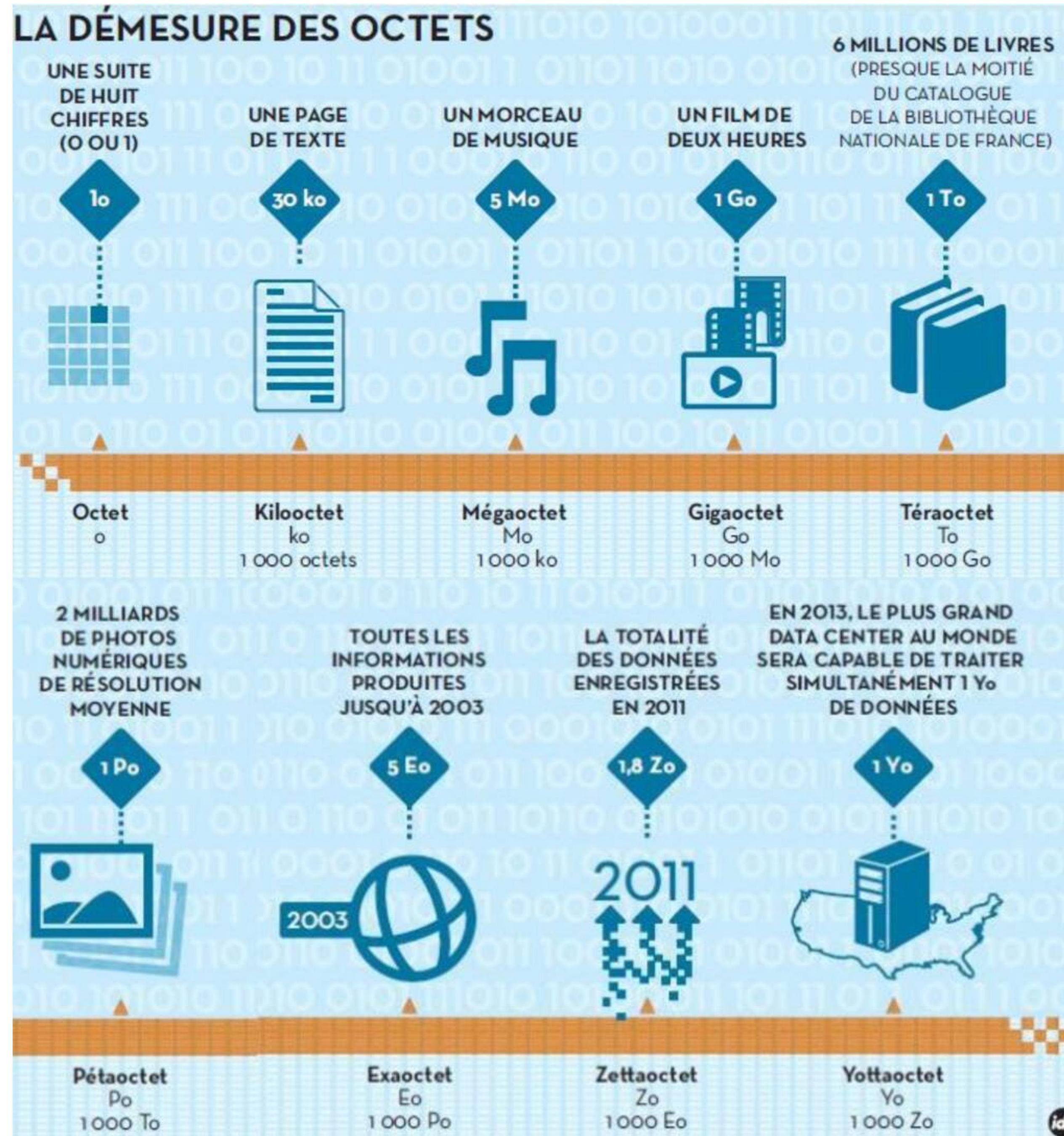


LA DÉMESURE DES OCTETS



“L’humanité produit autant d’informations en deux jours qu’elle ne l’a fait en deux millions d’années.”

(...)

Les chiffres donnent le tournis : chaque minute, environ 350 000 tweets, 15 millions de SMS et 200 millions de mails sont envoyés au niveau mondial ; pendant le même laps de temps, des dizaines d'heures de vidéos sont mises en ligne sur YouTube, des centaines de milliers de nouveaux fichiers sont archivés sur les serveurs de Facebook”

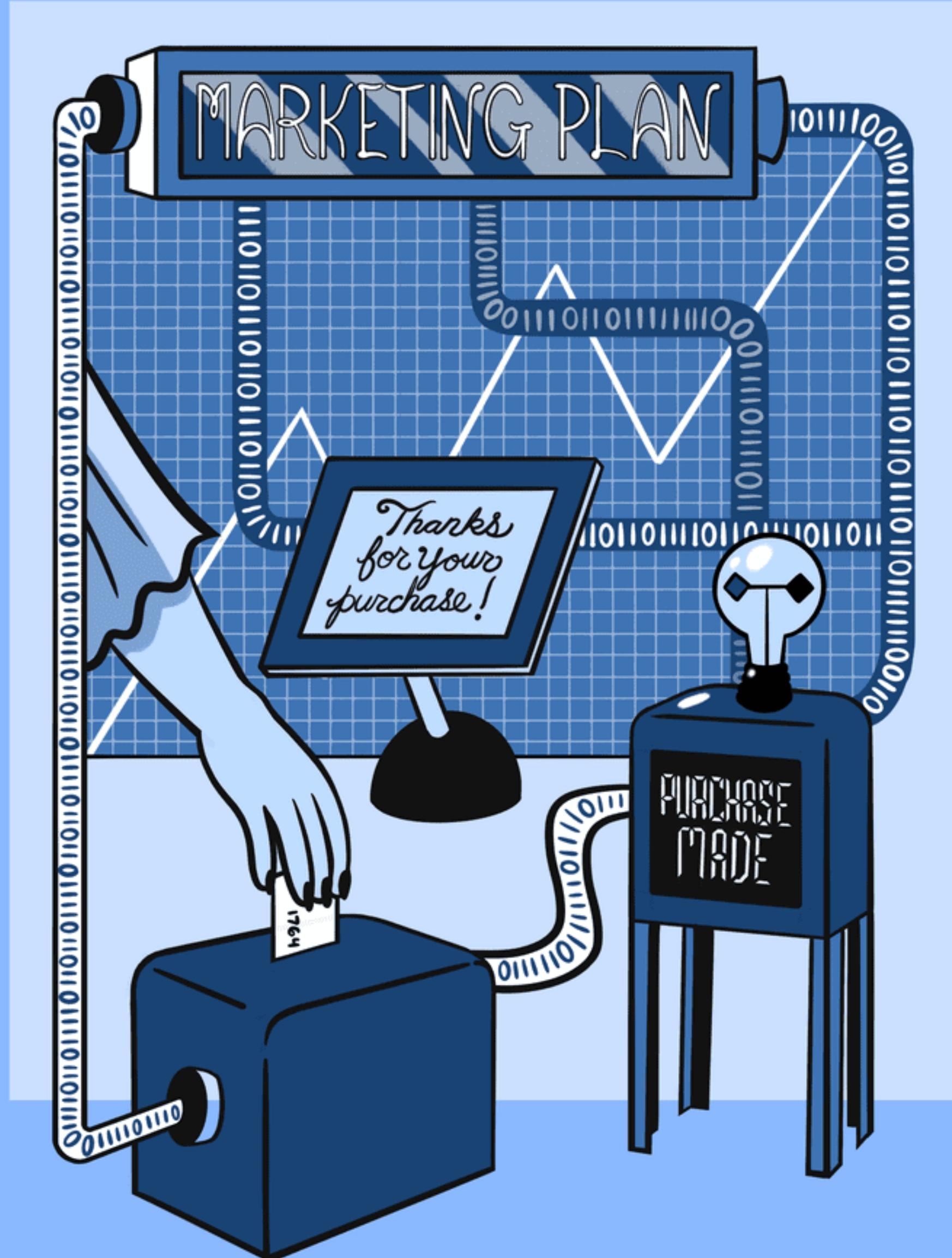
Et ça, c'était en 2012...

SAY BIG DATA



ONE MORE TIME

metacritic.com



Big Data

[*'big 'dā-tə*]

Large, diverse sets of information that grow at ever-increasing rates.

5 V's OF DATA



VOLUME

Amount of Data



VARIETY

Diversity of Data



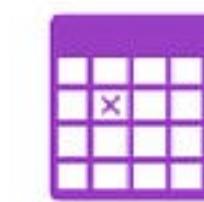
VELOCITY

Speed of
Data Generation



VALUE

Worth of Data

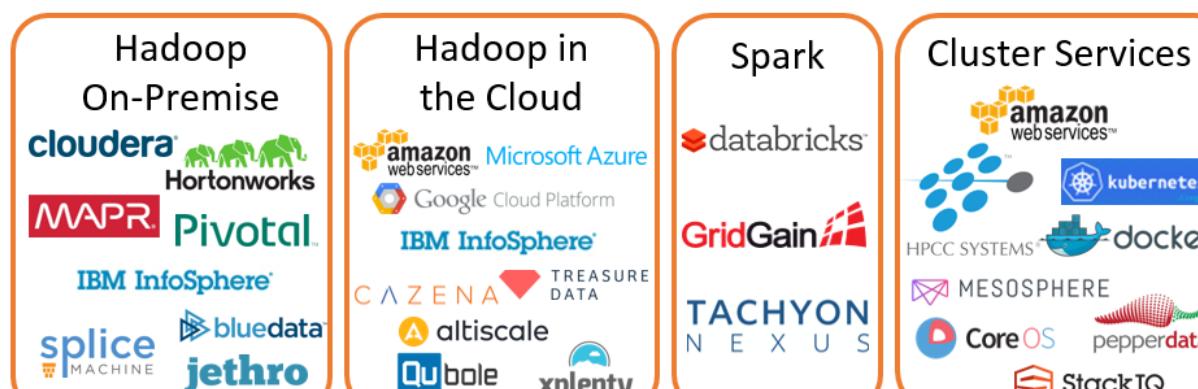


VERACITY

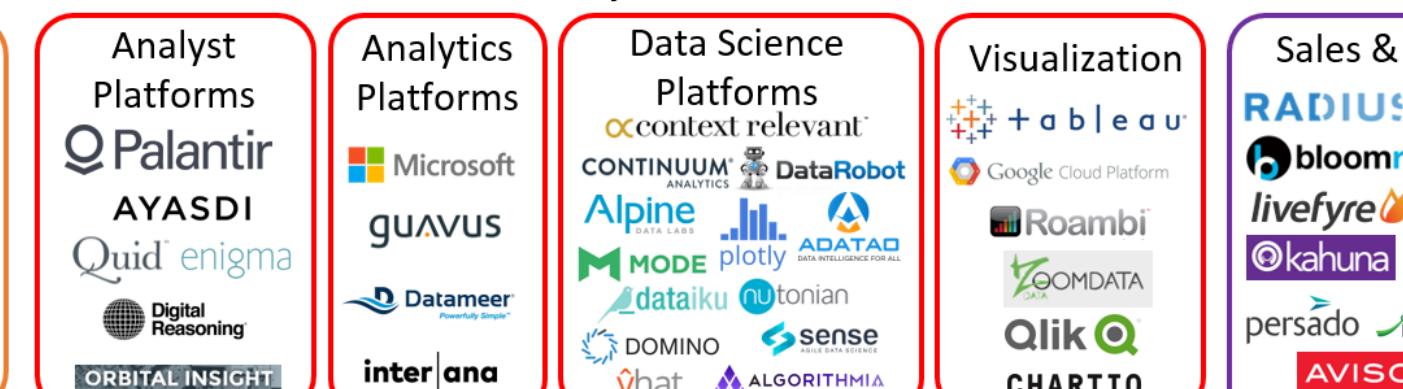
Accuracy of Data

Big Data Landscape 2016

Infrastructure



Analytics



Applications



NoSQL Databases



NewSQL Databases



BI Platforms



Statistical Computing



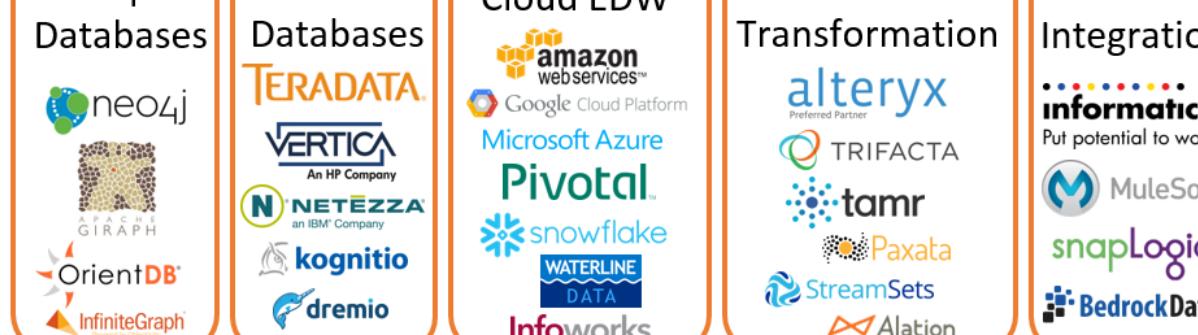
Log Analytics



Social Analytics



Graph Databases



MPP Databases



Cloud EDW



Data Transformation



Data Integration



Real-Time



Machine Learning



Speech & NLP



Horizontal AI



Publisher Tools



Govt/ Regulation



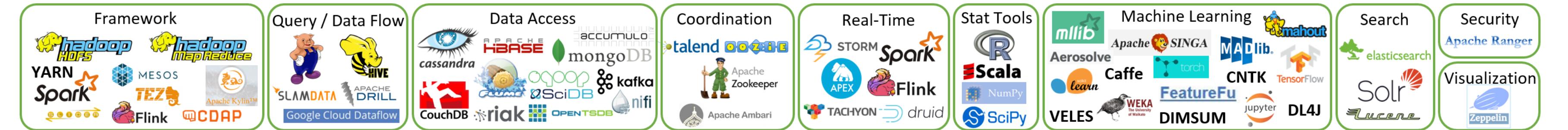
Finance



Cross-Infrastructure/Analytics



Open Source



Data Sources & APIs



THE 2023 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE



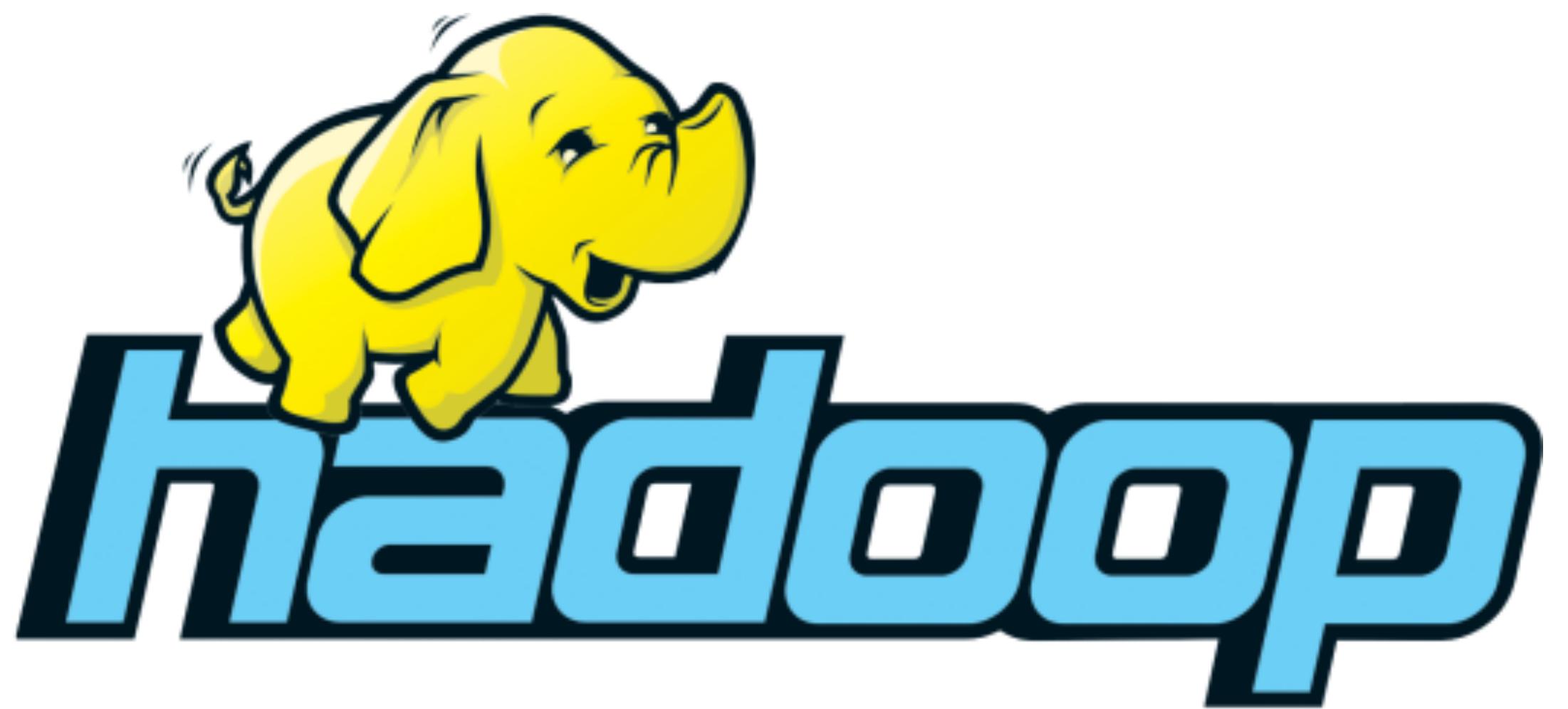
— OPEN SOURCE INFRASTRUCTURE —

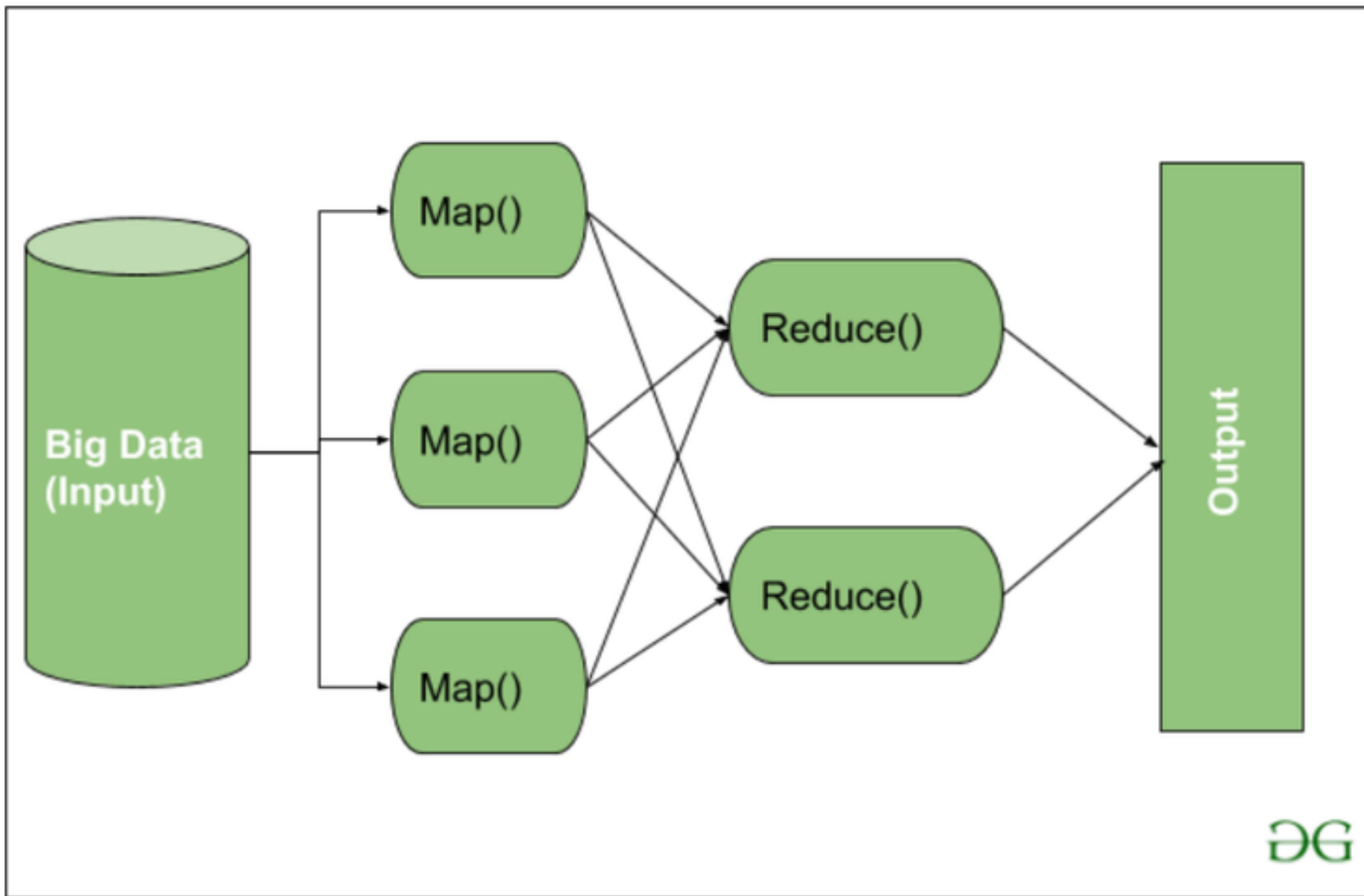


**10,000,000,000,000,000,
000,000,000,000,000,000,
000,000,000,000,000,000,
000,000,000,000,000,000,
000,000,000,000,000,000,
000,000,000 = 1 googol**

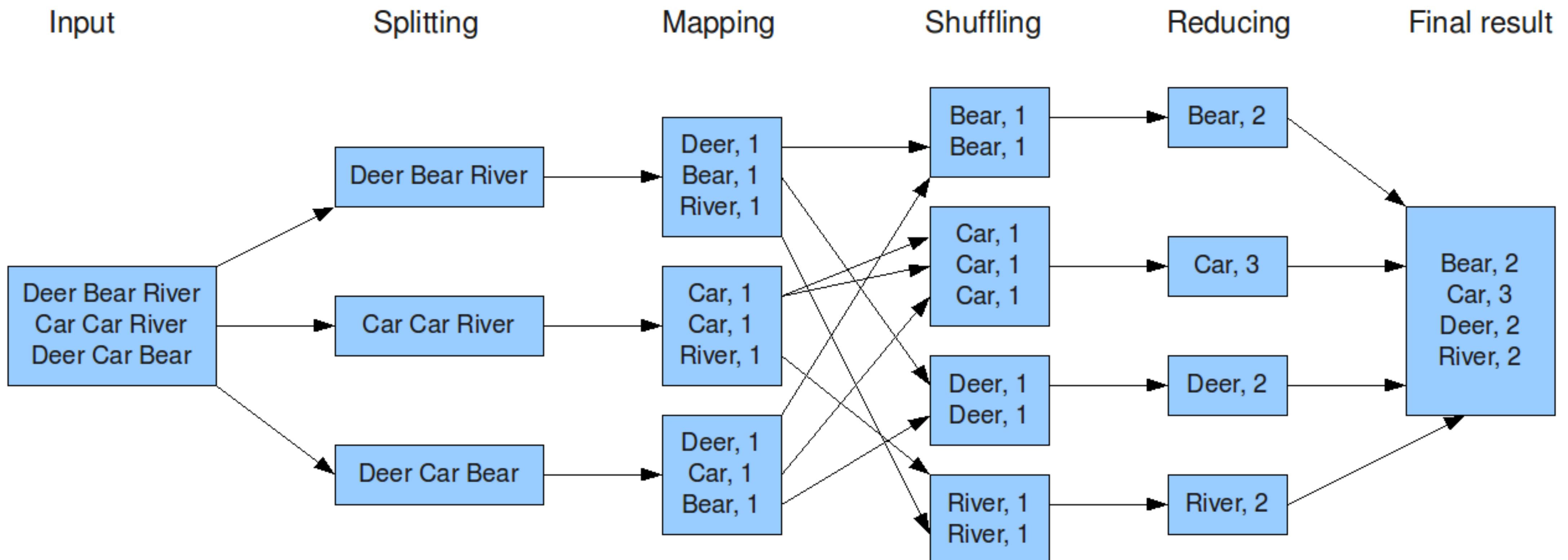


=>





The overall MapReduce word count process

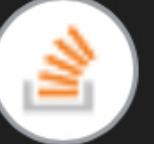




[CLÉMENT LEAFFRE](#) · 08h25, le 22 mars 2021

Le 10 mars, un important incendie ravageait les datacenters strasbourgeois d'OVH, leader français de l'hébergement de sites Internet. Depuis, une centaine de techniciens travaillent jour et nuit pour sauver ce qui peut l'être et redémarrer les serveurs. Mais des entreprises ont perdu des données capitales et l'image d'OVH risque d'en pâtir.



 Stack Overflow
<https://stackoverflow.com> › why-is...

Why is hadoop slow for a simple hello world job

Mar 1, 2019 — I am following the tutorial on the hadoop website: <https://hadoop.apache.org/docs/r3.1.2/hadoop-project-dist/hadoop-common/SingleCluster.html>. I ...

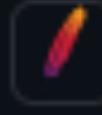


Apache Hadoop





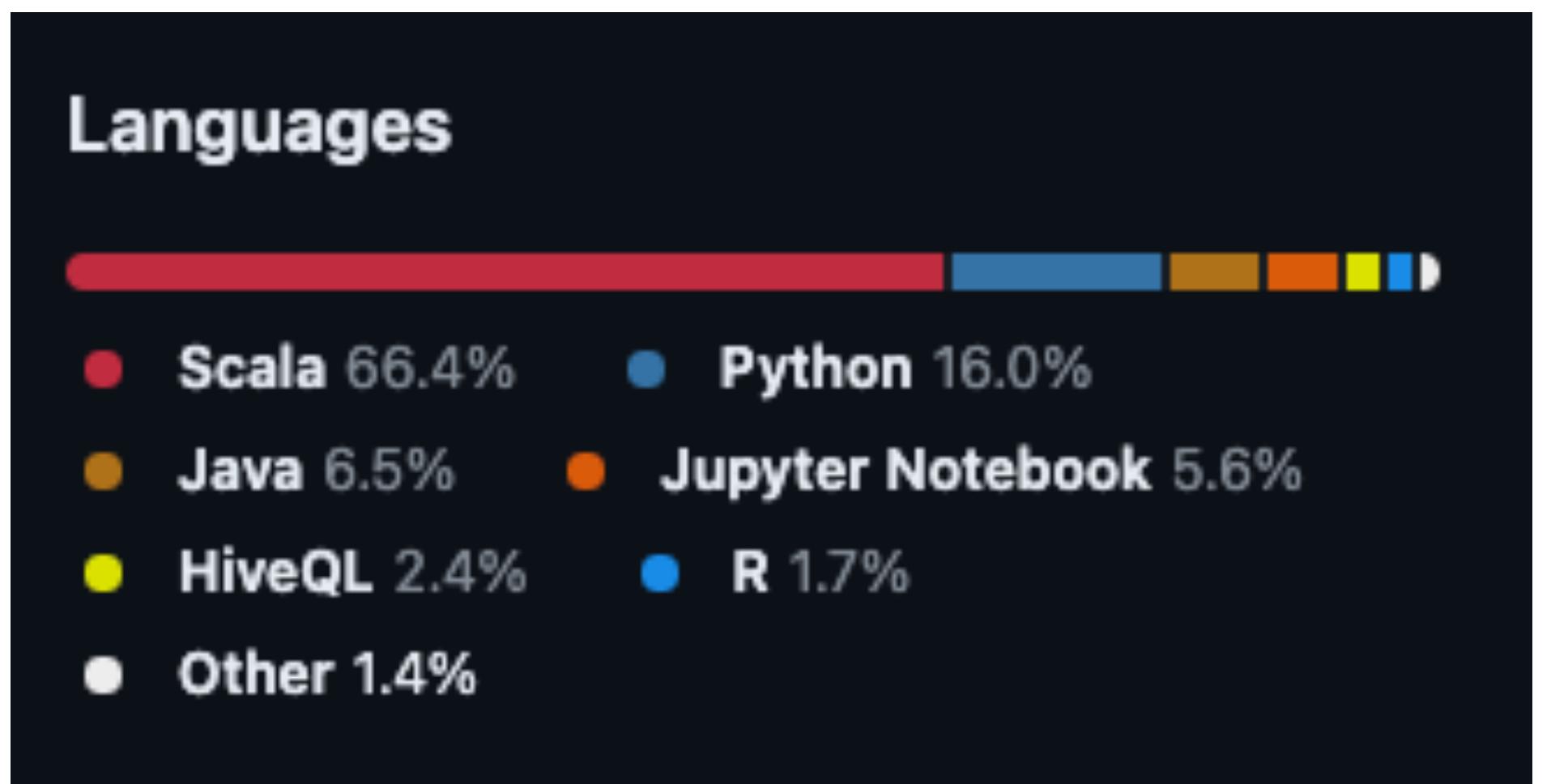
Spark

 apache/spark

Apache Spark - A unified analytics engine for large-scale data processing

python java r scala sql

Scala · ⭐ 38.1k · Updated 1 minute ago



Spark SQL and
DataFrames +
Datasets

Spark Streaming
(Structured
Streaming)

Machine Learning
MLlib

Graph
Processing
Graph X

Spark Core and Spark SQL Engine

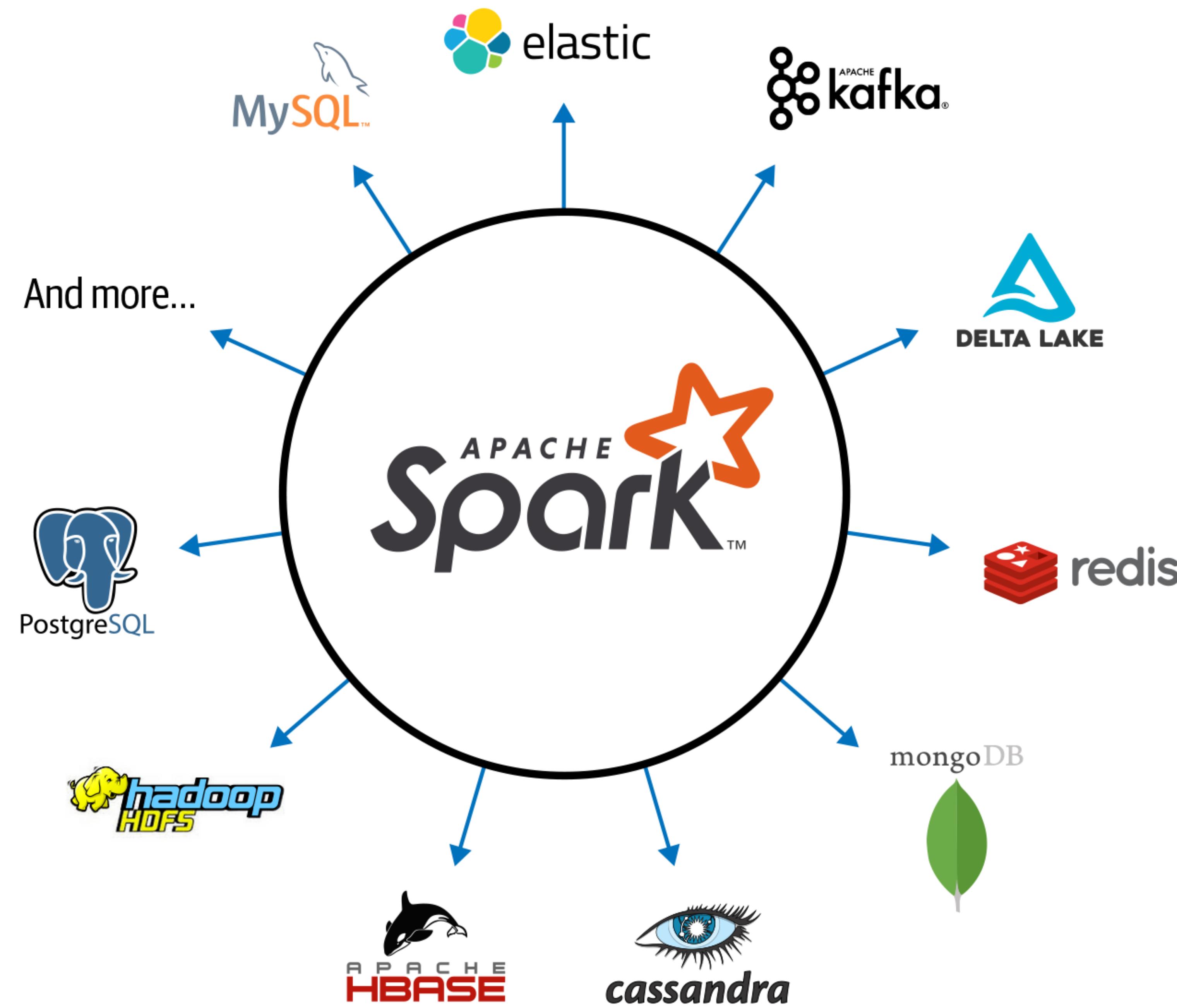
Scala

SQL

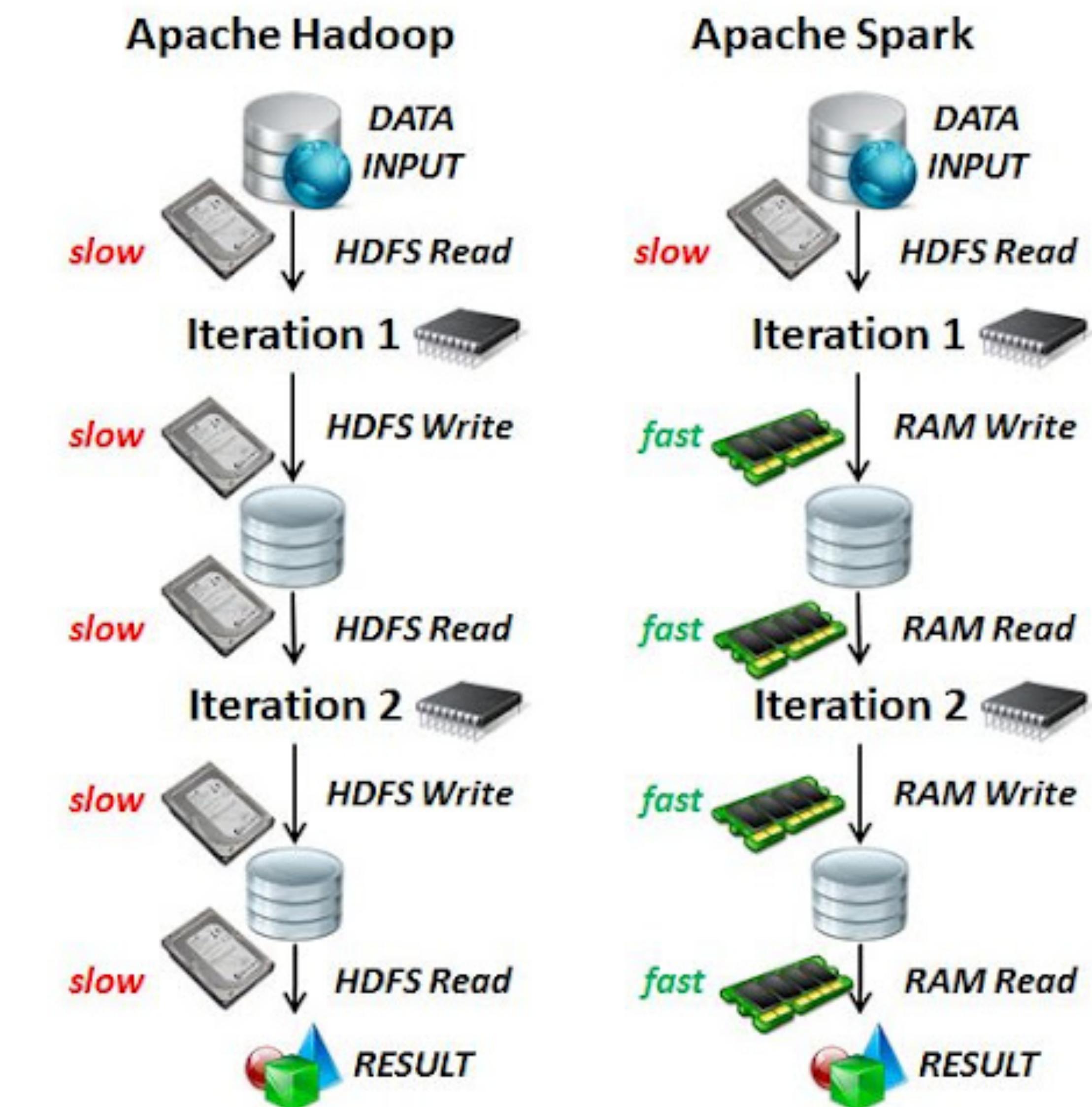
Python

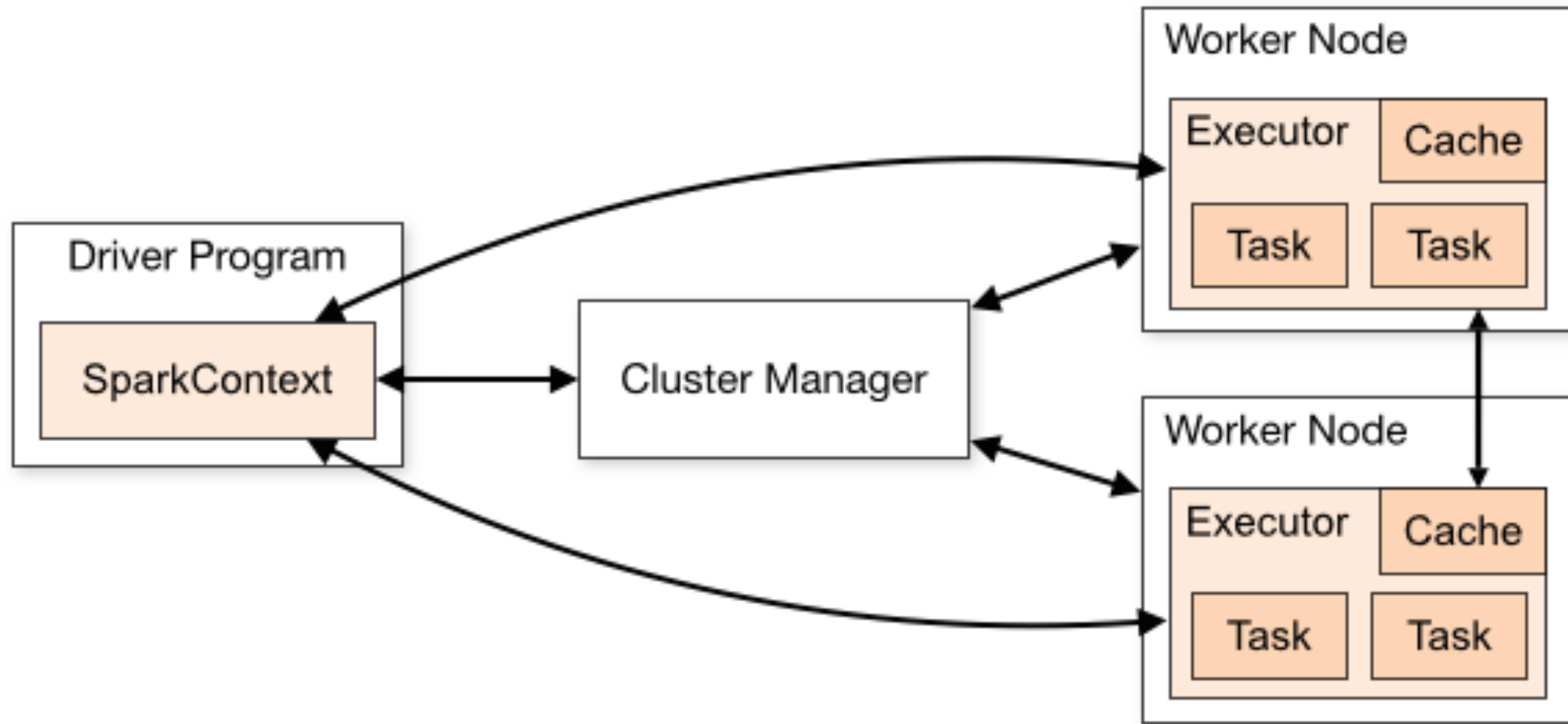
Java

R









Mode	Spark driver	Spark executor	Cluster manager
Local	Runs on a single JVM, like a laptop or single node	Runs on the same JVM as the driver	Runs on the same host
Standalone	Can run on any node in the cluster	Each node in the cluster will launch its own executor JVM	Can be allocated arbitrarily to any host in the cluster
YARN (client)	Runs on a client, not part of the cluster	YARN's NodeManager's container	YARN's Resource Manager works with YARN's Application Master to allocate the containers on NodeManagers for executors
YARN (cluster)	Runs with the YARN Application Master	Same as YARN client mode	Same as YARN client mode
Kubernetes	Runs in a Kubernetes pod	Each worker runs within its own pod	Kubernetes Master

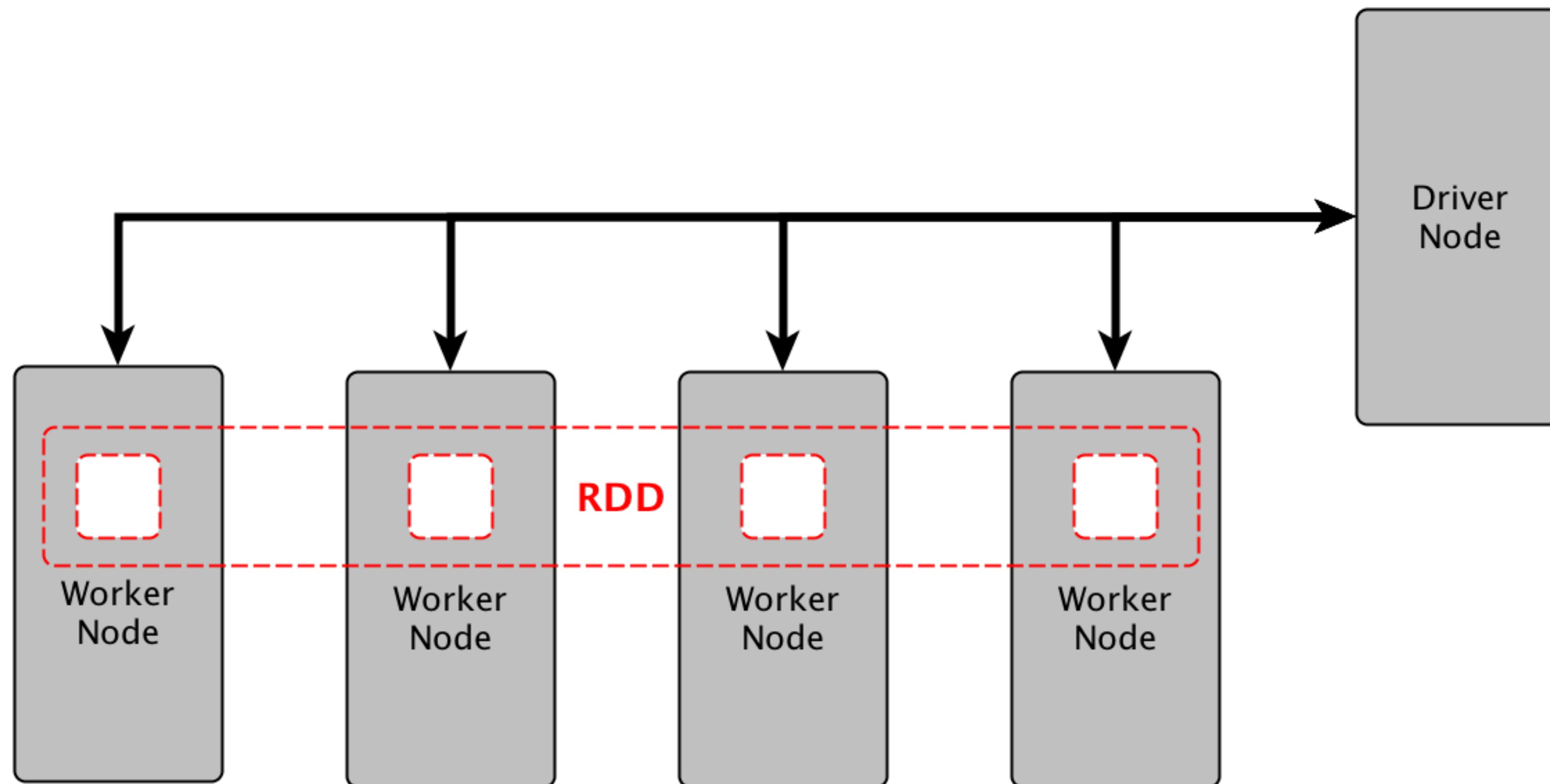
RDD

Resilient Distributed Dataset

Lazy Evaluation

Lazy evaluation (or call-by-name) is an evaluation strategy which delays the evaluation of an expression until its value is needed





 **Spark Operations**

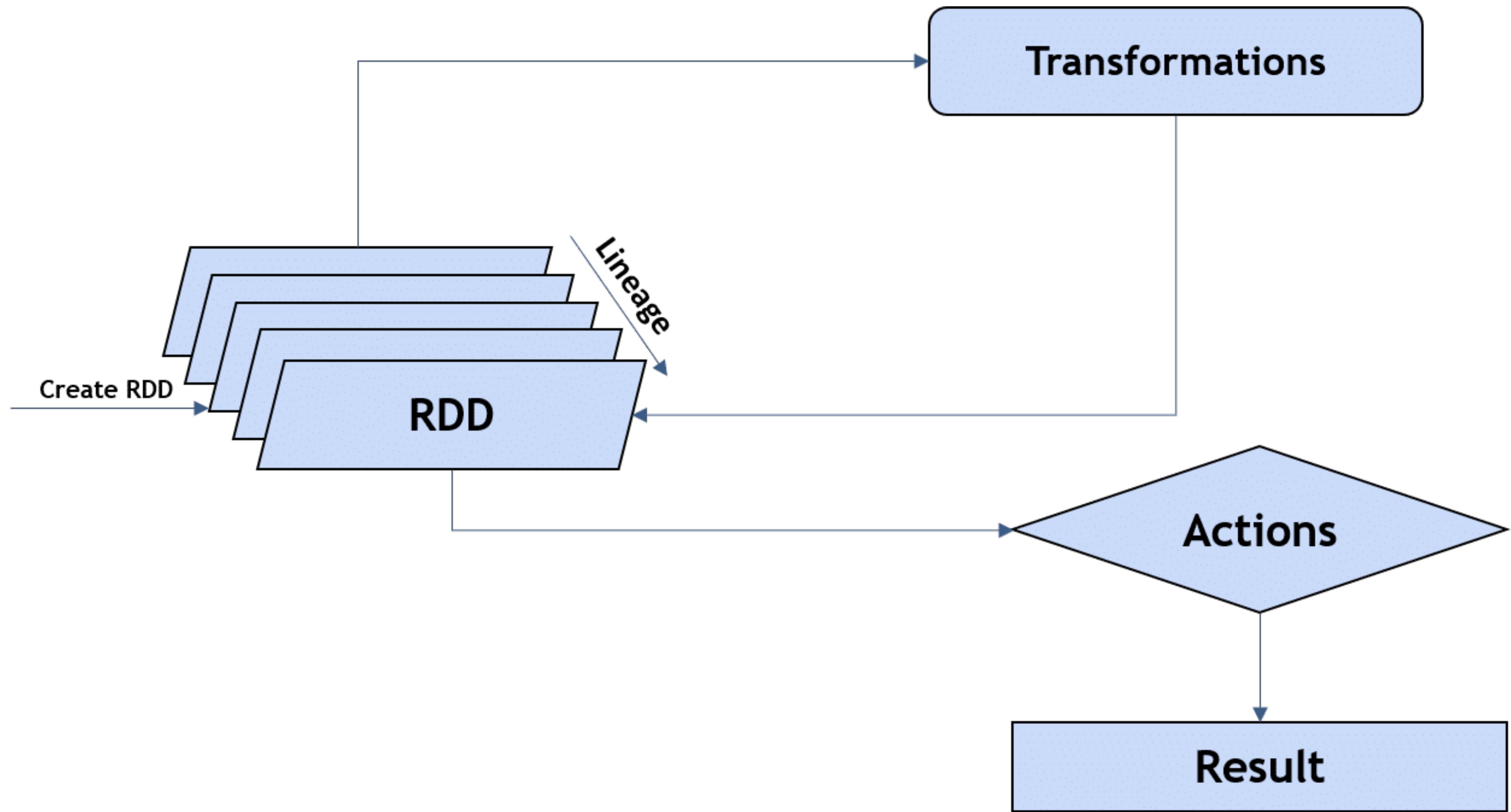
=

TRANSFORMATIONS

+



ACTIONS



RDD Operations

Transformations

Narrow

- map
- flatMap
- mapPartition
- filter
- sample
- union

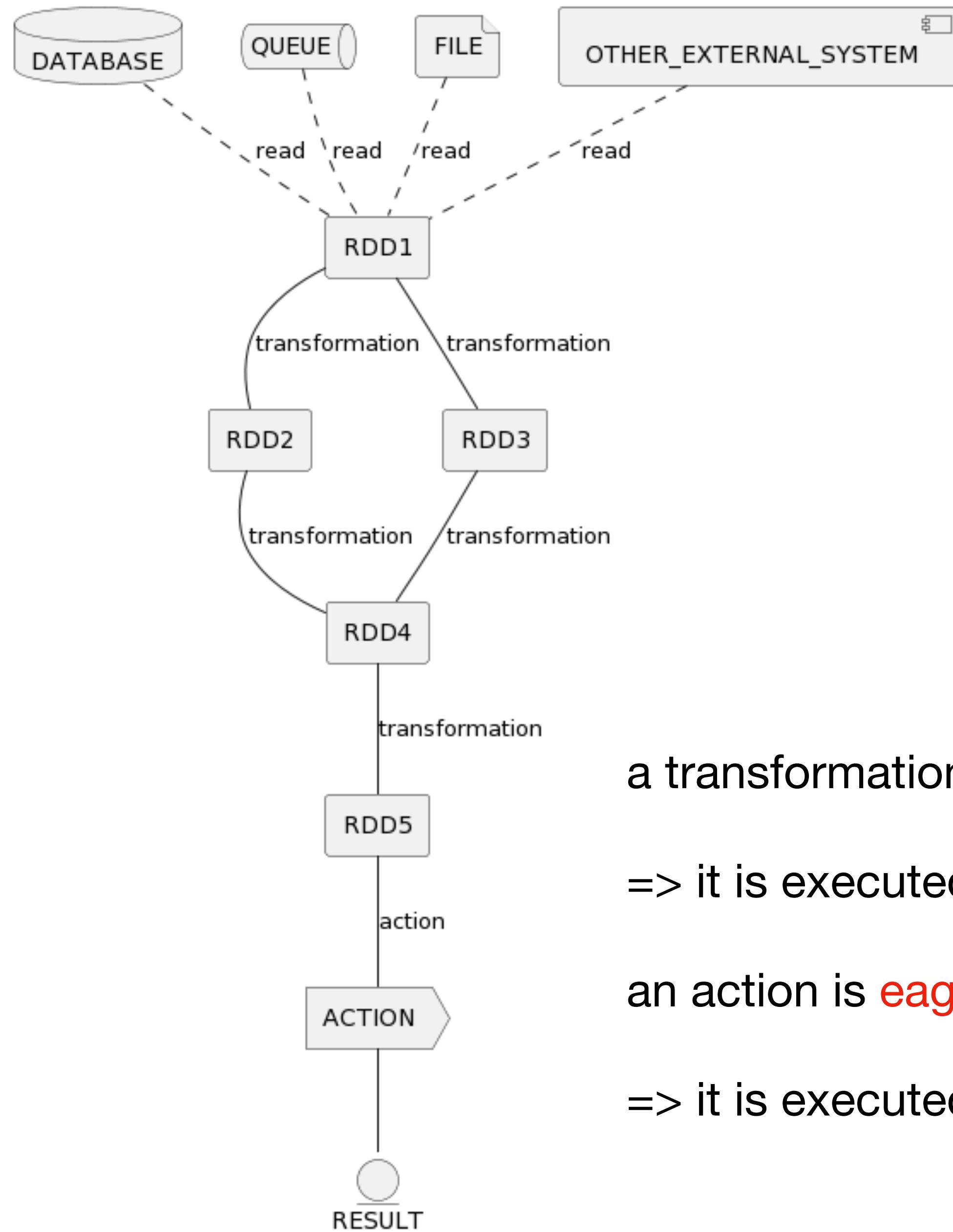
Wide

- intersection
- distinct
- reduceByKey
- groupByKey
- join
- cartesian
- repartition
- coalesce

Actions

- count
- collect
- take(n)
- countByValue
- top
- reduce
- fold
- aggregate
- foreach

Spark execution plan is a
Direct Acyclic Graph (DAG)



a transformation is **lazily** evaluated

=> it is executed when an action using it is called

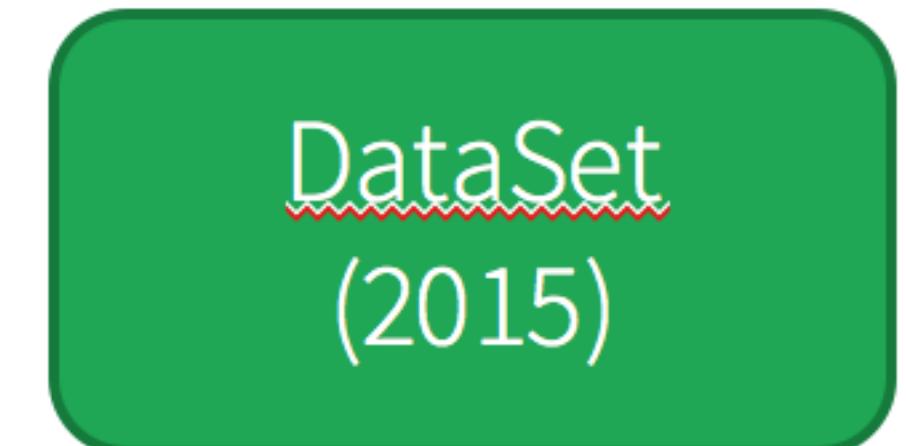
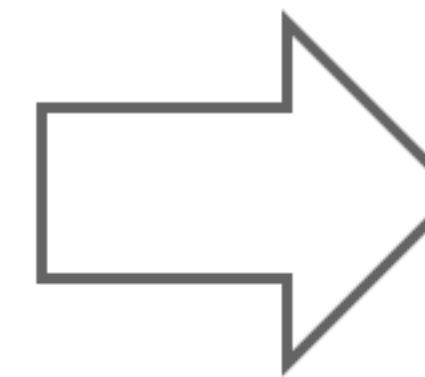
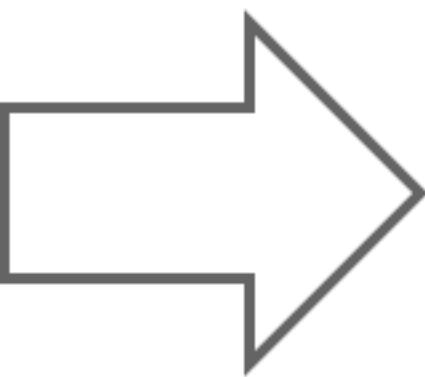
an action is **eagerly** evaluated

=> it is executed instantly



```
1 val schedules = spark.read.parquet("/input")
2
3 schedules
4 .filter(_.onlyInternationalFlights)
5 .map(_.makeFlightInfo)
6 .write.parquet("/output")
```

History of Spark APIs



Distribute collection
of JVM objects

Functional Operators (map,
filter, etc.)

Distribute collection
of Row objects

Expression-based operations
and UDFs

Logical plans and optimizer

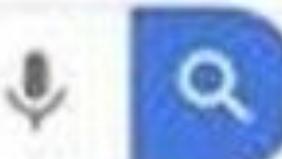
Fast/efficient internal
representations

Internally rows, externally
JVM objects

Almost the “Best of both
worlds”: type safe + fast

But slower than DF
Not as good for interactive
analysis, especially Python

when was php released



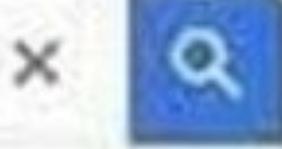
ALL IMAGES NEWS VIDEOS MAPS

PHP / Initial release date

June 8, 1995



How long does trash live?



ALL IMAGES VIDEOS NEWS MAPS

Trash / Lifespan

1000 years



Lower classifications and overview



Feature	Spark RDD	Spark DataFrame	Spark Dataset
Data representation	Immutable distributed collection of data	Structured data organized into named columns	Distributed collection of data with optional schema
Data processing	Fine-grained control	High-level abstraction	Ease of use and performance
Suitability	Developers who require precise control	Data analysts and SQL experts	Data professionals who need a balance of control and convenience
Key distinctions	Offers more control, but more complex	Offers more convenience, but less control	Offers a balance of control and convenience

