



GETTING STARTED WITH HADOOP AND AZURE HDINSIGHT



Niti Gupta



Niti Gupta



@nitigupta_16



@<https://www.linkedin.com/in/niti-gupta-33200086>

OUTLINE

- What is big data?
- Why Hadoop? ⁺ •
- What is HDInsight?
- Why HDInsight?
- Scenarios where you can use HDInsight.
- How can you get started?

○

FUN FACTS



Every human created about 1.7 mb of data per second in 2020

(Tech Jury)

+

- WhatsApp users exchange up to 65 billion messages daily

(Source: Connectiva Systems)

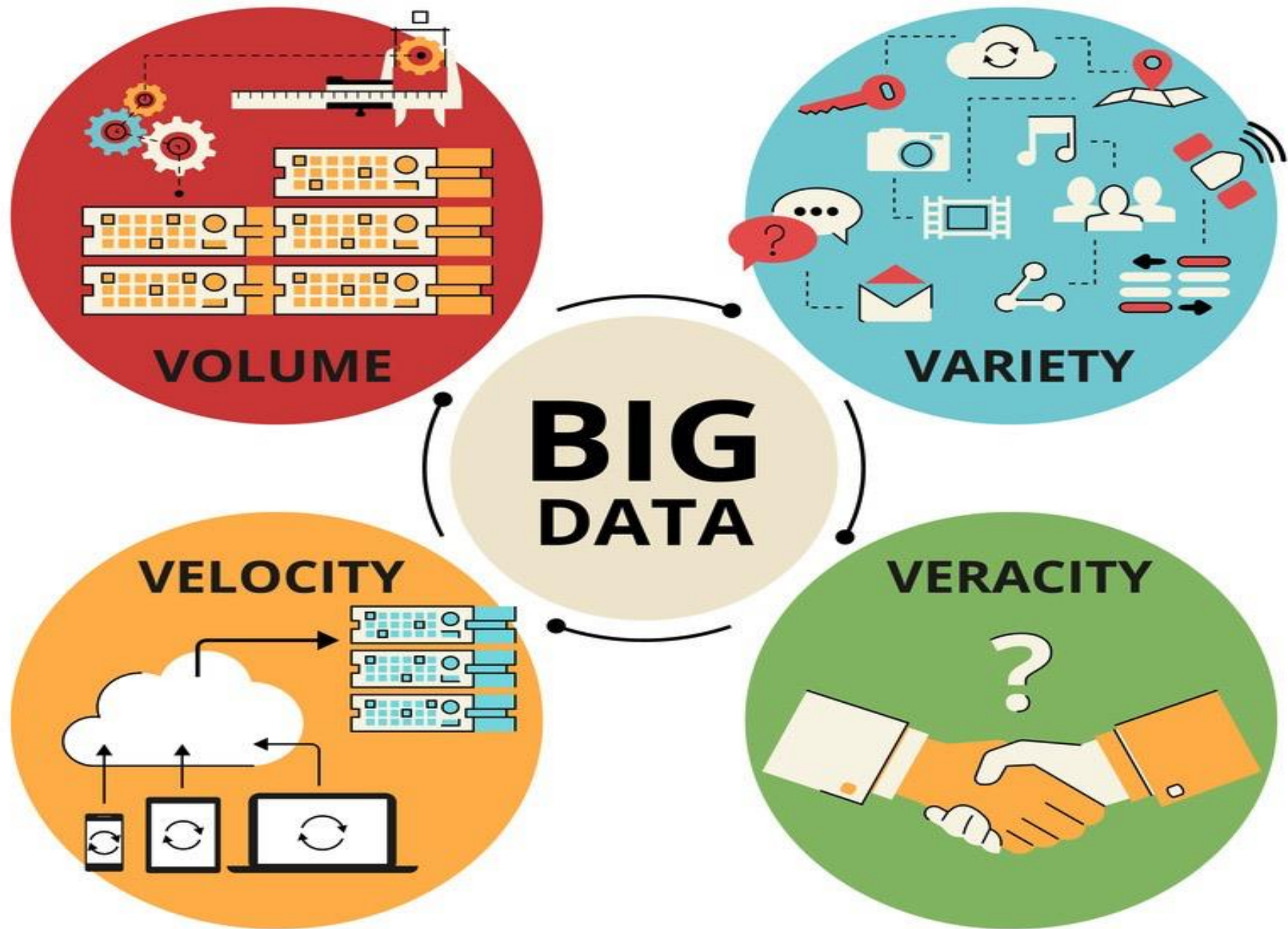
Every 2 days we create as much information as we did from the beginning of time until 2003

(CloudNine)

If you burned all of the data created in just one day onto DVDs, you could stack them on top of each other and reach the moon – twice

(CloudNine)





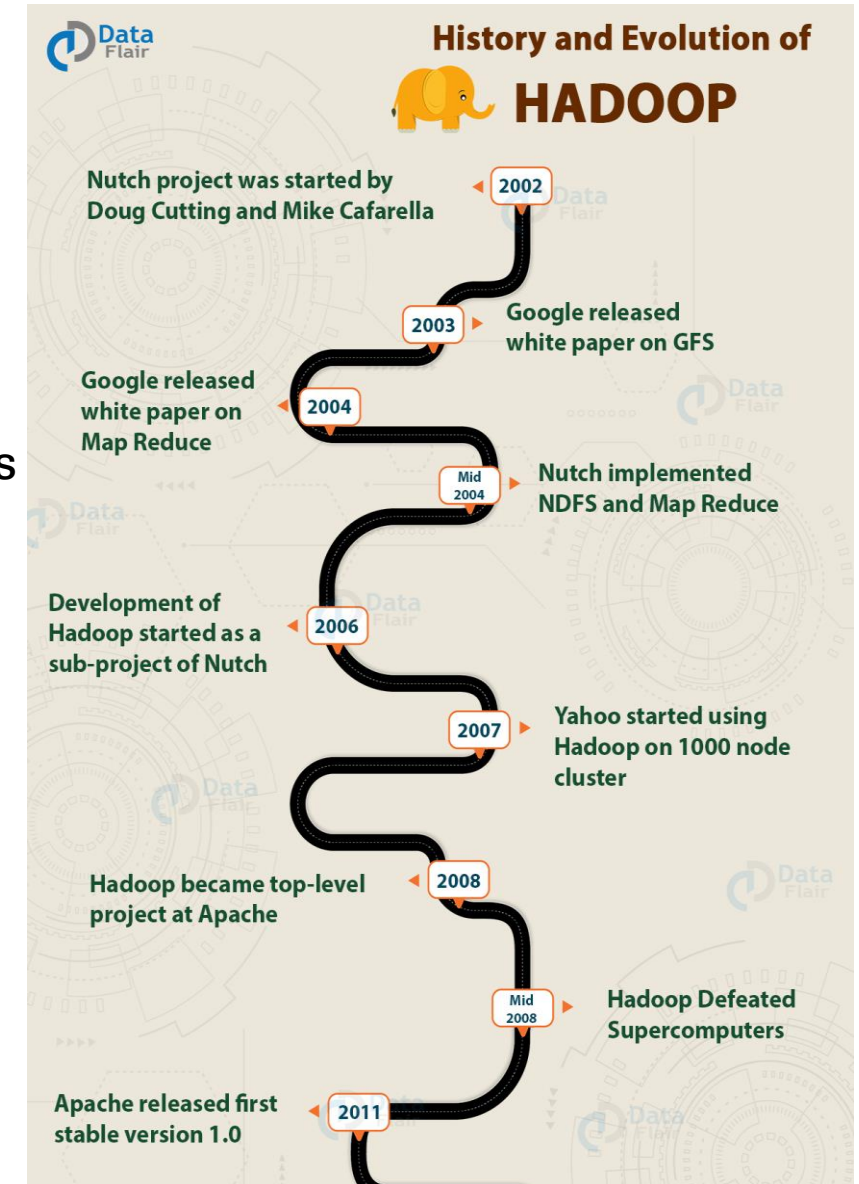
What is Hadoop?

Google's Solution

Google solved this problem using an algorithm called MapReduce. This algorithm divides the task into small parts and assigns them to many computers, and collects the results from them which when integrated, form the result dataset.

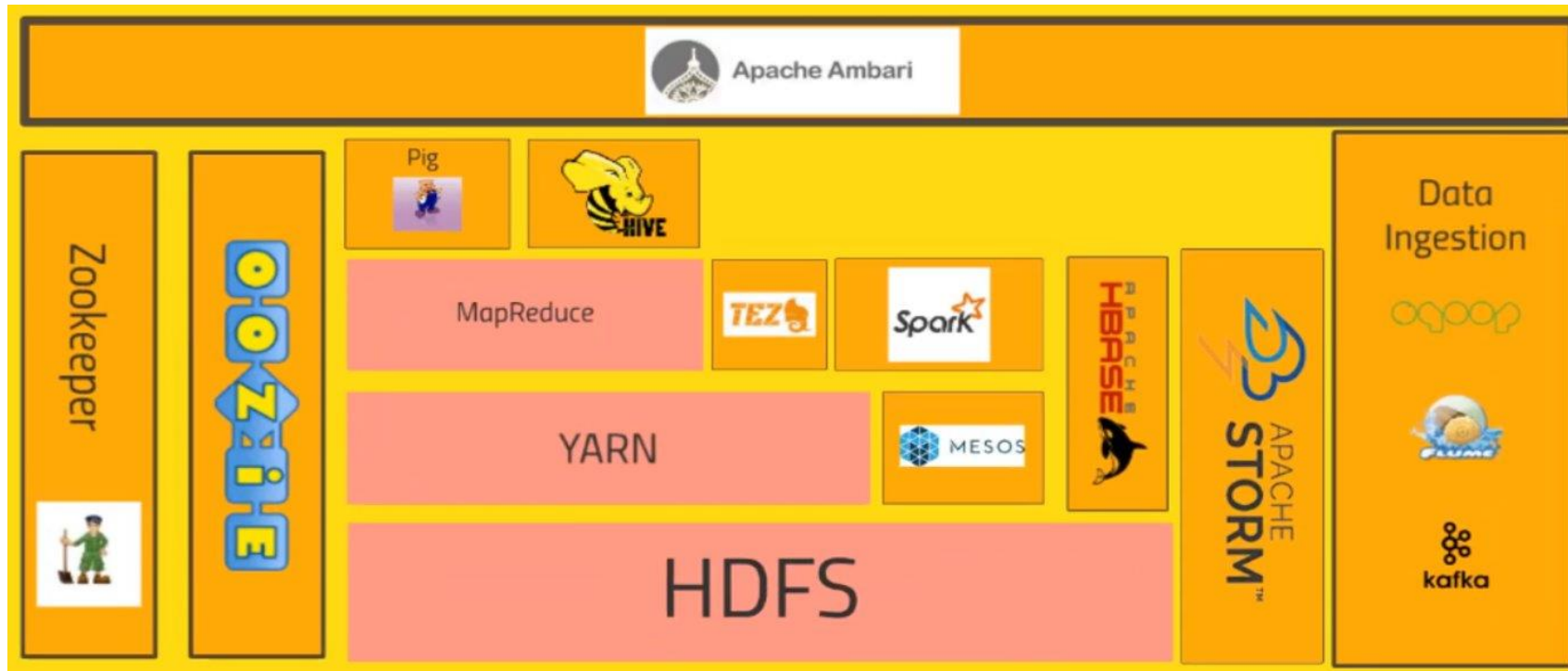
Using the solution provided by Google, Doug Cutting and his team developed an Open Source Project called **HADOOP**.

Hadoop runs applications using the MapReduce algorithm, where the data is processed in parallel with others.

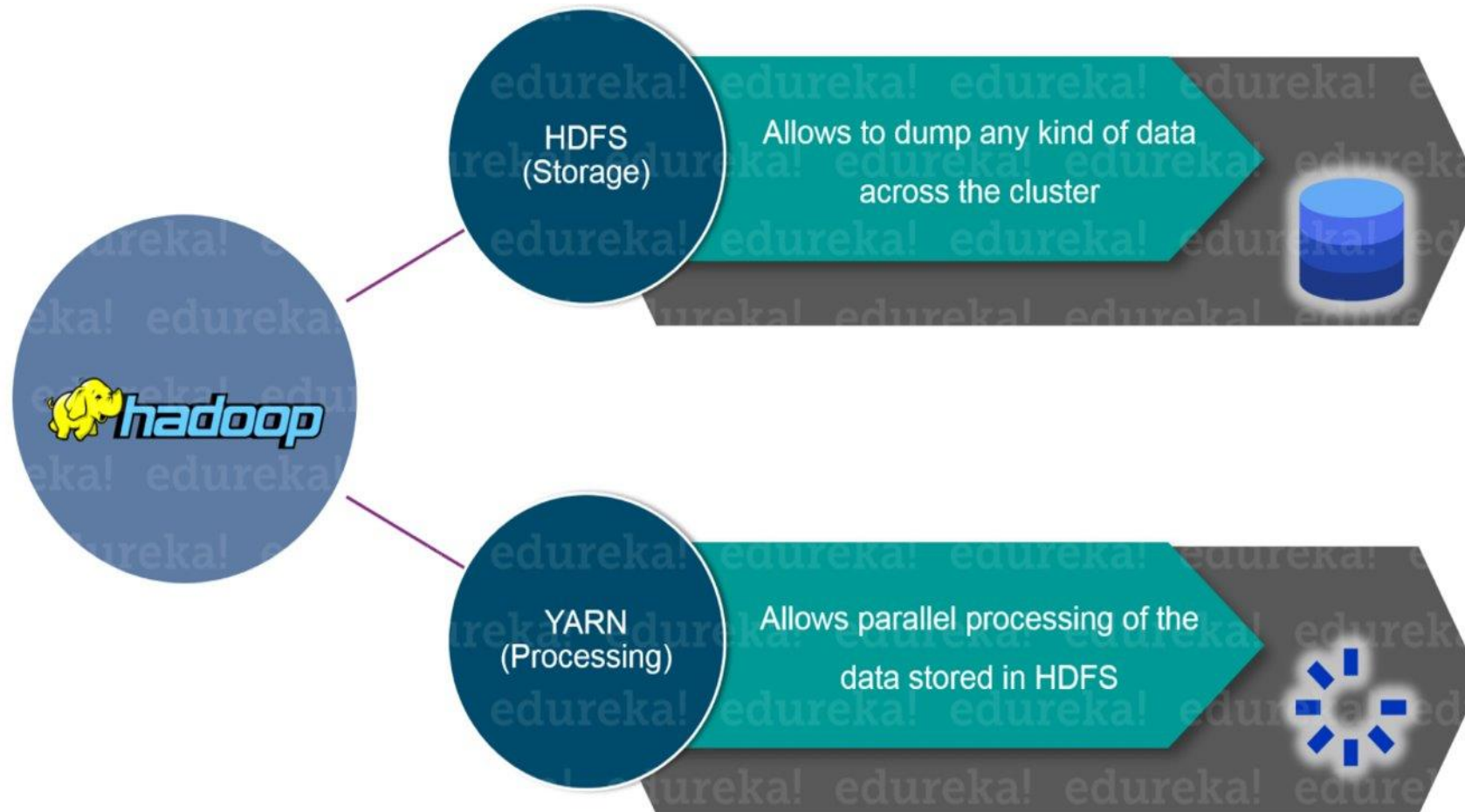


Discover Hadoop

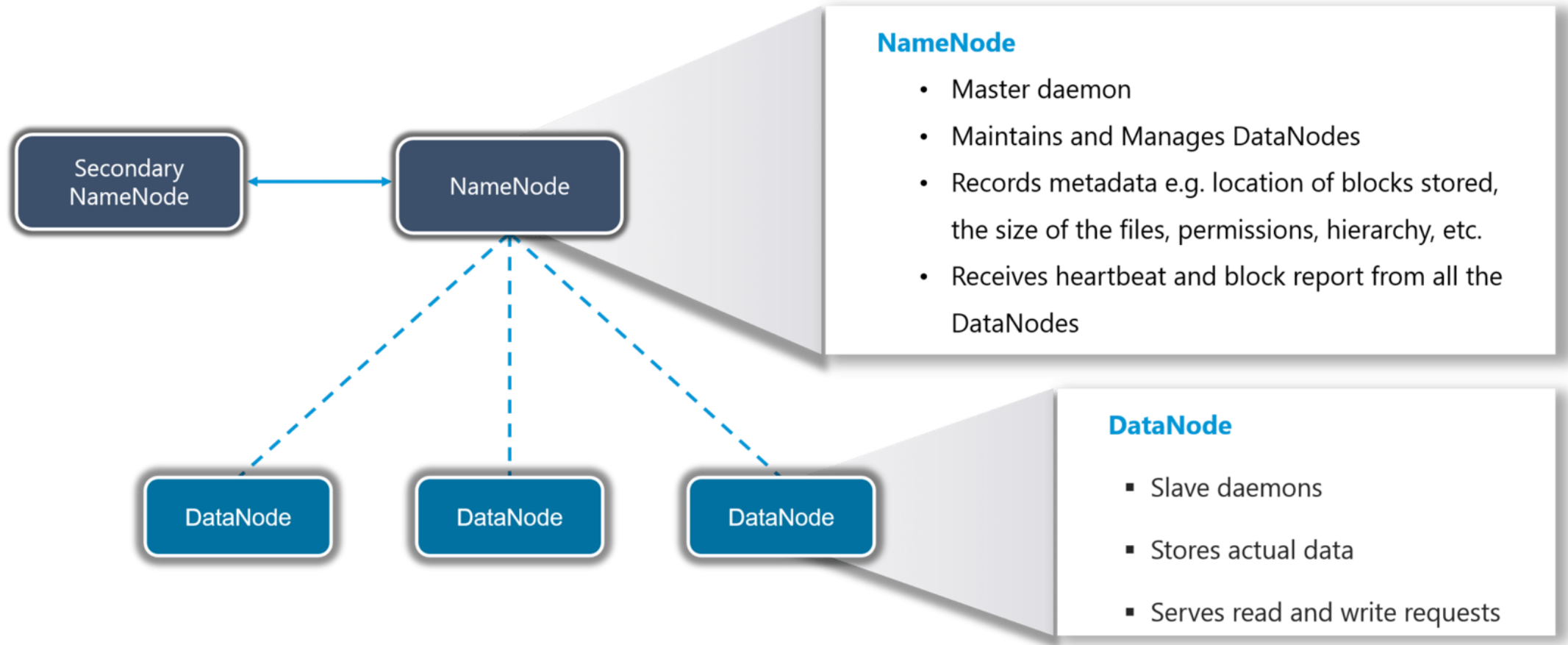
[Hadoop](#) is an open-source software framework for storing data and running applications on clusters of commodity hardware. It provides **massive storage** for any kind of data, **enormous processing power** and the ability to handle virtually **limitless concurrent tasks** or jobs.



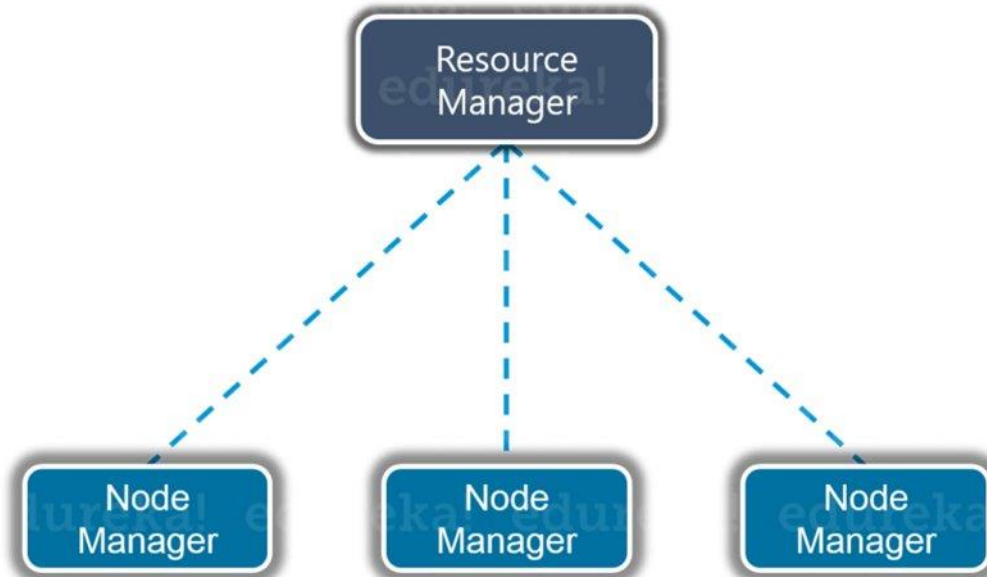
Hadoop Core Components



HDFS



YARN



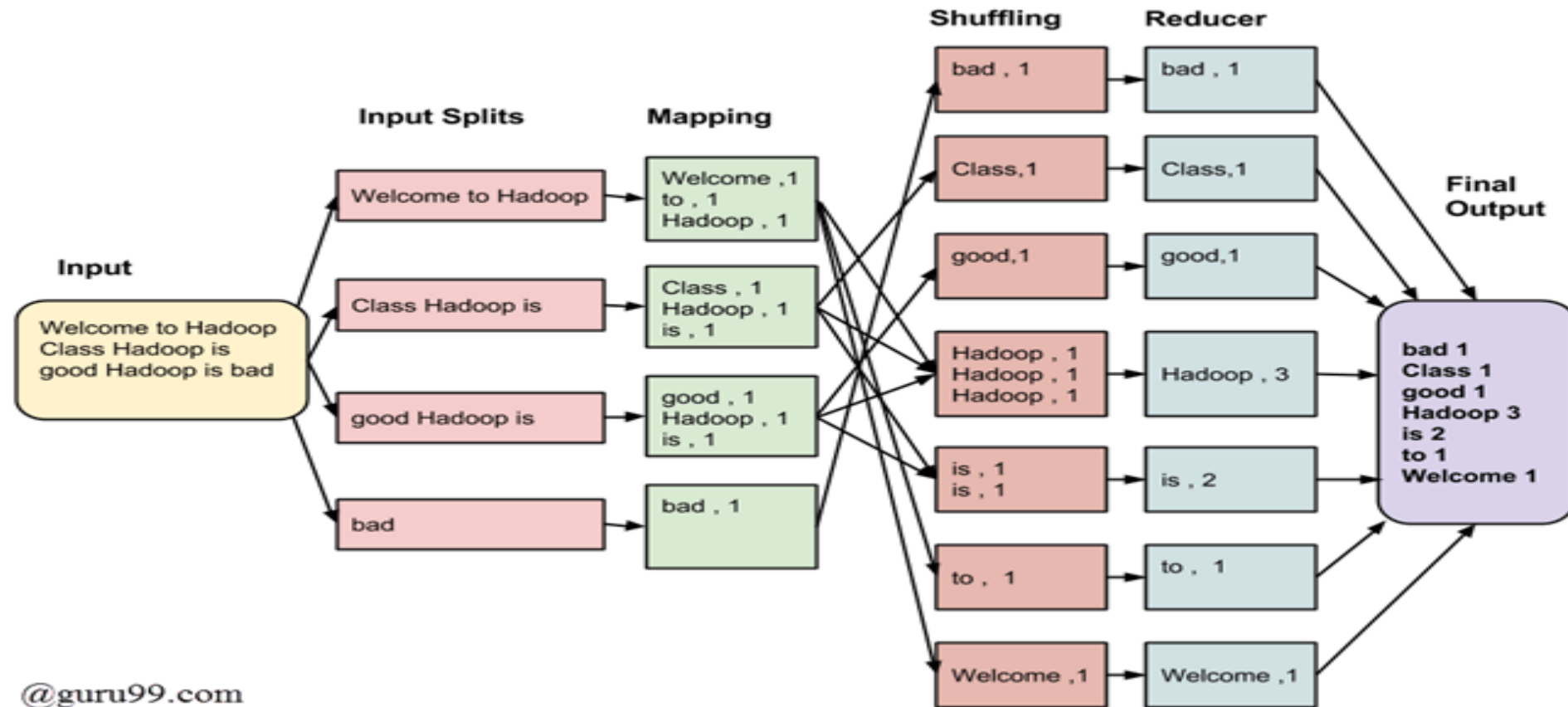
ResourceManager

- Receives the processing requests
- Passes the parts of requests to corresponding NodeManagers

NodeManagers

- Installed on every DataNode
- Responsible for execution of task on every single DataNode

MAP REDUCE



HDINSIGHT A WHAT?





Azure HDInsight



- Managed, full-spectrum, open-source analytics service in the cloud.
- Cloud distribution of Hadoop components.
- Core Components of Hadoop like HDFS, YARN, MapReduce are used to analyze batch data in Hadoop cluster of HDInsight.
- Integration with open-source frameworks such as Hadoop, Spark, Hive, LLAP, Kafka, Storm, R, and more supporting multiple programming languages.

FEATURES OF HDINSIGHT



Reduced
Costs

○



Secure &
Compliant



Highly
scalable



Optimized
Components

HDINSIGHT



ALL HADOOP USERS

memegenerator.net

HDInsight AZURE HDINSIGHT Advantages Over Hadoop

- Low Cost
- Automated Cluster Creation
- Monitoring
- Managed Hardware and Configuration
- Global Availability
- Security and Compliance
- Integration with other Azure Services like cosmos Db, blob storage, data factory
- Higher productive by using rich tools with Hadoop, Spark and Machine Learning

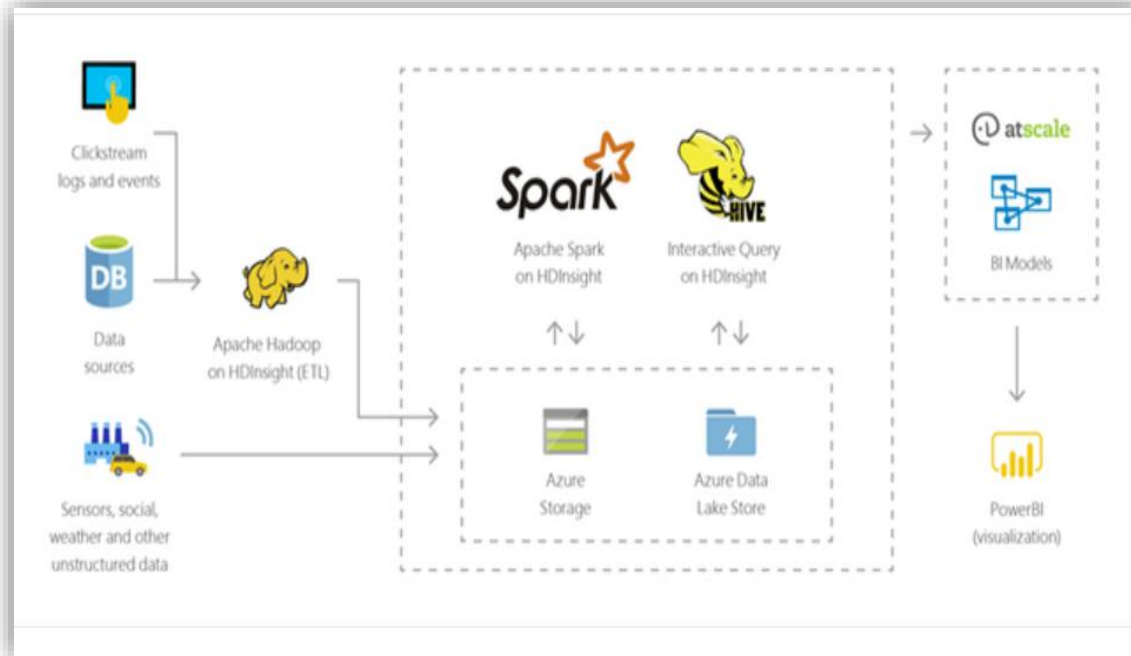
Cluster Type**Functionality**

Hadoop	Batch Query
HBase	Processing Schema free Data
Kafka	Streaming Platform
ML Services	Predictive Modelling
Spark	In-memory processing
Storm	Real-time event processing
Interactive Query	In-memory caching

Scenarios for using HDInsight

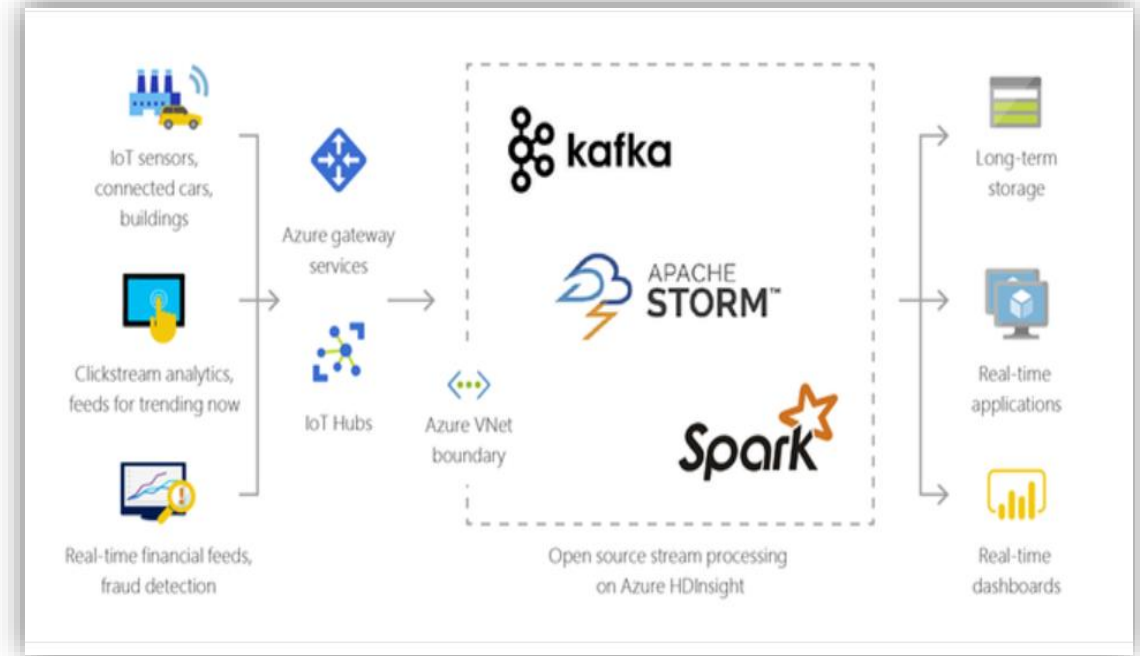
Batch Processing and data warehousing

- ETL at petabytes scale
- Interactive queries and models using BI capabilities



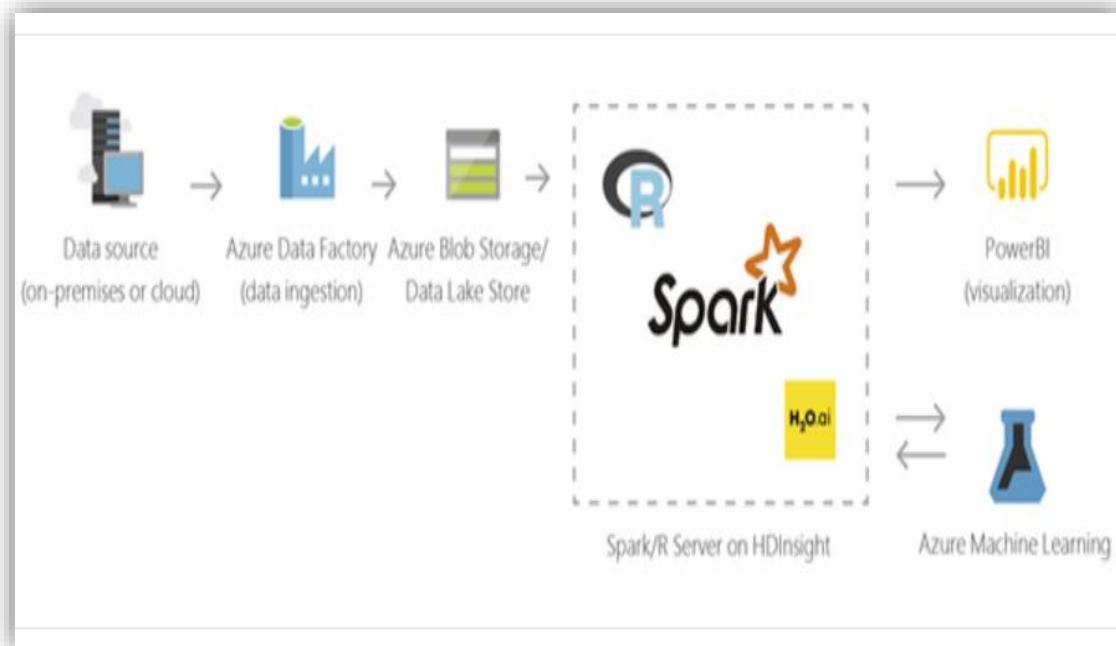
IoT

- Process streaming data in real-time received from different devices



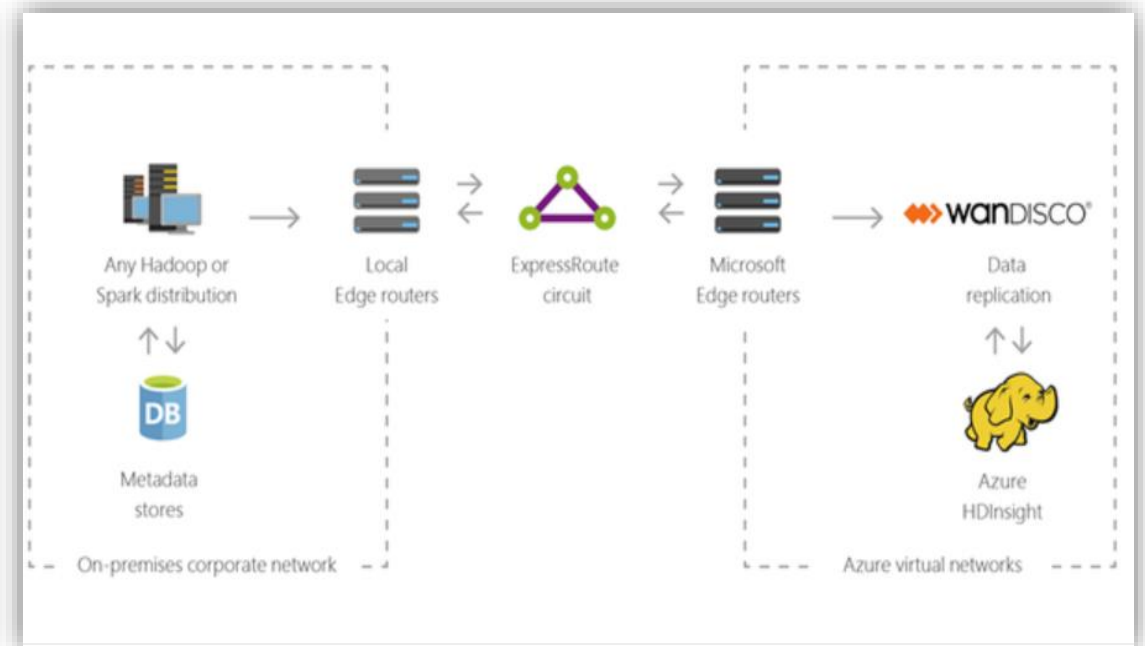
Data Science

- Extract critical insights from data
- Use Azure Machine learning to predict future trends



Hybrid

- Extend on-premise big data infrastructure to Azure to leverage cloud capabilities



-



[Home](#) > [Create a resource](#) > [Marketplace](#) >

Microsoft



Microsoft

♥ Add to Favorites



Plans

Reviews

With HDInsight, you can seamlessly process data of all types through Microsoft's modern data platform. Our platform provides simplicity, ease of management, and an open Enterprise-ready Big Data solution. HDInsight provides a platform for all of your Big Data needs including Batch, Interactive, No SQL and Streaming. It also comes with a strong eco-system of tools and developer environment.

Media

[See All](#)

Microsoft

Azure Service

Securely and Reliably update your devices with Device Update for IoT Hub.



Microsoft

Azure Service

Azure Front Door Standard/Premium (Preview) is security led, modern cloud CDN that provides static and



Microsoft

Azure Service

Azure VMware Solution (AVS) combines the VMware Software Defined Data Center (SDDC) with



Microsoft

Azure Service

scalable RESTful API with enterprise grade security, simple access control and auto SDK generation



DASHBOARD VIEW

21

Microsoft Azure (Preview)

Report a bug

Search resources, services, and docs (G+)

nitigupta@microsoft.com

MICROSOFT

Home > HDInsight_2022-02-07T13.17.04.138Z >

niti-hdinsight-spark

HDInsight cluster

Search (Ctrl+)

Delete

Refresh

Feedback

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Settings

Cluster size

Quota limits

SSH + Cluster login

Data Lake Storage Gen1

Storage accounts

Applications

Script actions

External metastores

Properties

Locks

Monitoring

Insights (Preview)

Alerts

Metrics

Diagnostic settings

Monitor integration

Logs (Preview)

Workbooks (Preview)

Automation

Tasks (preview)

Export template

Essentials

Resource group (move) : [Niti-TA-v2](#)

Status : Running

Location : East US

Subscription (move) : [NPMD Internal Azure Subscription 2](#)

Subscription ID : 3910e970-3d5a-4848-9697-0a8be1da0c8c

Tags (edit) : [Click here to add tags](#)

Learn More : [Documentation](#)

Cluster type, HDI version : Spark 2.4 (HDI 4.0)

URL : <https://niti-hdinsight-spark.azurehdinsight.net>

Cluster ID : fb548f42020449fbb83d05bc354a850e

Overview

Get started

Dashboards

Ambari home

Ambari views

Zeppelin notebook

Jupyter notebook

Spark history server

Yarn

Recommended features

Auto scale

Automatically increase or decrease the number of worker nodes based on a schedule or specific performance metrics.

Applications

Install third party applications.

Script actions

Customize Azure HDInsight clusters by using script actions.

Monitor integration

Monitor Azure HDInsight cluster with Azure Monitor logs.

Cluster size

Node type	Node size	Number of nodes
Head	E8 V3 (8 Cores, 64 GB RAM)	2
Worker	E8 V3 (8 Cores, 64 GB RAM)	1
Zookeeper	A2 v2 (2 Cores, 4 GB RAM)	3

KEY TAKEAWAYS

AZURE HDINSIGHT

- Data is growing tremendously
- Hadoop is an efficient⁺ distributed framework for parallelly processing
- HDInsight is a highly scalable and coherent cloud solution for big data processing and storage
- Seamless setup for azure HDInsight and numerous integration capabilities for moving from on-premise to cloud solution.

-



@nitigupta 16



nitigupta@microsoft.com



<https://www.linkedin.com/in/nitigupta-332000086>

Slides at: <https://github.com/nitigupta16/talks>

Why don't keyboards sleep?

○



@nitigupta_16

BECAUSE IT HAS TWO SHIFTS



Q & A



Thank you

