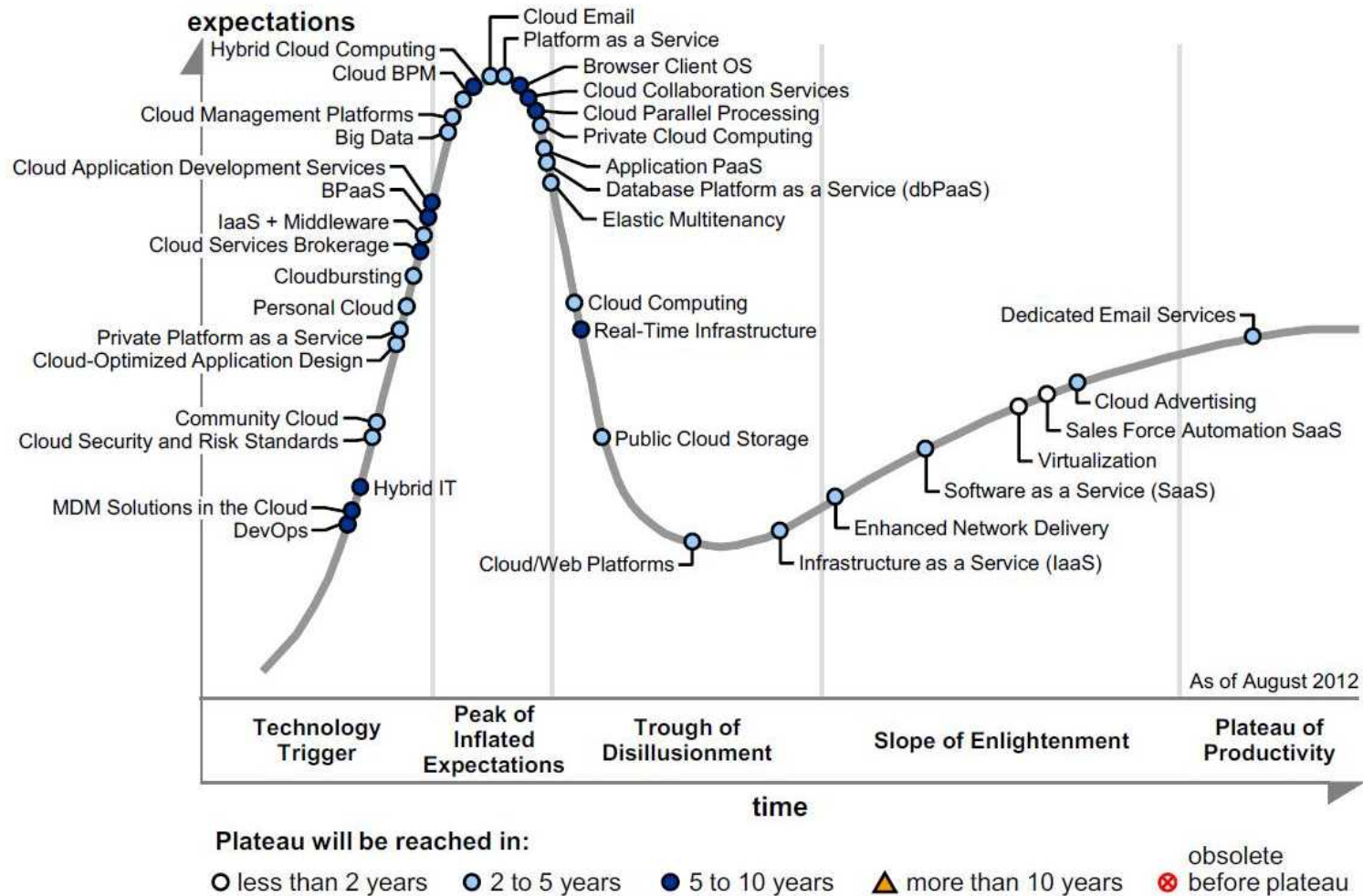# 빅데이터 분석과 상용 클라우드의 결합.

Big Data와 클라우드 컴퓨팅의 만남.

# First Question?

- Why We try to connect the two huge word

  'Big Data' & 'Cloud Computing'

  – Can you define 'Big Data' , 'Cloud Computing'?

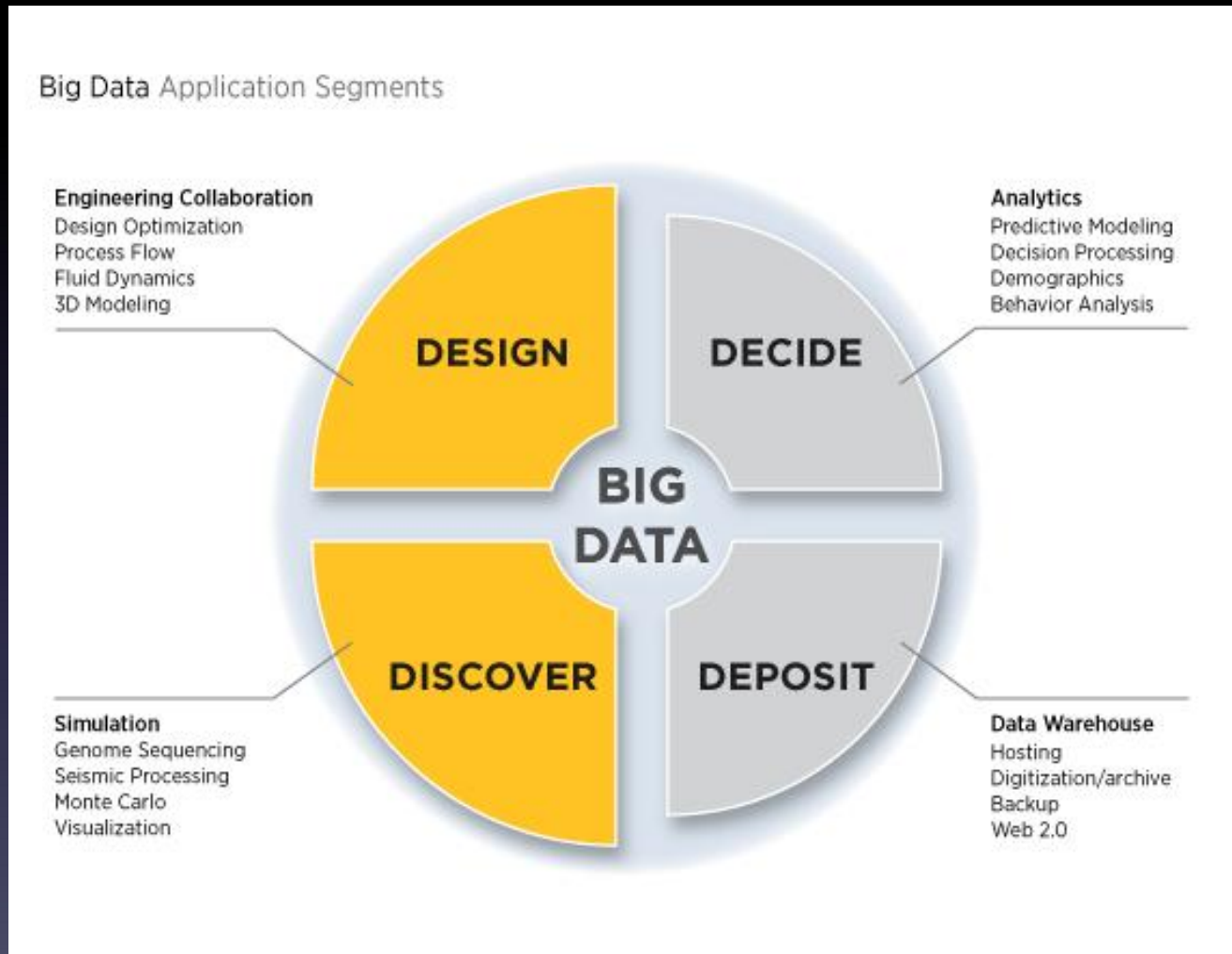  – isn't it just marketing words combining?
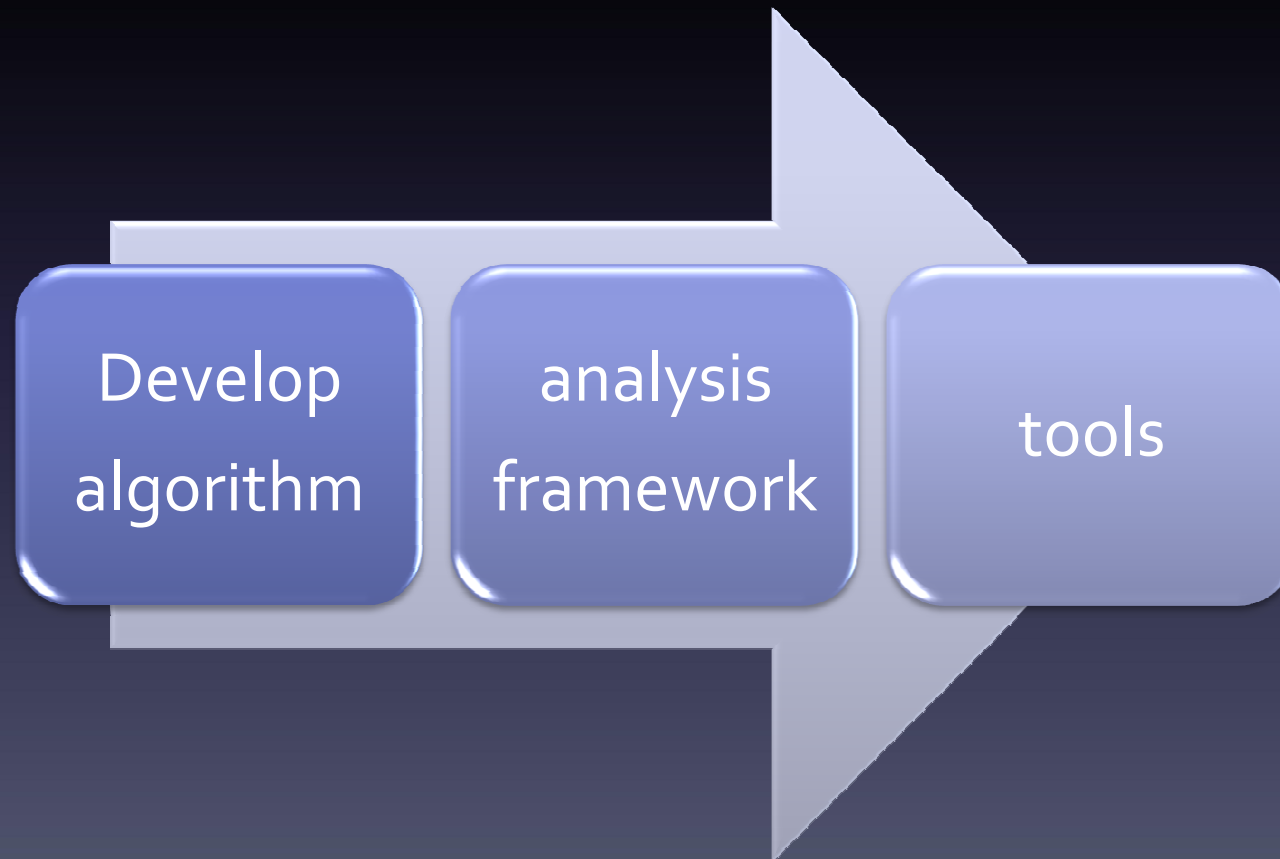
# Big Data hype & Fever?

Figure 1. Hype Cycle for Cloud Computing, 2012

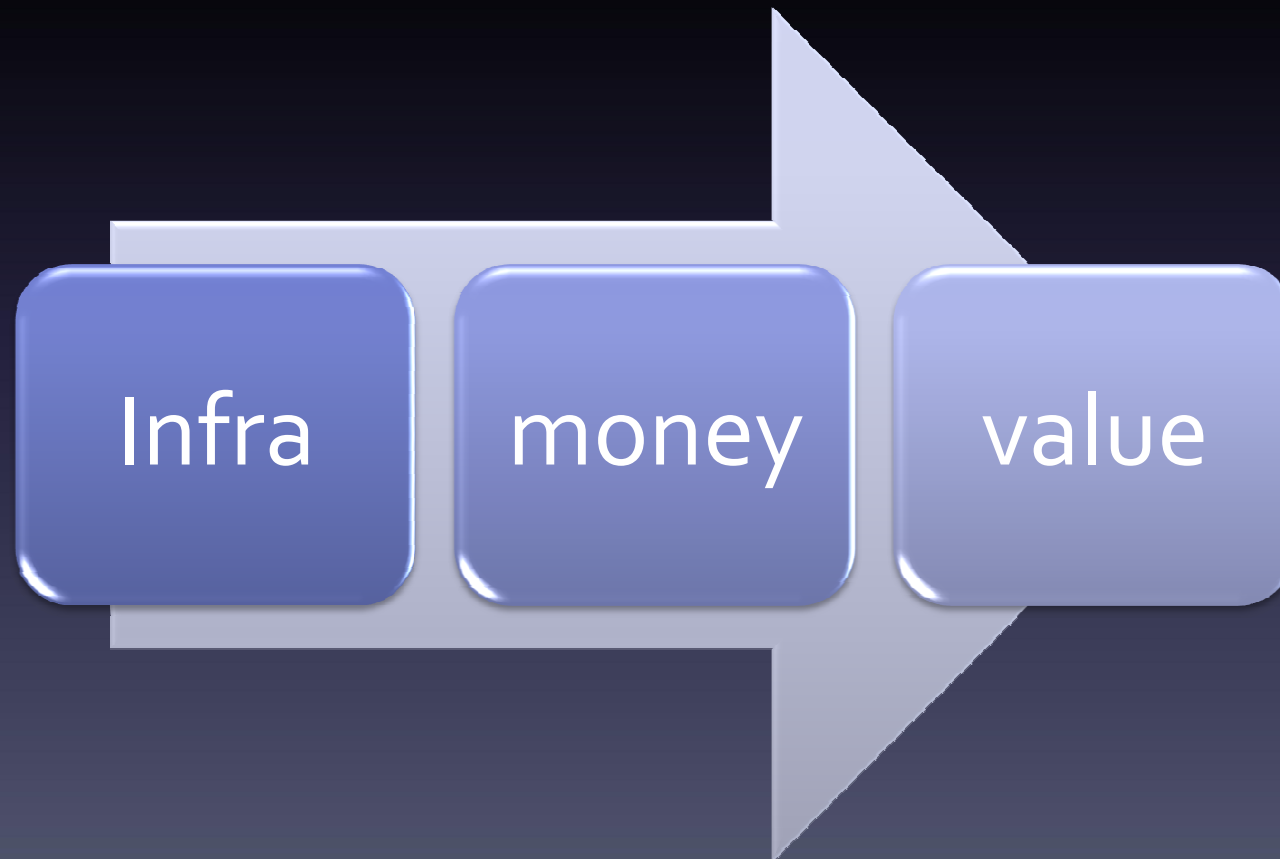# Big Data application segment



source: panasas.com

# Typical Decision for Data analysis.

- Things to decide when POC

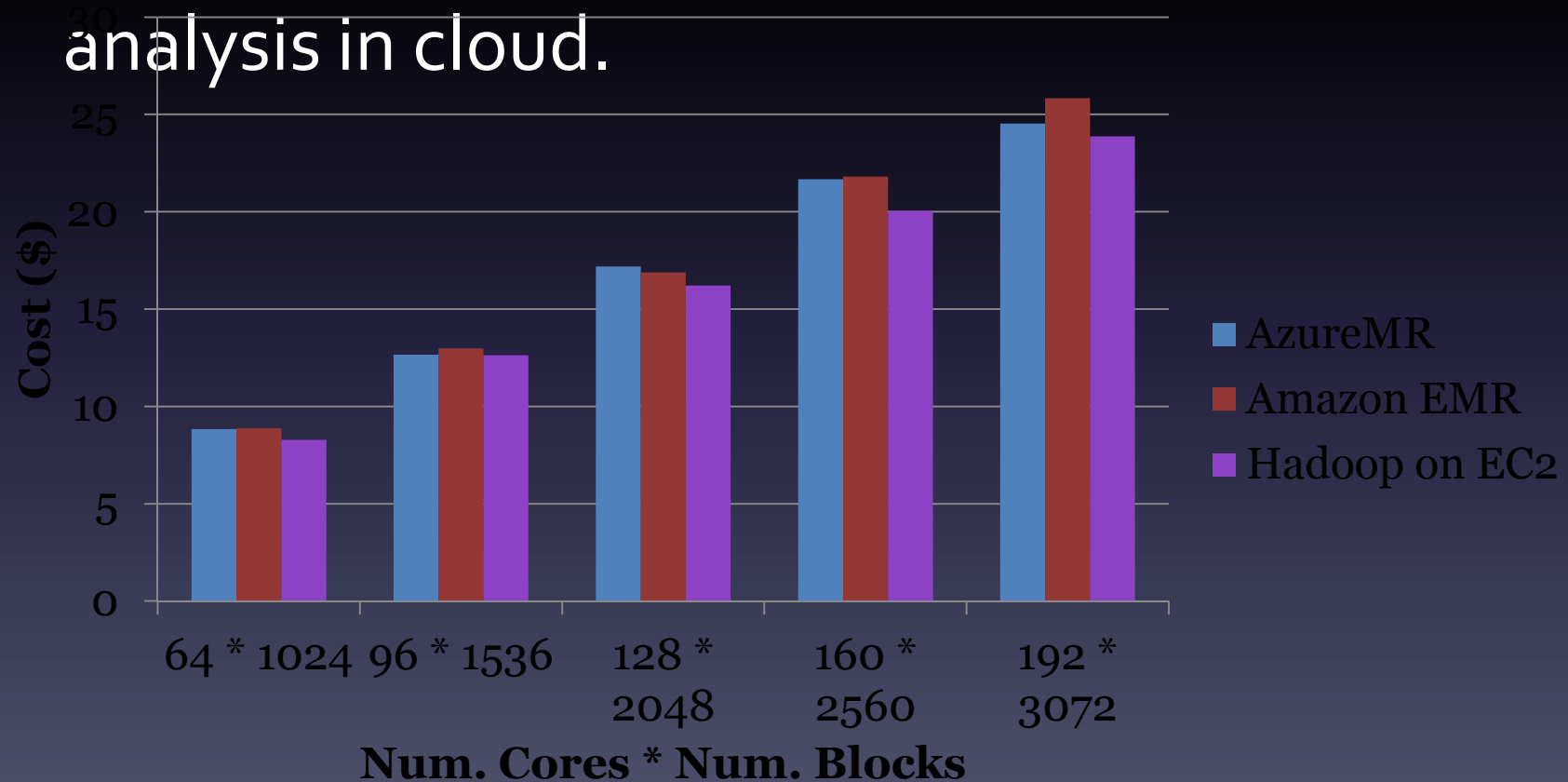# Typical Decision for Data analysis.

- things to decide when goes out.

# The example Case

- One Internet Fax company Case

  – they have 100,000 customer

  – the hosting cost of their machine is

    $24,000~30,000(including network cost)

  – Can you persuade the CEO of this company to invest

    a lot of money to build just 3 node hadoop cluster?

# How much?

- How much When you develop multi-node analysis in cloud.

# Time to implement

| Typical system | period |
|---|---|
| Meeting<br>-the server provider, network provider, IDC manager | 1month |
| Order<br>- Place order, you check importing schedule( always delayed) | 1month |
| Building<br>-installation Rack, Server, Service slide, networking (always missing/broken parts) | 1month |
| Testing<br>-check OS , Network, Disk<br>-install your application | 1month |

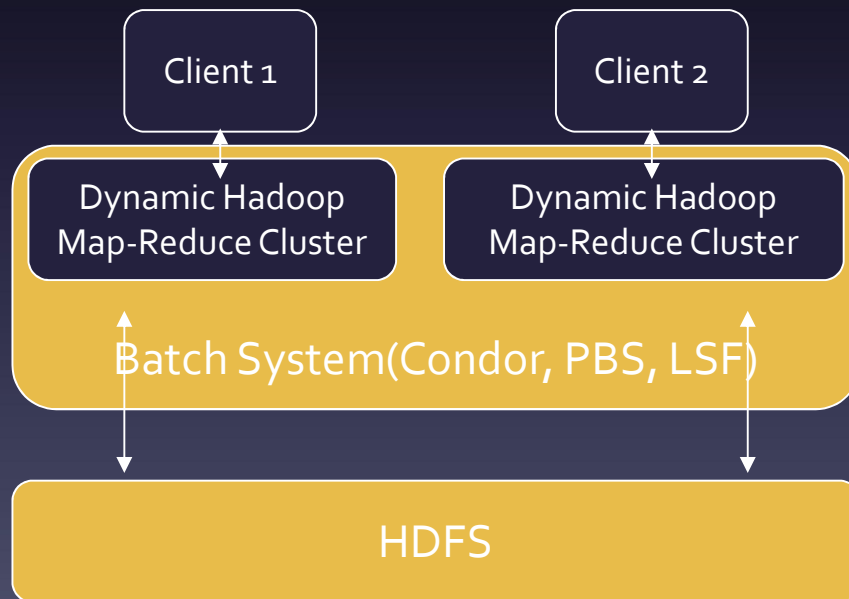| Cloud infra | period |
|---|---|
| Check the web site<br>-punch in your credit card number | 1hour |
| Order<br>- Make 256 virtual instance | 1~2 day |
| Building<br>- Waiting all instances coming up | 1~2 days |
| Testing<br>-install your application<br>-run your test | 1 month |
| Done<br>- Erase the instance | 2~4 hour |

# computing framework to cloud

- MapReduce  Framework (Bare-Metal Style)

  – Hadoop on EC2

  – Hadoop on google compute.

  – Hadoop on Azure instance.

# computing framework to cloud

- Pre-Configured Hadoop
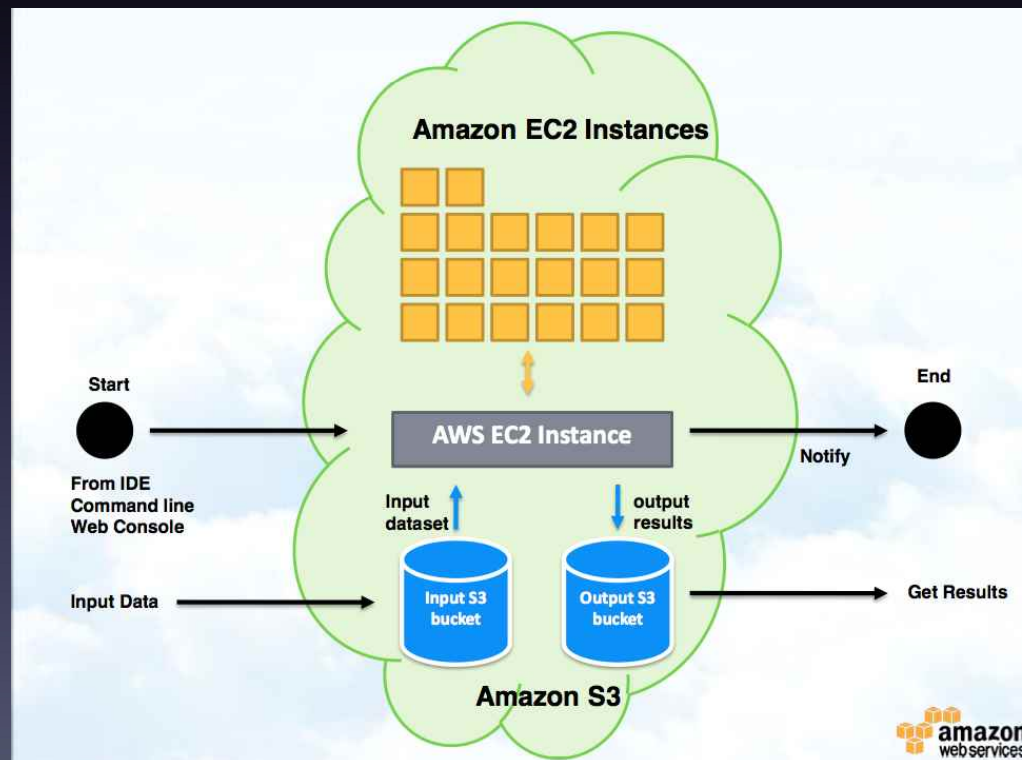
  – the 1st generation: HOD( Hadoop On Demand)

```
Client 1        Client 2
```

Dynamic Hadoop
Map-Reduce Cluster

Dynamic Hadoop
Map-Reduce Cluster

Batch System(Condor, PBS, LSF)

HDFS

Script Sample
hod allocate -d cluster_dir –n 16
hadoop --config ~/hod-
clusters/test jar /hadoop-
examples.jar wordcount /input
/path/to/output

# computing framework to cloud

- Pre-Configured Hadoop

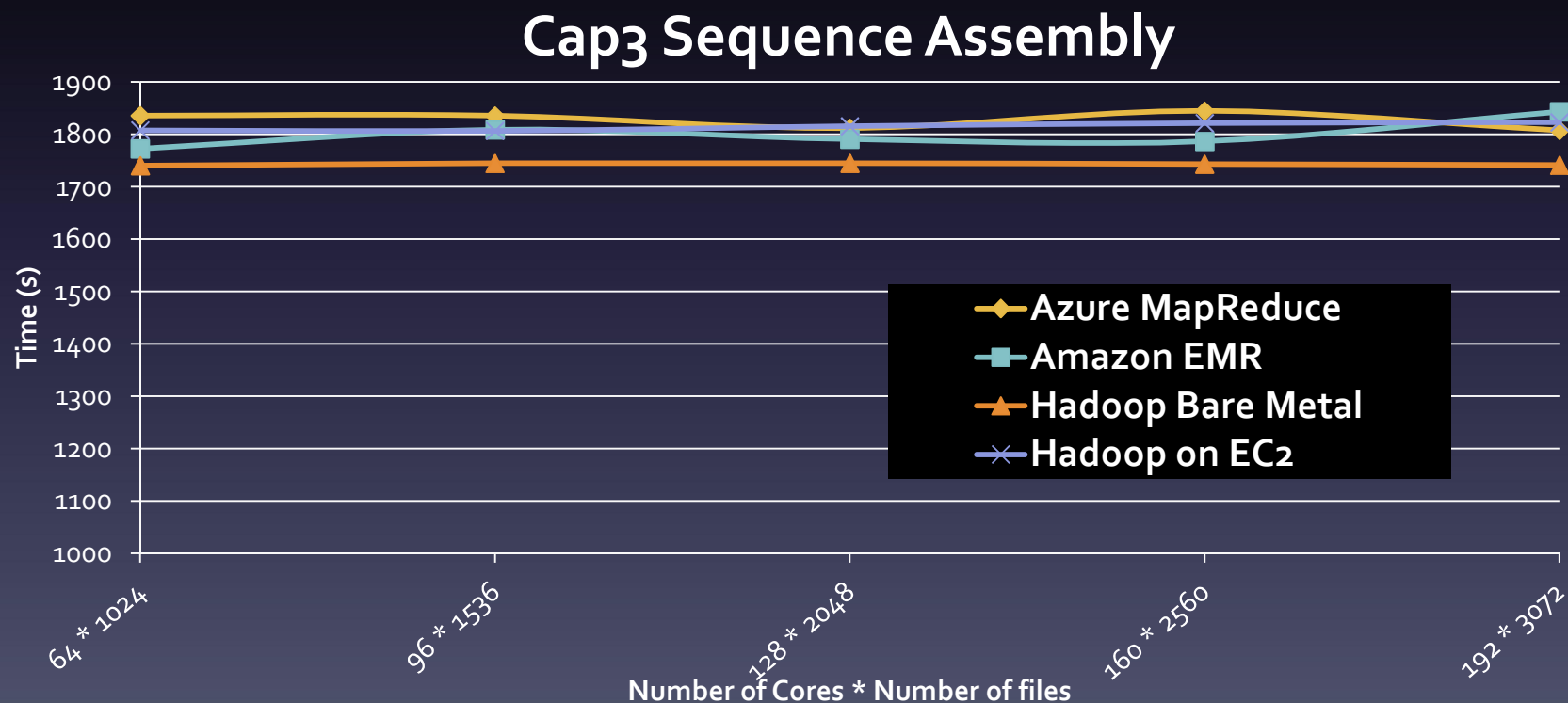  – the 2$^{nd}$ generation: EMR(Elastic MapReduce)

Amazon Elastic MapReduce =

Amazon EC2 + Hadoop



source: AWS

# computing framework to cloud

- Pre-Configured Hadoop

  – the 2$^{nd}$ generation: EMR(Elastic MapReduce)

**Cap3 Sequence Assembly**



Legend: Azure MapReduce, Amazon EMR, Hadoop Bare Metal, Hadoop on EC2

# computing framework to cloud

- Pre-configured hadoop – hadoopOnAzure

## Request a new cluster

**DNS name**

DNS name

mailboxpeak          **Available**

http://mailboxpeak.cloudapp.net

**Cluster size**

| ○ Small | ○ Medium | ◉ Large | ○ Extra large |
|---|---|---|---|
| 4 nodes | 8 nodes | 16 nodes | 32 nodes |
| 2 TB disk space | 4 TB disk space | 8 TB disk space | 16 TB disk space |
| **Available** | **Available** | **Available** | **Available** |

**Cluster login**

Username

campschurmann

Password

•••••••••••

Confirm Password

•••••••••••

**Request cluster**

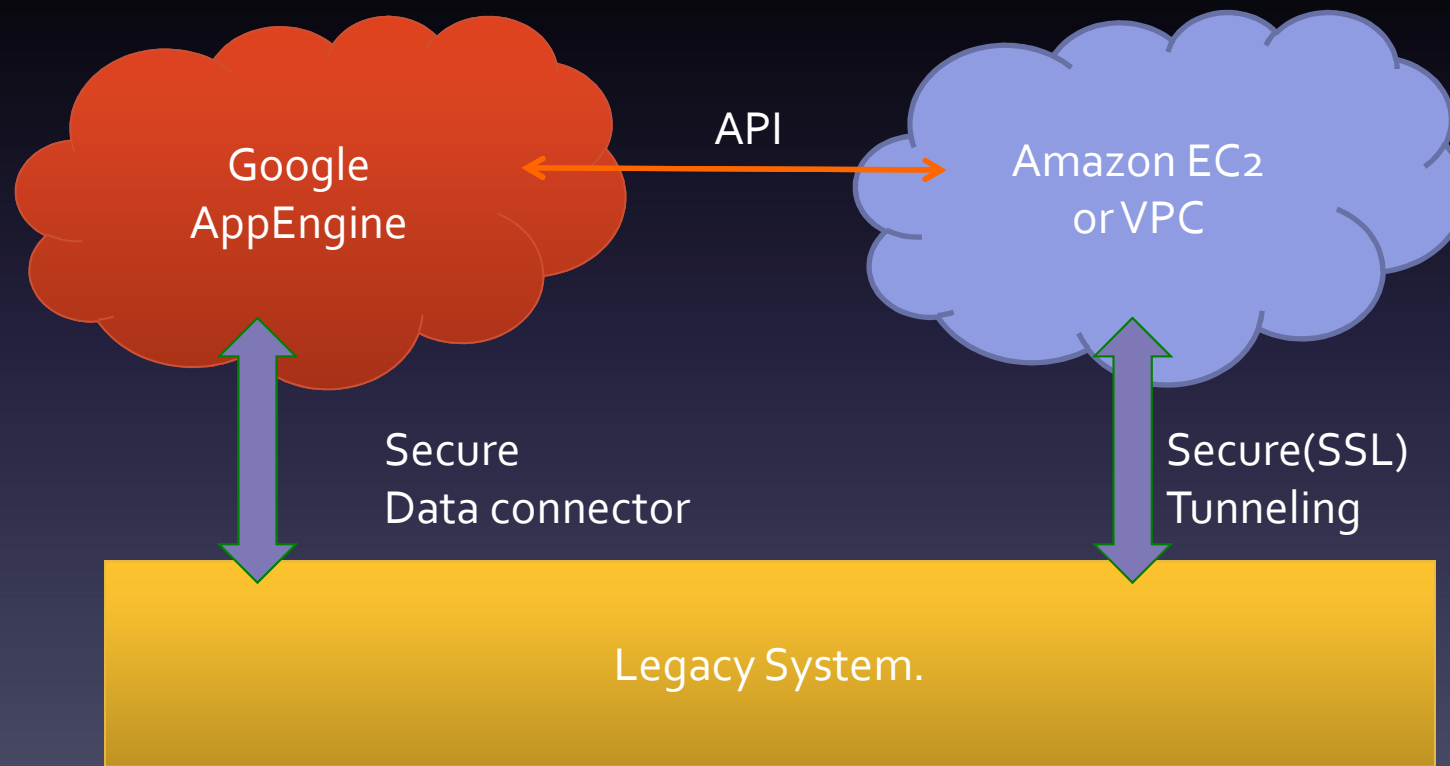# Computing framework to cloud

- EMR - Demo

# Computing framework to cloud

- Cloud only for MapReduce?

  – What about MPI(Message Passing Interface) based Clusters?

  – Most of Commercial Parallel Solutions developed in MPI library.

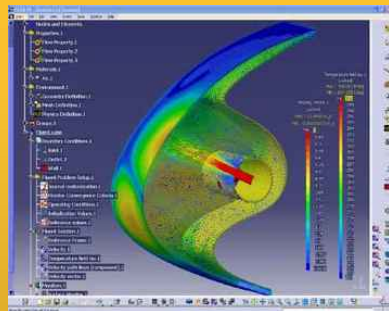# Computing framework to cloud

- Virtual Private Cloud

Google AppEngine

API

Amazon EC2 or VPC

Secure Data connector

Secure(SSL) Tunneling

Legacy System.

# Computing framework to cloud

- Cloudburst of Legacy System.

# Computing framework to cloud

- Cloudburst of Legacy System.
  - No limit in Extending your infra.

Legacy Solution
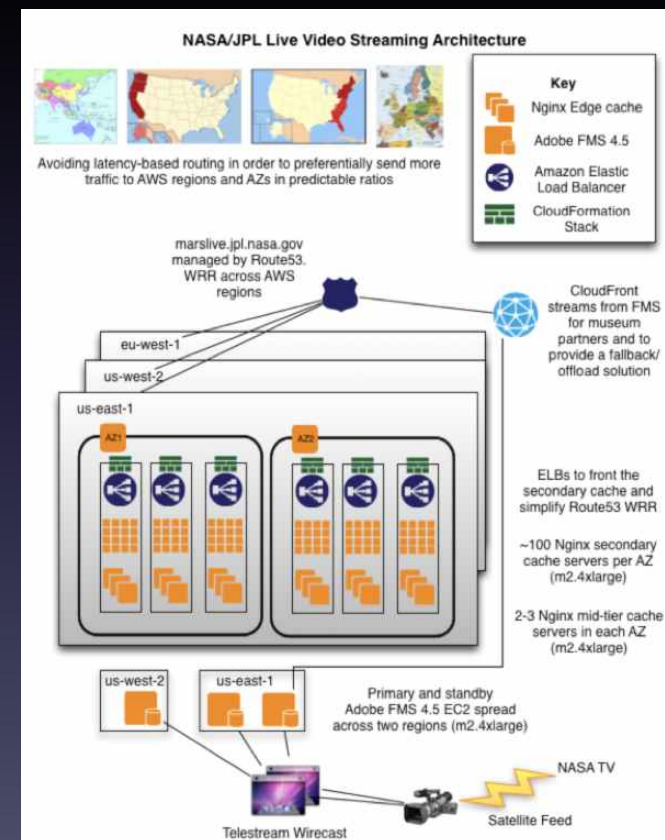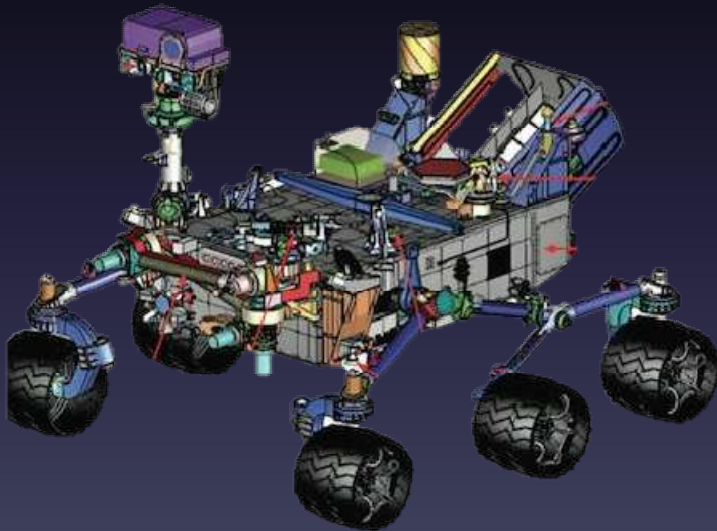
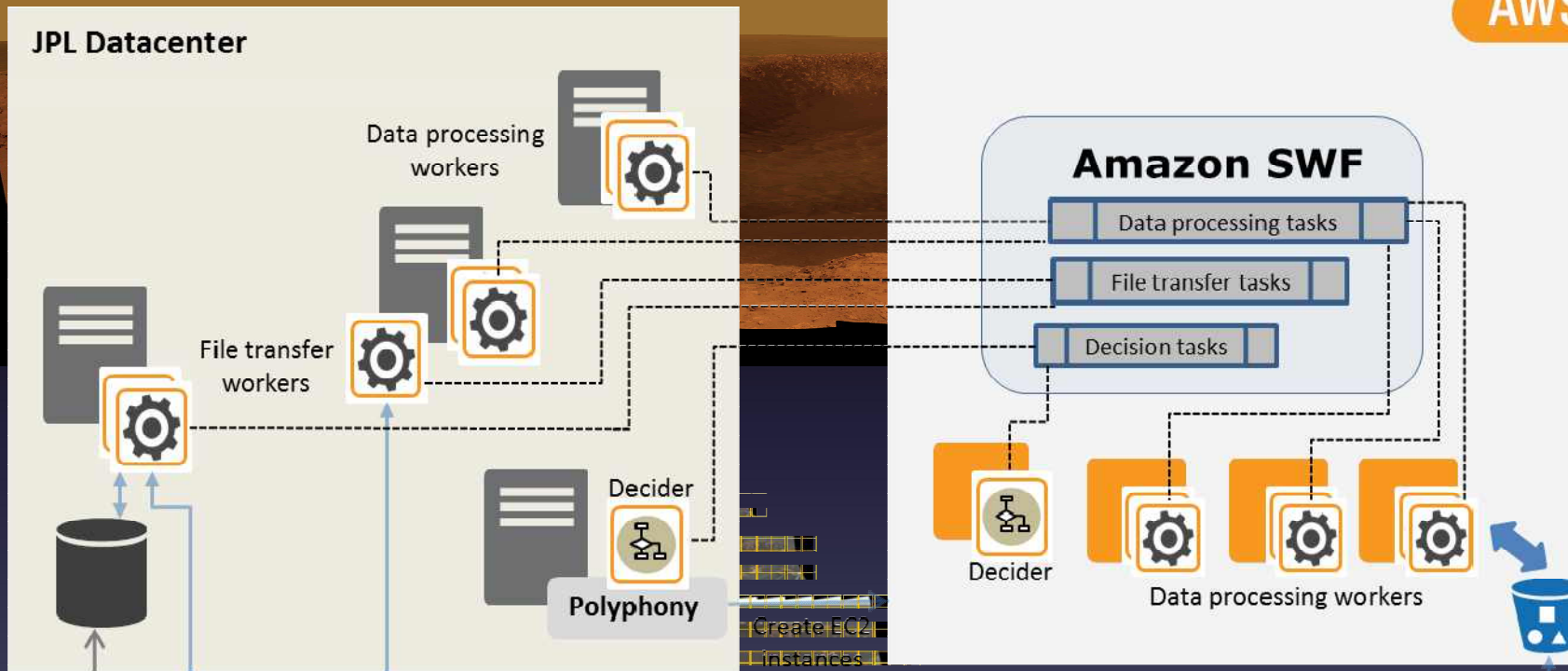Middle Ware (CloudSwitch)
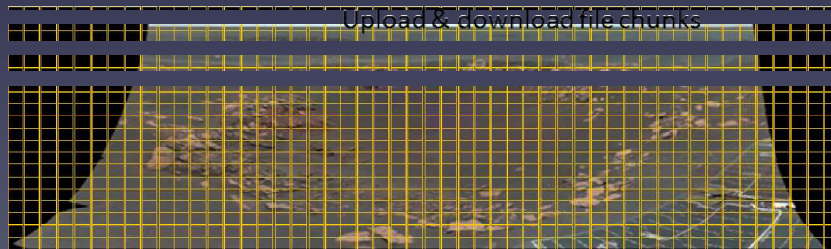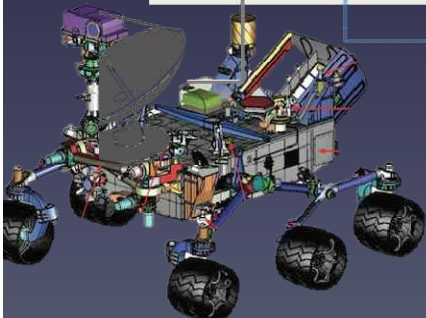
Local Cluster

Virtual Cluster

# Computing framework to Cloud

- Do you Remember the "Curiosity?"



NASA/JPL Live Video Streaming Architecture

# Computing framework to Cloud
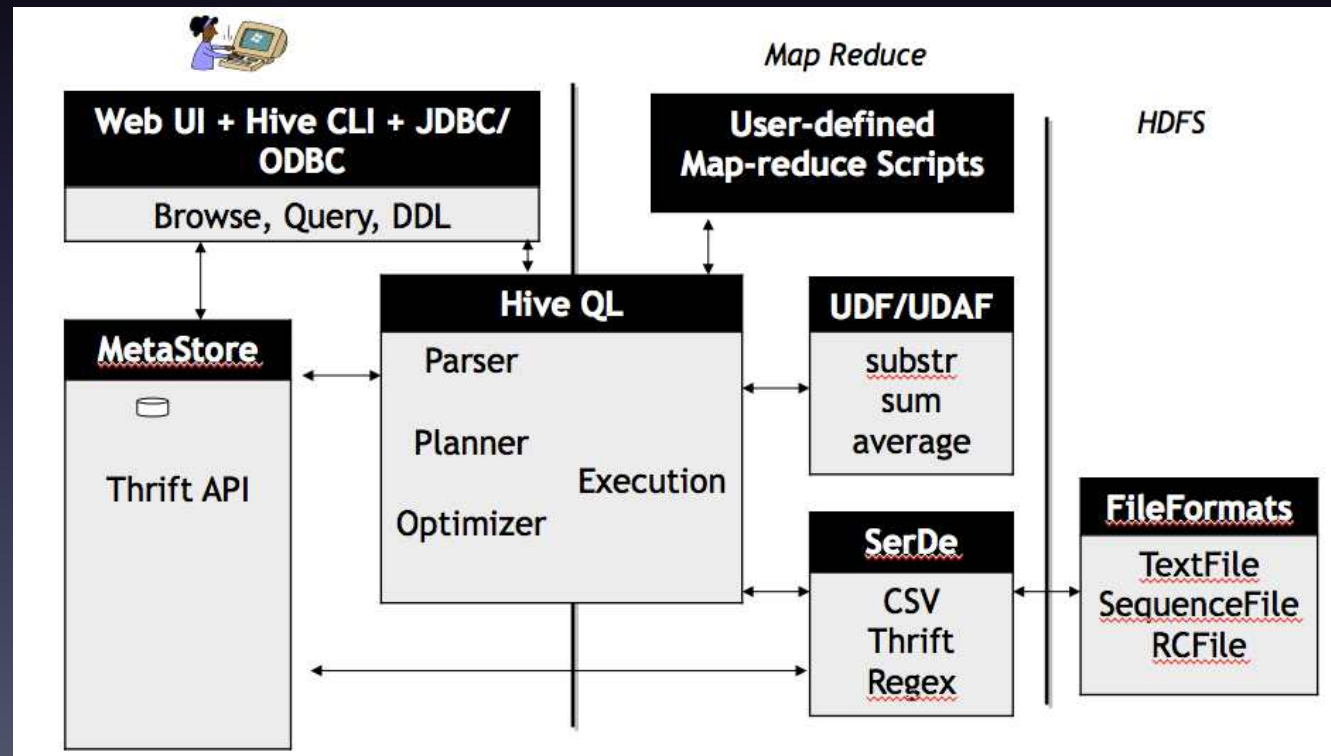
# Computing framework to cloud

- Query based Analytic Engine

    - Hive(hive.apache.org)

# Computing framework to cloud

- Hive on Cloud

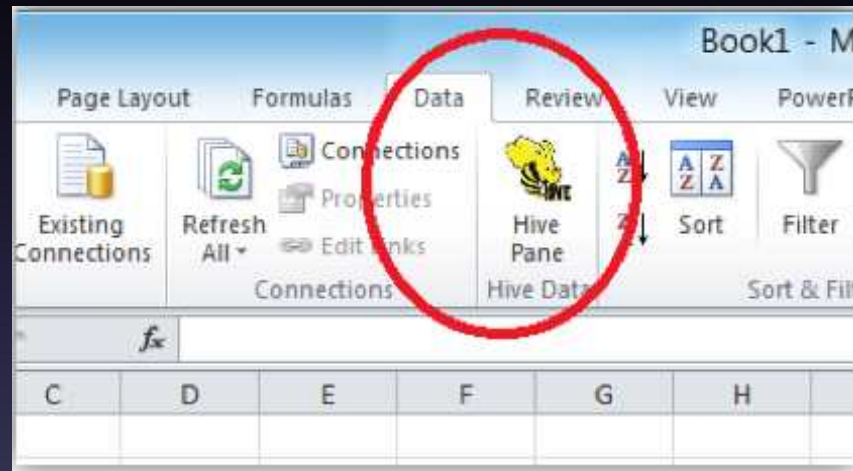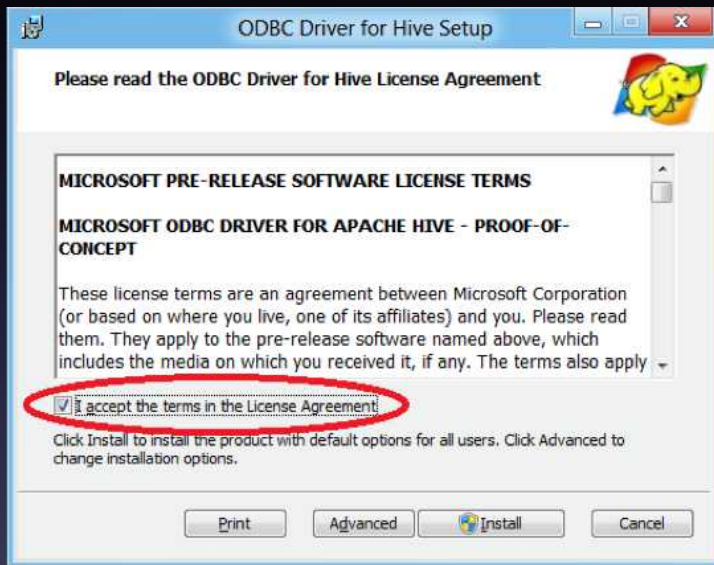  - Hive Job In EMR

  - Interactive mode

```
./elastic-mapreduce --create --name
"${JOB_NAME}"
    --hive-interactive --num-instances
${EMR_INSTANCES_NUM}
    --master-instance
${EMR_INSTANCES_TYPE} --alive
```

  - script mode

```
./elastic-mapreduce --create \
    --hive-script --args ${EMR_SCRIPT_PATH} \
    --args -
d,OUTPUT_PATH=${OUTPUT_LOCATION_S3} \
    --name "${JOB_NAME}" \
    --num-instances ${EMR_INSTANCES_NUM} \
    --instance-type ${EMR_INSTANCES_TYPE} \
    --credentials ${EMR_CREDENTIALS_FILE})
```

# BI(Business Intelligence) with HIVE

- EXCEL(most popular BI)

# BI(Business Intelligence) with HIVE

- Karmasphere BI