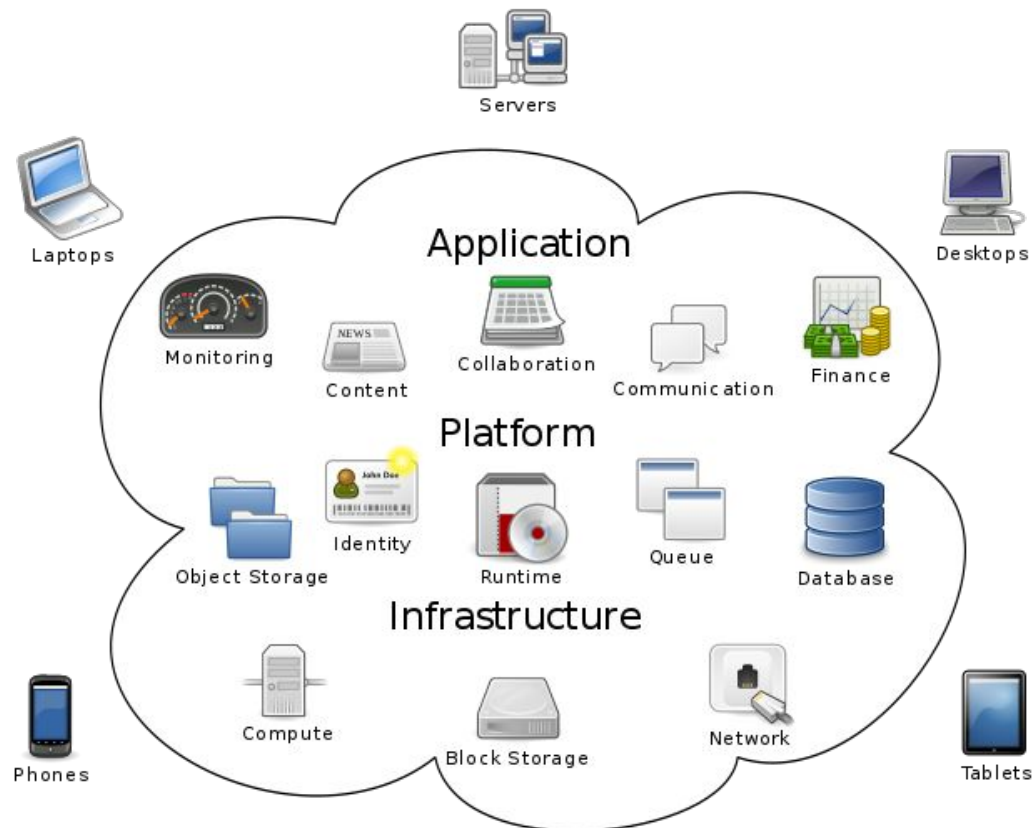


Cloud Computing, Big Data, and the Semantic Web

Dr. V. Raghava Mutharaju
IIIT-Delhi

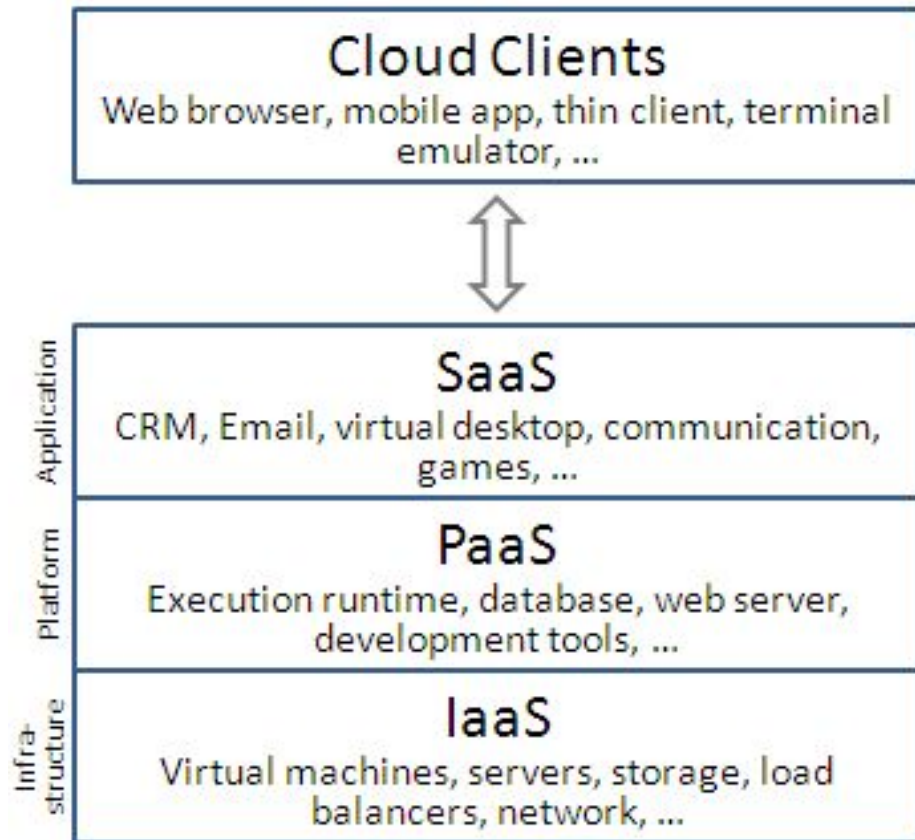


Cloud computing

Cloud Computing

- Laptops → Servers → Cloud
 - Maintenance of machines/software
 - Availability of machines
 - Need the machines/software for only a limited amount of time
 - Don't have the expertise to setup the software on networked machines
-
- Pay-as-you-go model

Cloud Computing



IAAS

- Infrastructure As A Service (IAAS)
 - Physical computing resources are available as a service in the form of VMs
 - Number of cores
 - RAM size
 - Disk

PAAS

- Platform As A Service (PAAS)
 - Platform/Application environment is provided as a service
 - Operating System
 - Database
 - Java runtime
 - .NET runtime

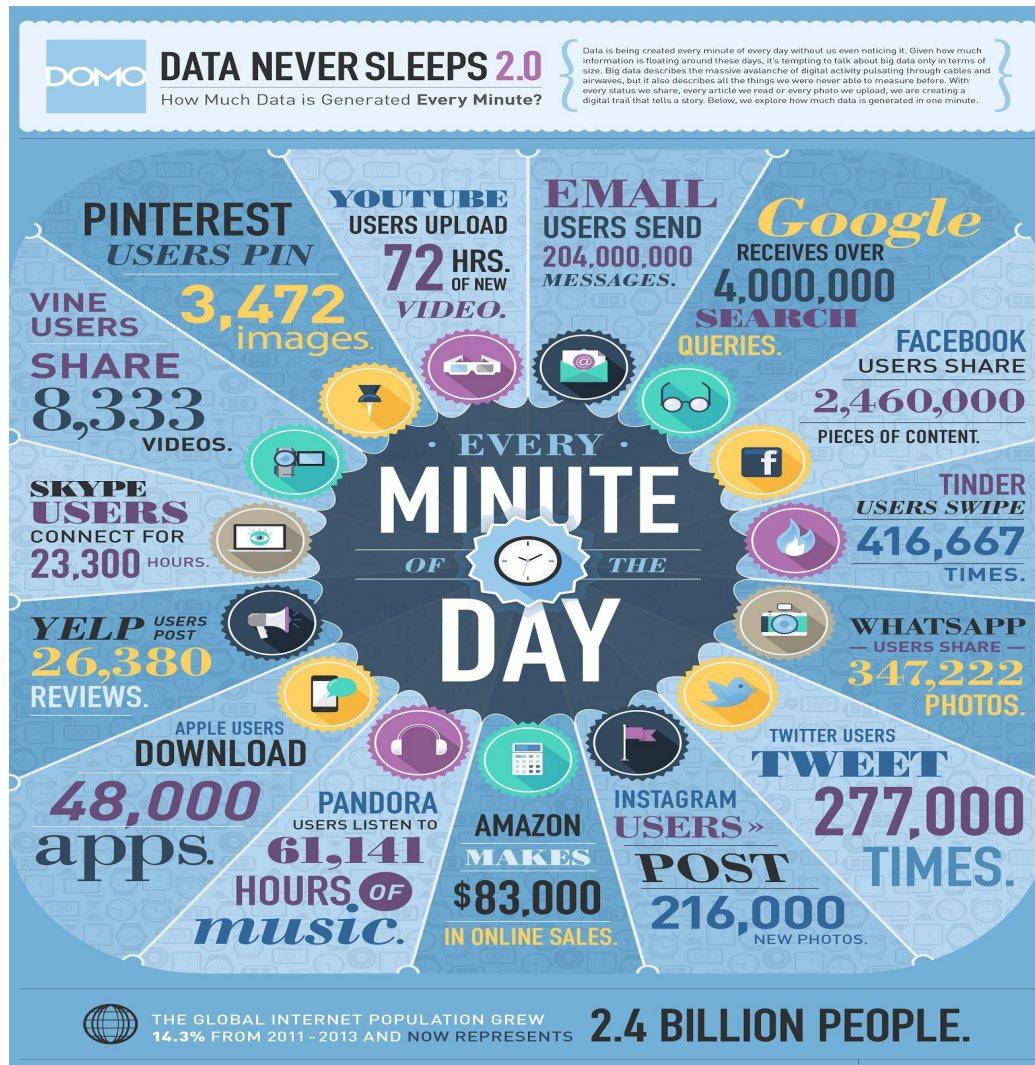
SAAS

- Software As A Service (SAAS)
 - Software is provided as a service
 - Office software
 - Messaging software
 - CAD software
 - Payroll, Accounting software
 - Content management software

IBM Cloud

- Combines PAAS with IAAS
- Demo

Big Data



VOLUME

Huge amount of data



VERACITY

Inconsistencies and uncertainty in data



VARIETY

Different formats of data from various sources



VELOCITY

High speed of accumulation of data



VALUE

Extract useful data



BIG DATA

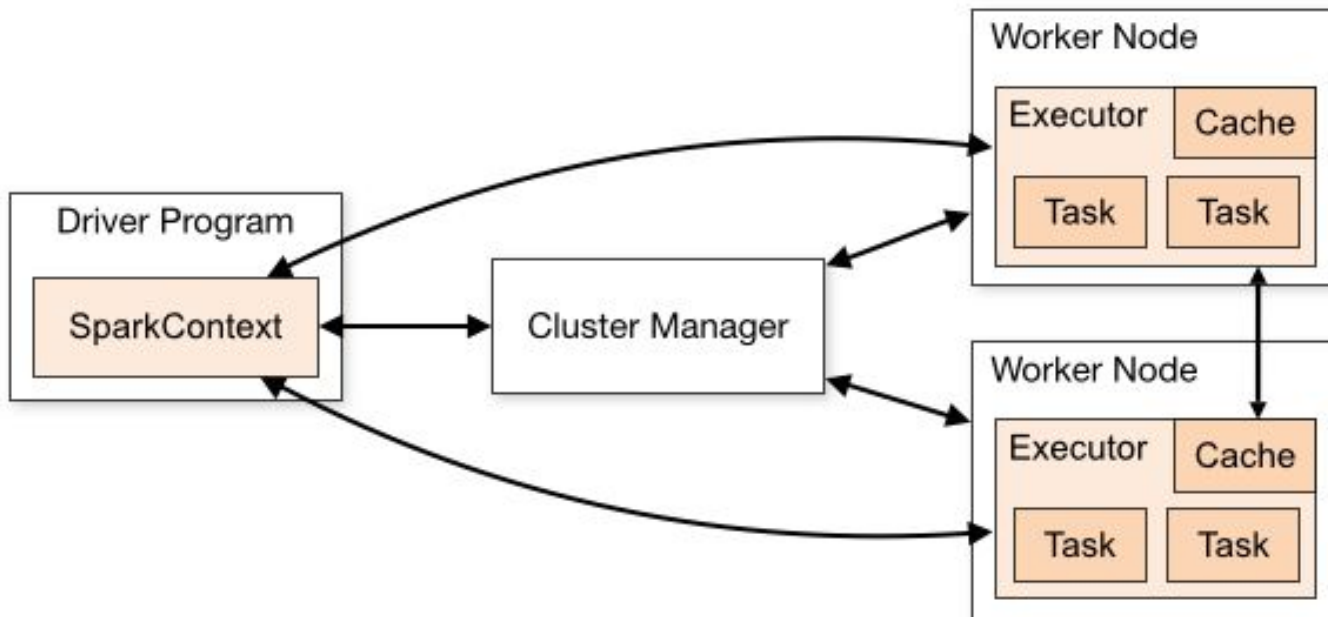
Apache Spark

- Handles the volume aspect very well
- Has also support for velocity
- Has support for Graph and Machine Learning algorithms that can be used to generate value from the data

Apache Spark

- Distribute the data and the computations across a cluster (group of machines)
- This increases the parallelism and reduces the load on each machine
 - Need more machines? Pay-as-you-go (add more machines depending on the load)
- Spark provides mechanism to
 - Capture the data in a form that is easily distributable
 - Perform computations that can be parallelized

Spark Cluster



Spark RDD

- Resilient Distributed Dataset (RDD)
 - It is the fundamental datastructure of Apache Spark
 - It is a immutable, fault tolerant collection of elements that can be operated in parallel
 - An RDD can be divided into logical partitions that can be distributed across the cluster
 - Elements of RDD can be text, numbers, or a combination of them

Transformations and Actions

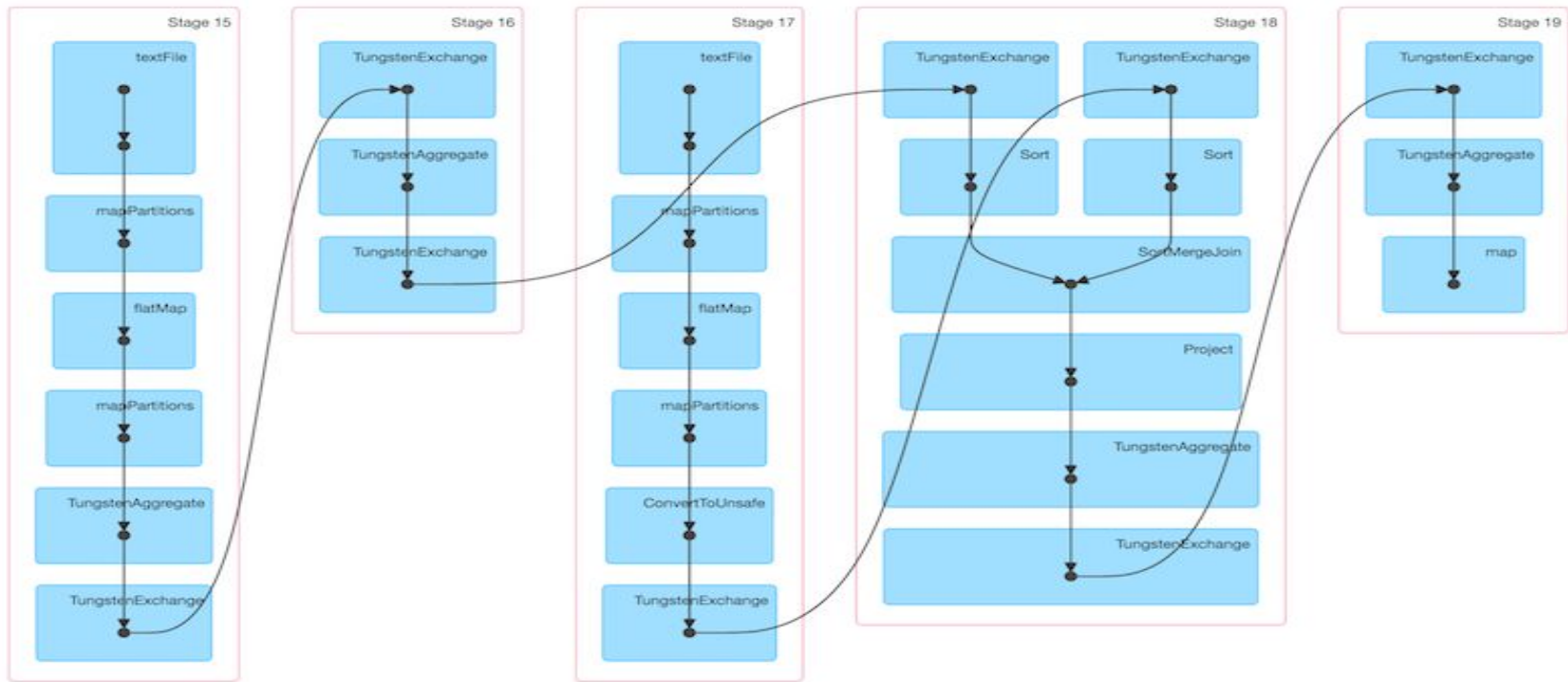
Transformations

map	join	union	distinct	repartition
mapPartitions	flatMap	intersection	pipe	coalesce
cartesian	cogroup	filter	sample	
sortByKey	groupByKey	reduceByKey	aggregateByKey	
mapPartitionsWithIndex		repartitionAndSortWithinPartitions		

Actions

reduce	take	collect	takeSample	count
takeOrdered	countByKey	first	foreach	saveAsTextFile
saveAsSequenceFile		saveAsObjectFile		

Directed Acyclic Graphs



Big Data Problems

- Parallelizing very large datasets (Protein sequencing)
- Estimate the runtime of Apache Spark jobs
- Automate the generation of code based on text descriptions and heuristics