



The Inside Scoop on Hadoop

Orion Gebremedhin

National Solutions Director – BI & Big Data , Neudesic LLC.
VTSP – Microsoft Corp.

Orion.Gebremedhin@Neudesic.COM

B-orgebr@Microsoft.com

[@OrionGM](#)





The Inside Scoop on Hadoop

Topics Covered

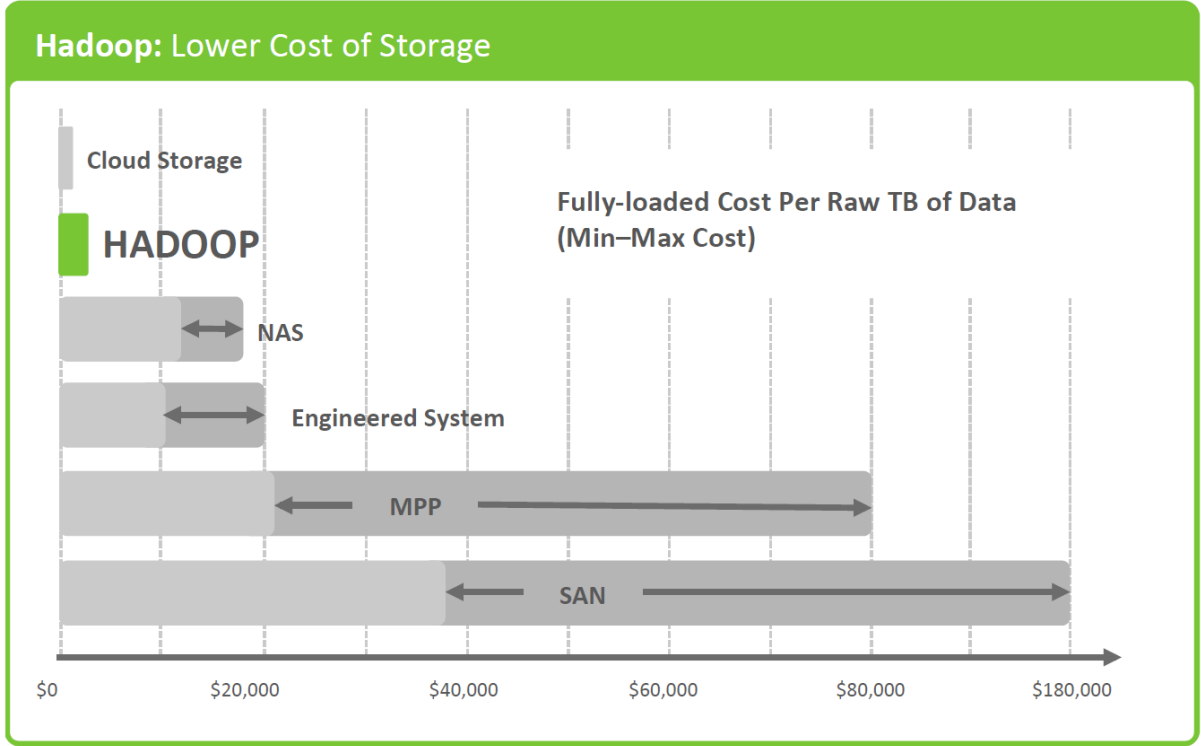


- Understanding Hadoop
- Big Data Solution Deployment Models
- Architecting the Modern Data Warehouse
- Summary



Understanding Hadoop

Big Data = Hadoop?



The Fundamentals of Hadoop

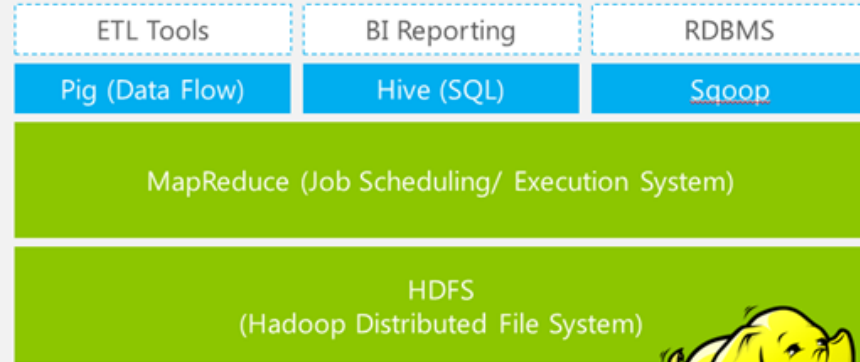


Hadoop evolved directly from commodity scientific supercomputing clusters developed in the 1990s.

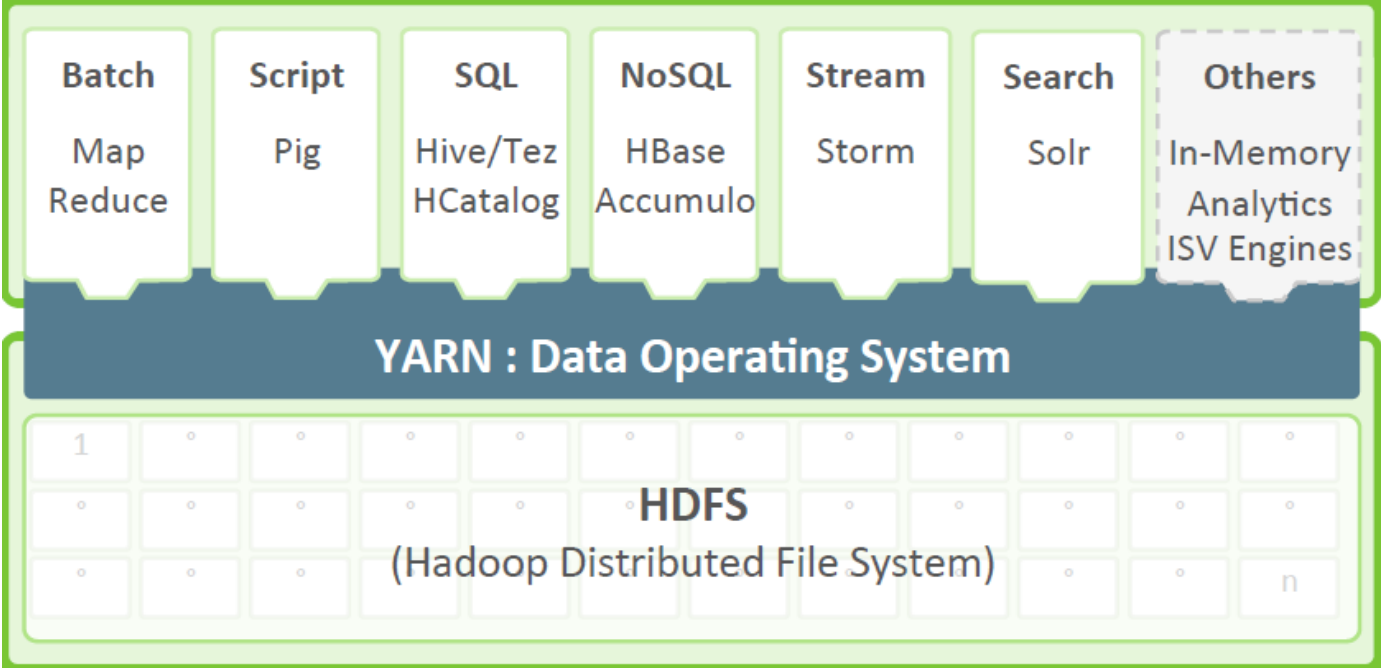
Hadoop consists of a parallel execution framework called

- Map/Reduce and
- Hadoop Distributed File System (HDFS).

The Hadoop Ecosystem



Latest Developments



HDFS



- Very high fault tolerance
- Can not be updated but corrections can be appended
- File blocks are replicated multiple times

Three types nodes:

Name Node (Directory)

Backup Node (checkpoint)

Data Node-actual data

MapReduce



- A programming framework for library and runtime. just like .NET
- **Map Function** - *Take a task and break it down into small tasks*
- **Reduce Function** - *Combine the partial answers and find the combined list*
- **Master (Job Tracker)**
 - Is where you submit a query. Manages the Task Trackers which do the actual Map or Reduce task.
- **Workers (Task Trackers)**
 - Do the work, just as each nodes in the cluster have a data node, they also have a task tracker

Basics of MapReduce



400 bills

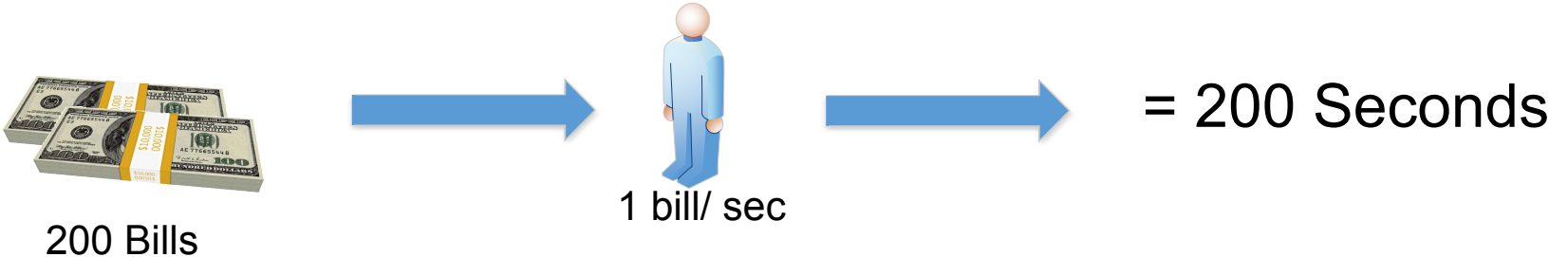
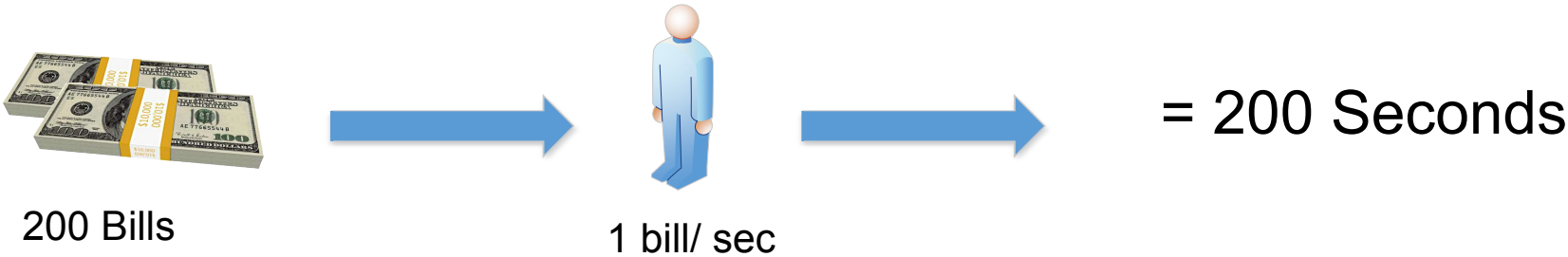


1 bill/ sec



= 400 Seconds

Basics of MapReduce



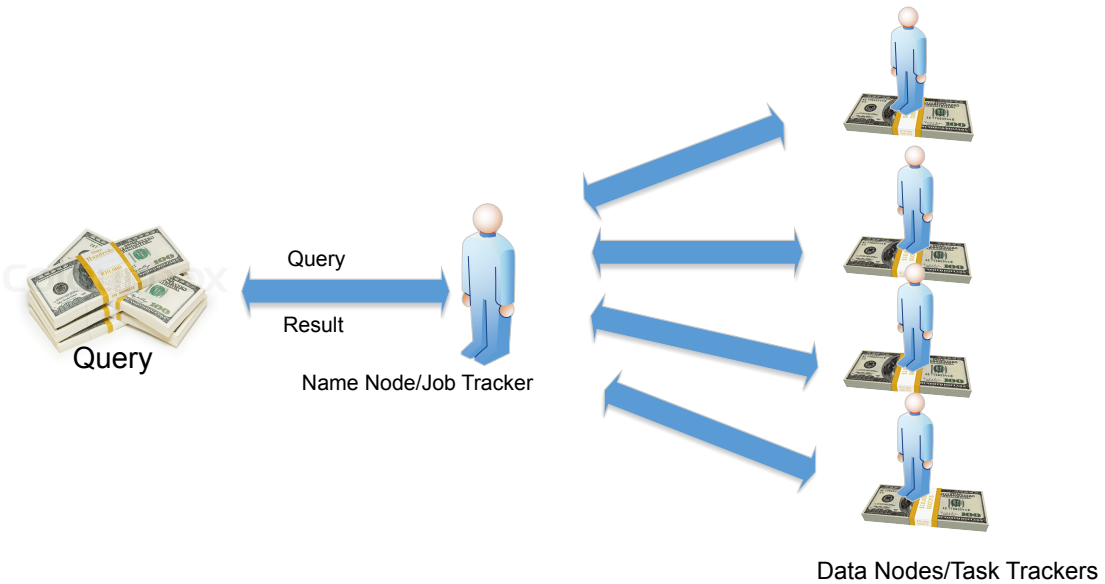
Total = 200 seconds

Basics of MapReduce



Total = 100 seconds

Basics of MapReduce



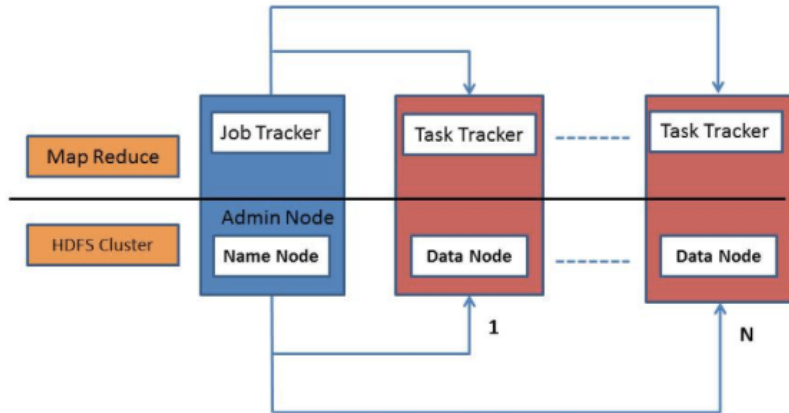
HDFS and MapReduce



The Main Node: runs the Job tracker and The name node controls the files.

Each node runs two processes: Task Tracker and Data Node

- HDFS – Hadoop Distributed File System (storage)
- MapReduce (processing)



Hive and Pig



MapReduce

- Java
- write many lines of code

Pig

- Mostly used by yahoo
- highly used for data processing
- Shares some constructs with SQL e.g. filtering, selecting, grouping, and ordering. But syntax is very different from sql.
- Is more Verbose
- Needs a lot of training for users with limited procedural programming background.
- Gives you more control over the flow of data.

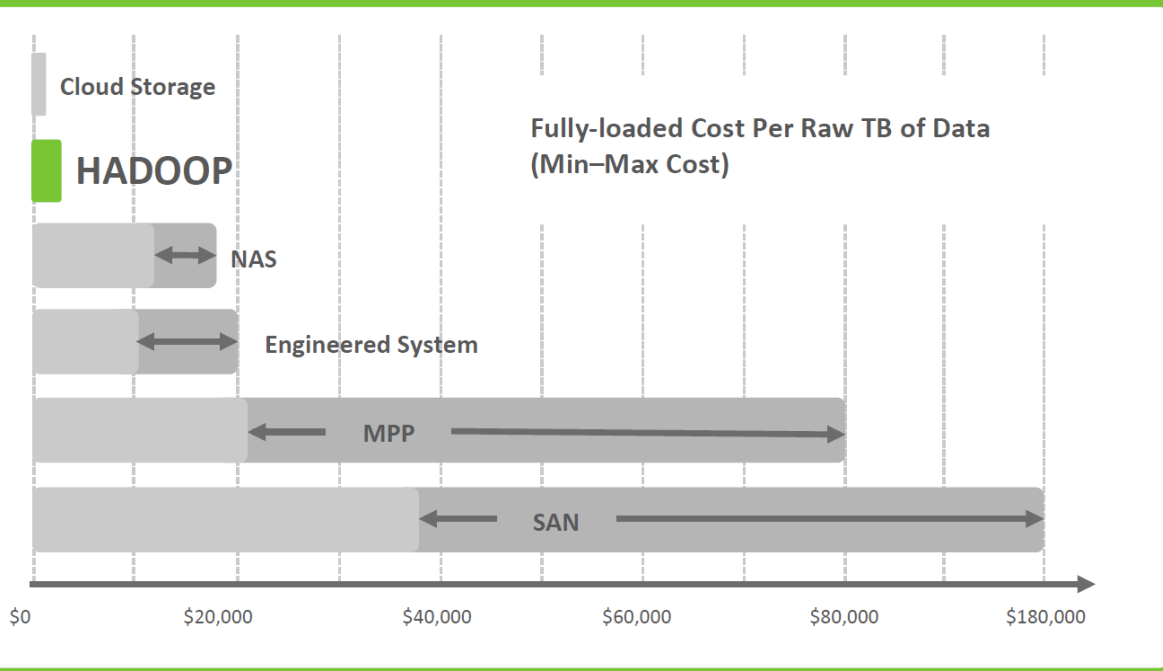
Hive

- Mostly used by Facebook for analytic purposes
- Used for analytics
- Relatively easier for developers with SQL experience.
- Less control over optimization of data flows compared to Pig

- Not as efficient as MapReduce
- Higher productivity for data scientists and developers



Hadoop: Lower Cost of Storage





Big Data Solution Deployment Models

Major Players in Big Data



- **Hortonworks**
- **Cloudera**
- **MapR**
- **Pentaho**
- **Amazon (AWS)**
- ...

Hortonworks



- June 2011 funded by \$23 million from Yahoo! and Benchmark Capital as an independent company
- Horton the Elephant - Horton Hears a Who!
- Employs contributors to project Apache Hadoop
- October 2011 partnered with Microsoft : Azure and Windows Server .
- Cloudera founded in October 2008...started the effort to be Microsoft Azure Certified in October 2014.



HDP User Interface



Navigation bar with icons for various tools and a user profile dropdown labeled 'hue'.

Configuration Check for misconfiguration Server Logs

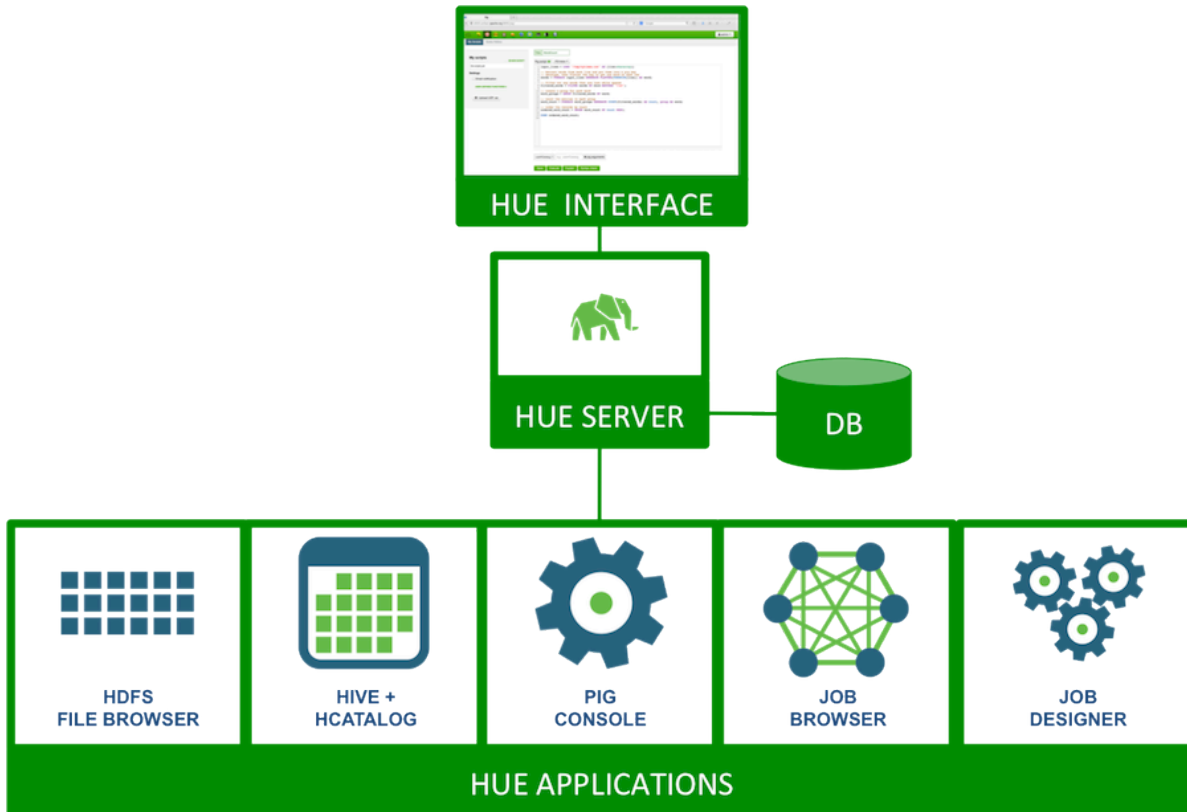
Hortonworks Sandbox 2.1

Leave Feedback

Component	Version	
Tutorials	2.0.005	<input type="button" value="Update"/>
Hue	2.3.1-385	
HDP	2.1.1	
Hadoop	2.4.0	
Pig	0.12.1	
Hive-Hcatalog	0.13.0	
Oozie	4.0.0	
Ambari	1.5.1	<input type="button" value="Enable"/>
HBase	0.98.0	
Knox	0.4.0	
Storm	0.9.1	
Falcon	0.5.0	
Sandbox Build	98e785a 18:26 04-21-14	



Copyright © 2013 The Apache Software Foundation.
Apache Hadoop, Hadoop, HDFS, HBase, Hive, Mahout, Pig, Zookeeper are trademarks of the Apache Software Foundation.
Hue and the Hue logo are trademarks of Cloudera, Inc. and licensed under the Apache 2 license. For more information: gethue.com



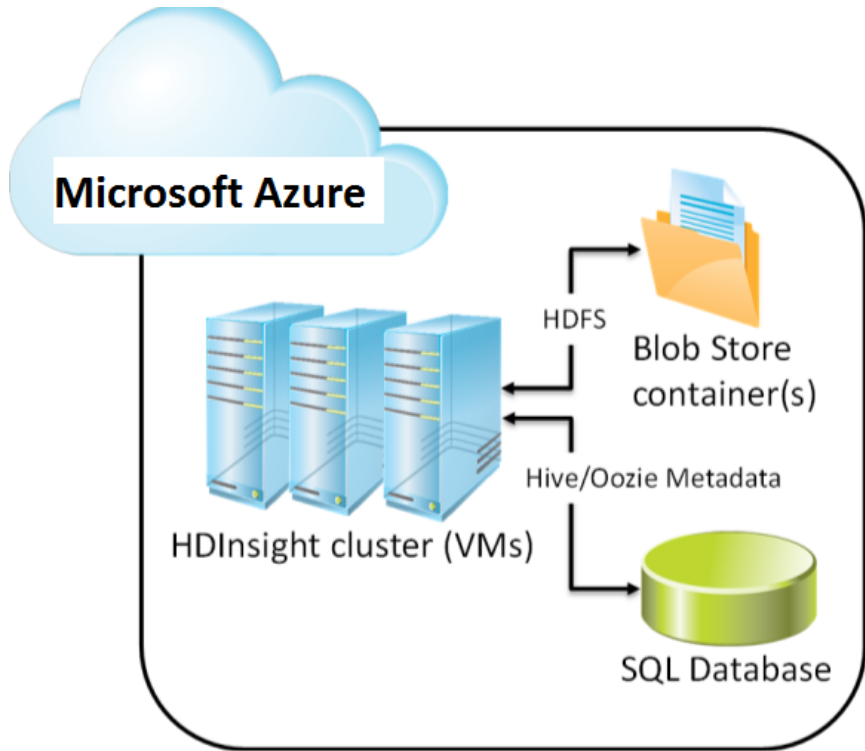
Deployment Models



- **On Premise Deployment**
 - Microsoft Analytics Platform System (APS)
 - Oracle Big Data Appliance
 - Hortonworks Data Platform (HDP)
 - Cloudera's CDH
 - Pivotal Data Computing Appliance (DCA)

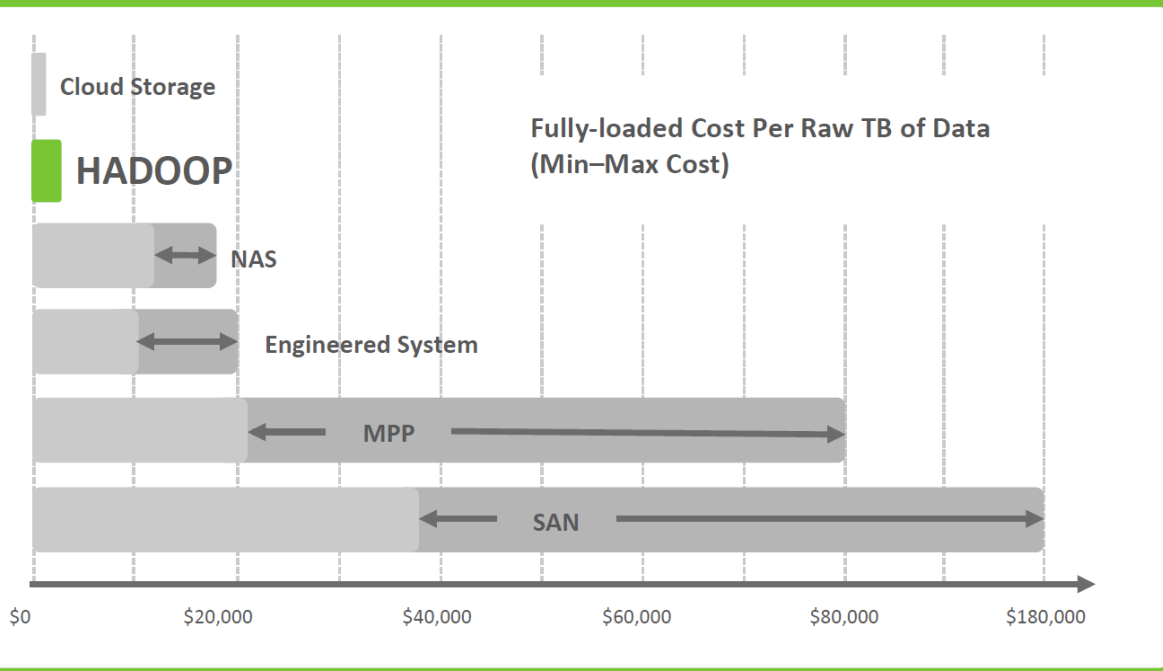
- **Big Data as a service**
 - HDInsight
 - Cloudera on AWS
 - Amazon RedShift
 - Amazon Elastic MapReduce

HDInsight: Hadoop As A Cloud Service





Hadoop: Lower Cost of Storage



HDInsight Versions



Microsoft Azure | Subscriptions | NeudesicHDInsightIncubation@outlo

hdinsight

NAME	STATUS	CLUSTER TYPE	SUBSCRIPTION NAME	LOCATION	VERS
neuhdinsightcluster	HDInsight Cluster Cqueued for Deleti...	Hadoop	Microsoft Azure Sponsorsh...	West US	
FridayLunch	* Accepted	Hadoop	Microsoft Azure Sponsorsh...	West US	

NEW HDINSIGHT CLUSTER

Cluster Details

CLUSTER NAME: *.azurehdinsight.net

CLUSTER TYPE:

HDINSIGHT VERSION [?]

- default (3.1)
- 3.1 (HDP 2.1, Hadoop 2.4)
- 3.0 (HDP 2.0, Hadoop 2.2)

HDInsight Version 3.1
Hortonworks Data Platform Version 2.1
Apache Hadoop Version 2.4

2 3 4 5



Architecting the Modern Data Warehouse

The ETL Automation Model



Hadoop: Data Warehouse Workload Optimization

Current Reality

EDW at capacity: some usage from low value workloads

Older transformed data archived, unavailable for ongoing exploration

Source data often discarded

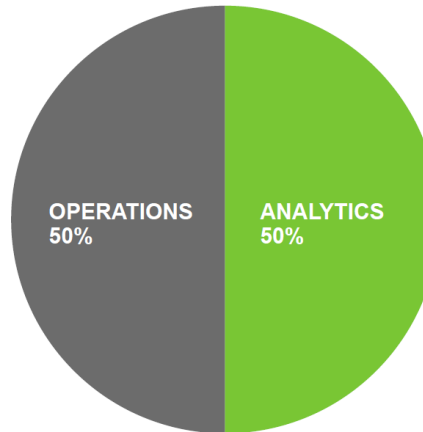
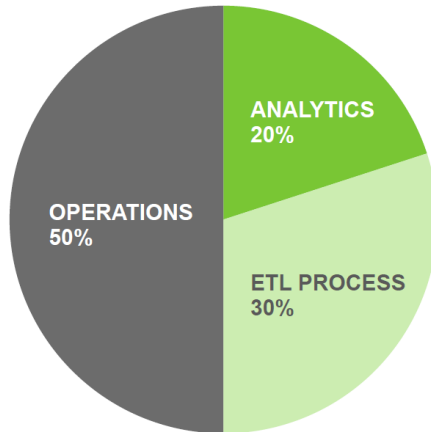


Augment with Hadoop

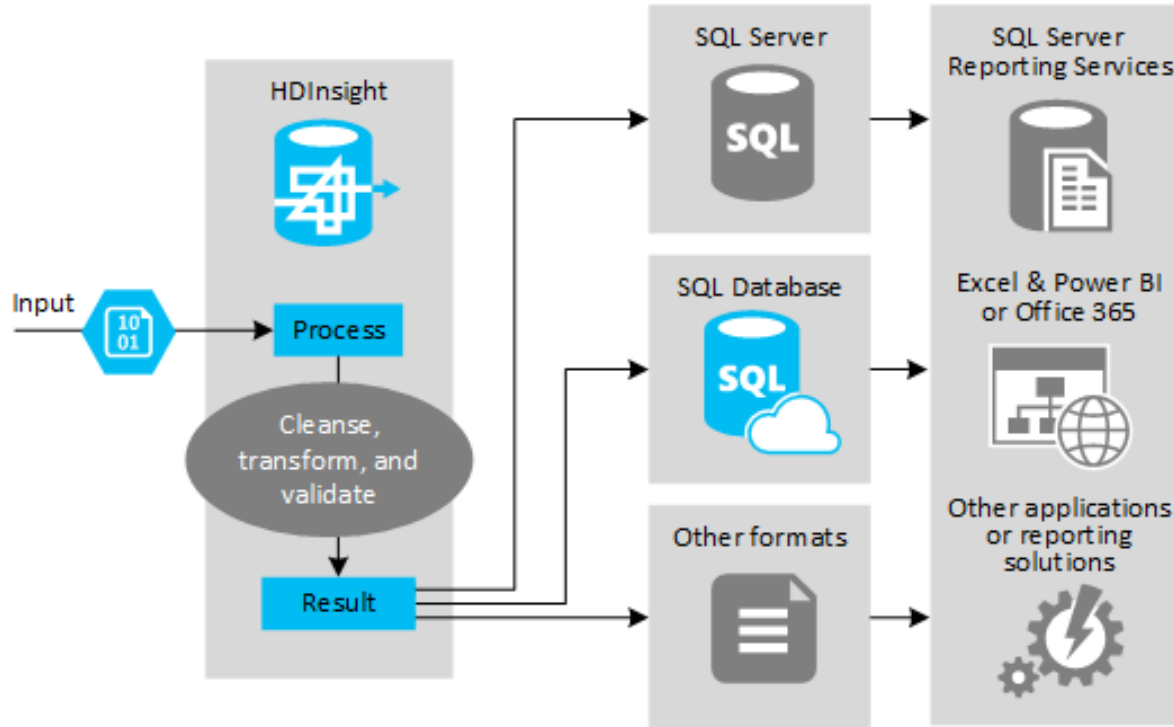
Free up EDW resources from low value tasks

Keep 100% of source data and historical data for ongoing exploration

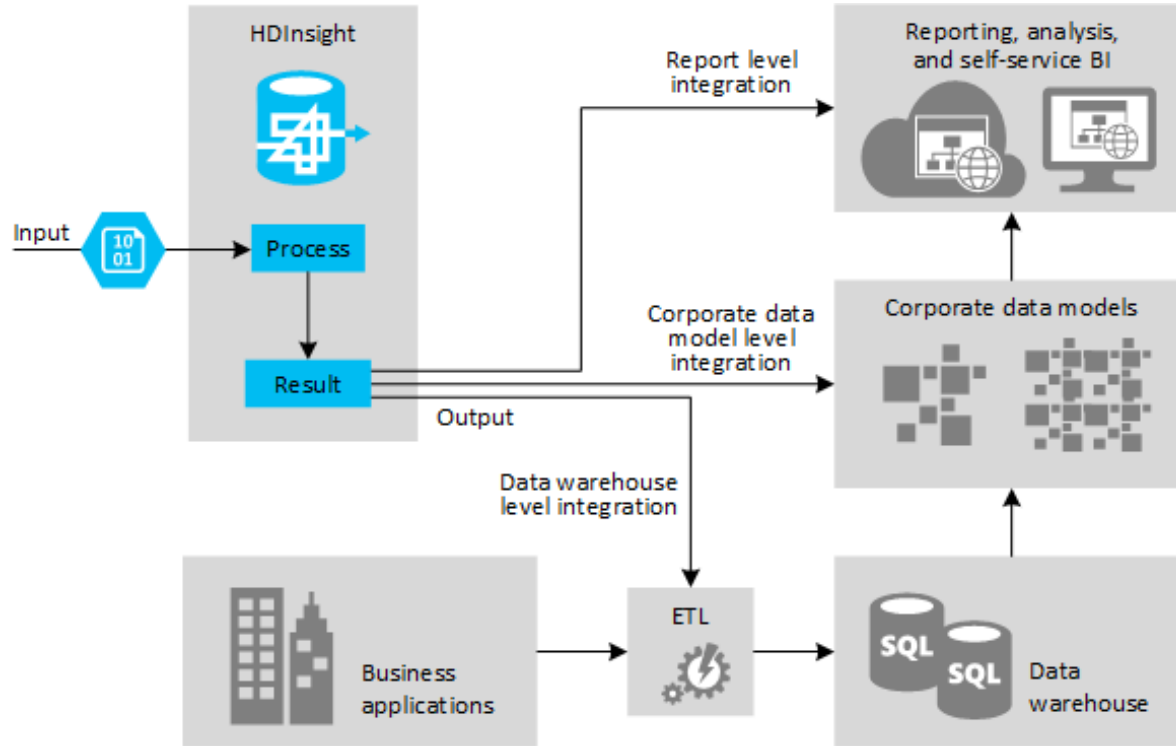
Mine data for value after loading it because of schema-on-read



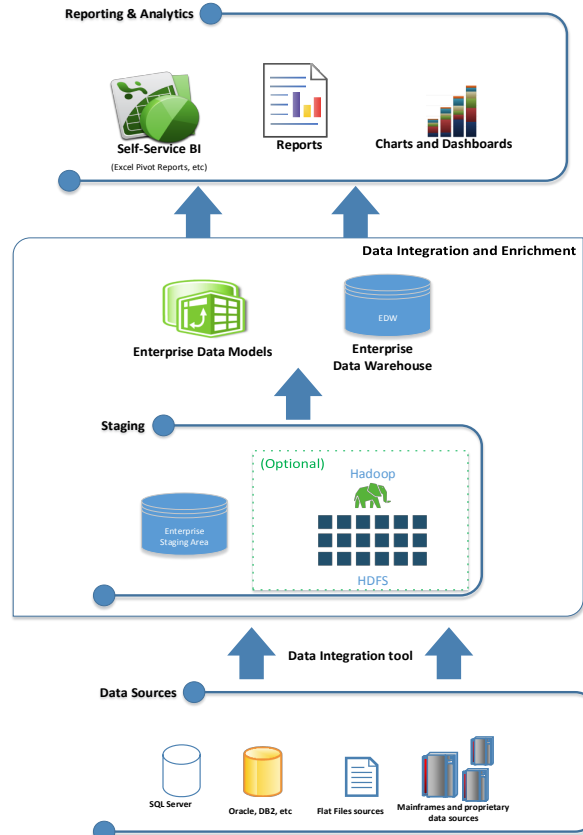
The ETL Automation Model



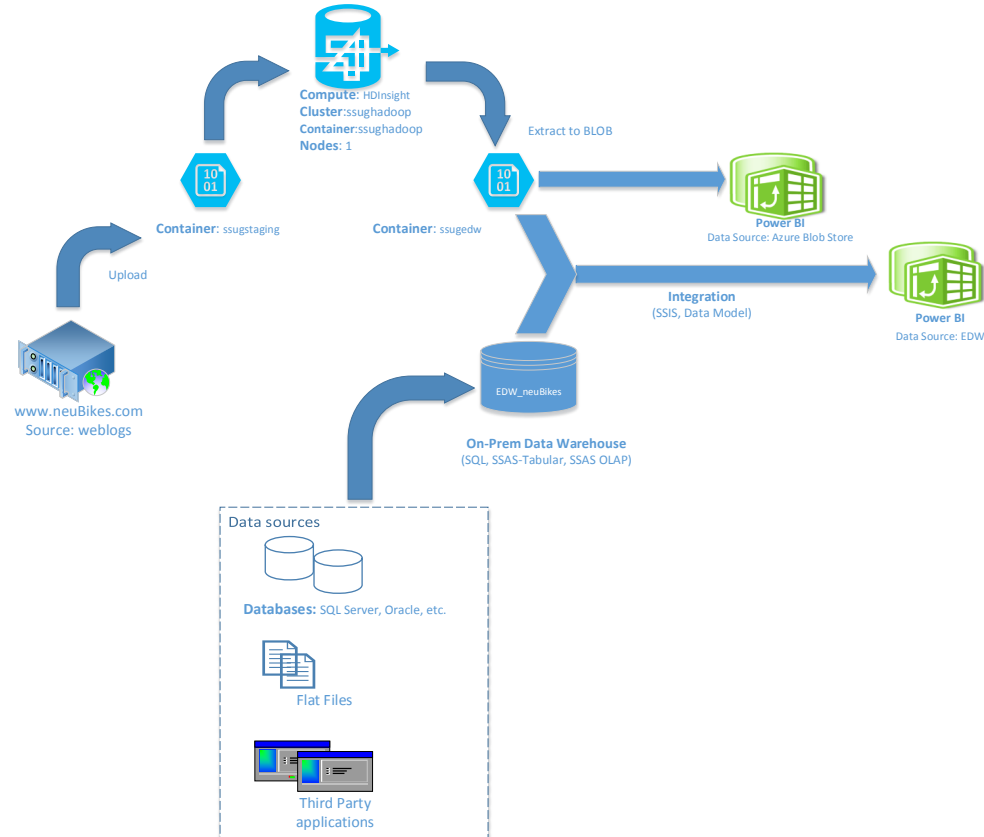
BI-Integration Model



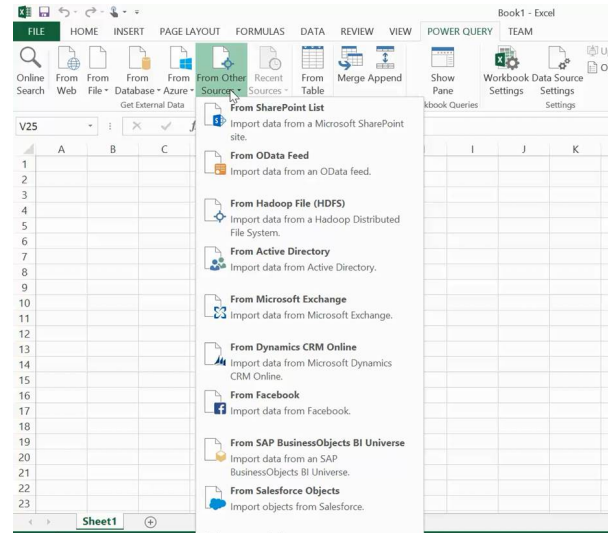
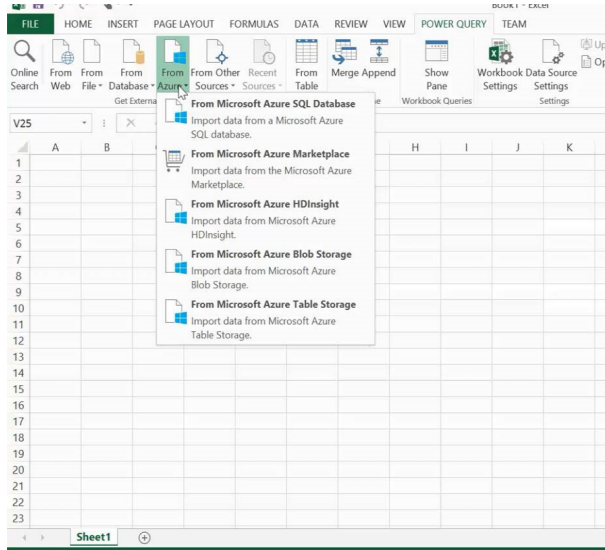
Hybrid BI-Integration Model



Hybrid BI-Integration Model



Hybrid BI-Integration Model



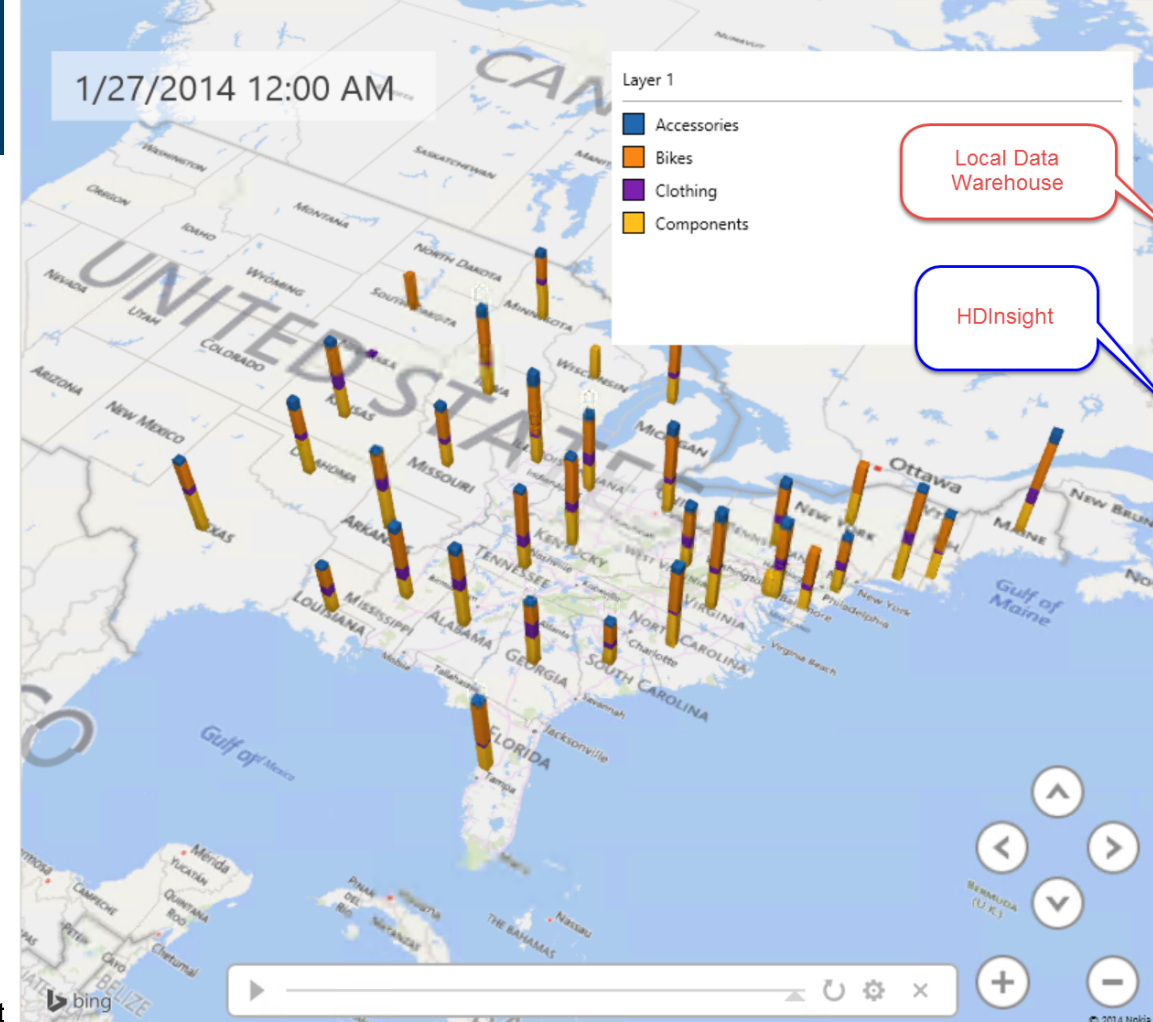
1/27/2014 12:00 AM

Layer 1

- Accessories
- Bikes
- Clothing
- Components

Local Data Warehouse

HDInsight



Layer 1

GEOGRAPHY

Map by State (State/Province) 100%

- DimDate
- DimLocations
- DimProduct
- FactProductWebHits
 - DateSK
 - IP
 - ProductID
- FactWebHitAggregation
 - CountOfWebHits
 - DateSK
 - ProductID

HEIGHT

ProductID (Count - Distinct)

CATEGORY

ProductCategoryName

TIME

FullDateAlternateKey (None)



Summary

Summary



- Understand your data growth to determine when to “Scale-Out”.
- Determine the right tool for the workload you have.
- Choose the right deployment of Big Data Solutions
- Hybridize, do not start from scratch!

Questions?



Questions and Discussion