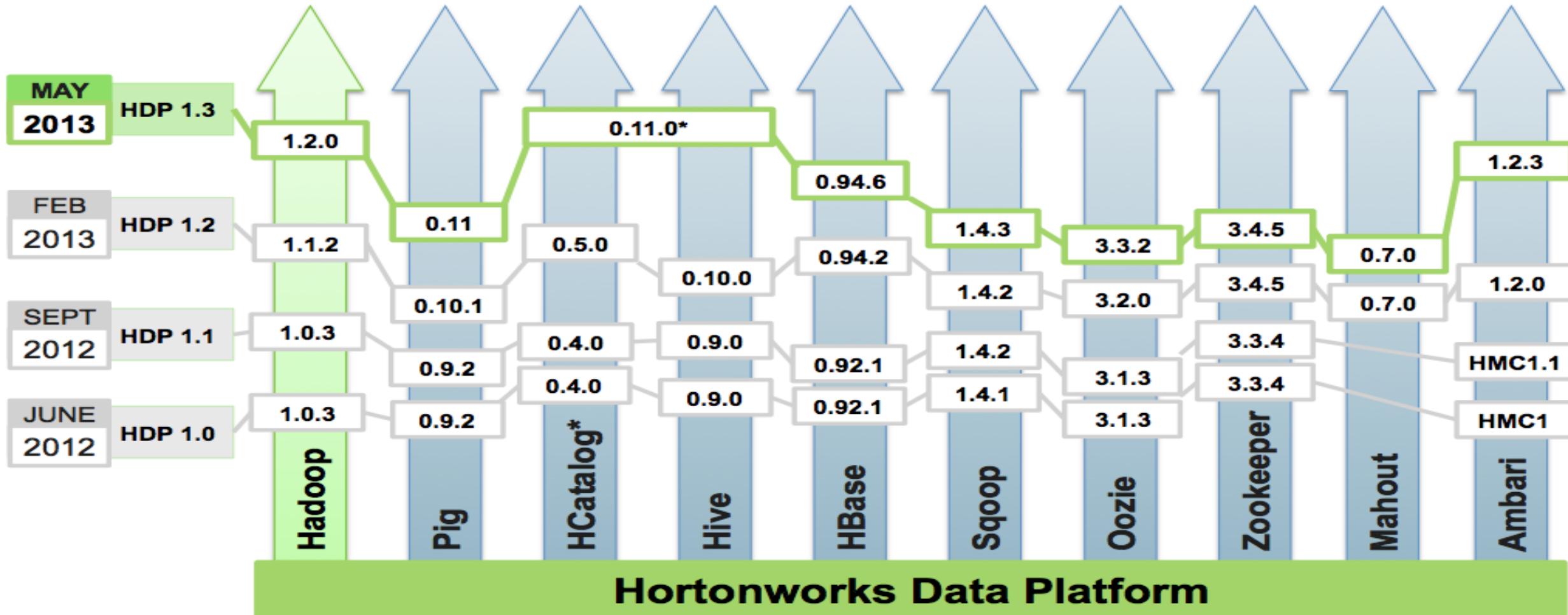
# Distribuições Hadoop



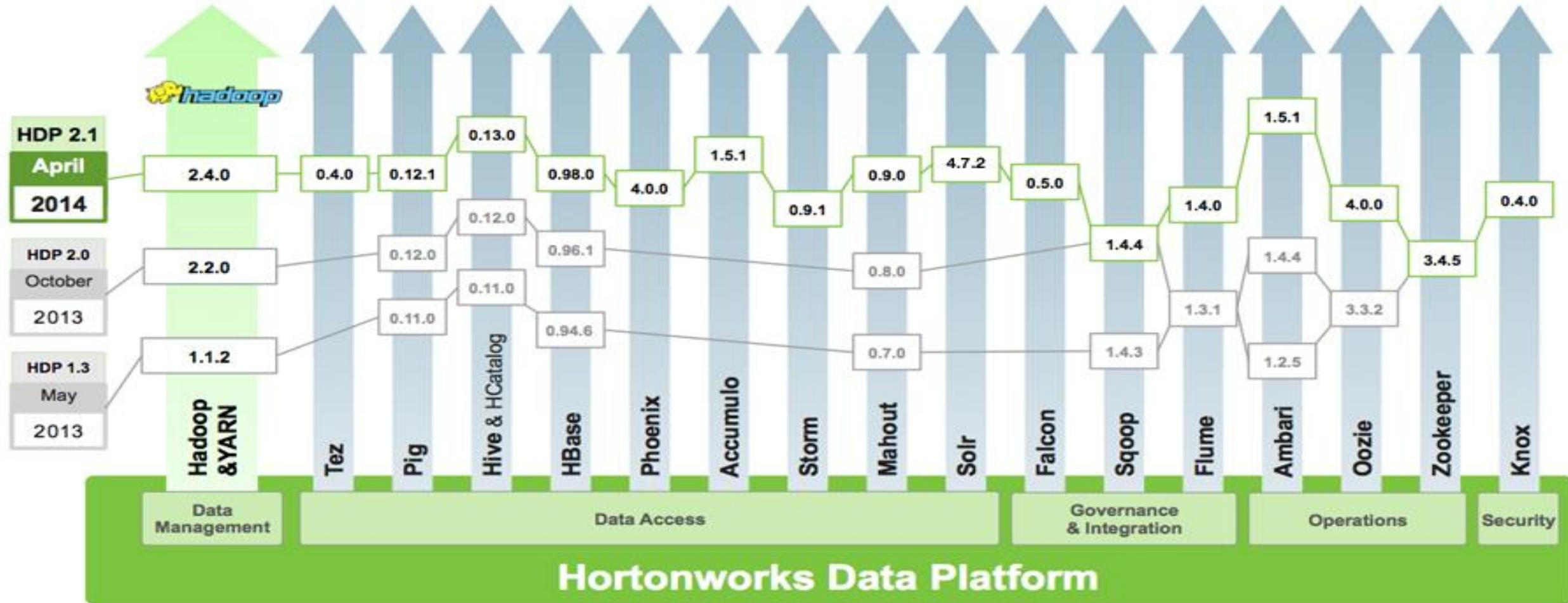
- *Hortonworks Data Platform (HDP)*  
<http://hortonworks.com>



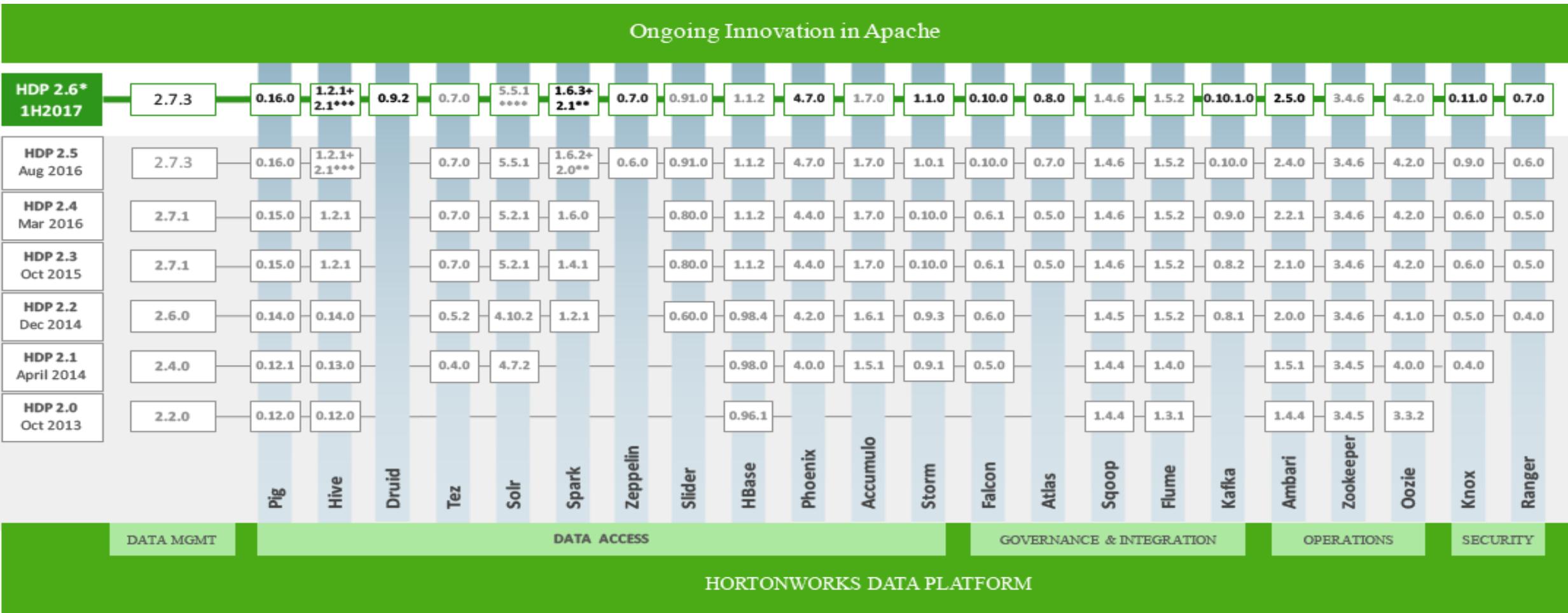
# HDP 1.3



# HDP 2.1



# HDP 2.6



\* HDP 2.6 – Shows current Apache branches being used. Final component version subject to change based on Apache release process.

\*\* Spark 1.6.3+ & Spark 2.1 – HDP 2.6 supports both Spark 1.6.3 and Spark 2.1 as GA.

\*\*\* Hive 2.1 is GA within HDP 2.6.

\*\*\*\* Apache Solr is available as an add-on product HDP Search.

DATA & ALLIANCESCAPE 2020

## INFRASTRUCTURE



ANALYTICS & MACHINE INTELLIGENCE –



#### APPLICATIONS – ENTERPRISE



A horizontal collage of logos from various fintech companies, categorized into five main sectors: Legal, RegTech & Compliance, Finance, Automation & RPA, and Security. The logos are arranged in a grid-like fashion, with some companies having multiple entries.

**ETL / DATA TRANSFORMATION**

- talend
- alteryx
- Paxata
- UNIFI
- dotform

**DATA INTEGRATION**

- pentaho
- MuleSoft
- Schibsted
- Zalando
- import.io
- innowise
- REEDSELL

**DATA GOVERNANCE**

- Informatica
- Alation
- collibra
- dremio
- INNOSTAT
- dataworld

**DATA QUALITY**

- SciPoint
- Alphawise
- MarkLogic
- FACTA
- OKERA
- dataopera
- MC
- DATA CARLO



<b>ADVERTISING</b>	<b>EDUCATION</b>	<b>REAL ESTATE</b>	<b>GOVT &amp; INTELLIGENCE</b>	<b>COMMERCE</b>	<b>FINANCE - LENDING</b>	<b>INSURANCE</b>
AppNexus  MediaMath	Unilink	REDFIN	Palantir	FAIRERE	STITCHFLX	ROOT
criteo  IAS	VTS	operator	OPENBIV	NowGood	affirm	Proclaimlife
ORACLE MARK	newton	Orchard	DODGINN	STANDARD	Monedo	ZEST
albert  gumgum	Declarra	iFinTech	MARK43	FINANCING	SIBA	Shift Technology
Opfer  theTradeDeck	KORET	GEOPRIFY	ANDURIL	AYASDI	BLOOMBERG BANC	CAPE
TAPAB		SPACE MAKER	FiscalNote	KENSIC	GURU	EvolutionIQ
			Quid  PRIMER	ADOSPAN	NUMERO	LoyaltyOne
				100Credit	Aptiv	EVOLV
					zestly	

A collage of logos for various AI and cloud computing companies, including Mgmt / Monitoring, Data Generation & Labelling, AI Ops, GPU DBs & Cloud, and AI Hardware.

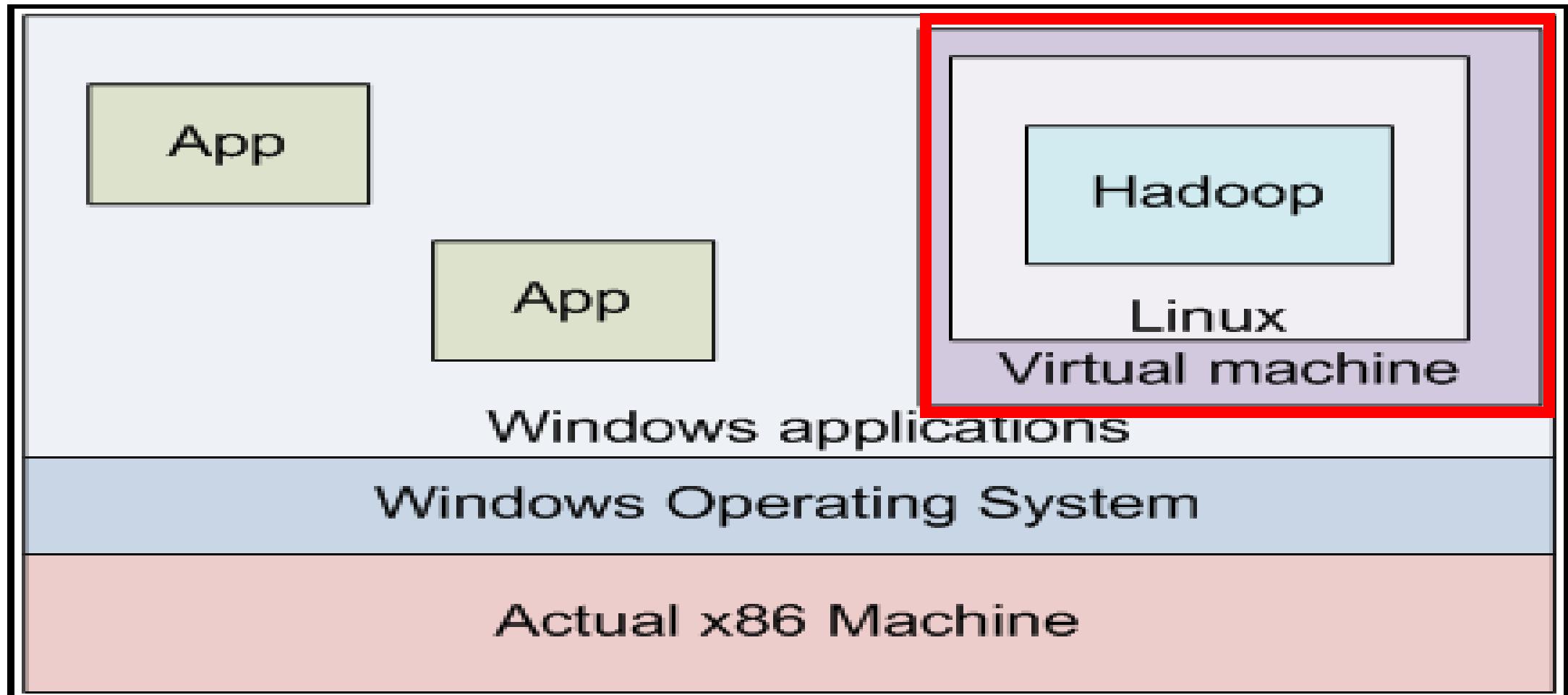


The banner is a horizontal strip containing logos for different companies. From left to right, it includes: DATA MARKETPLACES & DISCOVERY (aws Data Exchange, DAWEX, data.world, narrative); FINANCIAL & ECONOMIC DATA (Bloomberg, THOMSON REUTERS, DOW JONES, Quandl, SAP CAPITAL, ICB INSIGHTS, PLAID, QUALITY INSURANCE, AIRPORT FACILITY, ZENITH); and AIR / SPACE (Orbital Intel, TUMBLEWEED, Blue Origin, Blue Origin).



# Hortonworks Sandbox 2.1

- Usar uma Sandbox: Virtual Machine (VM) com tudo!



## Atividade 4

Dado o conjunto de programas exemplos disponível em:

```
$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar
```

Comente brevemente justificando e apresentando a aplicação de algum outro programa exemplo.

# Como faço a minha aplicação MR?

# MapReduce (MR) [by Google'04]

- **Modelo de programação (?) e uma implementação/framework**
  - Programação personalizada
  - Distribuição e Paralelização automática
  - Milhares de máquinas “comuns”
  - Balanceamento de carga
  - Otimização de rede e transferência de arquivos
  - Tolerância à falhas
  - Ambiente de desenvolvimento comum
    - Melhorias beneficiam todos usuários

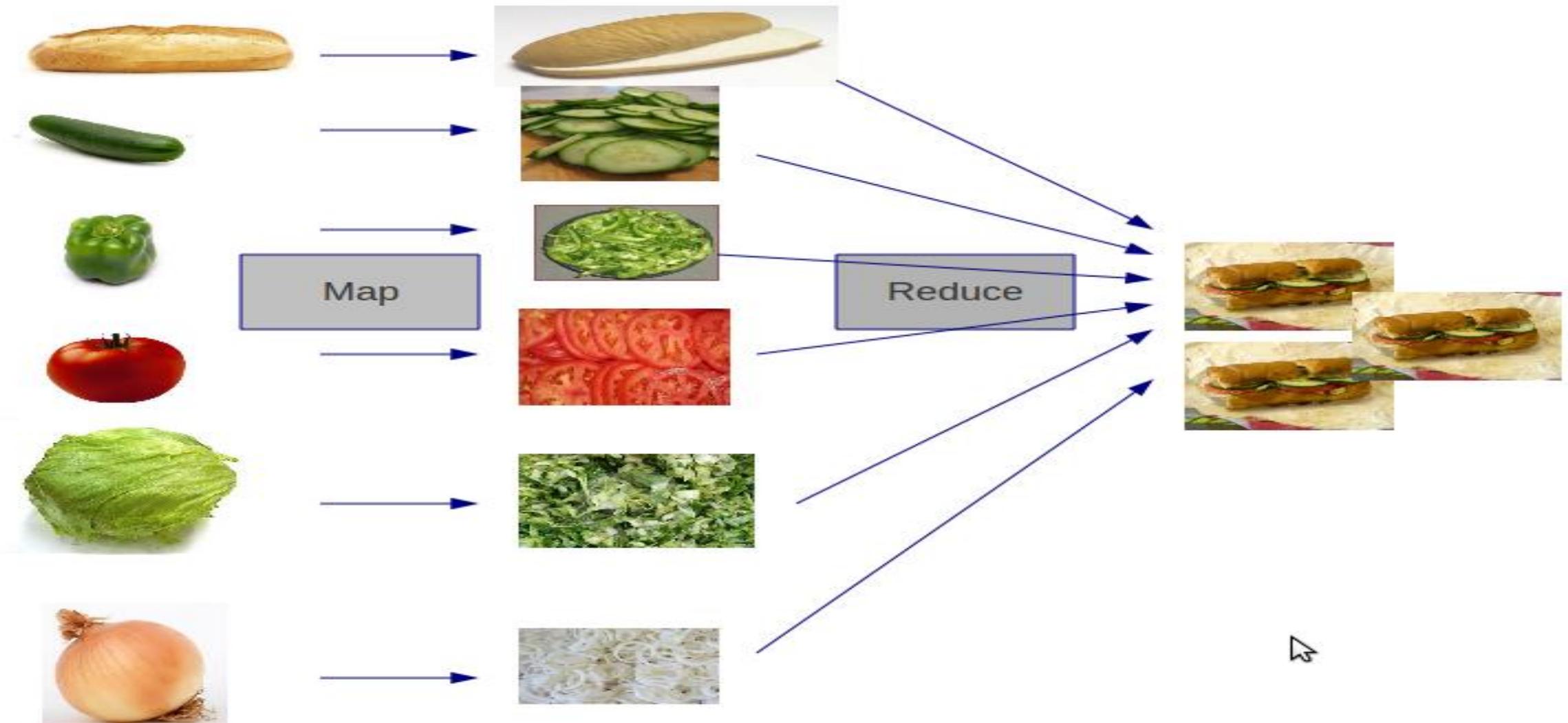
# Problema MR Típico

- Leia a entrada
- **Map**
  - identifique os registros (o dado relevante)
- Combine e Agrupe
- **Reduce**
  - agregue, some, resuma, filtre, transforme,...
- Escreva o resultado
- O esquema permanece o mesmo, altera-se o Map e o Reduce para atender cada problema.

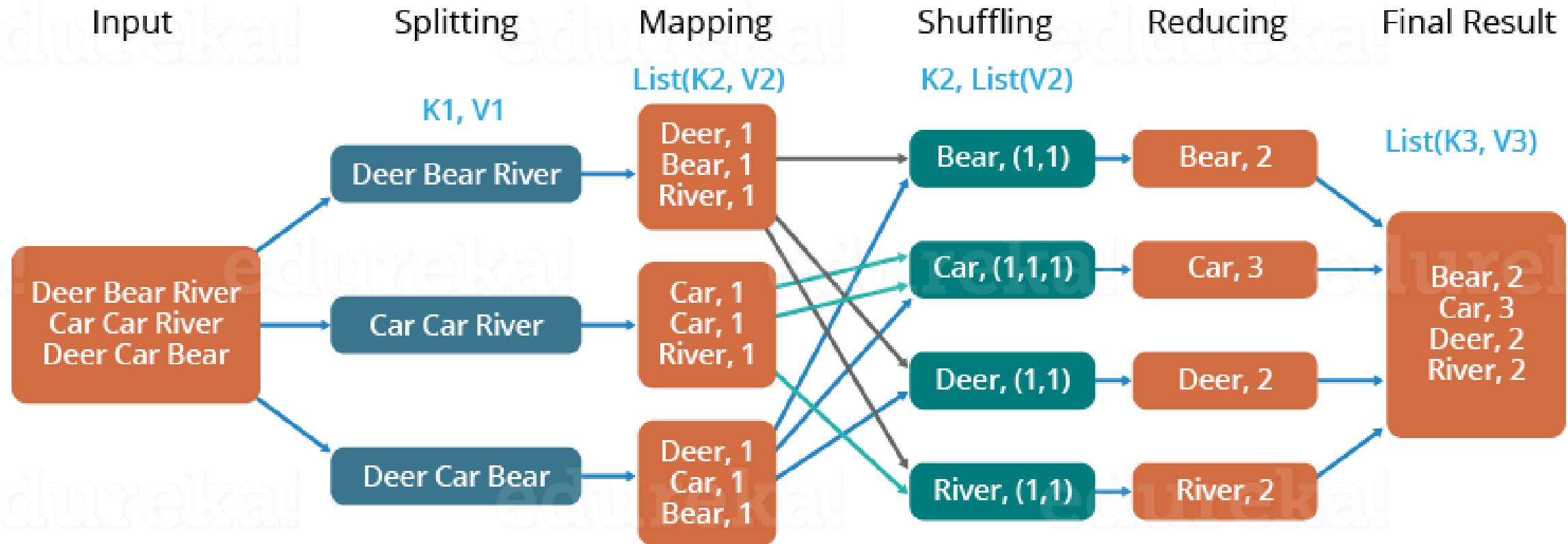
# Modelo de Programação

- Computação baseada em <chave, valor>
- **Map** (<chave, valor>)
  - Saída: um conjunto de pares <chave2, valor2>
- **Reduce** (<chave2, lista de valor2>)
  - Saída: um par <chave2, valor3>

# Aplicando Map e Reduce...



# Aplicando Map e Reduce – Contador de Palavras



# Algoritmo Map Reduce para Contar Palavras

```
map(String key, String value)
    // key: nome arquivo
    // value: conteúdo do arquivo
    for each word w in value:
        EmitIntermediate(w, "1");

reduce(String key, Iterator values)
    // key: uma palavra
    // values: uma lista de contadores
    result=0;
    for each v in values:
        result+= ParseInt(v);
    Emit(key, AsString(result));
```

# Função Map

```
public void map(LongWritable key, Text value,  
OutputCollector<Text, IntWritable> output, Reporter  
reporter) throws IOException {  
    String line = value.toString();  
    StringTokenizer tokenizer = new StringTokenizer(line);  
    while (tokenizer.hasMoreTokens()) {  
        word.set(tokenizer.nextToken());  
        output.collect(word, one);  
    }  
}
```

# Função Reduce

```
public void reduce(Text key, Iterator<IntWritable> values,  
OutputCollector<Text, IntWritable> output, Reporter  
reporter) throws IOException {  
    int sum = 0;  
    while (values.hasNext()) {  
        sum = sum + values.next().get();  
    }  
    output.collect(key, new IntWritable(sum));  
}
```

# Temperatura Máxima

Arquivo com cidades e temperaturas

Toronto, 20

New York, 22

New Jersey, 26

Toronto, 22

...

- Como é a função Map e Reduce para obter as temperaturas máximas das cidades?
  - Map (?, ?)
  - Reduce (?, ?)

# Temperatura Máxima

- **Map** (arquivo, conteúdo)

**Para cada cidade, emita <cidade, temp>**

(Toronto, 20) (Whitby, 25) (New York, 22) (Rome, 33)(Toronto, 18) (Whitby, 27) (New York, 32) (Rome, 37)(Toronto, 32) (Whitby, 20) (New York, 33) (Rome, 38)(Toronto, 22) (Whitby, 19) (New York, 20) (Rome, 31)(Toronto, 3 (Whitby, 22) (New York, 19) (Rome, 30)

- **Reduce** (cidade, valores)

**Emita o maior valor dentre os valores das temp.**

(Toronto, 32) (Whitby, 27) (New York, 33) (Rome, 38)

# Dados Geográficos

## Sistema de mapeamento de estradas

- **Objetivo**
  - Identificar as conexões das estradas
- **Map**
  - Cria pares de pontos conectados
  - (estrada, intercessão) ou (estrada, estrada)
- **Reduce**
  - Agrupa os pontos de conexão para cada estrada

# Grafo de Emails

- Contar quantos e-mails foram trocados entre cada par
  - Joao, Pedro: 30
  - Pedro, Ana: 14
  - Ana, Joao: 7
  - Daniel, Pedro: 43
- Como seria uma aplicação MapReduce?

# Grafo de Emails - INPUT

```
"_id" : ObjectId("4f2ad4c4d1e2d3f15a000000"),
"body" : "Here is our forecast\n\n",
"filename" : "1.",
"headers" : {
    "From" : "phillip.allen@enron.com",
    "Subject" : "Forecast Info",
    "X-bcc" : "",
    "To" : "tim.belden@enron.com",
    "X-Origin" : "Allen-P",
    "X-From" : "Phillip K Allen",
    "Date" : "Mon, 14 May 2001 16:39:00 -0700 (PDT)",
    "X-To" : "Tim Belden",
    "Message-ID" : "<18782981.1075855378110.JavaMail.evans@thyme>",
    "Content-Type" : "text/plain; charset=us-ascii",
    "Mime-Version" : "1.0"
```

- Map(?,?)
- Reduce(?,?)

# MR para Grafo de Emails

- Map(arquivo, conteúdo)
  - Para cada destinatário
    - Grave("remetente, destinatário", 1)
- Reduce(par, valores)
  - Para cada valor em valores
    - Some cada valor
  - Grave o resultado (par, soma)

# O que faz essa Aplicação MR?

```
Map (file, social.person records V1)
  for each social.person record in V1 do
    let Y be the person's age
    let N be the number of contacts the person has
    produce one output record (Y, (N,1))
```

```
Reduce ( age (in years) K2, list of (N,C) )
  for each input record (Y, (N,C)) do
    Accumulate in S the sum of N*C
    Accumulate in Cnew the sum of C
    let A be S/Cnew
    produce one output record (Y, (A,Cnew))
```

# Qual o Problema?

Average of the top 100 salaries for each unique department (key), if that average is larger than 100,000

```
public void reduce(String key, Iterator<Integer>  
salaries) {  
    int sum = 0; int i = 0;  
    while (salaries.hasNext() && i<100) {  
        sum += salaries.next();  
        i += 1; }  
    emit((i>0 && sum/i > 100000) ? sum/i : -1);  
}
```

[Csallner, C., Fegaras, L., & Li, C. (2011). [New Ideas Track: Testing MapReduce-Style Programs](#). In Proc. of the ESEC/FSE - SIGSOFT Symp. and the European Conf. on Foundations of Software Engineering (pp. 504–507). ACM.]

# Como implementar e implantar a sua aplicação MR? (hadoop-mapreduce-examples.jar)

- \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar **pi** 10 10
- \$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar **wordcount** /mysqld.log /saída

# Implementação e Implantação de uma Aplicação MR

- Gerar o jar com a aplicação MR no desktop \*
- Enviar o jar para a sandbox ou cluster Hadoop
  - Winscp, scp ou [filebrowser]! :)
- Executar na sandbox ou cluster

```
$hadoop jar seu.jar /mysql.log /saída
```

\* Seu ambiente de desenvolvimento (Sandbox) não tem JDK e fonte do Hadoop

# Gerar Jar

- Versão hadoop na sandbox? hadoop version
- <https://archive.apache.org/dist/hadoop/common/>
- hadoop[version].tar.gz (src não tem os jar)
- Netbeans, new project, java application
- Java: <https://tinyurl.com/AulaWC-java> (baseado em <http://tinyurl.com/htjlkic>)
  - Adequar o pacote (package)
  - Renomear arquivo para ficar igual a classe
- Adicionar bibliotecas (Project, Properties, Library, Add JAR)
  - share/hadoop/common/hadoop-common
  - share/hadoop/mapreduce/hadoop-mapreduce-client-core
- Project properties, sources, source/binary format
  - Versão Jdk na sandbox: java –version
- Netbeans, run, build project

# Executar a sua Aplicação MR

- Enviar o jar para a sandbox ou cluster Hadoop
  - Winscp, scp ou [filebrowser]! :)
- Executar na sandbox

```
$hadoop jar seu.jar /mysql.log /saída
```

# Problemas na Sandbox

- Processos param/congelam sem resposta
  - Espaço em disco
  - Apagar Logs
    - du
    - rm -rf ....
- Travar processo
  - Reiniciar a VM ou serviço via Ambari
- Ajuste de horário
  - yum install ntp; chkconfig ntpd on; service ntpd restart
  - rm /etc/localtime; ln -s /usr/share/zoneinfo/America/Sao\_Paulo /etc/localtime
  - date mmddHHMM; hwclock -w

# Próximas Aulas

- Fundamentos de Big Data
  - Big Data, Data Lake e Data Science
- Map Reduce e Hadoop
  - Utilização da Sandbox/VM
  - Personalização de aplicações Map Reduce
- Data Engineering (Gerenciamento/Ferramentas para Big Data)
  - NoSQL e NewSQL
  - Dados em movimento – Processamento de Streaming