



# PIG

is an open-source technology that

offers a high-level mechanism for parallel  
programming of MapReduce jobs

to be executed on Hadoop clusters



# PIG

IS A PART OF THE HADOOP ECOSYSTEM, IDEAL FOR  
USE WITH UNSTRUCTURED DATA WHERE

SCHEMAS ARE INCONSISTENT OR  
UNKNOWN



PIG

PIG DOESN'T REQUIRE DATA TO BE  
LOADED INTO TABLES FIRST

IT CAN OPERATE ON DATA AS SOON AS  
IT IS COPIED INTO HDFS



PIG VS



HIVE AND PIG SEEM TO HAVE SIMILAR  
APPLICATIONS AND FUNCTIONALITY

LEADING FOLKS TO WONDER  
HOW DO THEY STACK UP?



PIG VS



TO (OVER-)SIMPLIFY A BIT



DATA FROM  
TECH SYSTEMS

ETL  
→  
USING PIG



DATA  
WAREHOUSE

REPORTING  
→  
USING HIVE





PIG

IS GREAT FOR GETTING DATA **INTO** A  
DATA WAREHOUSE



PIG IS GREAT FOR GETTING DATA  
INTO A DATA WAREHOUSE

FROM POORLY STRUCTURED DATA  
SOURCES SUCH AS WEB LOGS



PIG IS GREAT FOR GETTING DATA  
INTO A DATA WAREHOUSE

FROM POORLY STRUCTURED DATA  
SOURCES SUCH AS WEB LOGS

THAT MAKES IT A GREAT  
COMPLEMENT TO HIVE



PIG IS GREAT FOR GETTING DATA  
INTO A DATA WAREHOUSE  
FROM POORLY STRUCTURED DATA SOURCES SUCH AS  
WEB LOGS

THAT MAKES IT A GREAT  
COMPLEMENT TO HIVE

WHICH IS GREAT FOR ANALYSING DATA  
ALREADY IN A DATA WAREHOUSE



PIG IS GREAT FOR GETTING DATA  
INTO A DATA WAREHOUSE  
FROM POORLY STRUCTURED DATA SOURCES SUCH AS  
WEB LOGS

THAT MAKES IT A GREAT COMPLEMENT TO HIVE WHICH IS GREAT FOR ANALYSING DATA  
ALREADY IN A DATA WAREHOUSE

PIG IS ALSO GREAT AT SOME NON-  
RELATIONAL DATA TRANSFORMATIONS



PIG IS ALSO GREAT AT SOME NON-  
RELATIONAL DATA TRANSFORMATIONS

SUCH AS CALCULATING THE CORRELATION  
OF 2 MILLION-ELEMENT VECTORS

THIS WOULD NOT BE A NATURAL  
OPERATION IN SQL OR HIVEQL



# PIG WORKS JUST FINE WITH HBASE

PIG CAN SERVE AS AN ABSTRACTION  
ATOP NON-RELATIONAL DATA

THIS MAKES IT MORE VERSATILE  
THAN HIVE

# HADOOP

is a distributed computing framework  
developed and maintained by

THE APACHE SOFTWARE FOUNDATION

written in Java

# HADOOP

## HDFS

A file system to  
manage the  
storage of data

## MapReduce

A framework to  
process data across  
multiple servers

# HADOOP

HDFS

MapReduce

Hadoop is an open  
source implementation  
of 2 proprietary  
technologies by Google

GFS  
MapReduce

# HADOOP

HDFS

MapReduce

A framework to  
process data across  
multiple servers

# HADOOP

HDFS

MapReduce

In Hadoop 2.0 the  
MapReduce block was  
broken into 2 parts

# HADOOP

HDFS

YARN

MapReduce

A framework  
to **run** the data  
processing  
task

A framework  
to **define** a data  
processing  
task

# HADOOP

# MapReduce

HDFS

YARN

A framework to  
**define** a data  
processing task

**MapReduce** is a way to parallelize a  
data processing task

# HADOOP

# MapReduce

HDFS

YARN

A framework to  
**define** a data  
processing task

Manages resources and memory  
across multiple nodes

# HADOOP

HDFS

YARN

MapReduce

Hadoop uses this to store data across multiple disks

# HADOOP

HDFS

YARN

MapReduce

YARN was introduced in Hadoop 2.0 to separately handle

the management of resources  
on the Hadoop cluster

# HADOOP

HDFS

YARN

MapReduce

YARN co-ordinates all the different  
MapReduce tasks running on the cluster

**THAT was Hadoop**

**LET US TALK ABOUT      HIVE**

**It is a data warehouse infrastructure**

**built on top of Hadoop**

**for providing data summarisation, query and analysis**

HIVE

HADOOP

HDFS

MapReduce

YARN

It is a data warehouse infrastructure

HIVE

HADOOP

HDFS

MapReduce

YARN

HIVE PROVIDES A SQL LIKE  
INTERFACE TO DATA IN HDFS

HIVE

HADOOP

HDFS

MapReduce

YARN

Traditional databases/closed-source  
data warehouses normally use SQL

# HIVE

HADOOP

HDFS

## MapReduce

YARN

HIVE WILL TRANSLATE THE QUERY  
INTO 1/MORE MAPREDUCE TASKS

# HIVE

HADOOP

HDFS

# MapReduce

YARN

THE MAPREDUCE TASKS WILL PROCESS THE DATA  
IN HDFS AND RETURN ANY RESULTS TO HIVE

# HIVE

## HADOOP

HDFS

MapReduce

YARN

HIVE STORES IT'S DATA  
AS FILES IN HDFS

# HIVE

## HADOOP

HDFS

MapReduce

YARN

THE HIVE DATA IN HDFS IS IN THE  
FORM OF FILES AND DIRECTORIES

HIVE

HIVE IS NOT THE SOLE OWNER OF THESE  
FILES

# HIVE

THESE FILES CAN BE USED AND  
MODIFIED BY OTHER CLIENTS

HIVE

HENCE, HIVE ENFORCES  
SCHEMA-ON-READ

# HIVE

IN SCHEMA-ON-READ

DURING LOAD/INSERT OPERATIONS, HIVE  
WILL JUST DUMP THE DATA INTO A FILE  
WITHOUT CHECKING THE SCHEMA

# HIVE

WHEN YOU READ THE DATA IN HIVE, IT  
WILL PARSE THE FILE AND TRY TO  
IMPOSE THE SCHEMA

# HIVE

**HIVE MAY NOT ALWAYS SUCCEED IN  
IMPOSING THE SCHEMA**

# HIVE

THE DATA IN HIVE IS MEANT TO BE  
USED FOR ANALYTICAL PURPOSES

# HIVE

HIVE IS OPTIMAL FOR PROCESSING  
REALLY LARGE-SCALE DATASETS

GIGABYTES/PETABYTES

# HIVE

HIVE IS OPTIMAL FOR PROCESSING REALLY LARGE-SCALE DATASETS

BUT HIVE IS NOT OPTIMAL IN SITUATIONS  
WHERE THE SCHEMA IS UNKNOWN,  
INCOMPLETE OR INCONSISTENT

# HIVE

BUT HIVE IS NOT OPTIMAL IN SITUATIONS WHERE THE SCHEMA IS UNKNOWN, INCOMPLETE OR INCONSISTENT

WHAT DO WE DO IF WE HAVE DATA WITH INCONSISTENT SCHEMA?



PIG COMES TO THE RESCUE



# PIG

is an open-source technology that

offers a high-level mechanism for parallel  
programming of MapReduce jobs

to be executed on Hadoop clusters



# PIG

IS A PART OF THE HADOOP ECOSYSTEM, IDEAL FOR  
USE WITH UNSTRUCTURED DATA WHERE

SCHEMAS ARE INCONSISTENT OR  
UNKNOWN



# PIG

**PIG DOESN'T REQUIRE DATA TO BE  
LOADED INTO TABLES FIRST**

**IT CAN OPERATE ON DATA AS SOON AS  
IT IS COPIED INTO HDFS**



# PIG

IT CAN OPERATE ON DATA AS  
SOON AS IT IS COPIED INTO HDFS



# PIG

```
2008-10-09 00:37:37.011] - 21384 12764 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:37:39.152] CHSvsr07 21384 12764 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:37:39.168] CHSvsr07 21384 12764 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:37:39.558] - 18568 18792 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:37:39.637] - 18900 20176 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:37:39.668] - 21572 17944 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:37:40.308] CHSvsr07 18900 20176 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:37:40.308] CHSvsr07 18900 20176 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:37:40.996] CHSvsr07 21572 17944 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:37:41.012] CHSvsr07 21572 17944 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:37:41.215] CHSvsr07 18568 18792 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:37:41.215] CHSvsr07 18568 18792 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:37:41.980] CHSvsr07 21384 12764 SW_SystemDispatcher::Init - 0 44951 "Successfully started Microsoft System Center Application virtualization Management Service."
2008-10-09 00:56:32.214] CHSvsr07 21384 12764 SW_SystemDispatcher::FinI - 0 44952 "Successfully shut down Microsoft System Center Application virtualization Management Service."
2008-10-09 00:56:32.292] CHSvsr07 21384 22340 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 00:56:32.573] CHSvsr07 18568 18792 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 00:56:32.573] CHSvsr07 18900 20176 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 00:56:32.573] CHSvsr07 21572 17944 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 00:58:37.375] - 1140 1164 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:58:41.250] CHSvsr07 1140 1164 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:58:41.265] CHSvsr07 1140 1164 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:58:41.765] - 1492 1496 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:58:41.781] - 1512 1516 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:58:41.796] - 1504 1508 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 00:58:41.234] CHSvsr07 1512 1516 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:58:41.234] CHSvsr07 1504 1508 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:58:43.234] CHSvsr07 1512 1516 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:58:43.234] CHSvsr07 1504 1508 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:58:43.359] CHSvsr07 1492 1496 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 00:58:43.359] CHSvsr07 1492 1496 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 00:58:46.328] CHSvsr07 1140 1164 SW_SystemDispatcher::Init - 0 44951 "Successfully started Microsoft System Center Application virtualization Management Service."
2008-10-09 01:15:25.292] CHSvsr07 1140 1164 SW_SystemDispatcher::FinI - 0 44952 "Successfully shut down Microsoft System Center Application virtualization Management Service."
2008-10-09 01:15:25.370] CHSvsr07 1140 1144 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 01:15:25.946] CHSvsr07 1492 1496 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 01:15:26.632] CHSvsr07 1512 1516 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 01:15:26.694] CHSvsr07 1504 1508 SW_MessageHandler::Close - - - 5 65535 "Shutdown complete."
2008-10-09 01:15:27.396] - 5192 5384 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 01:15:28.673] CHSvsr07 5192 5384 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 01:15:28.673] CHSvsr07 5192 5384 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 01:15:28.954] - 7432 6396 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 01:15:29.016] - 6400 7732 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 01:15:29.079] - 5732 6148 SW_MessageHandler::Open - - - 5 65535 "Initialization complete."
2008-10-09 01:15:29.904] CHSvsr07 6400 7732 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 01:15:29.904] CHSvsr07 6400 7732 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 01:15:30.278] CHSvsr07 7432 6396 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 01:15:30.294] CHSvsr07 7432 6396 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 01:15:30.356] CHSvsr07 5732 6148 SW_SQLDataConnection::Open - - - 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database: CHSSQL01, Error: 0x80004005, State: 1)."
2008-10-09 01:15:30.356] CHSvsr07 5732 6148 SW_SQLOutputHandler::Open - - - 1 44910 "Failed to connect to data source."
2008-10-09 01:15:31.385] CHSvsr07 5192 5384 SW_SystemDispatcher::Init - 0 44951 "Successfully started Microsoft System Center Application virtualization Management Service."
2008-10-09 01:15:56.785] CHSvsr07 7432 8120 SW_LicenseConduitLogger::LogMessage 1956118080 "Default Provider" "Justin Zarb" TECHDES.001/TECHDES.001.sft 0 40976 "License Setup"
2008-10-09 01:15:56.941] CHSvsr07 7432 8120 SW_RTSPHandler::HandleSetup 1956118080 "Default Provider" "Justin Zarb" TECHDES.001/TECHDES.001.sft 0 40960 "Session Setup"
2008-10-09 01:16:03.926] CHSvsr07 7432 5832 SW_LicenseConduitLogger::LogMessage 1956118080 "Default Provider" "Justin Zarb" TECHDES.001/TECHDES.001.sft 0 40977 "License Setup"
2008-10-09 01:16:03.926] CHSvsr07 7432 5832 SW_RTSPHandler::HandleTeardown 1956118080 "Default Provider" "Justin Zarb" TECHDES.001/TECHDES.001.sft 0 40961 "Session Teardown"
2008-10-09 01:20:32.766] CHSvsr07 6400 3444 SW_LicenseConduitLogger::LogMessage 1897464568 "Default Provider" "Justin Zarb" SCVMCON.001/SCVMCON.001.sft 0 40976 "License Setup"
2008-10-09 01:20:32.766] CHSvsr07 6400 3444 SW_RTSPHandler::HandleSetup 1897464568 "Default Provider" "Justin Zarb" SCVMCON.001/SCVMCON.001.sft 0 40960 "Session Setup"
2008-10-09 06:31:40.632] CHSvsr07 6400 8016 SW_LicenseConduitLogger::LogMessage 1897464568 "Default Provider" "Justin Zarb" SCVMCON.001_Package 0 40978 "License Abnormal"
2008-10-09 06:31:40.663] CHSvsr07 6400 6340 SW_SQLOutputHandler::HandleMessage - - - 1 44911 "Create record failed with error [1]."
2008-10-09 06:35:40.367] CHSvsr07 5732 3960 SW_DataConnectionPool::Reconnect - "Default Provider" "Justin Zarb" SCVMCON.001/SCVMCON.001.sft 2 41543 "Connection to data store on host CHSSQL01\8SQSTEPS failed"
2008-10-09 06:35:40.367] CHSvsr07 5732 3960 SW_ServerAuthenticationTask::Authorize - "Default Provider" "Justin Zarb" SCVMCON.001/SCVMCON.001.sft 2 41472 "Module execution failed."
2008-10-09 06:35:40.367] CHSvsr07 5732 5348 SW_SQLOutputHandler::HandleMessage - - - 1 44911 "Create record failed with error [1]."
2008-10-09 06:48:44.656] CHSvsr07 7432 6572 SW_DataConnectionPool::Reconnect - "Default Provider" "ecoe" - 2 41543 "Connection to data store on host CHSSQL01\8SQSTEPS failed"
2008-10-09 06:48:44.656] CHSvsr07 7432 6572 SW_ServerAuthenticationTask::Authorize - "Default Provider" "ecoe" - 2 41543 "Module execution failed."
2008-10-09 06:48:44.702] CHSvsr07 6400 5872 SW_DataConnectionPool::Reconnect - "Default Provider" "ecoe" - 2 41543 "Connection to data store on host CHSSQL01\8SQSTEPS failed"
2008-10-09 06:48:44.702] CHSvsr07 7432 7864 SW_SQLOutputHandler::HandleMessage - - - 1 44911 "Create record failed with error [1]."
2008-10-09 06:48:44.702] CHSvsr07 6400 5872 SW_ServerAuthenticationTask::Authorize - "Default Provider" "ecoe" - 2 41472 "Module execution failed."
2008-10-09 08:33:07.803] CHSvsr07 5732 6508 SW_LicenseConduitLogger::LogMessage 231148567 "Default Provider" Swelshman SIMSNET.005/SIMSNET.005_7.sft 0 40976 "License Session"
2008-10-09 08:33:07.818] CHSvsr07 5732 6508 SW_RTSPHandler::HandleSetup 231148567 "Default Provider" Swelshman SIMSNET.005/SIMSNET.005_7.sft 0 40960 "Session Setup"
2008-10-09 08:39:02.202] CHSvsr07 5732 3960 SW_LicenseConduitLogger::LogMessage 2076584603 "Default Provider" YWESTWOOD SIMSNET.005/SIMSNET.005_7.sft 0 40976 "License Session"
2008-10-09 08:39:02.218] CHSvsr07 5732 3960 SW_RTSPHandler::HandleSetup 2076584603 "Default Provider" YWESTWOOD SIMSNET.005/SIMSNET.005_7.sft 0 40960 "Session Setup"
2008-10-09 08:44:58.679] CHSvsr07 7432 6492 SW_LicenseConduitLogger::LogMessage 1008015806 "Default Provider" DLANEY ADRDR811.001/ADRDR811.001.sft 0 40976 "License Session"
2008-10-09 08:44:58.726] CHSvsr07 7432 6492 SW_RTSPHandler::HandleSetup 1008015806 "Default Provider" DLANEY ADRDR811.001/ADRDR811.001.sft 0 40960 "Session Setup"
2008-10-09 08:45:03.412] CHSvsr07 7432 1528 SW_LicenseConduitLogger::LogMessage 1008015806 "Default Provider" DLANEY ADRDR811.001/ADRDR811.001.sft 0 40977 "License Session"
```

CONSIDER  
SERVER LOG DATA  
WHICH HAS NO  
SCHEMA



CONSIDER THE DATA OF SERVER LOGS WHICH HAS  
NO SCHEMA

# THIS SERVER DATA CAN BE TRANSFORMED INTO DATA WITH PROPER SCHEMA IN PIG

THIS SERVER DATA CAN BE TAKEN INTO DATA WITH PROPER SECURITY

TIME- STAMP	SERVER NAME	PID	TID	Module	SESSION ID	PROVIDE	MESSA GE
----------------	----------------	-----	-----	--------	---------------	---------	-------------



CONSIDER THE DATA OF SERVER LOGS WHICH HAS  
NO SCHEMA

# EVERY ENTRY IN THE LOG IS SPLIT INTO VARIOUS COLUMNS

**ANSWER** The answer is **100**.

TIME- STAMP	SERVER NAME	PID	TID	Module	SESSION ID	PROVIDE	MESSA GE
----------------	----------------	-----	-----	--------	---------------	---------	-------------



# CONSIDER THE DATA OF SERVER LOGS WHICH HAS NO SCHEMA

# DIFFERENT ERROR MESSAGES WILL HAVE DIFFERENT STRUCTURES

# DIFFERENT ERROR MESSAGES HAVE DIFFERENT STYLES

**TIME- SERVER PID TID Module SESSION ID PROVIDE MESSA STAMP NAME**



# CONSIDER THE DATA OF SERVER LOGS WHICH HAS NO SCHEMA

THIS IS VERY SIMPLE COMPARED TO THE COMPLEXITY WE WILL SEE IN WEBSITE LOGS

```
2008-10-09 00:37:37.011 - 21384 12764 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:37:39.152 CHSVSR07 21384 12764 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database:
2008-10-09 00:37:39.168 CHSVSR07 21384 12764 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:37:39.558 - 18568 18792 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:37:39.637 - 18900 18792 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database:
2008-10-09 00:37:39.668 CHSVSR07 18900 18792 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:37:40.308 CHSVSR07 18900 18792 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database:
2008-10-09 00:37:40.308 CHSVSR07 18900 18792 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:37:40.996 CHSVSR07 21572 18792 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:37:41.012 CHSVSR07 21572 18792 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database:
2008-10-09 00:37:41.215 CHSVSR07 18568 18792 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:37:41.215 CHSVSR07 18568 18792 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, Database:
2008-10-09 00:37:41.980 CHSVSR07 18568 18792 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:56:32.214 CHSVSR07 1140 1164 SW_SystemDispatcher::Init --- 0 44951 "Successfully started Microsoft System Center Application virtualization Management"
2008-10-09 00:56:32.292 CHSVSR07 1140 1164 SW_SystemDispatcher::Fini --- 0 44952 "Successfully shut down Microsoft System Center Application virtualization Management"
2008-10-09 00:56:32.573 CHSVSR07 1140 1164 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:56:32.573 CHSVSR07 1140 1164 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 00:58:37.375 CHSVSR07 1140 1164 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:58:41.250 CHSVSR07 1140 1164 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 00:58:41.265 CHSVSR07 1140 1164 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:58:41.765 - 1492 1496 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:58:41.781 - 1512 1516 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:58:41.796 - 1504 1508 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 00:58:43.234 CHSVSR07 1512 1516 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 00:58:43.234 CHSVSR07 1504 1508 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 00:58:43.234 CHSVSR07 1512 1516 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:58:43.234 CHSVSR07 1504 1508 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:58:43.359 CHSVSR07 1492 1496 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 00:58:43.359 CHSVSR07 1492 1496 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 00:58:46.328 CHSVSR07 1140 1164 SW_SystemDispatcher::Init --- 0 44951 "Successfully started Microsoft System Center Application virtualization Management"
2008-10-09 01:15:25.292 CHSVSR07 1140 1164 SW_SystemDispatcher::Fini --- 0 44952 "Successfully shut down Microsoft System Center Application virtualization Management"
2008-10-09 01:15:25.370 CHSVSR07 1140 1164 SW_MessageHandler::Close --- 5 65535 "Shutdown complete."
2008-10-09 01:15:25.940 CHSVSR07 1492 1496 SW_MessageHandler::Close --- 5 65535 "Shutdown complete."
2008-10-09 01:15:26.630 CHSVSR07 1512 1516 SW_MessageHandler::Close --- 5 65535 "Shutdown complete."
2008-10-09 01:15:26.694 CHSVSR07 1504 1508 SW_MessageHandler::Close --- 5 65535 "Shutdown complete."
2008-10-09 01:15:27.394 - 5192 5184 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 01:15:28.673 CHSVSR07 5192 5184 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 01:15:28.673 CHSVSR07 5192 5184 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 01:15:28.954 - 7432 6796 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 01:15:29.016 - 6400 732 5184 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 01:15:29.079 - 5732 648 5184 SW_MessageHandler::Open --- 5 65535 "Initialization complete."
2008-10-09 01:15:29.904 CHSVSR07 6400 732 5184 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 01:15:29.904 CHSVSR07 6400 732 5184 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 01:15:30.278 CHSVSR07 7432 639 5184 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 01:15:30.294 CHSVSR07 7432 639 5184 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 01:15:30.354 CHSVSR07 5732 614 5184 SW_SQLDataConnection::Open --- 2 41494 "Failed to establish a connection to the data source: (Server: CHSSQL01, database: 5"
2008-10-09 01:15:30.354 CHSVSR07 5732 614 5184 SW_SQLOutputHandler::Open --- 1 44910 "Failed to connect to data source."
2008-10-09 01:15:31.385 CHSVSR07 5192 5184 SW_LicenseConduitLogger::LogMessage 1956118080 "default provider" "Justin Zarb" TECHDES_001/TECHDES_001.sft 0 40960 "License Se
2008-10-09 01:15:36.941 CHSVSR07 7432 6120 SW_RTSPHandler::HandleSetup 1956118080 Default Provider "Justin Zarb" TECHDES_001/TECHDES_001.sft 0 40960 "Session Setup"
2008-10-09 01:16:03.920 CHSVSR07 7432 6120 SW_LicenseConduitLogger::LogMessage 1956118080 "Default Provider" "Justin Zarb" TECHDES_001/TECHDES_001.sft 0 40977 "License Se
2008-10-09 01:16:03.920 CHSVSR07 7432 6120 SW_RTSPHandler::HandleTeardown 1956118080 "Default Provider" "Justin Zarb" TECHDES_001/TECHDES_001.sft 0 40961 "Session Teardown"
2008-10-09 01:20:32.768 CHSVSR07 6400 6444 SW_LicenseConduitLogger::LogMessage 1897464568 "Default Provider" "Justin Zarb" SCVMCON_001/SCVMCON_001.sft 0 40976 "License"
2008-10-09 01:20:32.768 CHSVSR07 6400 6444 SW_RTSPHandler::HandleSetup 1897464568 "Default Provider" "Justin Zarb" SCVMCON_001/SCVMCON_001.sft 0 40960 "Session Setup"
2008-10-09 06:31:40.632 CHSVSR07 6400 6031 SW_LicenseConduitLogger::LogMessage 1897464568 "Default Provider" "Justin Zarb" SCVMCON_001_Package 0 40978 "License Abnormal"
2008-10-09 06:31:40.661 CHSVSR07 6400 6341 SW_SQLOutputHandler::HandleMessage --- 1 44911 "Create record failed with error [1]."
2008-10-09 06:35:40.367 CHSVSR07 5732 6196 SW_DataConnectionPool::Reconnect - "Default Provider" "Justin Zarb" SCVMCON_001/SCVMCON_001.sft 2 41543 "Connection to data
2008-10-09 06:35:40.367 CHSVSR07 5732 6196 SW_ServerAuthenticationTask::Authorize - "Default Provider" "Justin Zarb" SCVMCON_001/SCVMCON_001.sft 2 41472 "Module executi
2008-10-09 06:35:40.367 CHSVSR07 5732 6134 SW_SQLOutputHandler::HandleMessage --- 1 44911 "Create record failed with error [1]."
2008-10-09 06:48:44.656 CHSVSR07 7432 6571 SW_DataConnectionPool::Reconnect - "Default Provider" "ecoe" - 2 41543 "Connection to data store on host CHSSQL01\BSQSTEPS fail
2008-10-09 06:48:44.656 CHSVSR07 7432 6571 SW_ServerAuthenticationTask::Authorize - "Default Provider" "ecoe" - 2 41472 "Module execution failed."
2008-10-09 06:48:44.703 CHSVSR07 6400 687 SW_DataConnectionPool::Reconnect - "Default Provider" "ecoe" - 2 41543 "Connection to data store on host CHSSQL01\BSQSTEPS fail
2008-10-09 06:48:44.656 CHSVSR07 7432 686 SW_SQLOutputHandler::HandleMessage --- 1 44911 "Create record failed with error [1]."
2008-10-09 06:48:44.703 CHSVSR07 6400 687 SW_ServerAuthenticationTask::Authorize - "Default Provider" "ecoe" - 2 41472 "Module execution failed."
2008-10-09 08:33:07.801 CHSVSR07 5732 6503 SW_LicenseConduitLogger::LogMessage 231148567 "default provider" Swelshman SIMSNET_005/SIMSNET_005_7.sft 0 40976 "License Sessi
2008-10-09 08:33:07.815 CHSVSR07 5732 6503 SW_RTSPHandler::HandleSetup 231148567 "Default Provider" Swelshman SIMSNET_005/SIMSNET_005_7.sft 0 40960 "Session Setup"
2008-10-09 08:39:02.203 CHSVSR07 5732 6196 SW_LicenseConduitLogger::LogMessage 2076584603 "Default Provider" YWESTWOOD SIMSNET_005/SIMSNET_005_7.sft 0 40976 "License Sess
2008-10-09 08:39:02.218 CHSVSR07 5732 6196 SW_RTSPHandler::HandleSetup 2076584603 "Default Provider" YWESTWOOD SIMSNET_005/SIMSNET_005_7.sft 0 40960 "Session Setup"
2008-10-09 08:44:58.679 CHSVSR07 432 6491 SW_LicenseconduitLogger::LogMessage 1008015806 "Default Provider" DLANEY ADRDR811.001/ADRDR811.001.sft 0 40976 "License Session"
2008-10-09 08:44:58.726 CHSVSR07 7432 6491 SW_RTSPHandler::HandleSetup 1008015806 "Default Provider" DLANEY ADRDR811.001/ADRDR811.001.sft 0 40960 "Session Setup"
2008-10-09 08:45:01.412 CHSVSR07 7432 6521 SW_LicenseconduitLogger::LogMessage 1008015806 "Default Provider" DLANEY ADRDR811.001/ADRDR811.001.sft 0 40977 "License session
```

TIME-STAMP	SERVER NAME	PID	TID	Module	SESSION ID	PROVIDE	MESSEGE
------------	-------------	-----	-----	--------	------------	---------	---------



# PIG

COMES WITH A LANGUAGE, **PIG LATIN**, FOR  
DOING THIS KIND OF DATA PROCESSING



DON'T WORRY, THIS IS NOT  
THE SYNTAX OF PIG LATIN!



# PIG LATIN

IS A PROCEDURAL DATA FLOW LANGUAGE

IT IMPLIES THAT DATA FROM ONE OR MORE INPUTS

CAN BE PROCESSED, READ,



# PIG LATIN

IS A PROCEDURAL **DATA FLOW LANGUAGE**

IT IMPLIES THAT **DATA FROM ONE OR MORE INPUTS**  
**CAN BE PROCESSED, READ,**

**AND THEN STORED TO ONE OR MORE**  
**OUTPUTS IN PARALLEL**



# PIG LATIN

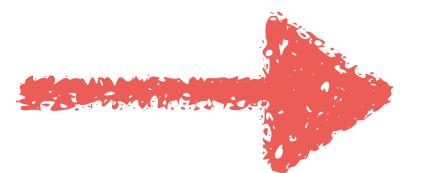
IS A PROCEDURAL DATA FLOW LANGUAGE

→ THERE ARE NO **IF** STATEMENTS AND **FOR** LOOPS LIKE WE SEE IN OBJECT ORIENTED OR FUNCTIONAL LANGUAGES



# PIG LATIN

IS A PROCEDURAL DATA FLOW LANGUAGE



IT DOES HAVE MOST OF THE USUAL DATA  
PROCESSING CONCEPTS THAT SQL HAS

THE SYNTAX WILL BE PRETTY DIFFERENT THOUGH



# PIG LATIN

IT DOES HAVE MOST OF THE USUAL DATA  
PROCESSING CONCEPTS THAT SQL HAS

THE SYNTAX WILL BE PRETTY DIFFERENT THOUGH

## SQL

```
SELECT function(column_name)  
FROM table GROUP BY column_name;
```

## PIG LATIN

```
FOREACH (GROUP table BY column_name)  
GENERATE function(table.column_name);
```

WE WILL COVER ALL OF IT IN DETAIL LATER



# PIG LATIN

LOOKS VERY SIMILAR TO SQL

BUT FUNDAMENTALLY THEY ARE VERY DIFFERENT



# PIG LATIN

SQL IS A QUERY LANGUAGE  
IT IS USED IN QUERYING

PIG LATIN IS A DATA-FLOW LANGUAGE  
IT FOCUSES ON DATA FLOW



# PIG LATIN

IN SQL, YOU DON'T WORRY ABOUT HOW THE QUERIES  
ARE EXECUTED

IN PIG LATIN, THE USER DESCRIBES EXACTLY  
HOW TO PROCESS THE DATA



# PIG LATIN

IN SQL, YOU DON'T WORRY ABOUT HOW THE QUERIES  
ARE EXECUTED

IN PIG LATIN CAN BE USED AS A FACTLY  
SQL SUBSTITUTE



# PIG LATIN

PIG LATIN CAN BE USED AS A SQL SUBSTITUTE

BUT, PIG LATIN VARIES TO A GREAT EXTENT AND  
THUS IT TAKES SOME TIME TO MASTER PIG

HIVEQL SHOULD BE USED IF THE PURPOSE  
OF DATA EXTRACTION IS ANALYSIS



# PIG LATIN

IN SQL, PIG IS DESIGNED WITH A LONG SERIES  
OF DATA OPERATIONS IN MIND

IN PIG LATIN, THE USER DESCRIBES EXACTLY  
HOW TO PROCESS THE DATA



# PIG LATIN

PIG IS DESIGNED WITH A LONG SERIES OF DATA OPERATIONS IN MIND

A PRIMARY USE CASE OF PIG IS  
TRADITIONAL EXTRACT TRANSFORM LOAD(ETL)  
DATA PIPELINES



# PIG LATIN

A PRIMARY USE CASE OF PIG IS  
TRADITIONAL EXTRACT TRANSFORM LOAD(ETL) DATA PIPELINES

WEBSITE DATA LOGS, SERVER DATA LOGS ALL  
CONTAIN HUMUNGOUS AMOUNT OF SUPER  
VALUABLE DATA IN UNSTRUCTURED FORMATS



# PIG LATIN

WEBSITE DATA LOGS, SERVER DATA LOGS ALL CONTAIN HUMUNGOUS AMOUNT OF SUPER VALUABLE DATA IN UNSTRUCTURED FORMATS

## PIG IS USED

- TO CLEAN THESE RECORDS WITH INCONSISTENT AND UNKNOWN SCHEMA
- TO PRECOMPUTE COMMON AGGREGATES BEFORE PUTTING INTO DATA WAREHOUSE



# PIG

is an open-source technology that

offers a high-level mechanism for parallel  
programming of MapReduce jobs

to be executed on Hadoop clusters



# PIG

is an open-source technology that

offers a high-level mechanism for parallel  
programming of MapReduce jobs

to be executed on Hadoop clusters



PIG

HADOOP

HDFS

MapReduce

YARN

PIG RUNS ON HADOOP



# PIG

## HADOOP

HDFS

MapReduce

YARN

PIG READS INPUT FILES FROM HDFS,  
USE HDFS TO STORE INTERMEDIATE FILES  
AND WRITES ITS OUTPUT TO HDFS



# PIG

## HADOOP

HDFS

MapReduce

YARN

PIG MANAGES DECOMPOSING THE  
OPERATIONS INTO MAPREDUCE JOBS



# PIG

## HADOOP

HDFS

# MapReduce

YARN

PIG PROVIDES SEVERAL ADVANTAGES  
OVER USING MAPREDUCE DIRECTLY



# PIG

## HADOOP

HDFS

## MapReduce

YARN

PIG PROVIDES SOME COMPLEX , NON-TRIVIAL  
IMPLEMENTATIONS OF STANDARD DATA OPERATIONS  
WHICH GIVES GREAT RETURNS IN EFFICIENCY

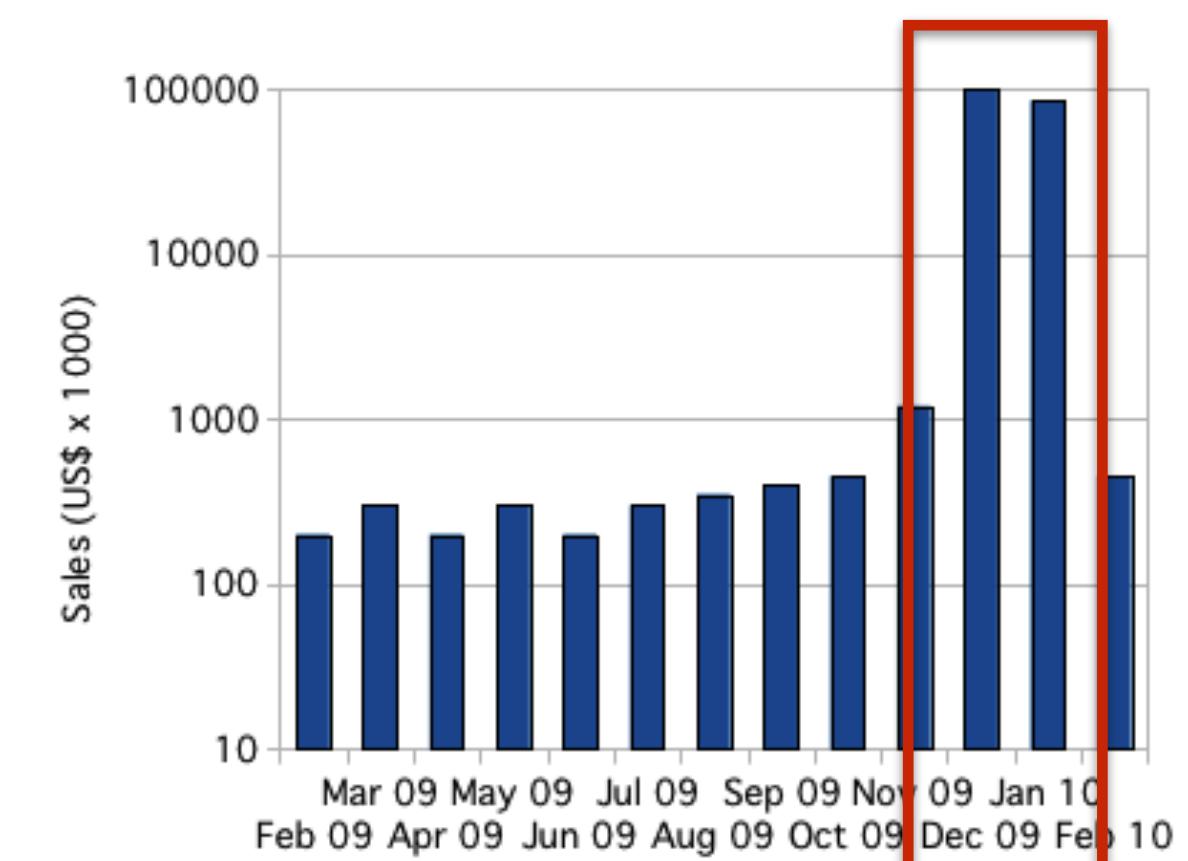


# PIG

PIG PROVIDES SOME COMPLEX , NON-TRIVIAL IMPLEMENTATIONS OF STANDARD DATA OPERATIONS WHICH GIVES GREAT RETURNS IN EFFICIENCY

CONSIDER THE CASE, WHERE A JOIN FUNCTION IS TO BE PERFORMED ON A SKEWED DATA SET

LET US ASSUME THE JOIN IS DONE ON MONTH KEY



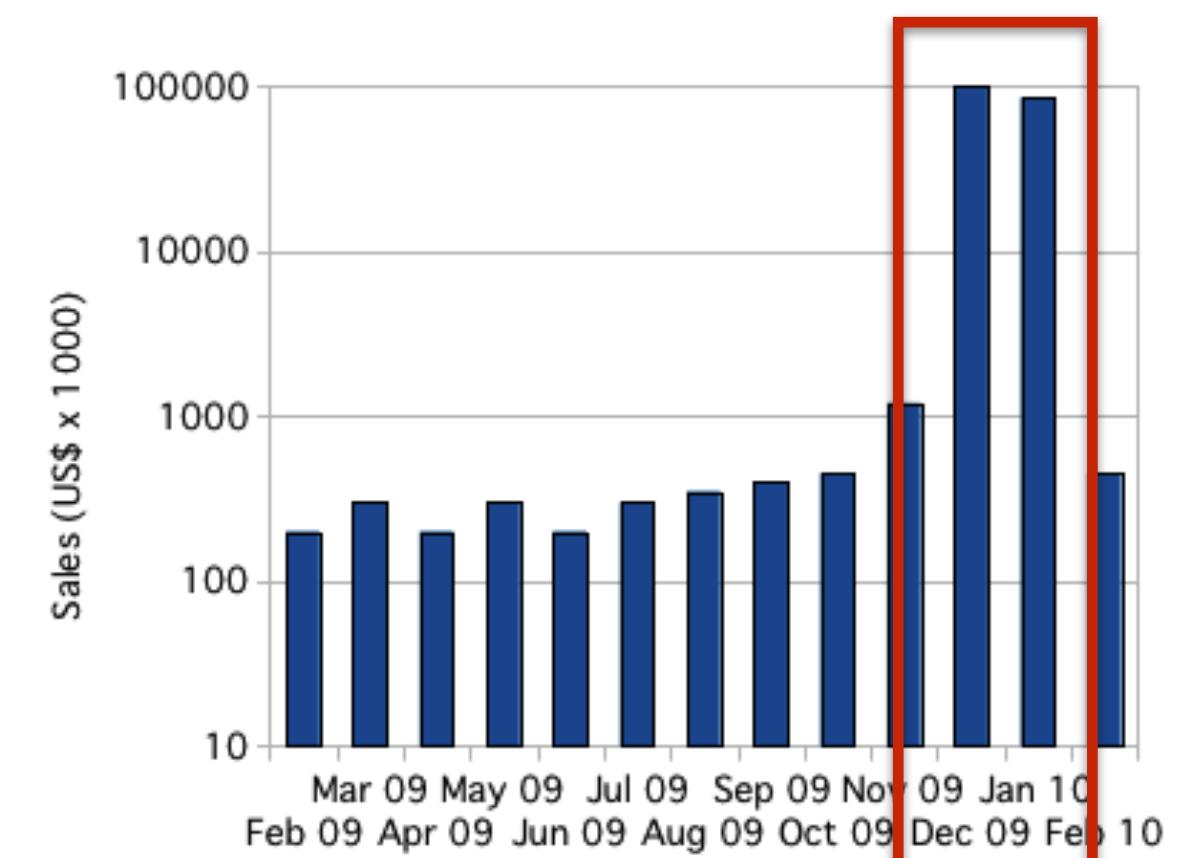


# PIG

PIG PROVIDES SOME COMPLEX , NON-TRIVIAL IMPLEMENTATIONS OF STANDARD DATA OPERATIONS WHICH GIVES GREAT RETURNS IN EFFICIENCY

LET US ASSUME THE JOIN IS DONE ON MONTH KEY

SOME REDUCERS WILL GET 10 OR MORE TIMES THE DATA THAN THE OTHER REDUCERS



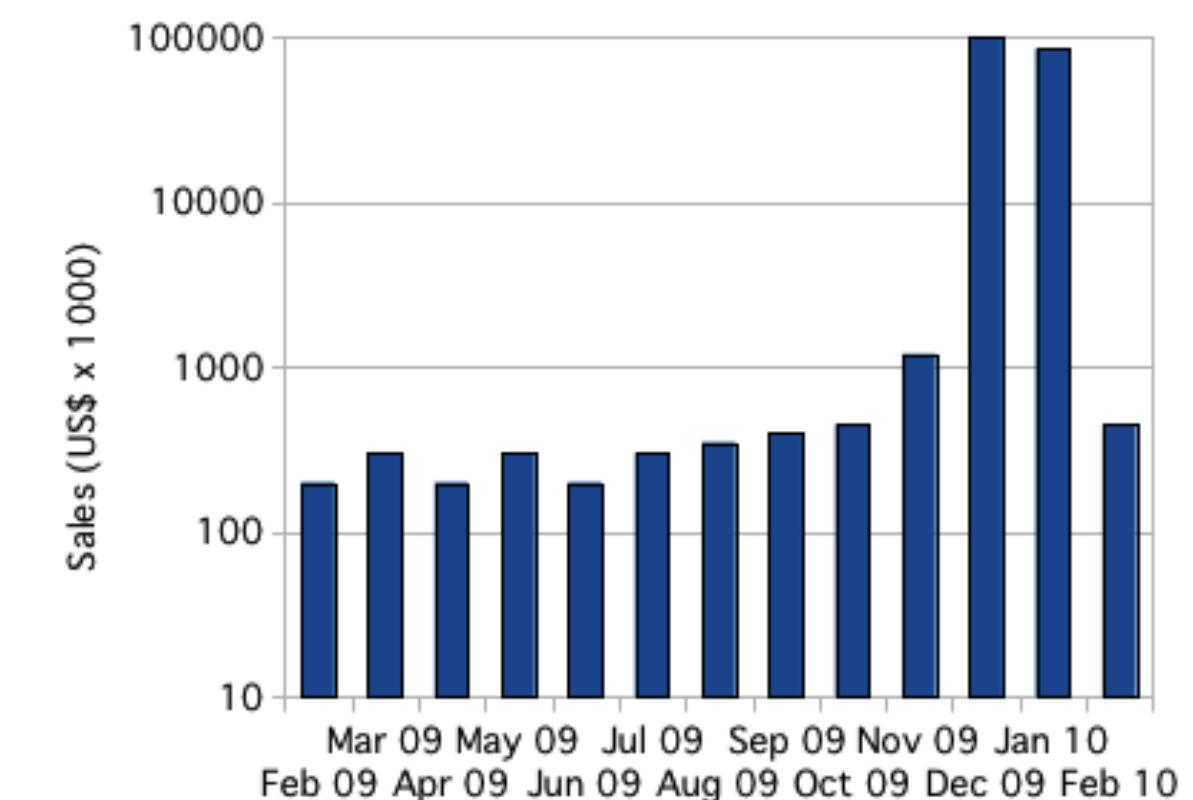


# PIG

PIG PROVIDES SOME COMPLEX , NON-TRIVIAL IMPLEMENTATIONS OF STANDARD DATA OPERATIONS WHICH GIVES GREAT RETURNS IN EFFICIENCY

LET US ASSUME THE JOIN IS DONE ON MONTH KEY

PIG HAS JOIN AND ORDER BY OPERATIONS THAT WILL REBALANCE THE REDUCERS





# PIG

PIG PROVIDES SOME COMPLEX , NON-TRIVIAL IMPLEMENTATIONS OF STANDARD DATA OPERATIONS WHICH GIVES GREAT RETURNS IN EFFICIENCY

WRITING ALL THIS IN MAPREDUCE OR JAVA WILL BE VERY TIME CONSUMING



# PIG

## HADOOP

HDFS

## MapReduce

YARN

PIG CAN ALSO ANALYSE THE PIG LATIN  
SCRIPT AND UNDERSTAND THE DATA FLOW



# PIG

## HADOOP

HDFS

## MapReduce

YARN

IT CAN DO EARLY ERROR CHECKING AND  
OPTIMISATIONS WHICH MAPREDUCE CAN'T



# PIG OVER MAPREDUCE

PIG PROVIDES SOME COMPLEX, NON-TRIVIAL  
IMPLEMENTATIONS OF STANDARD DATA OPERATIONS WHICH  
GIVES GREAT RETURNS IN EFFICIENCY

IT CAN DO EARLY ERROR CHECKING AND  
OPTIMISATIONS WHICH MAPREDUCE CAN'T



# PIG OVER MAPREDUCE

PIG PROVIDES SOME COMPLEX , NON-TRIVIAL  
IMPLEMENTATIONS OF STANDARD DATA OPERATIONS WHICH  
GIVES GREAT RETURNS IN EFFICIENCY

## WHERE CAN THIS BE USEFUL?

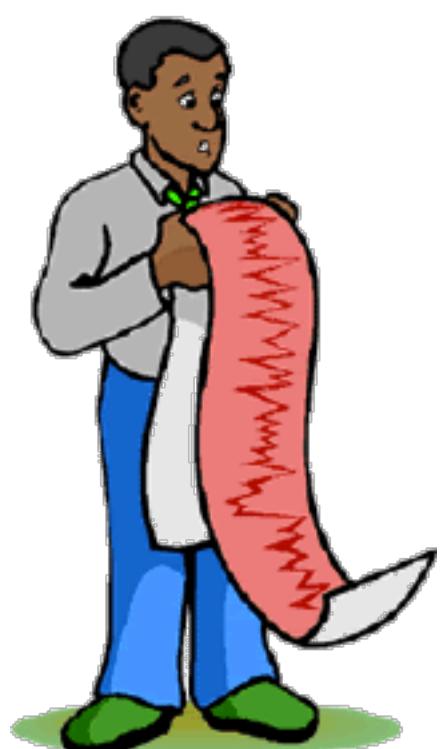
IT CAN DO EARLY ERROR CHECKING AND  
OPTIMISATIONS WHICH MAPREDUCE CAN'T



# PIG

SUPPOSE YOU ARE A RESEARCHER WORKING IN THE FIELD OF SEISMOLOGY

YOU ARE TRYING TO MAKE HIGH-PRECISION AND SENSITIVE DETECTORS



TO PREDICT EARTHQUAKES

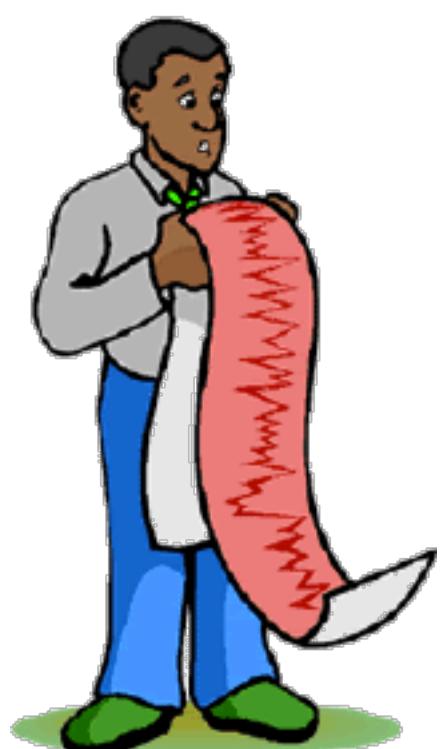




# PIG

SUPPOSE YOU ARE A RESEARCHER WORKING IN THE FIELD OF SEISMOLOGY

YOU ARE TRYING TO MAKE HIGH-PRECISION AND SENSITIVE DETECTORS



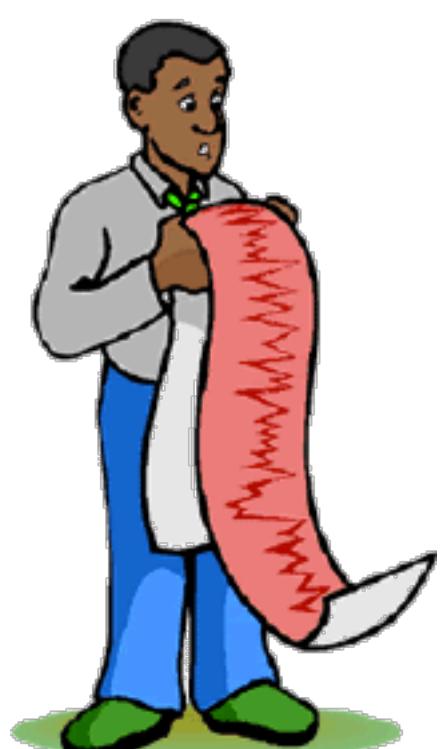
TO DO THIS, THE RESEARCHER HAS  
TO WORK WITH HUGE SEISMIC DATA



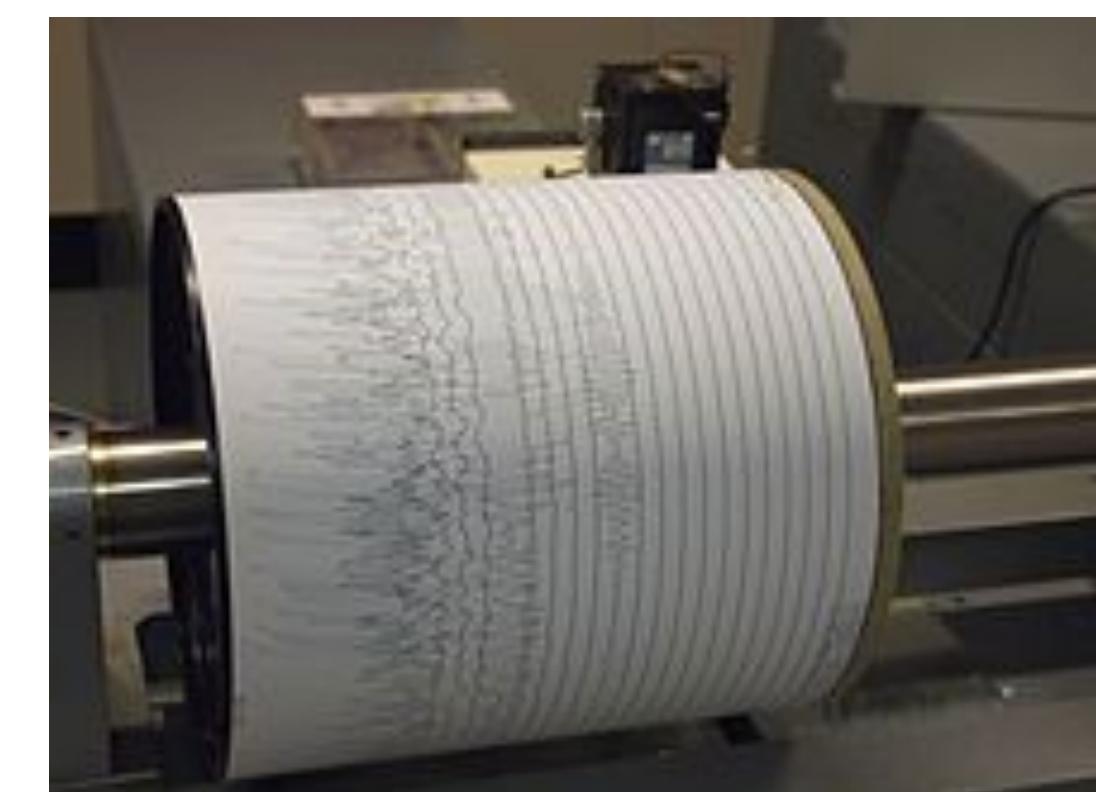
# PIG

SUPPOSE YOU ARE A RESEARCHER WORKING IN THE FIELD OF SEISMOLOGY

YOU ARE TRYING TO MAKE HIGH-PRECISION AND SENSITIVE DETECTORS



THE SEISMIC DATA THAT THEY  
STUDY IS CALLED SEISMOGRAM



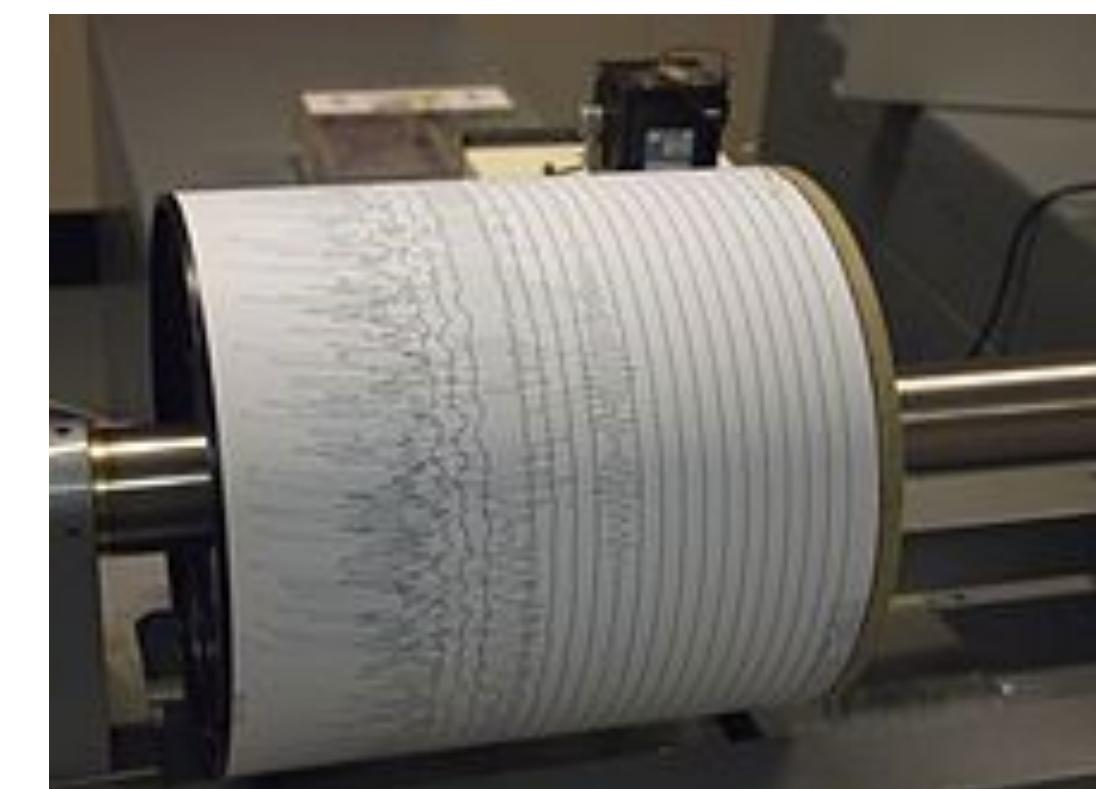


# PIG

SUPPOSE YOU ARE A RESEARCHER WORKING IN THE FIELD OF SEISMOLOGY  
YOU ARE TRYING TO **MAKE HIGH-PRECISION AND SENSITIVE DETECTORS**



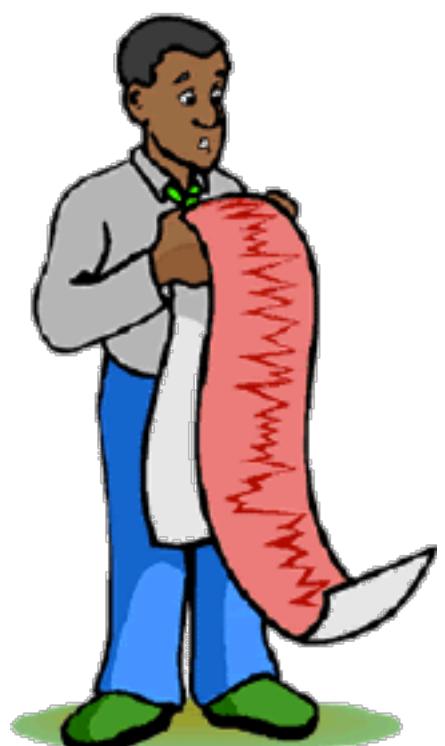
THE SEISMOGRAMS ARE DIGITAL  
TIME SERIES OF GROUND MOTION  
RECORDED BY **SEISMOMETERS**  
**INSTALLED AT THE SEISMIC STATIONS**





# PIG

SUPPOSE YOU ARE A RESEARCHER WORKING IN THE FIELD OF SEISMOLOGY  
YOU ARE TRYING TO **MAKE HIGH-PRECISION AND SENSITIVE DETECTORS**



A TECHNIQUE CALLED WAVEFORM  
**CROSS CORRELATION** IS APPLIED ON  
THESE SEISMOGRAMS

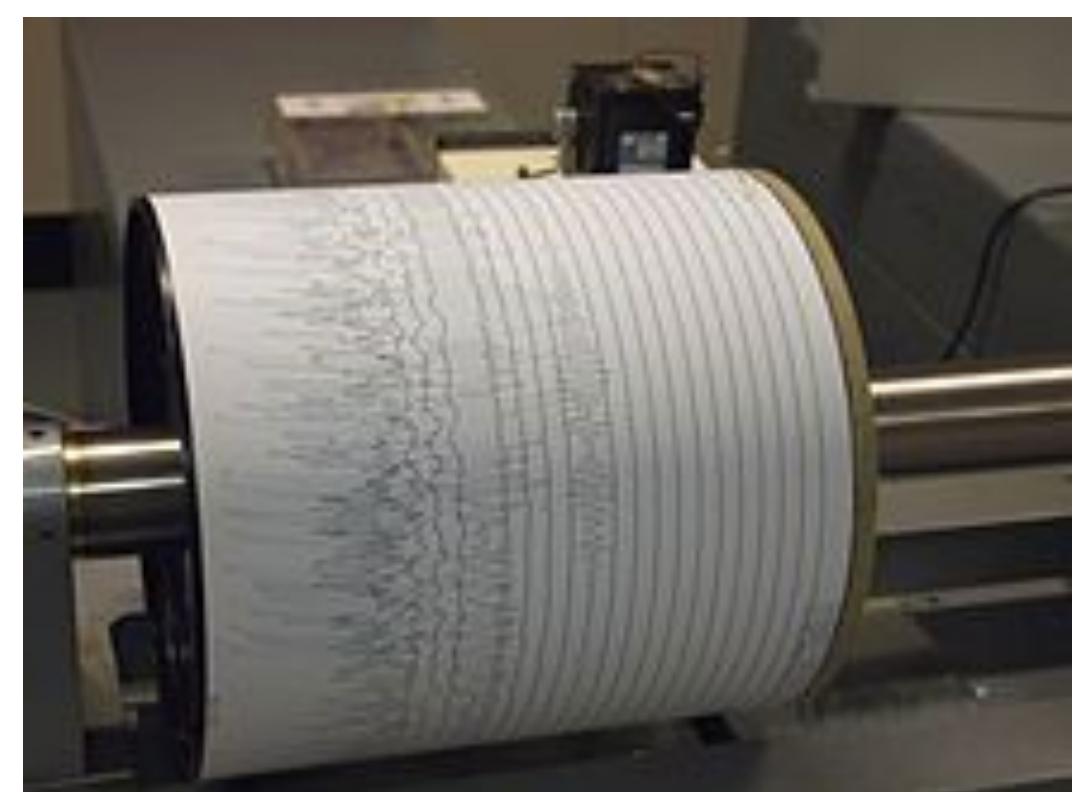
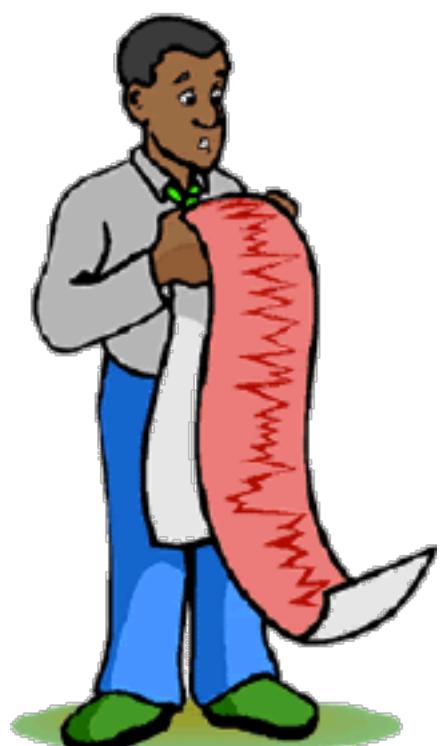


# PIG

SUPPOSE YOU ARE A RESEARCHER WORKING IN THE FIELD OF SEISMOLOGY

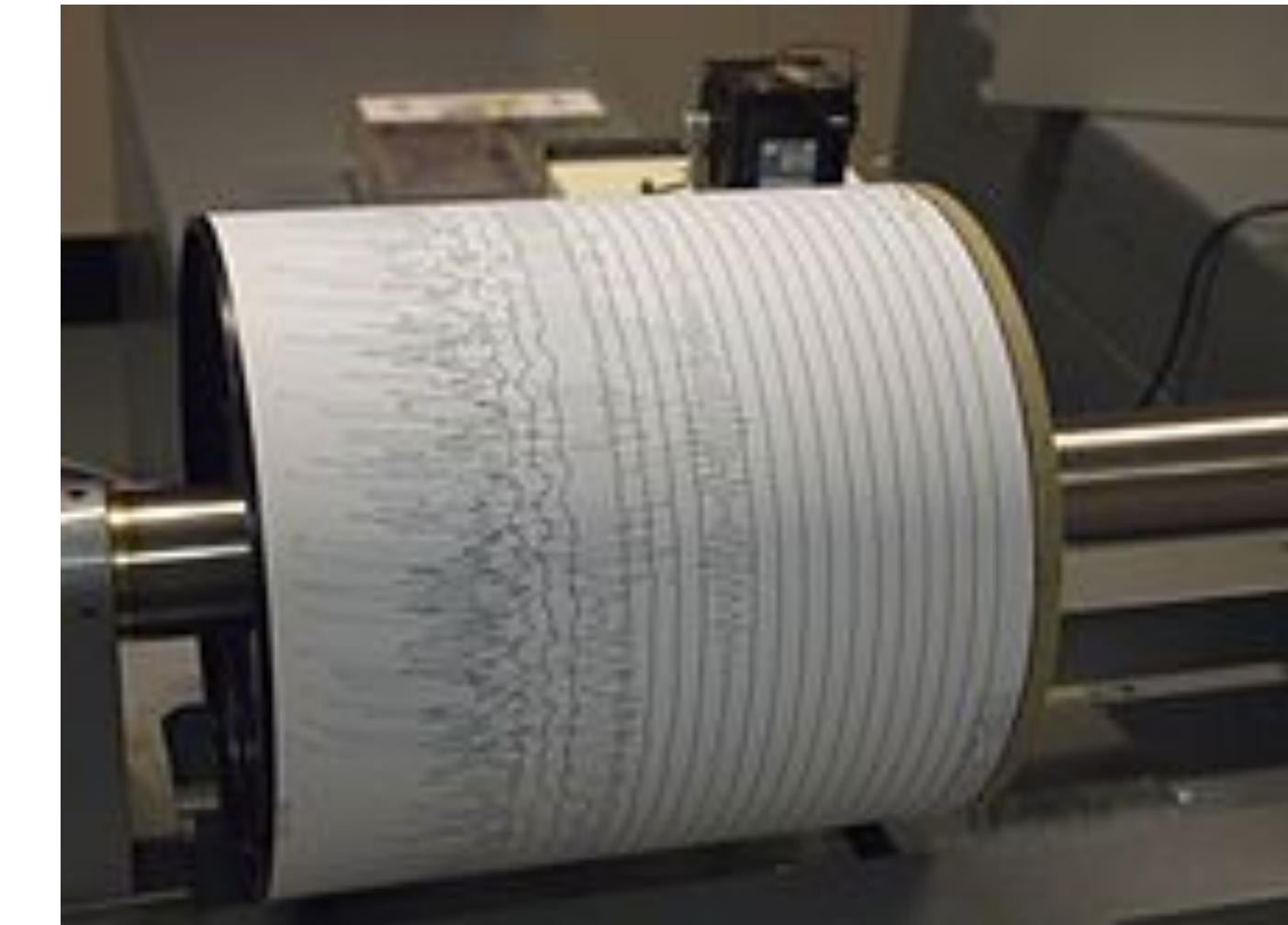
A TECHNIQUE CALLED WAVEFORM CROSS CORRELATION IS  
APPLIED ON THESE SEISMOGRAMS

THE DATASET HAS OVER 300 MILLION  
SEISMOGRAMS



THE DATASET HAS OVER 300  
MILLION SEISMOGRAMS

AND

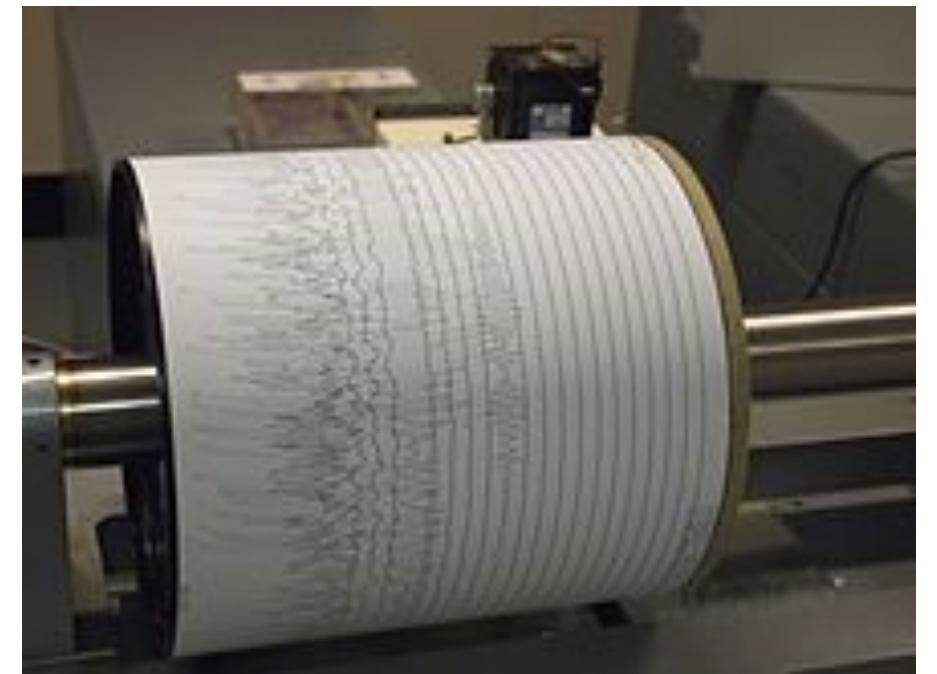


CORRELATIONS BETWEEN VECTORS WITH  
MILLIONS OF ELEMENTS NEED TO BE FOUND

THE DATASET HAS OVER 300 MILLION  
SEISMOGRAMS

AND

CORRELATIONS BETWEEN VECTORS WITH MILLIONS  
OF ELEMENTS NEED TO BE FOUND

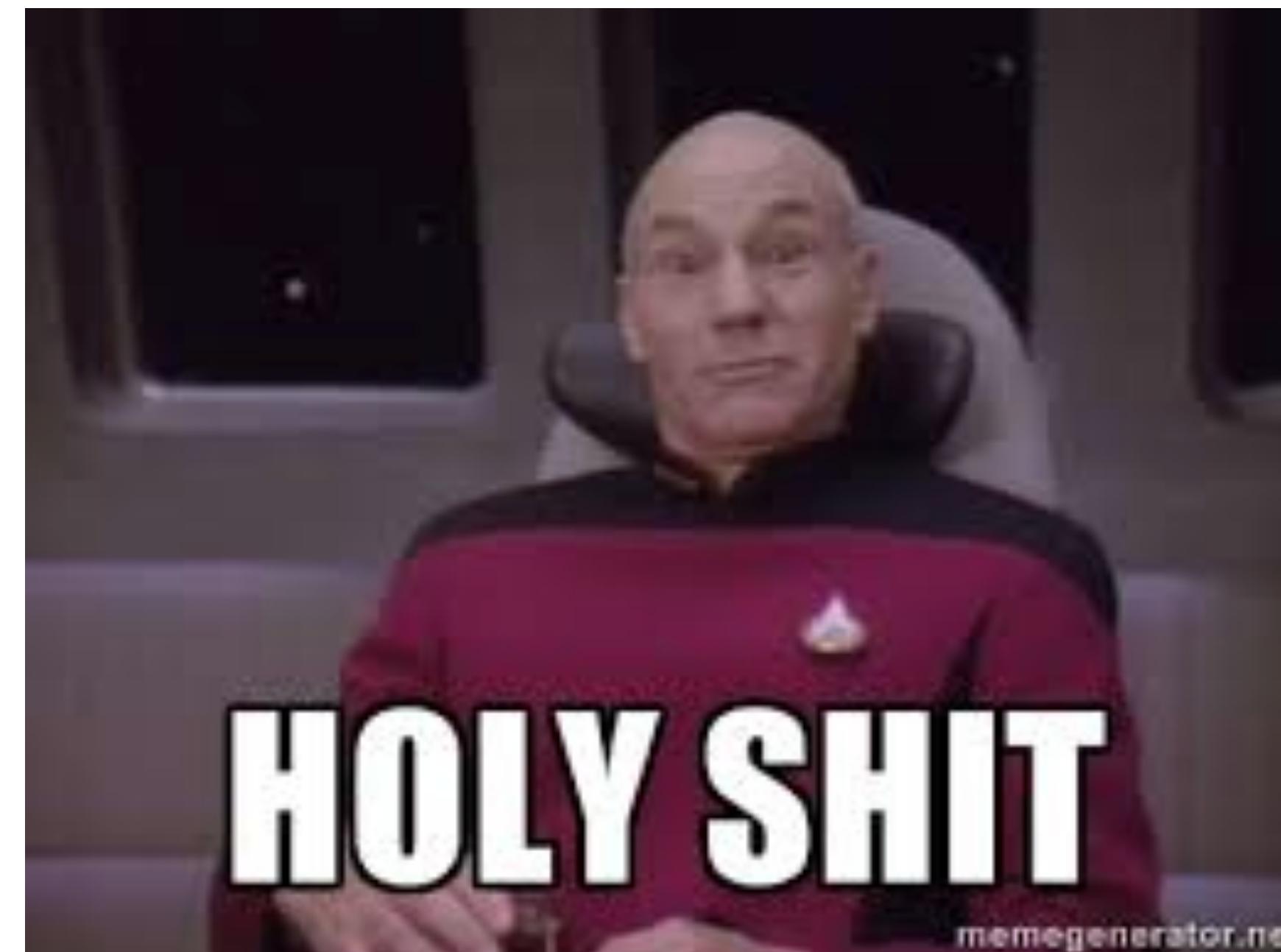


THE COMPLEXITY  
WILL BE  
 $O(N \times N)$   
 $=O(300M \times 300M)$

THE FIRST ATTEMPT INVOLVED 4 SERVERS WITH  
44 CORES AND 613 GB OF RAM

50 TB OF WAVEFORM DATA MANAGED BY A 2-HEAD HITACHI FILE  
SERVER WAS PROCESSED IN

42 DAYS

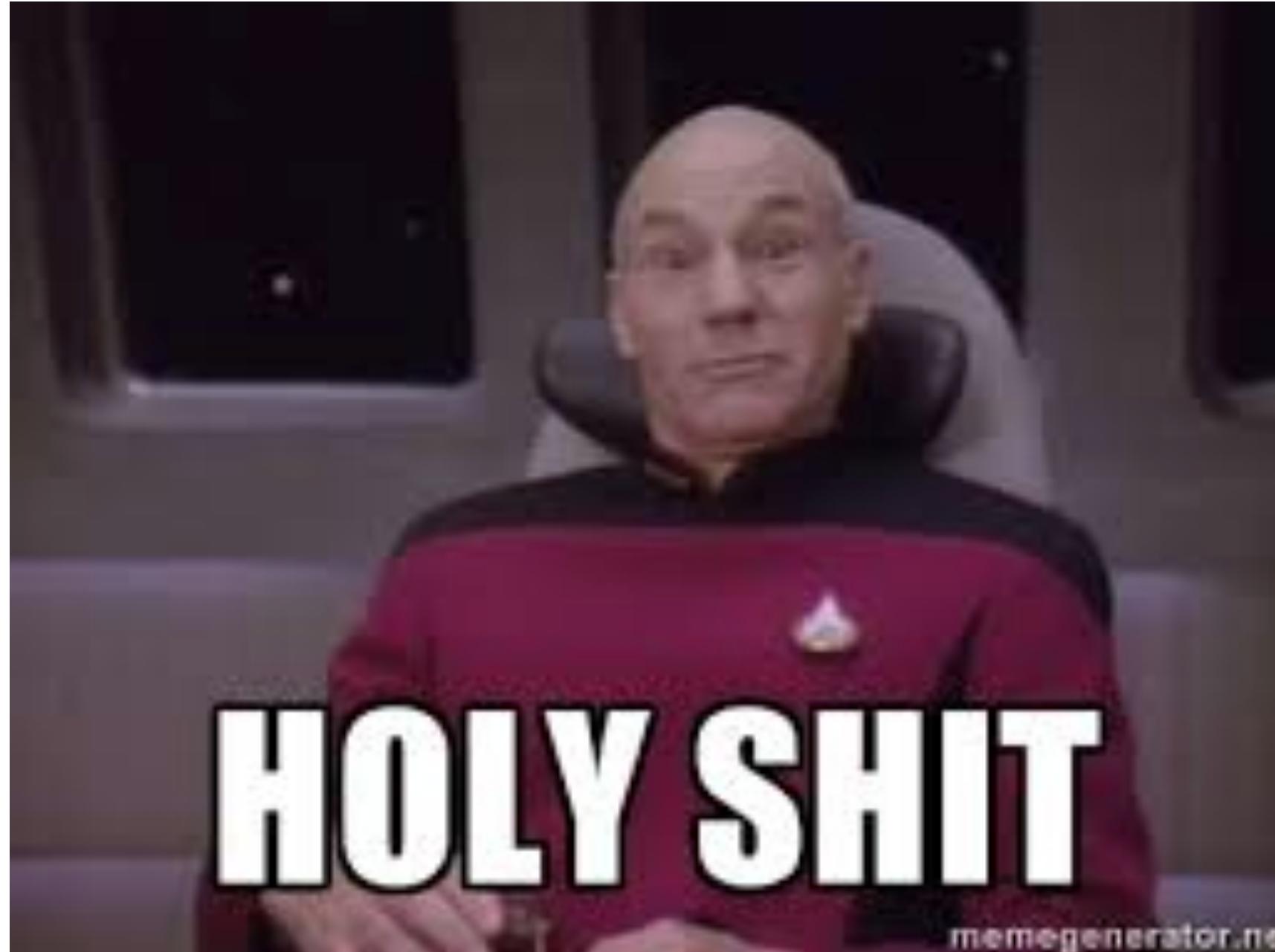


**CLEARLY, HADOOP COULD HELP**

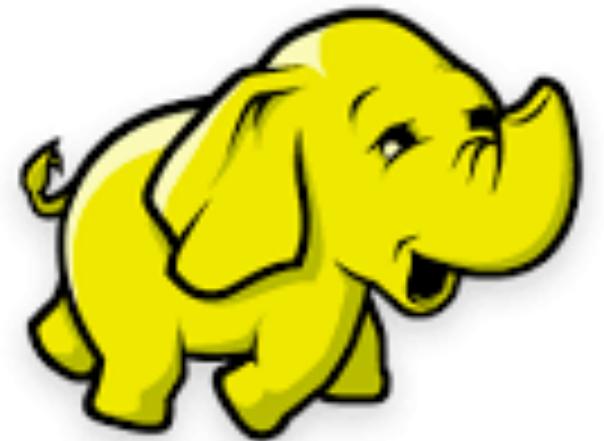


**TO THE RESCUE**

**42 DAYS**

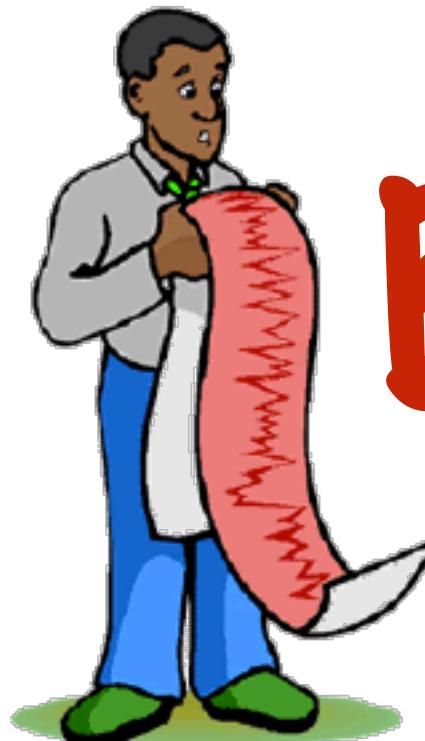


**CONTINUES**

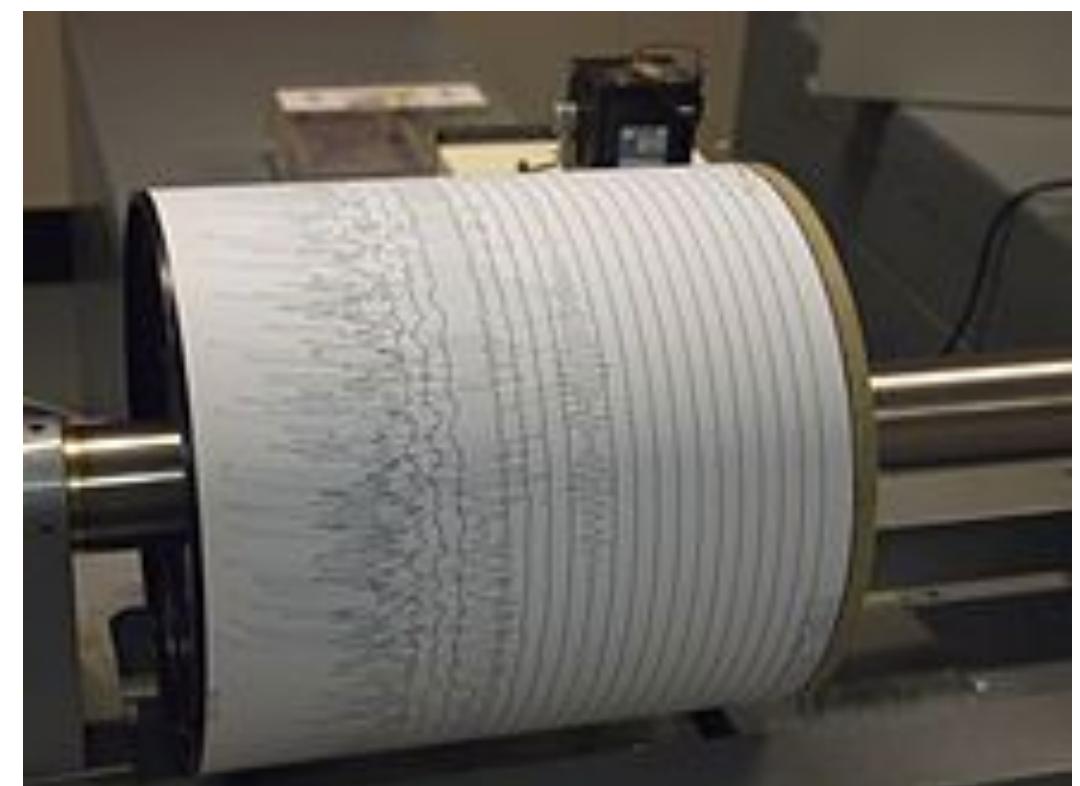


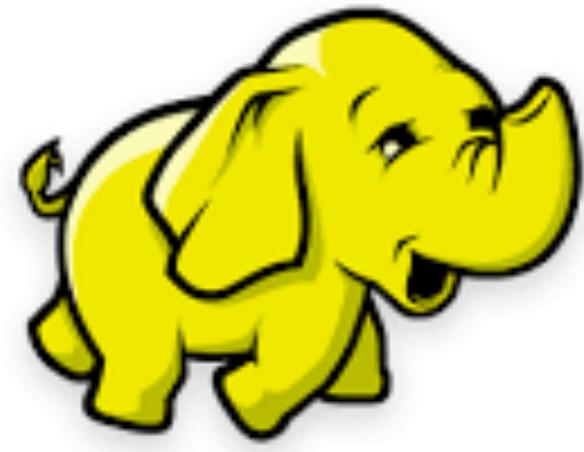
## TO THE RESCUE

THE RESEARCHER CAN RUN A SERIES OF  
MAPREDUCE JOBS ON A HADOOP CLUSTER



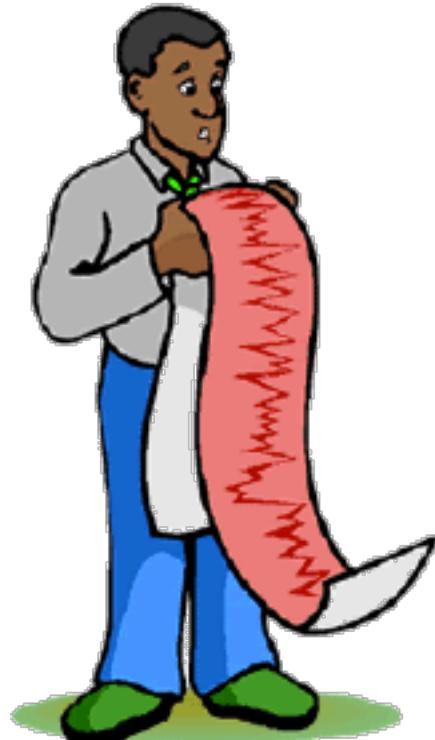
BUT WILL THAT BE ENOUGH?  
THE COMPLEXITY WAS  $O(300M \times 300M)$



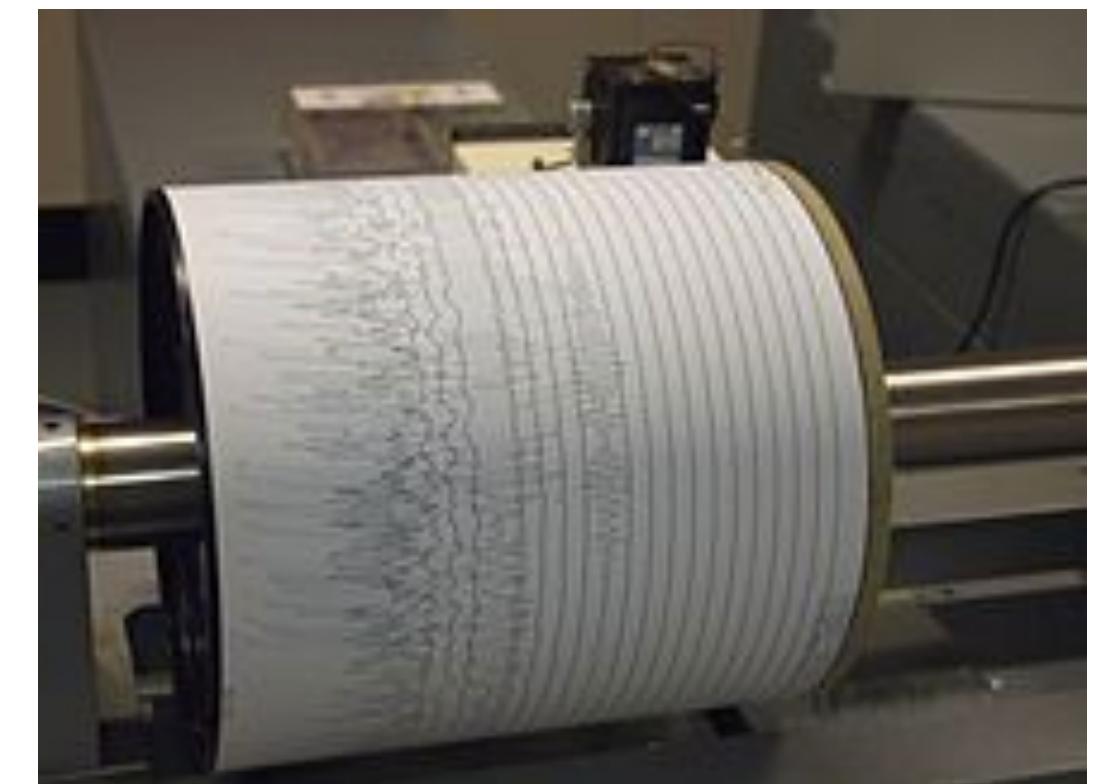


COMES TO THE RESCUE

THE RESEARCHERS FURTHER  
OPTIMISED THESE MAPREDUCE JOBS



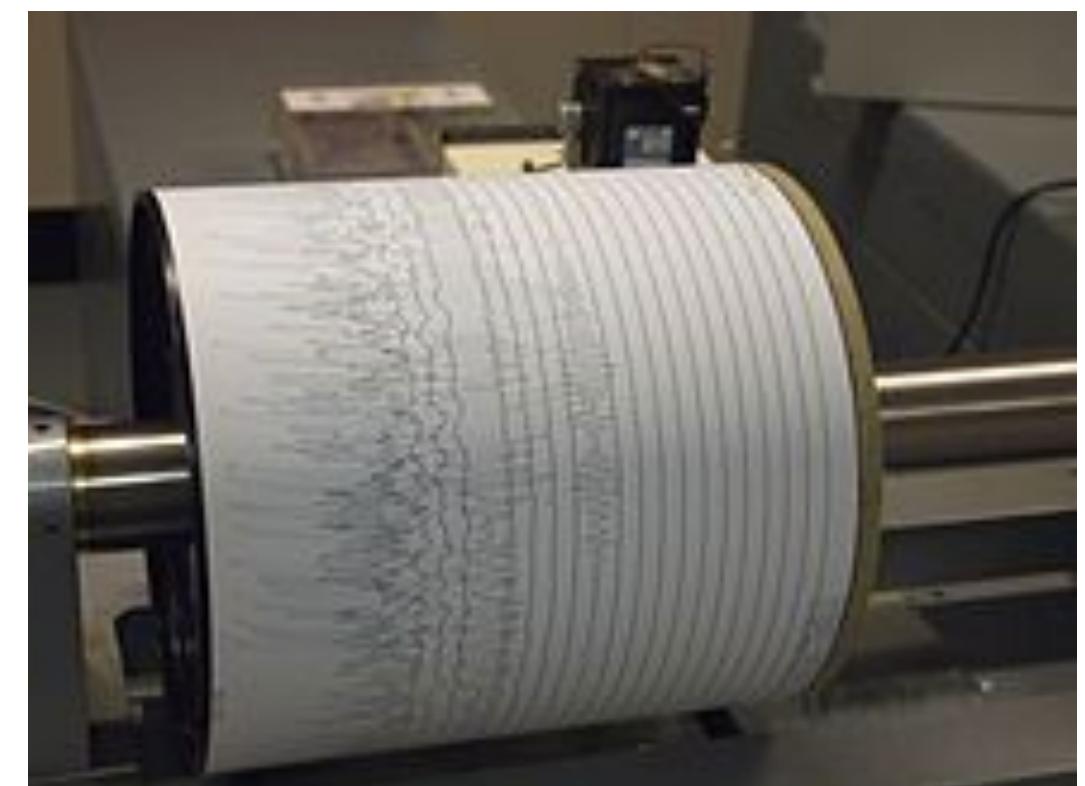
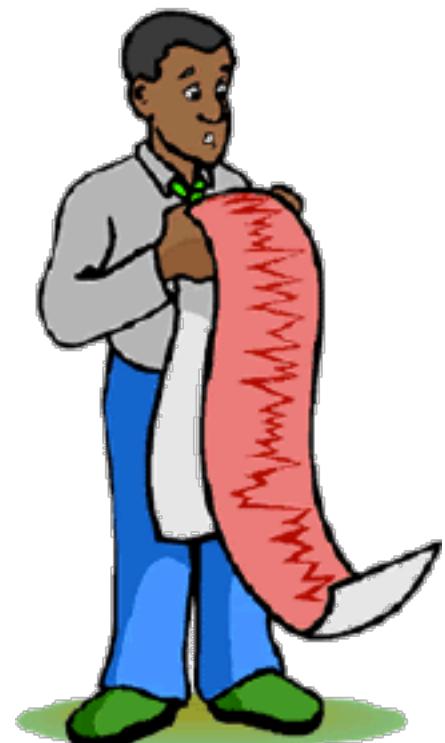
USING





COMES TO THE RESCUE

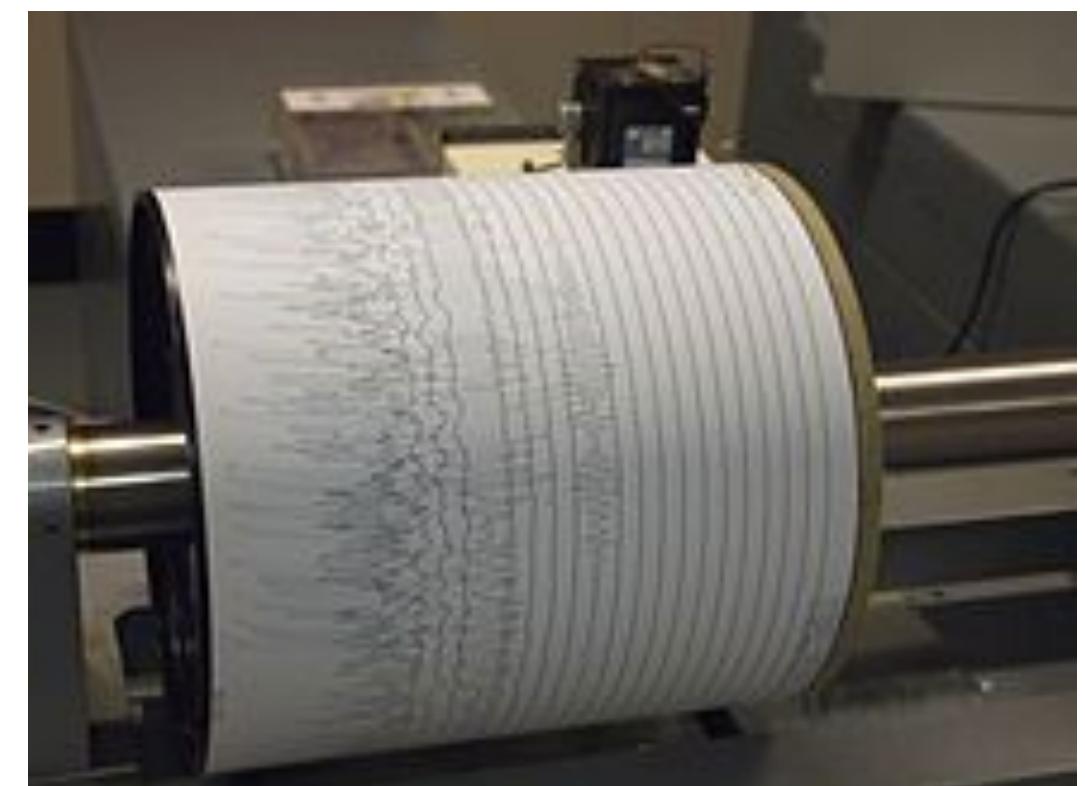
PIG DATA FLOWS AND THEIR IMPLEMENTATIONS  
ARE AT CONCEPTUALLY HIGHER LEVEL THAN THE  
SERIES OF MAPREDUCE JOBS

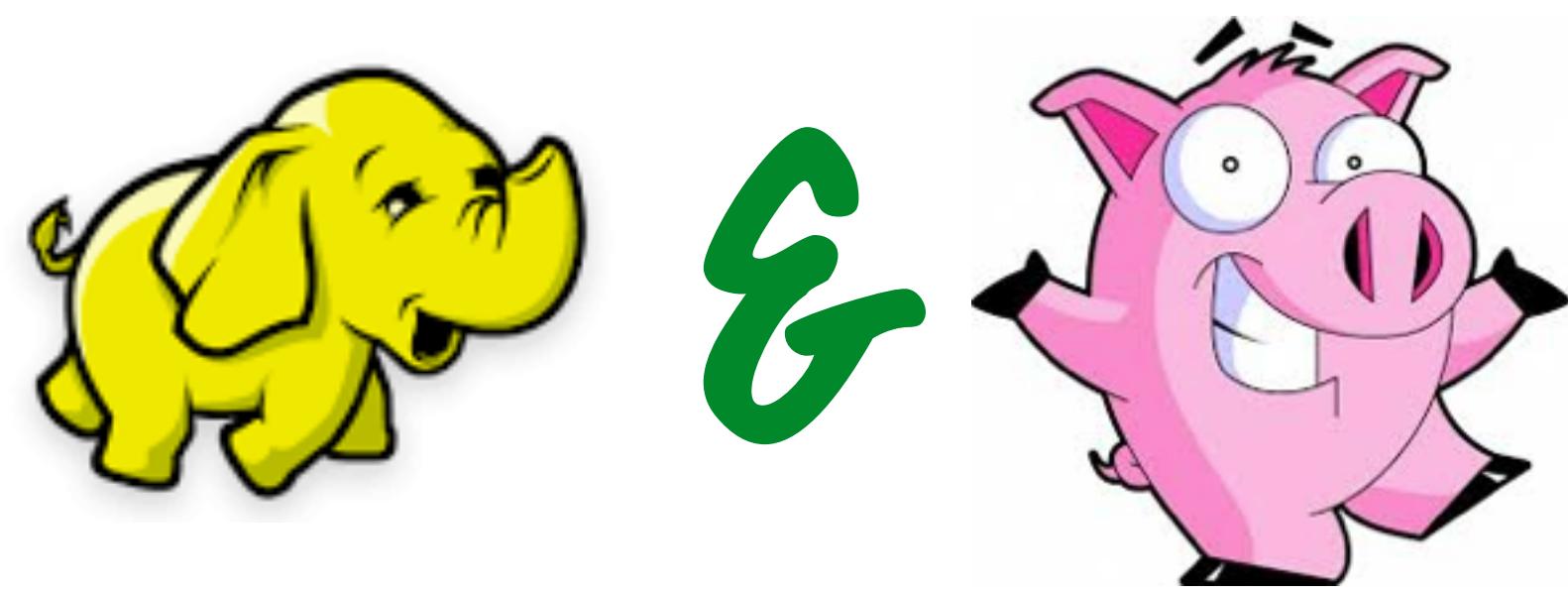




COMES TO THE RESCUE

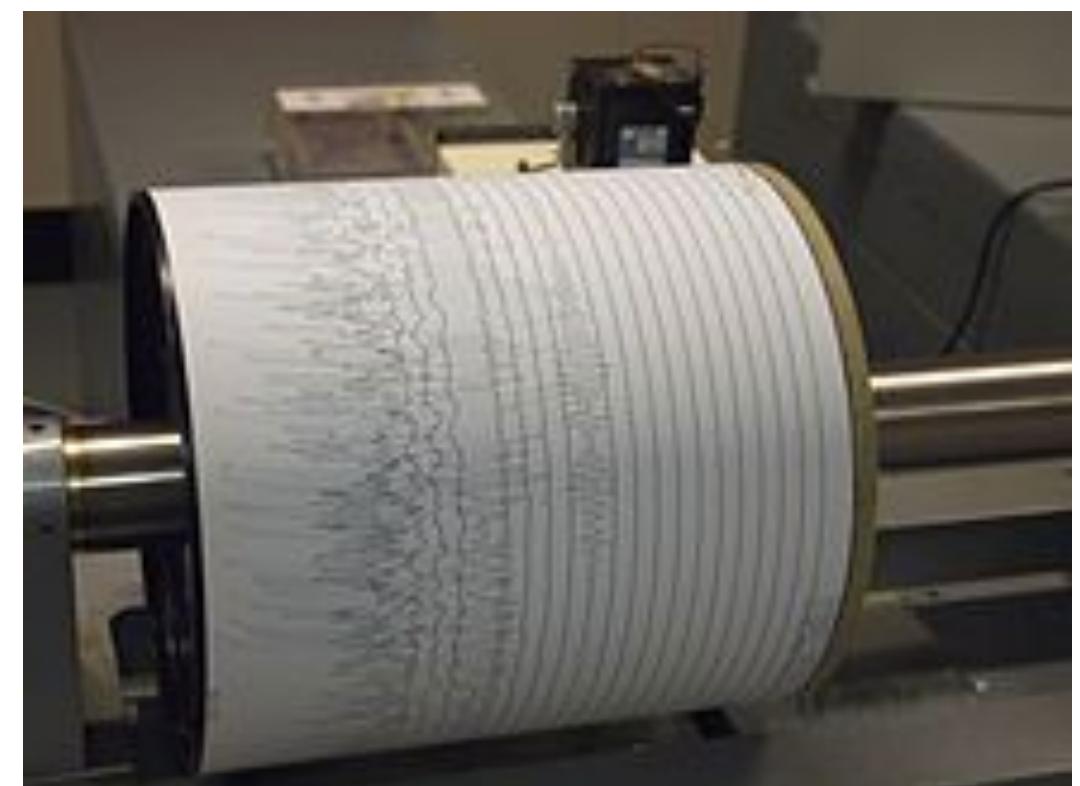
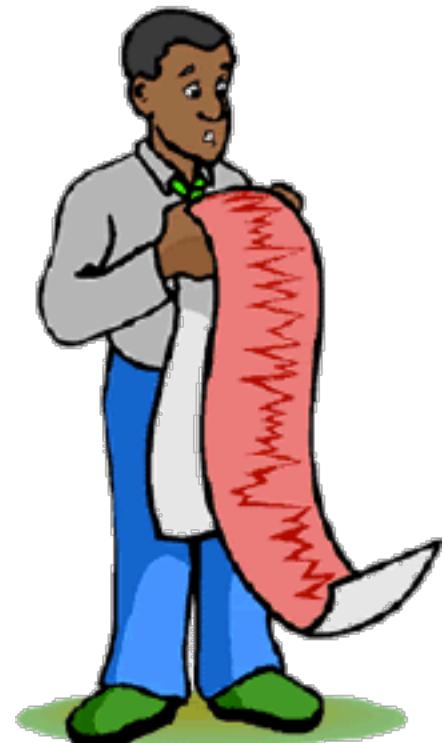
THE RESEARCHER USED PIG FOR THE  
IMPLEMENTATION OF THE CORRELATOR TO SPEED  
UP THE CALCULATION TIME





COMES TO THE RESCUE

**IMPROVED THE PERFORMANCE OF THE SYSTEM  
BY A FACTOR OF 19 ON THE GIVEN DATASET**





# PIG

CAN EXECUTE ITS HADOOP JOBS NOT  
JUST ON MAPREDUCE BUT ALSO ON

**TEZ**  APACHE TEZ

 APACHE SPARK



PIG &



HBASE

**HBASE IS NON-RELATIONAL  
DISTRIBUTED DATABASE**

**IT RUNS ON TOP OF HADOOP**

**AND IS WRITTEN IN JAVA**



PIG &



HBASE

**HBASE IS COLUMN-ORIENTED KEY-VALUE DATA STORE**

PRODUCT ID	CATEGORY	PRODUCT_NAME	PRICE	SIZE
1	Book	Bible	\$7	
2	Mobile	Moto G4	\$300	
3	T-Shirt	Mynta	\$8	M



PRODUCT ID	ATTRIBUTE	ATTRIBUTE VALUE
1	CATEGORY	Book
2	CATEGORY	Mobile
3	CATEGORY	T-Shirt
1	PRODUCT_NAME	Bible
2	PRODUCT_NAME	Moto G4
3	PRODUCT_NAME	Mynta
1	PRICE	\$7
2	PRICE	\$300
3	PRICE	\$8
3	SIZE	M



PIG &



HBASE

PRODUCT_ID	CATEGORY	PRODUCT_NAME	PRICE	SIZE
1	Book	Bible	\$7	
2	Mobile	Moto G4	\$300	
3	T-Shirt	Myntra	\$8	M



PRODUCT_ID	ATTRIBUTE	ATTRIBUTE_VALUE
1	CATEGORY	Book
2	CATEGORY	Mobile
3	CATEGORY	T-Shirt
1	PRODUCT_NAME	Bible
2	PRODUCT_NAME	Moto G4
3	PRODUCT_NAME	Myntra
1	PRICE	\$7
2	PRICE	\$300
3	PRICE	\$8
3	SIZE	M

HBASE IS HANDY WHEN THE DATA CONTAINS

TOO MANY (AND DYNAMICALLY CHANGING!) COLUMNS



PIG &



HBASE

PRODUCT_ID	CATEGORY	PRODUCT_NAME	PRICE	SIZE
1	Book	Bible	\$7	
2	Mobile	Moto G4	\$300	
3	T-Shirt	Myntra	\$8	M



PRODUCT_ID	ATTRIBUTE	ATTRIBUTE_VALUE
1	CATEGORY	Book
2	CATEGORY	Mobile
3	CATEGORY	T-Shirt
1	PRODUCT_NAME	Bible
2	PRODUCT_NAME	Moto G4
3	PRODUCT_NAME	Myntra
1	PRICE	\$7
2	PRICE	\$300
3	PRICE	\$8
3	SIZE	M

HBASE IS HANDY WHEN THESE ATTRIBUTES ARE DYNAMIC  
CONSIDER THE ATTRIBUTE SIZE HERE WHICH APPLIES ONLY TO CERTAIN PRODUCT CATEGORY



PIG &



HBASE

**HBASE IS WELL-SUITED FOR FASTER READ  
AND WRITE OPERATIONS ON LARGE  
DATASETS WITH HIGH THROUGHPUT AND  
LOW INPUT/OUTPUT LATENCY**



PIG &



HBASE

**HBASE IS USED BY SEVERAL  
MAJOR TECH FIRMS**





PIG &



HBASE

**PIG** WAS NAMED BECAUSE ITS CREATORS FELT IT HAD  
SOME PORCINE ATTRIBUTES

**PIGS ARE OMNIVOROUS**

**PIG CAN OPERATE ON ANY TYPE OF DATA**

**PIG CAN OPERATE ON DATA THAT LIES IN HBASE**



PIG &



HBASE

PIG PROVIDES HBASESTORAGE (LOAD  
FUNCTION) TO READ DATA FROM  
AND WRITE DATA TO HBASE TABLES

IT IS NOT AS FAST AS LOADING DATA FROM HDFS



PIG VS





PIG VS



HIVE AND PIG SEEM TO HAVE SIMILAR  
APPLICATIONS AND FUNCTIONALITY

LEADING FOLKS TO WONDER  
HOW DO THEY STACK UP?



PIG VS



TO (OVER-)SIMPLIFY A BIT



DATA FROM  
TECH SYSTEMS

ETL  
→  
USING PIG



DATA  
WAREHOUSE

REPORTING  
→  
USING HIVE





DATA FROM  
TECH SYSTEMS

ETL  
USING PIG



DATA  
WAREHOUSE

REPORTING  
USING HIVE





DATA FROM  
TECH SYSTEMS

ETL

USING PIG



DATA  
WAREHOUSE

REPORTING

USING HIVE



PIG IS USED BY  
DEVELOPERS

HIVE IS USED BY  
ANALYSTS





DATA FROM  
TECH SYSTEMS

ETL

USING PIG



DATA  
WAREHOUSE

REPORTING

USING HIVE



PIG HAS A  
PROCEDURAL DATA  
FLOW LANGUAGE

HIVE HAS A DECLARATIVE  
SQL LIKE LANGUAGE



PIG VS



TO (OVER-)SIMPLIFY A BIT



DATA FROM  
TECH SYSTEMS

ETL  
→  
USING PIG



DATA  
WAREHOUSE

REPORTING  
→  
USING HIVE

