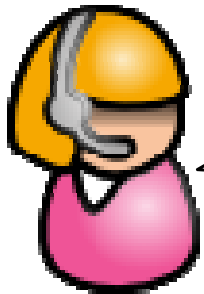




Hadoop只支援用**Java**開發嘛？
Is Hadoop only support Java ?

總不能全部都重新設計吧？如何與舊系統相容？

Can Hadoop work with existing software ?

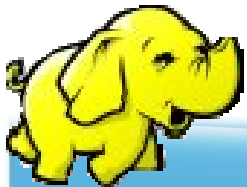


可以跟資料庫結合嘛？

Can Hadoop work with Databases ?

開發者們有聽到大家的需求.....

Yes, we hear the feedback of developers ...





Top

Common

Chukwa

HBase

HDFS

Hive

MapReduce

Pig

ZooKeeper

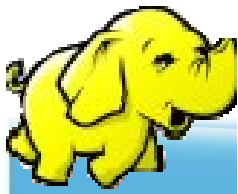
About

- Welcome
- Who We Are?
- Mailing Lists

Welcome to Apache Hadoop!

二、Hadoop 相關子專案

- **Hadoop Common:** The common utilities that support the other Hadoop subprojects.
- **HDFS:** A distributed file system that provides high throughput access to application data.
- **MapReduce:** A software framework for distributed processing of large data sets on compute clusters.



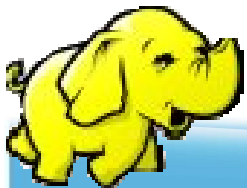
Hadoop 相關子專案

- **Chukwa**: A data collection system for managing large distributed systems.
- **HBase**: A scalable, distributed database that supports structured data storage for large tables.
- **Hive**: A data warehouse infrastructure that provides data summarization and ad hoc querying.
- **Pig**: A high-level data-flow language and execution framework for parallel computation.
- **ZooKeeper**: A high-performance coordination service for distributed applications.

Hadoop 生態系 (Ecosystem)

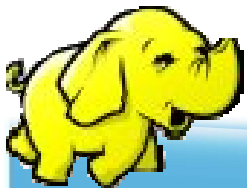
Pig	Chukwa	Hive	HBase
MapReduce		HDFS	ZooKeeper
Hadoop Core (Hadoop Common)			Avro

Source: *Hadoop: The Definitive Guide*



Avro

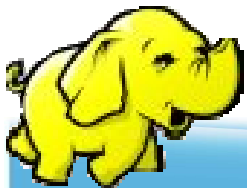
- Avro is a **data serialization system**.
- It provides:
 - *Rich data structures.*
 - *A compact, fast, binary data format.*
 - *A container file, to store persistent data.*
 - *Remote procedure call (RPC).*
 - *Simple integration with dynamic languages.*
- For more detail, please check the official document:
<http://avro.apache.org/docs/current/>



Zoo Keeper

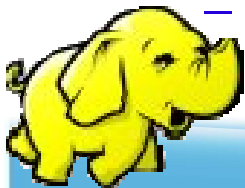


- <http://hadoop.apache.org/zookeeper/>
- ZooKeeper is a **centralized service** for maintaining configuration information, naming, providing distributed synchronization, and providing group services. All of these kinds of services are used in some form or another by distributed applications.



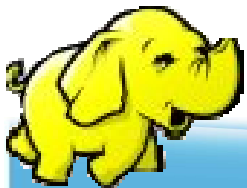
Pig

- <http://hadoop.apache.org/pig/>
- Pig is a platform for analyzing large data sets that consists of a **high-level language** for expressing data analysis programs, coupled with infrastructure for evaluating these programs.
- Pig's infrastructure layer consists of a **compiler** that produces sequences of **Map-Reduce programs**
- Pig's language layer currently consists of a textual language called **Pig Latin**, which has the following key properties:
 - Ease of programming
 - Optimization opportunities
 - Extensibility



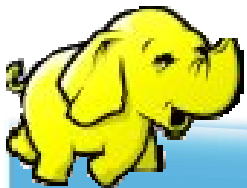
Hive

- <http://hadoop.apache.org/hive/>
- Hive is a **data warehouse** infrastructure built on top of Hadoop that provides tools to enable easy **data summarization**, **adhoc querying** and analysis of large datasets data stored in Hadoop files.
- **Hive QL** is based on SQL and enables users familiar with SQL to query this data.



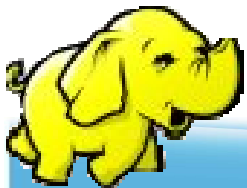
HBase

- HBase is a distributed **column-oriented database** built on top of HDFS.
- A distributed data store that can scale horizontally to 1,000s of commodity servers and **petabytes** of indexed storage.
- Designed to operate on top of the Hadoop distributed file system (**HDFS**) or Kosmos File System (**KFS**, aka Cloudstore) for scalability, fault tolerance, and high availability.
- Integrated into the Hadoop **map-reduce** platform and paradigm.



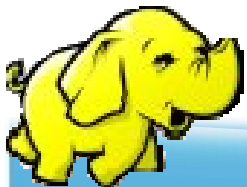
Chukwa

- <http://hadoop.apache.org/chukwa/>
- Chukwa is an open source **data collection system** for monitoring large distributed systems.
- built on top of HDFS and Map/Reduce framework
- includes a flexible and powerful toolkit for displaying, monitoring and analyzing res make the best use of the collected data.



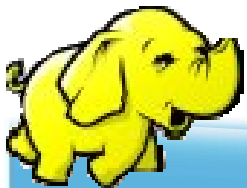
Mahout

- <http://mahout.apache.org/>
- Mahout is a scalable **machine learning libraries**.
- implemented on top of Apache Hadoop using the map/reduce paradigm.
- Mahout currently has
 - Collaborative Filtering
 - User and Item based recommenders
 - **K-Means, Fuzzy K-Means**
 - Mean Shift clustering
 - More ...



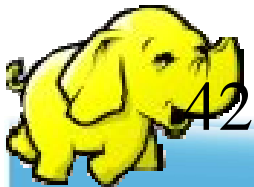
Hadoop 只支援 Java 嗎？

- Although the Hadoop framework is implemented in JavaTM, **Map/Reduce applications need not be written in Java.**
- **Hadoop Streaming** is a utility which allows users to create and run jobs with any executables (e.g. shell utilities) as the mapper and/or the reducer.
- **Hadoop Pipes** is a SWIG-compatible C++ API to implement Map/Reduce applications (non JNITM based).



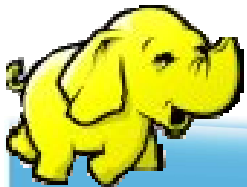
Hadoop Pipes (C++, Python)

- Hadoop Pipes allows C++ code to use Hadoop DFS and map/reduce.
- The C++ interface is "swigable" so that interfaces can be generated for python and other scripting languages.
- For more detail, check the API Document of org.apache.hadoop.mapred.pipes
- You can also find example code at
 - hadoop-*/src/examples/pipes
- About the pipes C++ WordCount example code:
 - <http://wiki.apache.org/hadoop/C++WordCount>



Hadoop Streaming

- Hadoop Streaming is a utility which allows users to create and run Map-Reduce jobs **with any executables (e.g. Unix shell utilities)** as the mapper and/or the reducer.
- It's useful when you need to run **existing program** written in shell script, perl script or even PHP.
- Note: both the **mapper** and the **reducer** are executables that read the input from **STDIN** (line by line) and emit the output to **STDOUT**.
- For more detail, check the official document of Hadoop Streaming



Running Hadoop Streaming

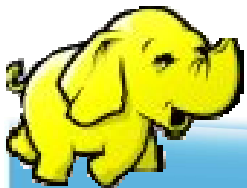
```
jazz@hadoop:~$ hadoop jar hadoop-streaming.jar -help
```

```
Usage: $HADOOP_HOME/bin/hadoop [--config dir] jar \  
      $HADOOP_HOME/hadoop-streaming.jar [options]
```

Options:

```
-input      <path>          DFS input file(s) for the Map step  
-output     <path>          DFS output directory for the Reduce step  
-mapper     <cmd|JavaClassName>    The streaming command to run  
-combiner   <JavaClassName> Combiner has to be a Java class  
-reducer    <cmd|JavaClassName>    The streaming command to run  
-file       <file>          File/dir to be shipped in the Job jar file  
-dfs        <h:p>|local  Optional. Override DFS configuration  
-jt         <h:p>|local  Optional. Override JobTracker configuration  
-additionalconfspec specfile  Optional.  
-inputformat  
    TextInputFormat(default) | SequenceFileAsTextInputFormat | JavaClassNam  
    e Optional.  
-outputformat TextOutputFormat(default) | JavaClassName  Optional.
```

... More ...

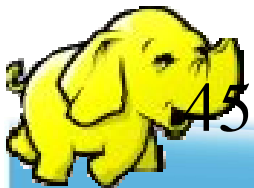


Hadoop Streaming 範例 (1)

```
hadoop:~$ hadoop fs -rmr input output
```

```
hadoop:~$ hadoop fs -put /etc/hadoop/conf  
input
```

```
hadoop:~$ hadoop jar hadoop-streaming.jar  
-input input -output output -mapper  
/bin/cat -reducer /usr/bin/wc
```



Hadoop Streaming 範例 (2)

```
hadoop:~$ echo "sed -e \"s/ /\n/g\" | grep ." >  
streamingMapper.sh  
hadoop:~$ echo "uniq -c | awk '{print \$2 \"\t\" \$1}'" >  
> streamingReducer.sh  
hadoop:~$ chmod a+x streamingMapper.sh  
hadoop:~$ chmod a+x streamingReducer.sh  
hadoop:~$ hadoop fs -put /etc/hadoop/conf input  
hadoop:~$ hadoop jar hadoop-streaming.jar -input input  
-output output -mapper streamingMapper.sh -reducer  
streamingReducer.sh -file streamingMapper.sh  
-file streamingReducer.sh
```

