# HOW MUCH DOES THIS DRIVE COST?

3 TB

# HOW MUCH DOES THIS DRIVE COST?

That's $0.03 per GB!

| TODAY | $69 |
| --- | --- |
| 2010 | $270 |
| 2005 | $3,720 |

Insight: Disk Space is FREE!

| 1980 | $1,312,500,000 |
| --- | --- |

# IT'S NOT JUST ABOUT COST!

How long does it take to read **3 TB of data**?

3 TB

# IT'S NOT JUST ABOUT COST!

How long does it take to read **3 TB of data**?

3 TB

**4.17 Hours**

# IT'S NOT JUST ABOUT COST!

How long does it take

What happens if you add more disks?

# HOW LONG DOES IT TAKE TO READ A 3 TB FILE?

**1 disk** — 4.17 hr

**100 disks** — 2.5 min

**1000 disks** — 15 sec

# HOW LONG DOES IT TAKE TO READ A 3 TB FILE?

| | |
|---|---|
| 1 disk | 4.17 hr |

Insight: More Disks are FASTER!

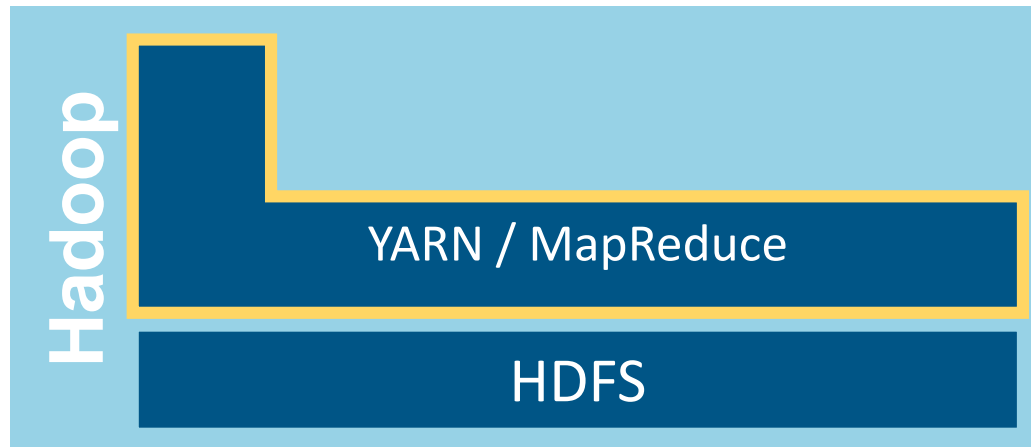| | |
|---|---|
| 1000 disks | 15 sec |

# Hadoop is a Storage Platform

**Hadoop**

**HDFS**

- Distributed Storage Performs Great

- Data is Replicated

- Reasonable Cost

- Sits on the OS File System

# Hadoop is a Processing Platform



- MapReduce/YARN
- Distributed Processing
- Data Locality
- Usually Java

# Apache Pig
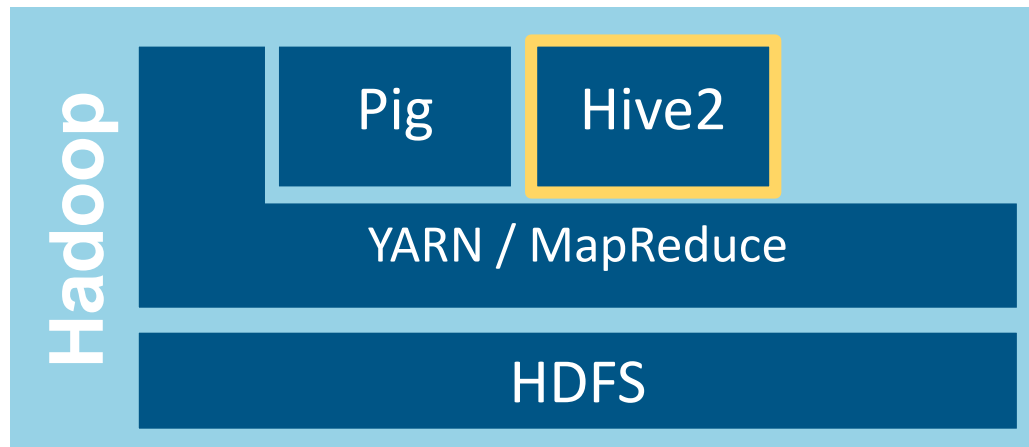
Hadoop

Pig

YARN / MapReduce

HDFS

- Scripting Language

- Higher level than programming Java MapReduce

- Pig Latin scripts are converted to MapReduce jobs

- Great for joining data

- Great for transforming data

# Apache Pig: Example Program

- Distributed Processing

```
people = LOAD '/user/training/customers' AS (cust_id, name);
orders = LOAD '/user/training/orders' AS (ord_id, cust_id, cost);
groups = GROUP orders BY cust_id;
totals = FOREACH groups GENERATE group, SUM(orders.cost) AS t;
result = JOIN totals BY group, people BY cust_id;
DUMP result;
```

Hadoop

# Apache Hive



Hadoop

Pig | Hive2

YARN / MapReduce

HDFS

- SQL on Hadoop
- Similar to traditional SQL
- Reduces development time
- Enables BI on Hadoop
- Schema-on-Read
- You choose underlying file format
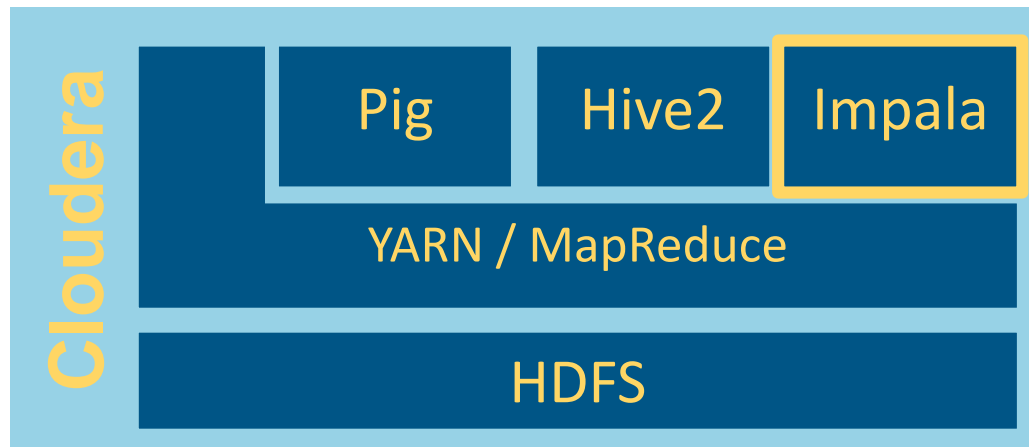
# Apache Hive

- SQL on Hadoop

**Hadoop**

```
SELECT zipcode, SUM(cost) AS total
FROM customers
JOIN orders
ON (customers.cust_id = orders.cust_id)
WHERE zipcode LIKE '63%'
GROUP BY zipcode
ORDER BY total DESC;
```

# Apache Impala is a SQL Engine

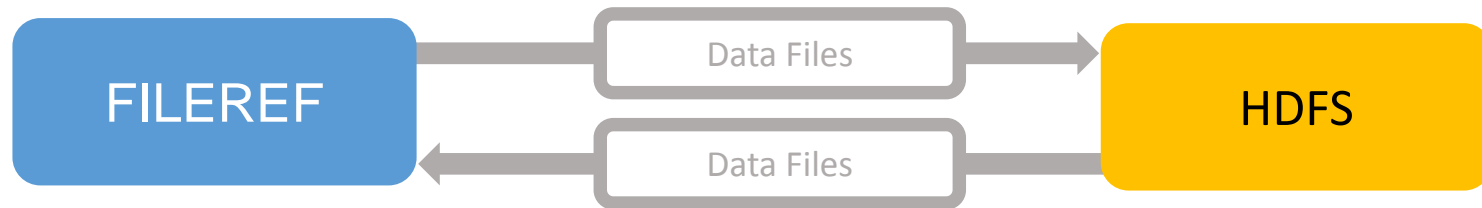| Cloudera | Pig | Hive2 | Impala |
| --- | --- | --- | --- |
| | YARN / MapReduce | | |
| | HDFS | | |

- High-performance SQL engine

- Handles concurrency well

- Does not rely on MapReduce

- Supports a dialect of SQL very similar to Hive's

- 100% open source

- Apache License

How can SAS Interact with Hadoop?

# Using Base SAS 9.4 with Hadoop

#1  FILEREF  → Data Files →  HDFS
          ← Data Files ←

# SAS FILENAME Statement for Hadoop

# SAS FILENAME Statement for Hadoop

```
options set=SAS_HADOOP_CONFIG_PATH="\\sashq\cdh45p1";
options set=SAS_HADOOP_JAR_PATH="\\sashq\cdh45";

FILENAME hdp1 hadoop 'test.txt';

/* Write file to HDFS */
data _null_;
    file hdp1;
    put ' Test Test Test';
run;


/* Read file from HDFS */
data test;
    infile hdp1;
    input textline $15.;
run;
```

# Using Base SAS 9.4 with Hadoop

#1  FILEREF  →  Data Files  →  HDFS
         ←  Data Files  ←

#2  PROC Hadoop  →  MapReduce + HDFS commands  →  Hadoop

# Hadoop Procedure

SAS

Hadoop

Pig | Hive2 | Impala

YARN / MapReduce

HDFS

- Submit HDFS commands
- Submit MapReduce Jobs
- Submit Pig Latin programs

# How Do I Submit HDFS Commands?

```
filename cfg 'C:\Hadoop_cfg\cdh57.xml';

/* Copy war_and_peace.txt to HDFS. */
/* Copy moby_dick.txt    to HDFS. */
proc hadoop options=cfg username="sasxjb" verbose;
    HDFS mkdir='/user/sasxjb/Books';
    HDFS COPYFROMLOCAL="C:\Hadoop_data\moby_dick.txt"
             OUT='/user/sasxjb/Books/moby_dick.txt';
    HDFS COPYFROMLOCAL="C:\Hadoop_data\war_and_peace.txt"
             OUT='/user/sasxjb/Books/war_and_peace.txt';
run;
```
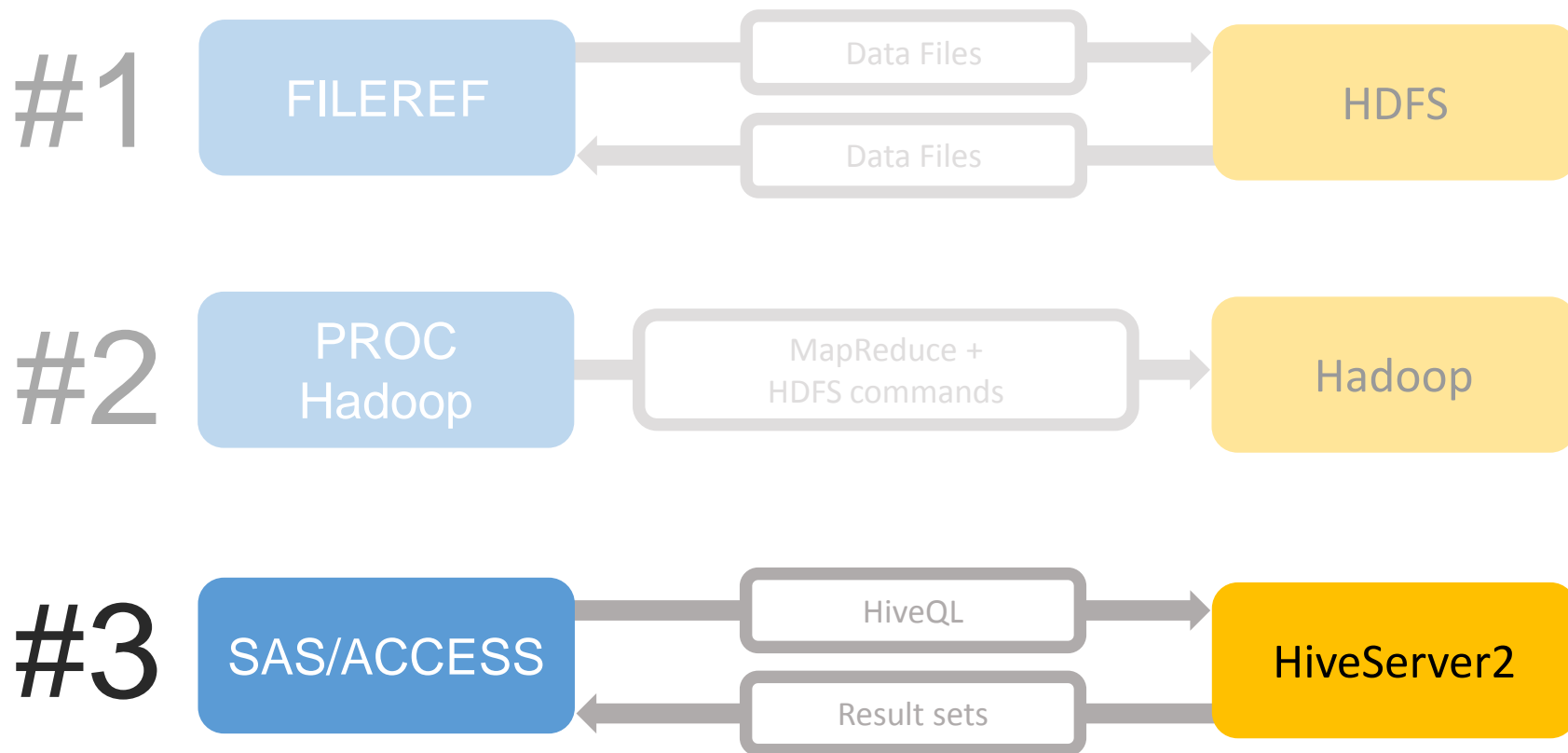
# How Do I Submit MapReduce Jobs?

```
filename cfg 'C:\Hadoop_cfg\cdh57.xml';

proc hadoop options=cfg user="sasxjb" verbose;
   mapreduce input='/user/sasxjb/Books/moby_dick.txt'
      output='/user/sasxjb/outBook'
      jar='C:\Hadoop_examples\hadoop-examples-1.2.0.1.3-96.jar'
      outputkey="org.apache.hadoop.io.Text"
      outputvalue="org.apache.hadoop.io.IntWritable"
      reduce="org.apache.hadoop.examples.WordCount$IntSumReducer"
      combine="org.apache.hadoop.examples.WordCount$IntSumReducer"
      map="org.apache.hadoop.examples.WordCount$TokenizerMapper";
run;
```
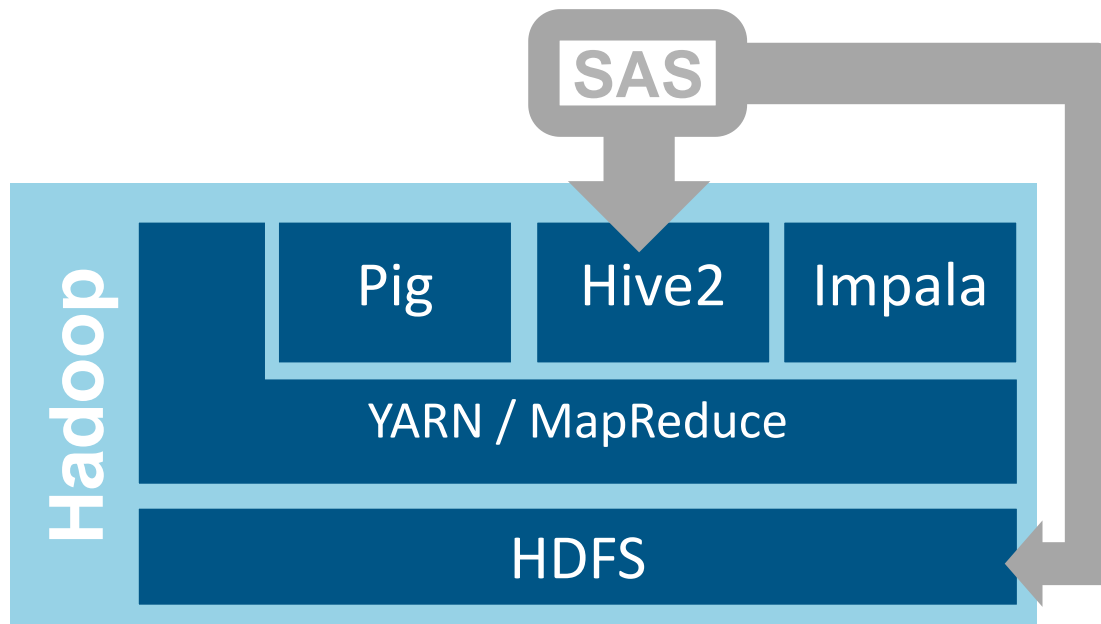
# How Do I Submit Pig Latin Programs?

```
filename cfg 'C:\Hadoop_cfg\cdh57.xml';

proc hadoop options=cfg username="sasxjb" verbose;
     pig code=pigcode ;
run;
```

# Using Base SAS 9.4 with Hadoop



**#1** FILEREF → Data Files → HDFS
HDFS ← Data Files ← FILEREF

**#2** PROC Hadoop → MapReduce + HDFS commands → Hadoop

**#3** SAS/ACCESS → HiveQL → HiveServer2
HiveServer2 → Result sets → SAS/ACCESS

# SAS/ACCESS Interface to Hadoop



- Connects via JDBC

- Makes Hive tables look like SAS data sets

- Bulk loads directly to HDFS

- Can read directly from HDFS

# How Does SAS/ACCESS Talk to Hadoop?

```
proc sql;
    select count(*) from mycdh.customer_dim
        where loyalty_program='Chocolate Club';
run;
```

**?**

# How Does SAS/ACCESS Talk to Hadoop?

```
proc sql;
    select count(*) from mycdh.customer_dim
        where loyalty_program='Chocolate Club';
run;
```

```
select COUNT(*) from `CUSTOMER_DIM` TXT_1
WHERE TXT_1.`loyalty_program` = 'Chocolate Club'
```

**SAS Generated This SQL**

# How Does SAS/ACCESS Talk to Hadoop?
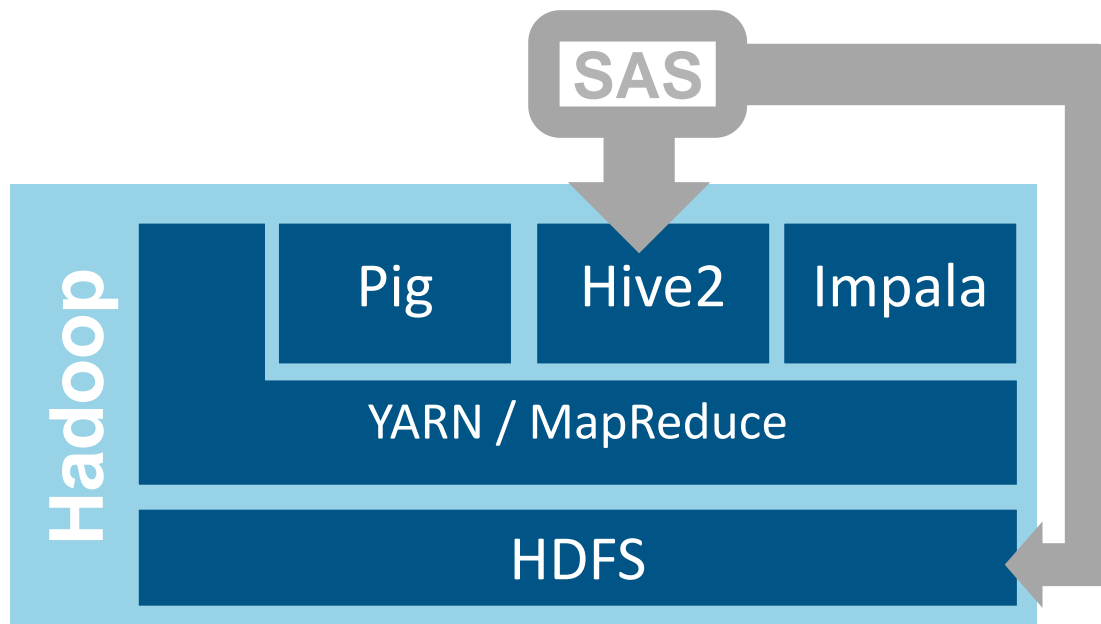
```
proc sql;
    select count(*) from mycdh.customer_dim
        where loyalty_program='Chocolate Club';
run;
```

```
OPTIONS SASTRACE=',,,d' SASTRACELOC=SASLOG NOSTSUFFIX;
```

```
select COUNT(*) from `CUSTOMER_DIM` TXT_1
WHERE TXT_1.`loyalty_program` = 'Chocolate Club'
```

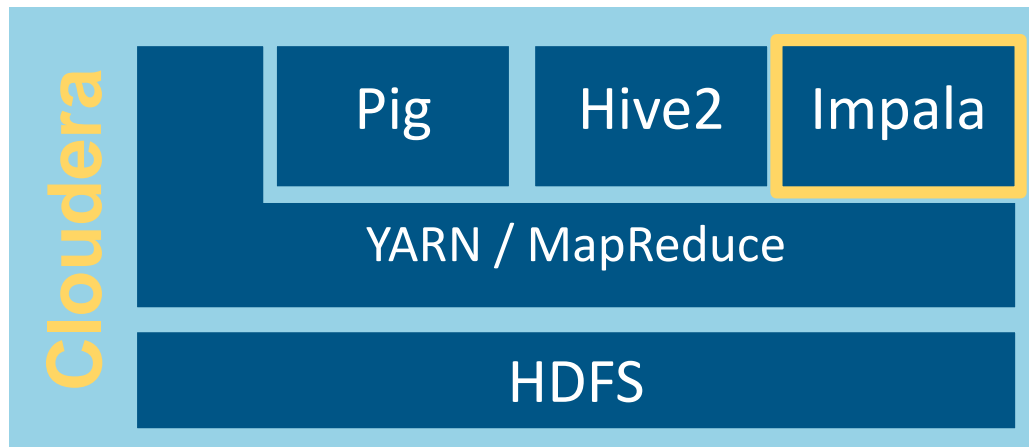**SAS Generated This SQL**

# SAS/ACCESS Interface to Hadoop



- **Generates HiveQL**
- Connects via JDBC
- Makes Hive tables look like SAS data sets
- Bulk loads directly to HDFS
- Can read directly from HDFS

# We Can Write Our Own HiveQL!

```
proc sql;
    connect to hadoop (server=quickstart
                          user=cloudera);
    execute (create table store_cnt
                    row format delimited
                    fields terminated by '\001`
                    stored as parquet
                 as
                    select customer_rk, count(*) as tot
                      from order_fact
                    group by customer_rk) by hadoop;
quit;
```
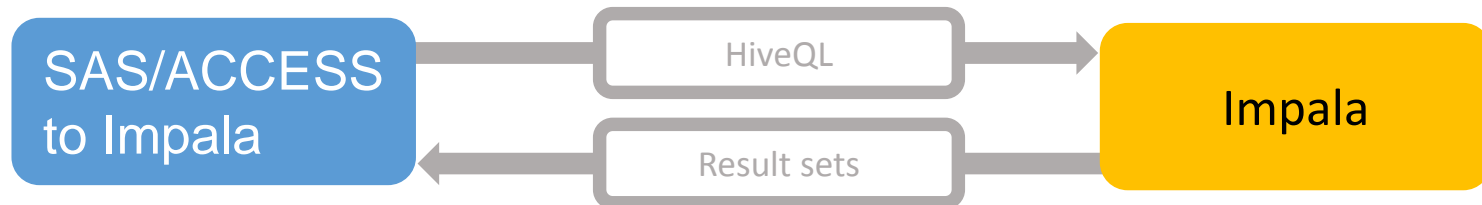
**Explicit Pass-Through**

## What about Apache Impala?

Cloudera

| Pig | Hive2 | Impala |

YARN / MapReduce

HDFS

SAS/ACCESS Interface to Impala:

- Connects via ODBC

- Makes Hive tables look like SAS data sets

- Bulk loads directly to HDFS

#4 SAS/ACCESS to Impala → HiveQL → Impala

Result sets ←

# What is SAS In-Database Code Accelerator?

```
proc ds2 indb=yes;
    thread tpgm / overwrite=yes;
        method run();
            set hdplib.intable;
            output;
    end;
    endthread;
    run;
    data hdplib.outdata

(overwrite=yes);
        dcl thread tpgm hdpdata;
        method run();
            set from hdpdata;
        end;
    enddata;
    run;
quit;
```
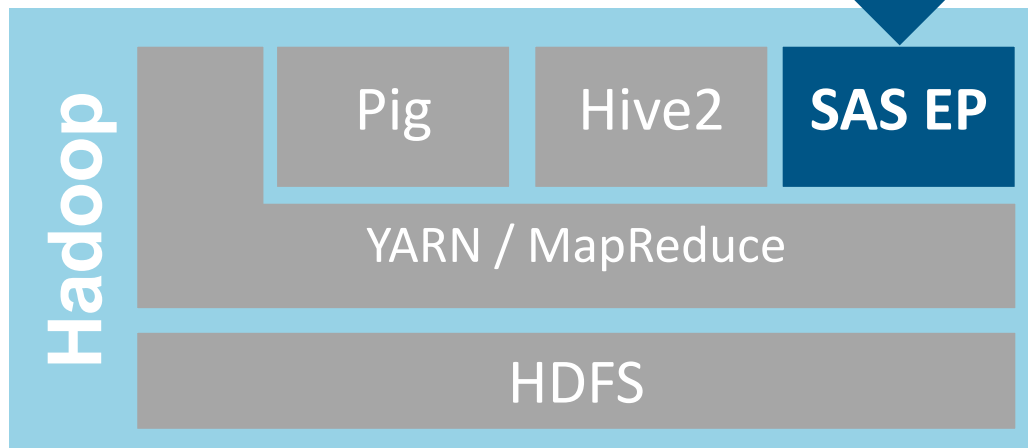
SAS In-Database Code Accelerators let you run SAS code inside Hadoop. With this you get:

- DS2 processing (modern DATA Step)

- More Data Types

- Code Packages

- More Programming Structures

- Parallel Database Operations

- Thread Programs Run Inside Database

# In-Database Code Accelerator Runs in Hadoop
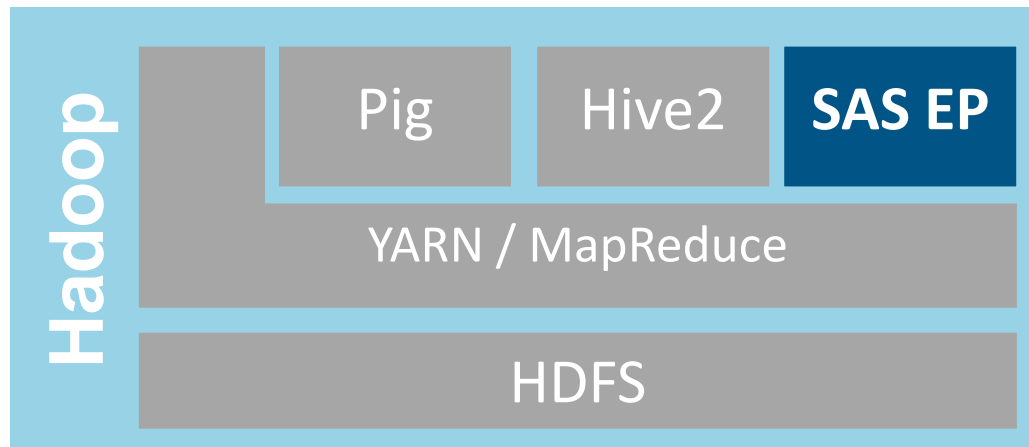
```
proc ds2 indb=yes;
    thread tpgm / overwrite=yes;
        method run();
            set hdplib.intable;
            output;
    end;
    endthread;
    run;
    data hdplib.outdata
```

```
(overwrite=yes);
        dcl thread tpgm hdpdata;
        method run();
            set from hdpdata;
        end;
    enddata;
    run;
quit;
```

**Hadoop**

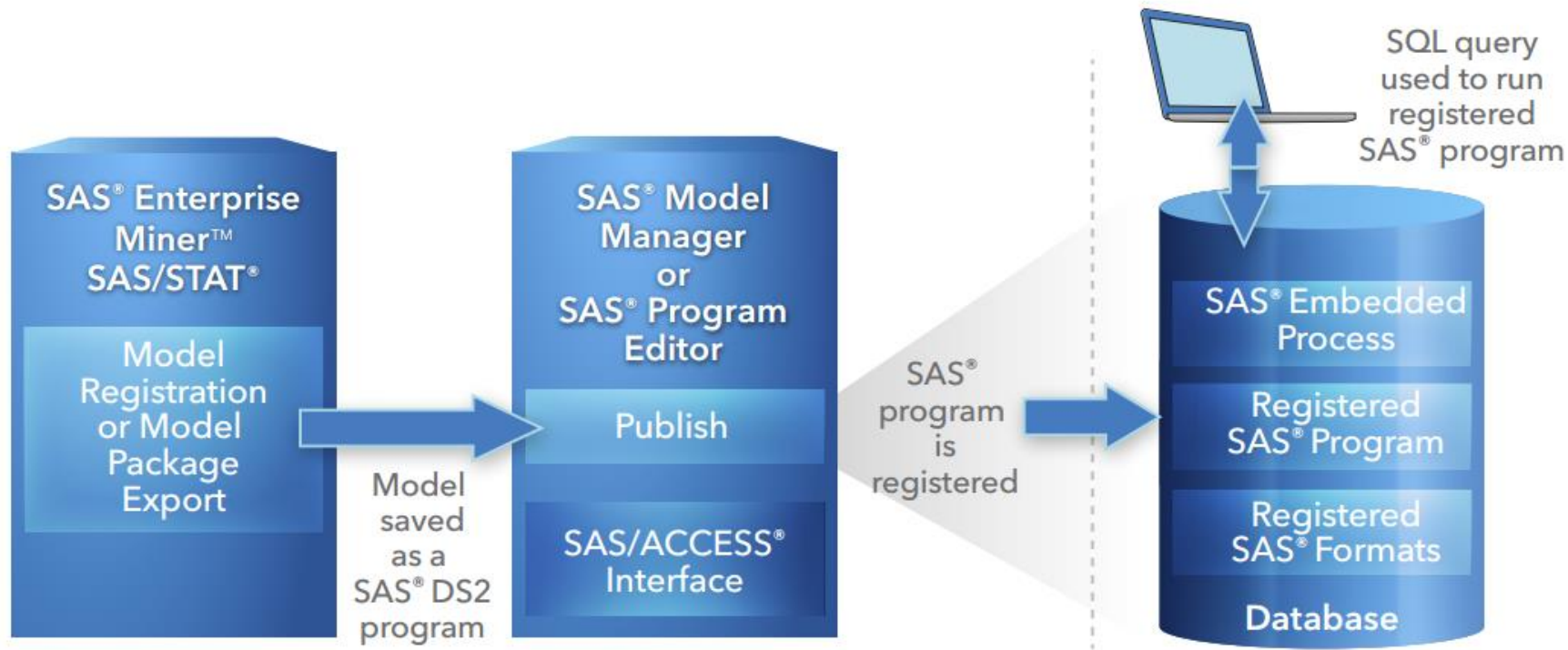| Pig | Hive2 | **SAS EP** |

YARN / MapReduce

HDFS

# What is SAS In-Database Scoring Accelerator?



SAS In-Database Scoring Accelerator lets you score models inside the cluster. With this you get:
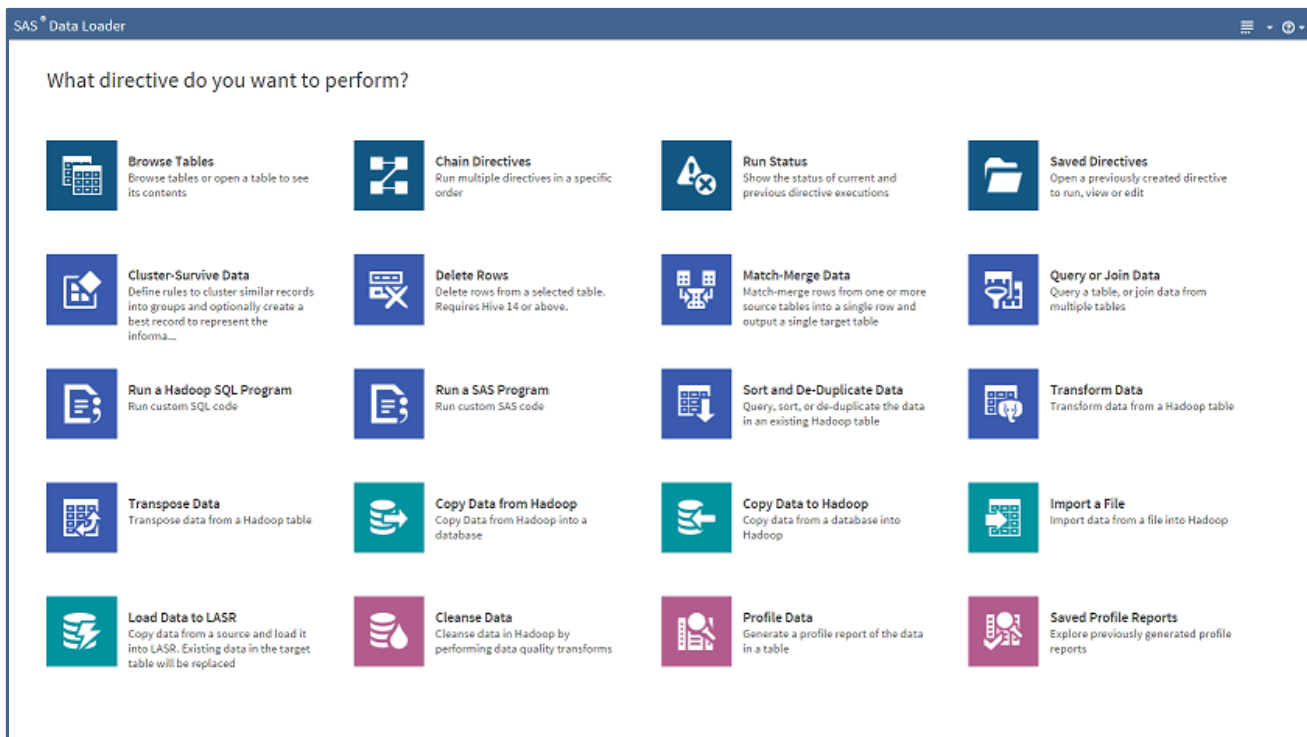
- Uses the SAS Embedded Process

- Faster Scoring

- Less data movement – score data where it lives

- Uses fewer resources

# What does the Scoring Process Look like?

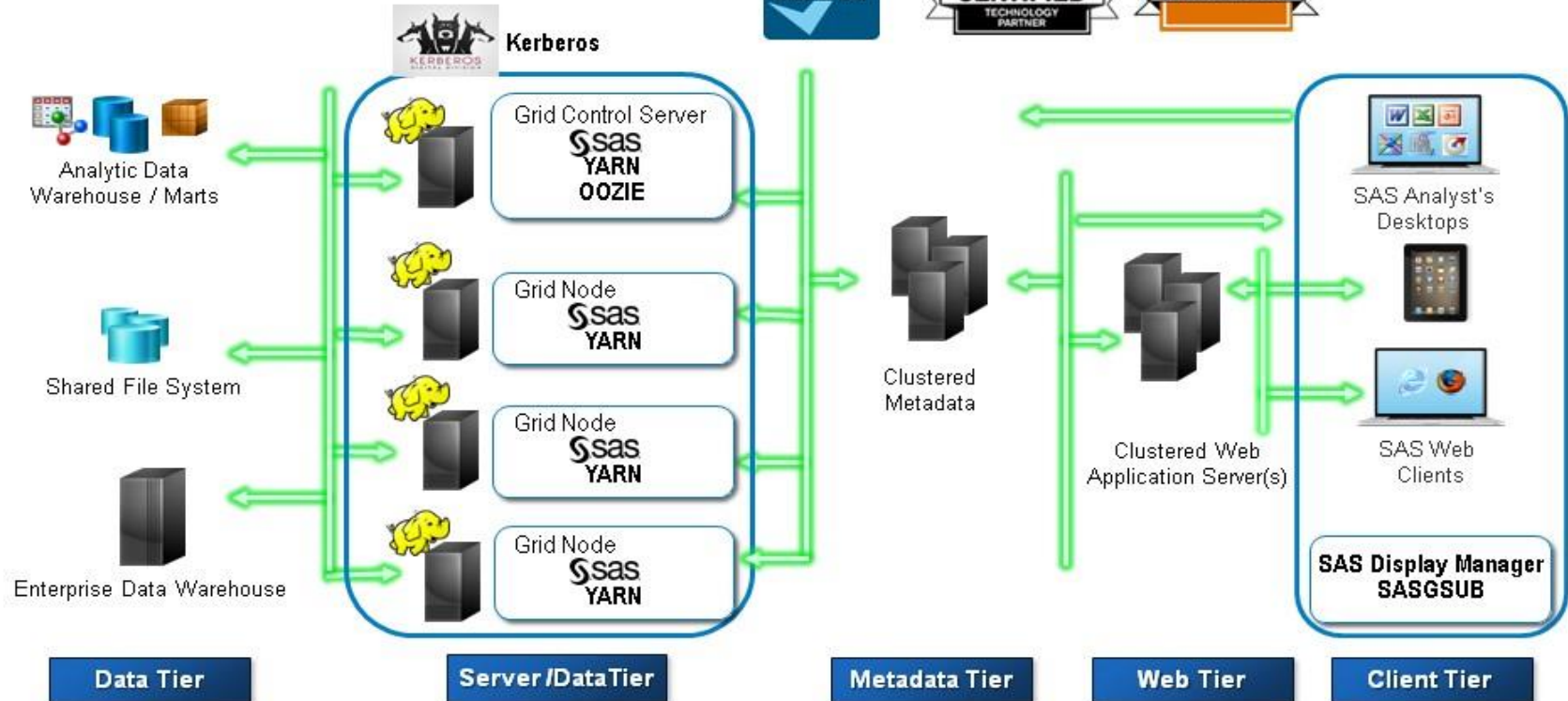# Data Loader for Hadoop – Self Service Big Data



- Easy to use UI
- Query Data
- Manage Data
- Transform Data
- Run Custom Code
- Move Data

SAS Grid Manager for Hadoop
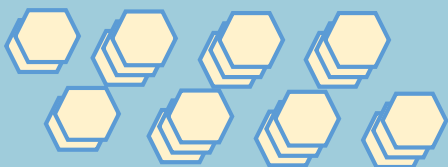
# What is SAS Grid Manager for Hadoop?

# Servers

# SAS Viya + Hadoop Architecture

## Cloud Analytic Services (CAS)
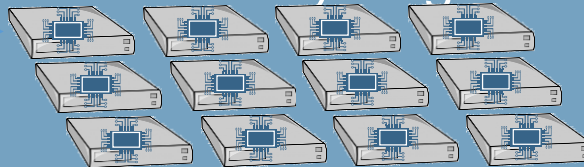
**Microservices**

**In-Memory Engine**

HDFS as infrastructure

SAS Data Connect Accelerator for Hadoop

SAS Data Connector to Hadoop

Hadoop

# Feel Free to Contact Me!

Jeff.Bailey@sas.com

http://www.linkedin.com/in/jeffreydbailey

https://github.com/Jeff-Bailey/SAS13341_SAS_Hadoop

ANALYTICS EXPERIENCE
2016

#analyticsx