#### Hadoop 程式設計

# 四、寫Code環境準備

4.A: Console 端

編譯與執行

4.B:透過 Eclipse 開發與測試運算



#### Java之編譯與執行

#### 1. 編譯

- 2. 封裝
  - lacktriangle jar  $_{\Delta}$  -cvf  $_{\Delta}$  MyJar.jar  $_{\Delta}$  -C  $_{\Delta}$  MyJava  $_{\Delta}$  .
- 3. 執行
- · 所在的執行目錄為Hadoop\_Home
- ./MyJava = 編譯後程式碼目錄
- · My jar. jar = 封裝後的編譯檔

- · 先放些文件檔到HDFS上的input目錄
- ./input; ./ouput = hdfs的輸入、輸出目錄

#### WordCount1練習(I)

- 1. cd \$HADOOP\_HOME; mkdir input\_local
- 2. echo "I like NCHC Cloud Course." > input\_local/input1
- 3. echo "I like nchc Cloud Course, and we enjoy this crouse." > input\_local/input2
- 4. bin/hadoop dfs -put input\_local input
- 5. bin/hadoop dfs -ls input



#### WordCount1練習(II)

- 1. 編輯WordCount.java
  <a href="http://trac.nchc.org.tw/cloud/attachment/wiki/jazz/Hadoop\_Lab6/WordCount.java?format=raw">http://trac.nchc.org.tw/cloud/attachment/wiki/jazz/Hadoop\_Lab6/WordCount.java?format=raw</a>
- 2. mkdir MyJava
- 3. javac -classpath hadoop-\*-core.jar -d MyJava WordCount.java
- 4. jar -cvf wordcount.jar -C MyJava •
- 5. bin/hadoop jar wordcount.jar WordCount input/ output/
- · 所在的執行目錄為Hadoop\_Home (因為hadoop-\*-core.jar )
- · javac編譯時需要classpath,但hadoop jar時不用
- · wordcount. jar = 封裝後的編譯檔,但執行時需告知class name
- · Hadoop進行運算時,只有 input 檔要放到hdfs上,以便hadoop分析運算,執行檔(wordcount.jar)不需上傳,也不需每個node都放,程式的 載之交由java處理

#### WordCount1練習(III)

```
|waue@vPro:/opt/hadoop$ bin/hadoop dfs -put input input
waue@vPro:/opt/hadoop$ mkdir MvJava
waue@vPro:/opt/hadoop$ javac -classpath hadoop-*-core.jar -d MyJava WordCount.java
waue@vPro:/opt/hadoop$ jar -cvf wordcount.jar -C MyJava .
新增 manifest
|新增:WordCount.class (讀=1516)(寫=740)(壓縮 51%)
新增:WordCount$Reduce.class (讀=1591)(寫=642)(壓縮 59%)
新增:WordCount$Map.class (讀=1918)(寫=795)(壓縮 58%)
waue@vPro:/opt/hadoop$ bin/hadoop jar wordcount.jar WordCount input/ output/
09/03/22 11:39:01 WARN mapred.JobClient: Use GenericOptionsParser for parsing the argu
ments. Applications should implement Tool for the same.
09/03/22 11:39:01 INFO mapred.FileInputFormat: Total input paths to process : 1
09/03/22 11:39:01 INFO mapred.FileInputFormat: Total input paths to process : 1
09/03/22 11:39:02 INFO mapred.JobClient: Running job: job_200903201526_0007
09/03/22 11:39:03 INFO mapred.JobClient: map 0% reduce 0%
09/03/22 11:39:08 INFO mapred.JobClient: map 100% reduce 0%
09/03/22 11:39:15 INFO mapred.JobClient: Job complete: job_200903201526_0007
09/03/22 11:39:15 INFO mapred.JobClient: Counters: 16
09/03/22 11:39:15 INFO mapred.JobClient:
                                          File Systems
09/03/22 11:39:15 INFO mapred.JobClient:
                                            HDFS bytes read=320950
09/03/22 11:39:15 INFO mapred.JobClient:
                                            HDFS bytes written=130568
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Local bytes read=168448
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Local bytes written=336932
09/03/22 11:39:15 INFO mapred.JobClient:
                                          Job Counters
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Launched reduce tasks=1
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Launched map tasks=1
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Data-local map tasks=1
09/03/22 11:39:15 INFO mapred.JobClient:
                                          Map-Reduce Framework
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Reduce input groups=9284
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Combine output records=18568
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Map input records=7868
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Reduce output records=9284
Map output bytes=445846
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Map input bytes=320950
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Combine input records=47227
Map output records=37943
09/03/22 11:39:15 INFO mapred.JobClient:
                                            Reduce input records=9284
waue@vPro:/opt/hadoop$
```

#### WordCount1 練習(IV)

```
|waue@vPro:/opt/hadoop$ bin/hadoop dfs -cat output/part-00000
Cloud 2
Course,
Course.
NCHC
and
course.
enjoy
like
nchc
this
We
```

#### **BTW** ...

- ●雖然Hadoop框架是用Java實作,但 Map/Reduce應用程序則不一定要用Java 來寫
- Hadoop Streaming :
  - ◆執行作業的工具,使用者可以用其他語言 (如:PHP) 套用到Hadoop的mapper和 reducer
- Hadoop Pipes : C++ API



#### Hadoop 程式設計

# 四、寫Code環境準備

4.A: Console 端編譯與執行

4.B:透過 Eclipse 開發與測試運算



#### Requirements

- Hadoop 0.20.0 up
- •Java 1.6
- Eclipse 3.3 up
- Hadoop Eclipse Plugin 0.20.0 up

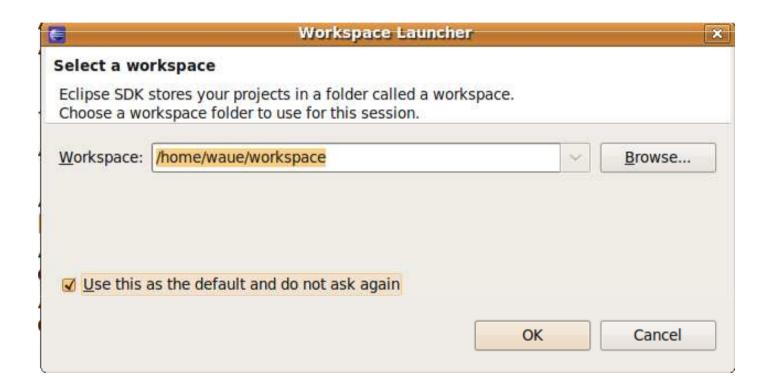


#### 安裝Hadoop Eclipse Plugin

- Hadoop Eclipse Plugin 0.20.0
  - ◆From \$Hadoop\_0.20.0\_home/contrib/eclipse-plugin/hadoop-0.20.0-eclipse-plugin.jar
- Hadoop Eclipse Plugin 0.20.1
  - ◆ Compiler needed
  - ◆ Or download from <a href="http://hadoop-eclipse-plugin.googlecode.com/files/hadoop-0.20.1-eclipse-plugin.jar">http://hadoop-eclipse-plugin.googlecode.com/files/hadoop-0.20.1-eclipse-plugin.jar</a>
- copy to \$Eclipse\_home/plugins/

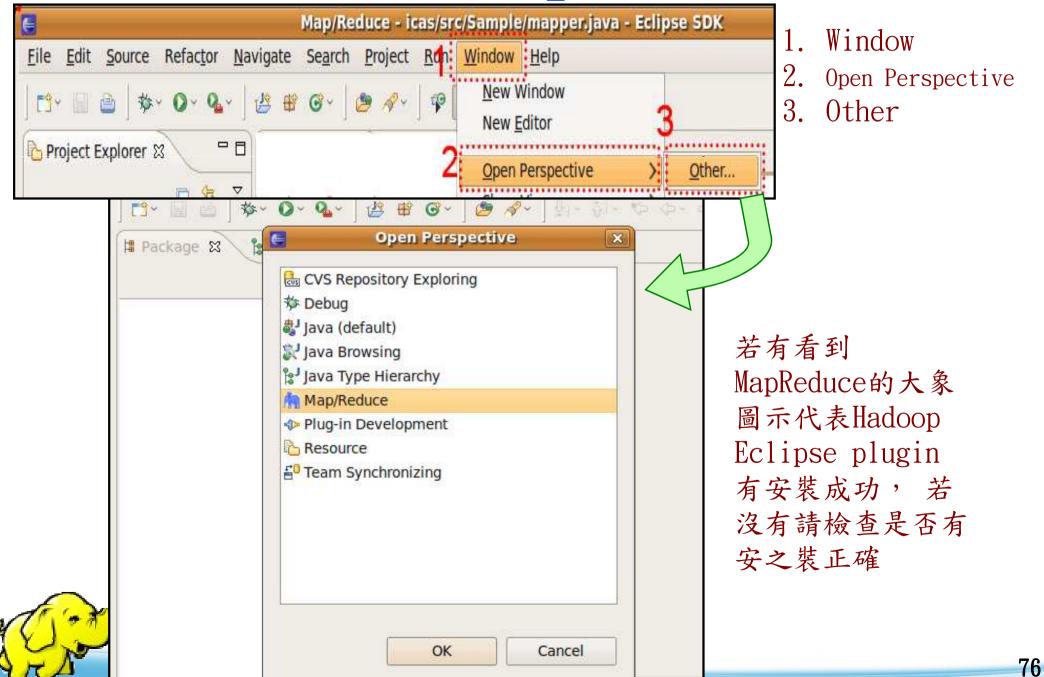


# 1 打開Eclipse, 設定專案目錄

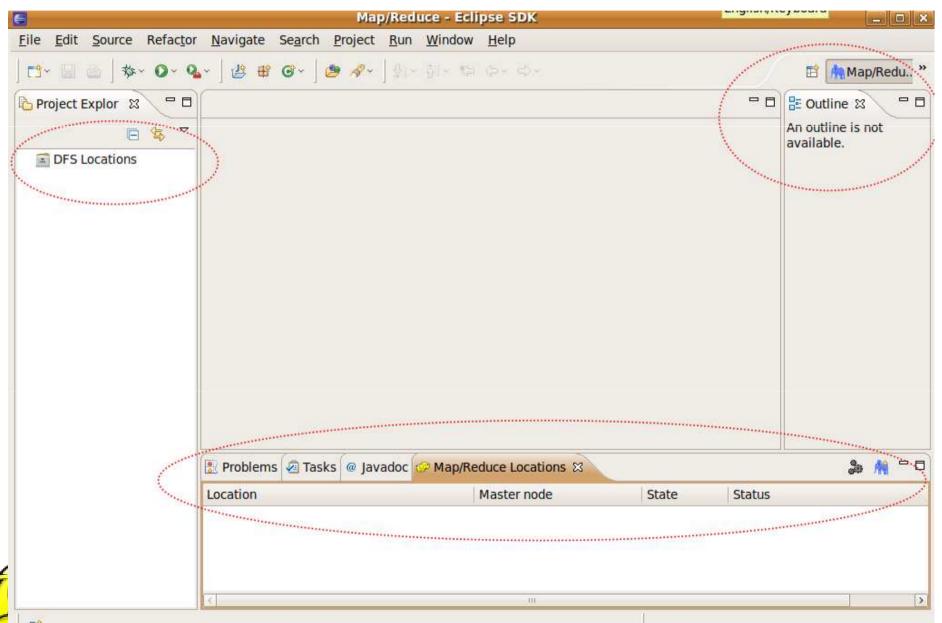




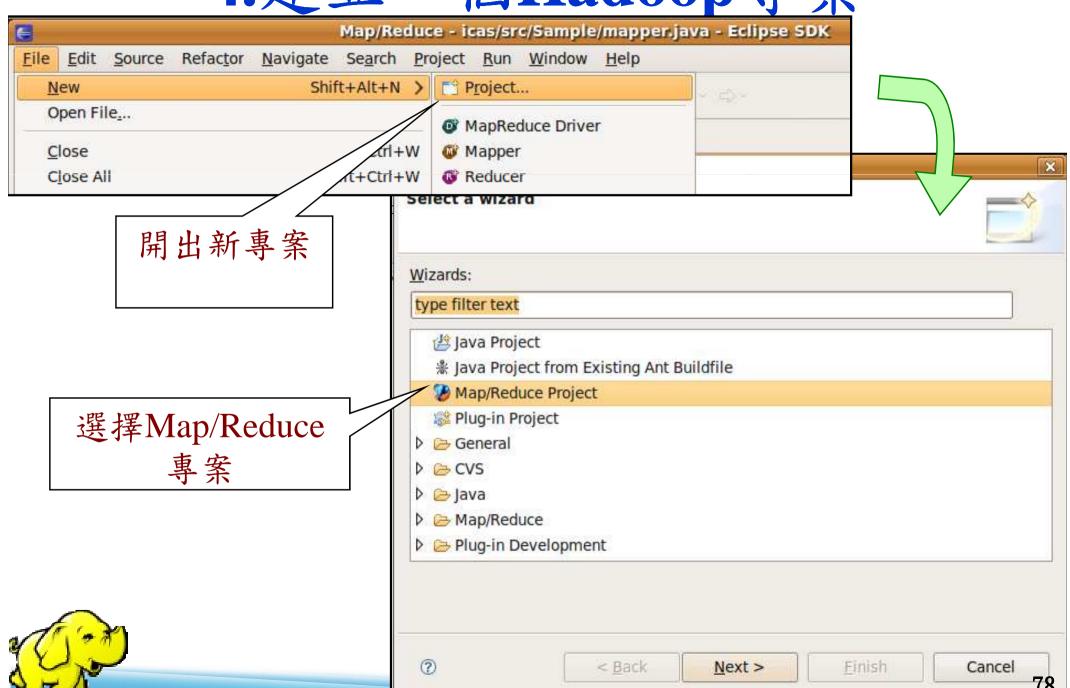
#### 2. 使用Hadoop mode視野



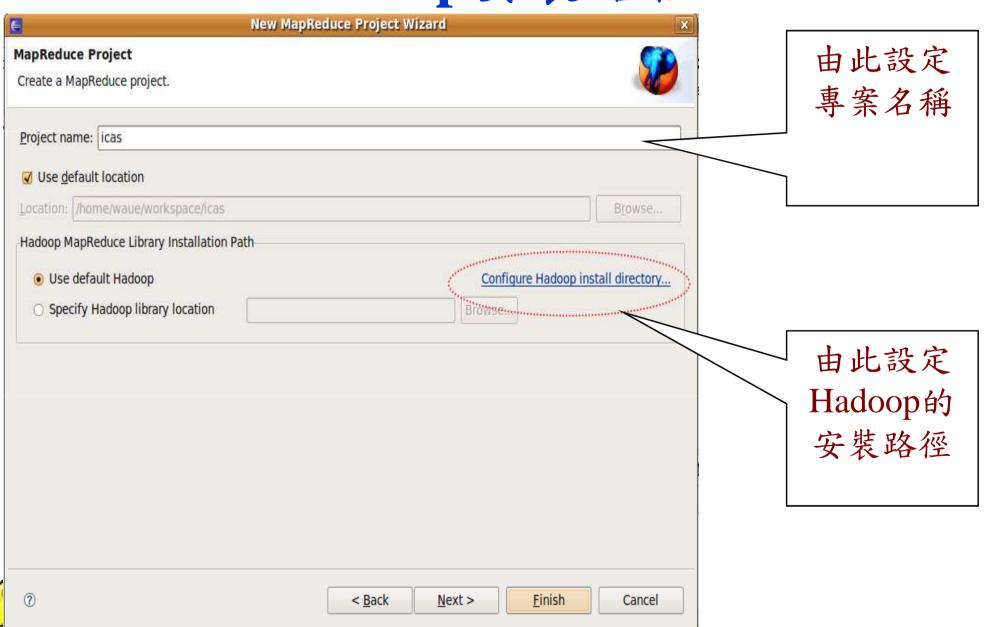
# 3. 使用Hadoop視野,主畫面將出現三個功能



4.建立一個Hadoop專案



# 4-1. 輸入專案名稱並點選設定 Hadoop安裝路徑

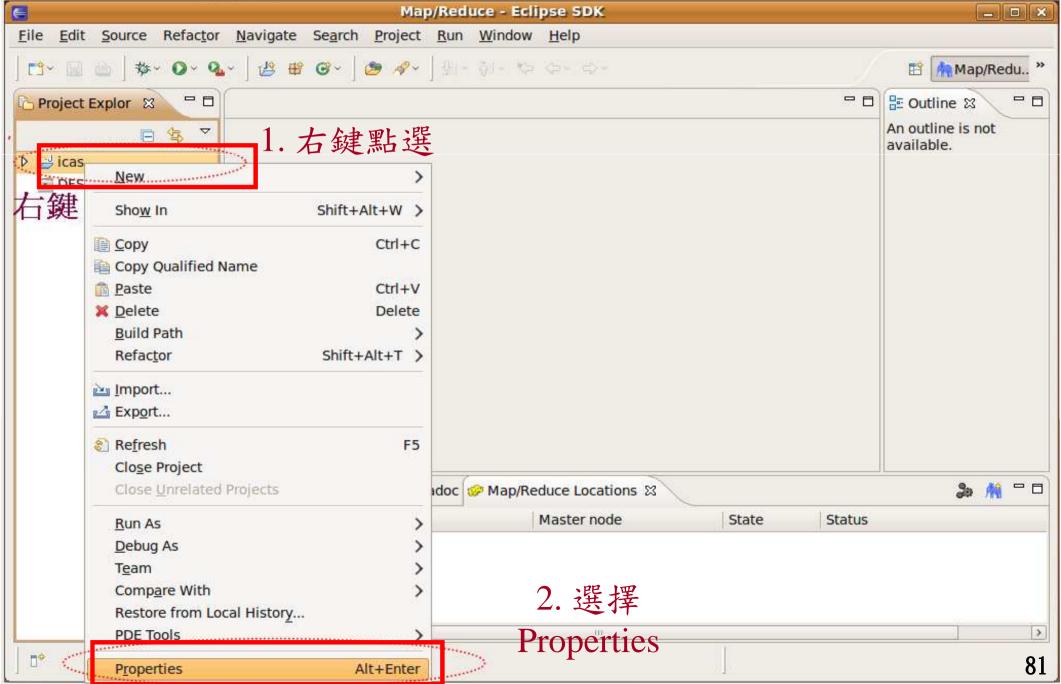


### 4-1-1. 填入Hadoop安裝路徑





# 5. 設定Hadoop專案細節



#### 5-1. 設定原始碼與文件路徑

選擇 Java Build Path

(?)

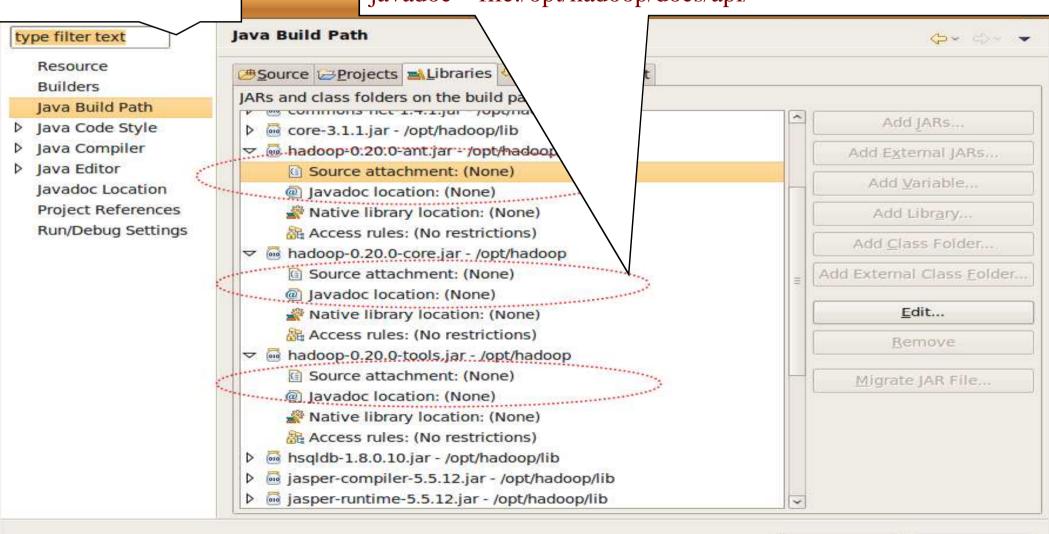
以下請輸入正確的Hadoop原始碼與API文件檔路徑,如

OK

Cancel

source : /opt/hadoop/src/

javadoc: file:/opt/hadoop/docs/api/

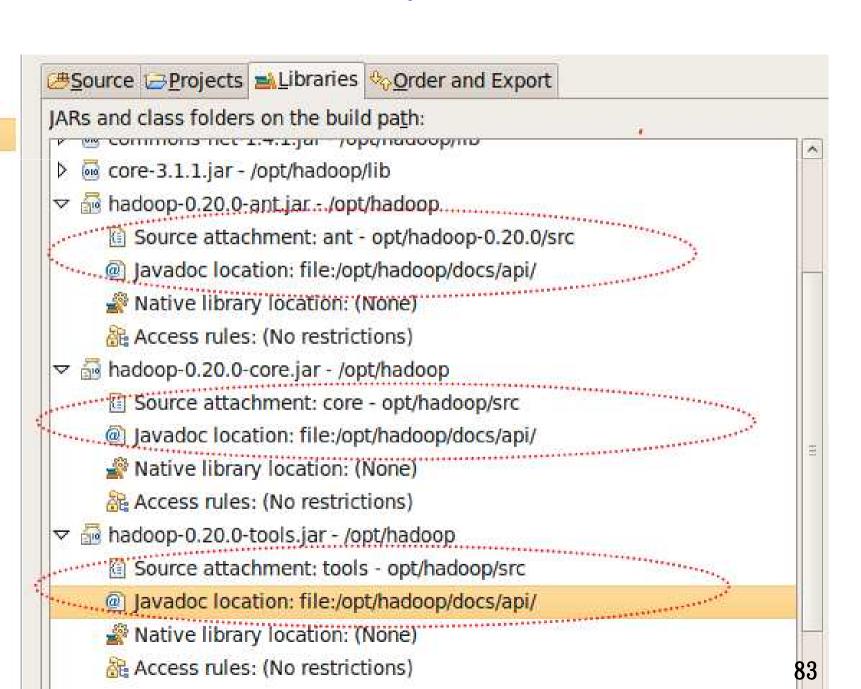


#### 5-1-1. 完成圖

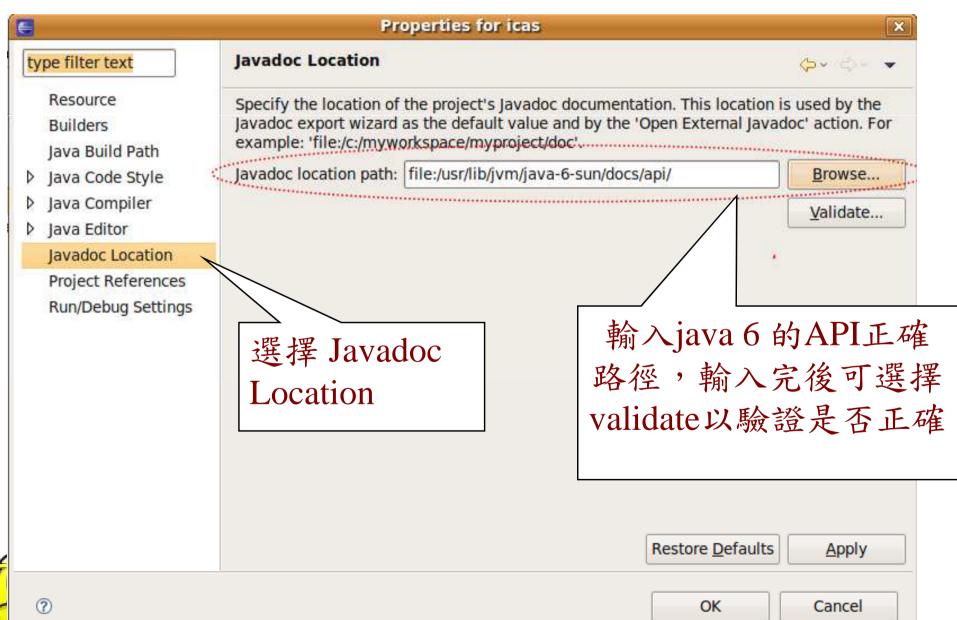
Resource Builders

#### Java Build Path

- Dava Code Style
- Java Compiler
- Java Editor Javadoc Location Project References Run/Debug Settings



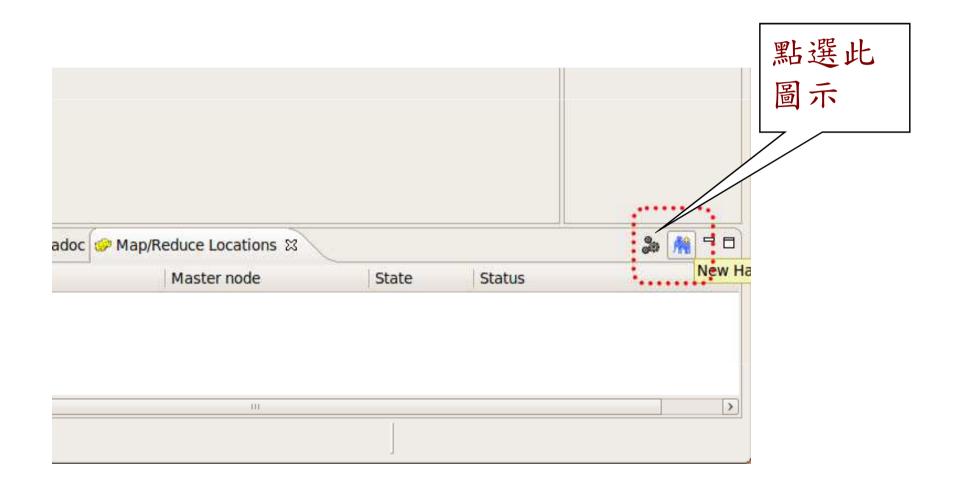
# 5-2. 設定java doc的完整路徑





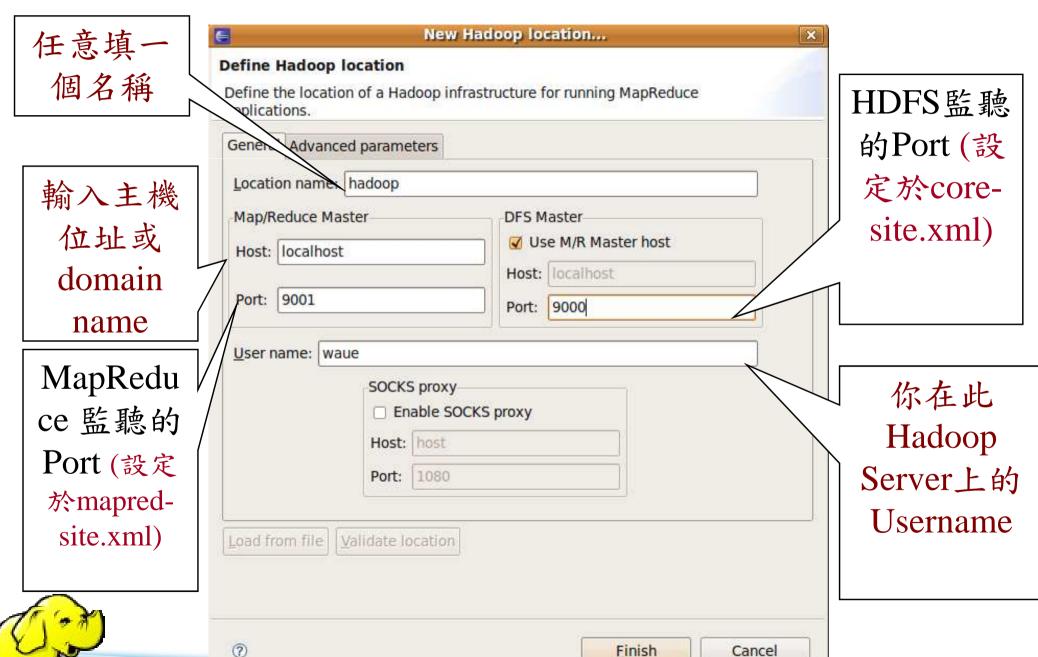
84

# 6. 連結Hadoop Server與Eclipse





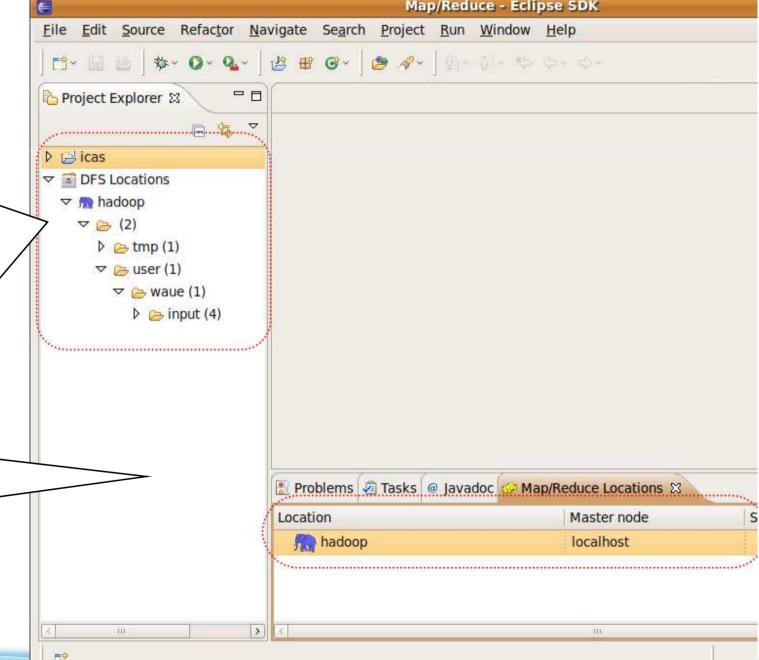
#### 6-1. 設定你要連接的Hadoop主機



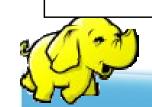
86

#### 6-2 若正確設定則可得到以下畫面

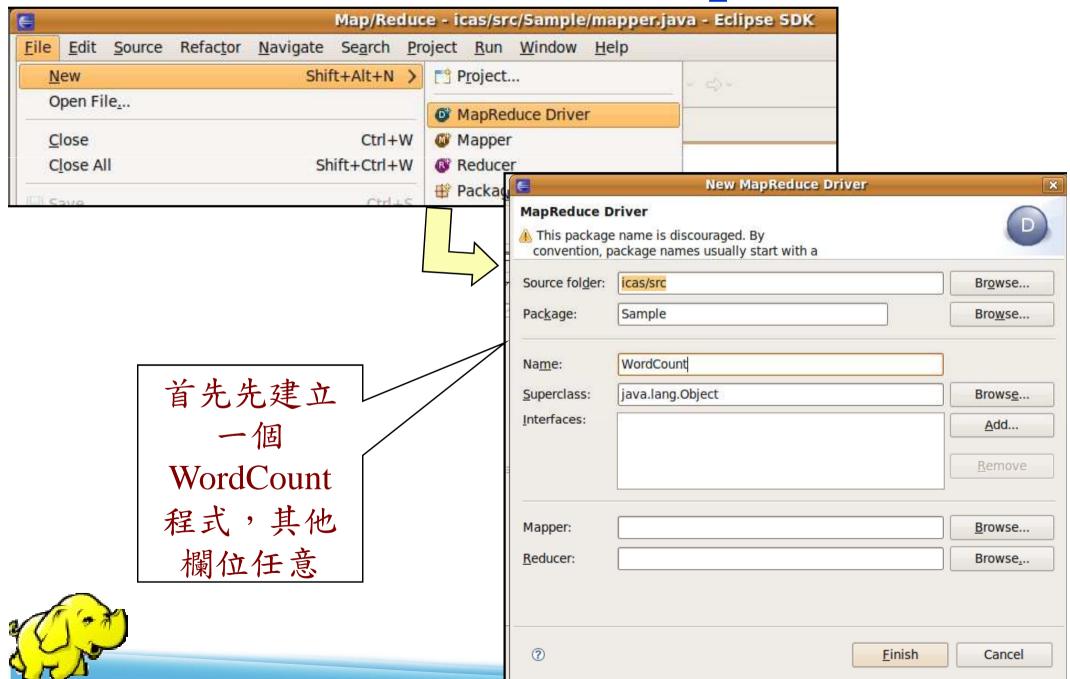
HDFS的資訊, 可直接於此操 作檢視、新增、 上傳、刪除等 命令



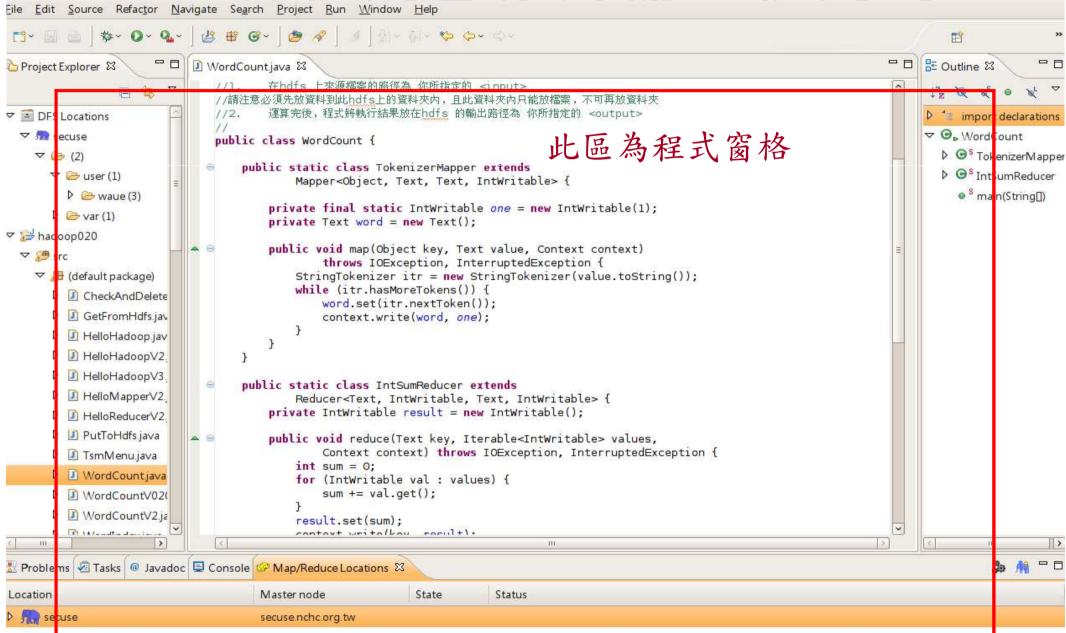
若有Job運作, 可於此視窗 檢視



# 7. 新增一個Hadoop程式

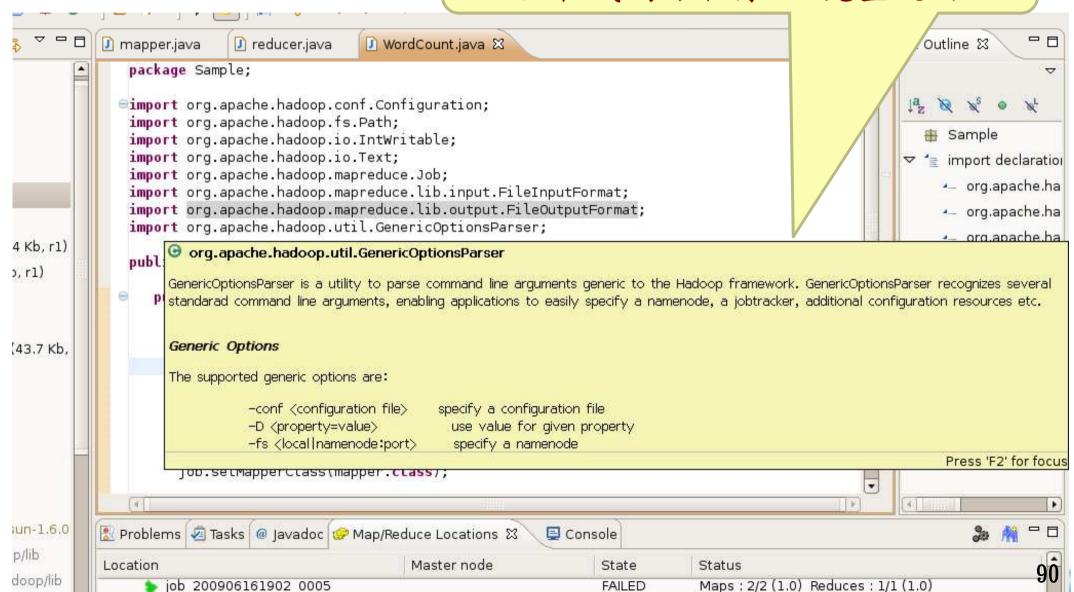


#### 7.1 於程式窗格內輸入程式碼



#### 7.2 補充

若之前doc部份設定正確,則滑鼠移 至程式碼可取得API完整說明



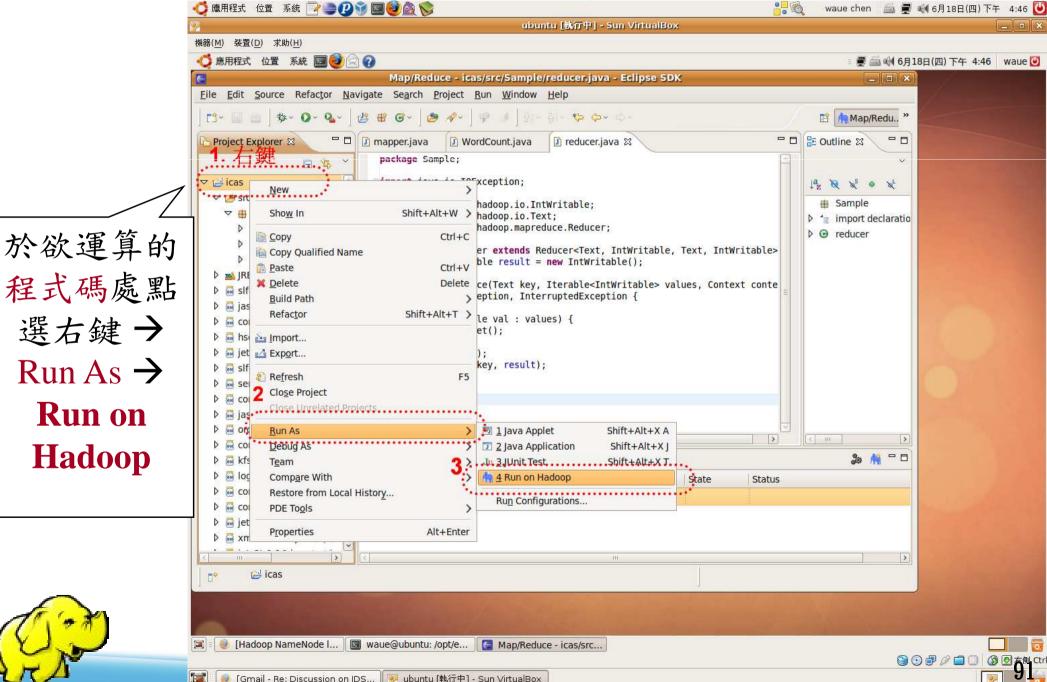
#### 8. 運作

Run on

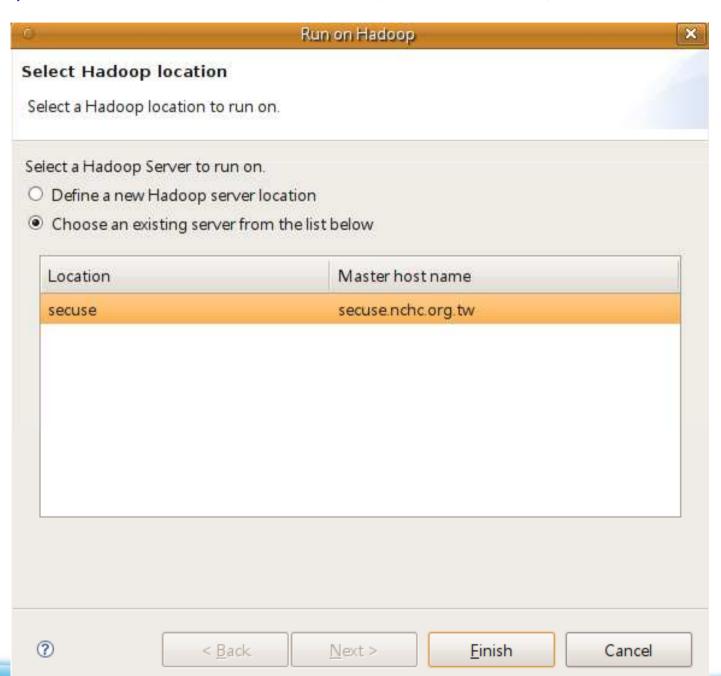
**Hadoop** 

[Gmail - Re: Discussion on IDS...

🥦 ubuntu [執行中] - Sun VirtualBox

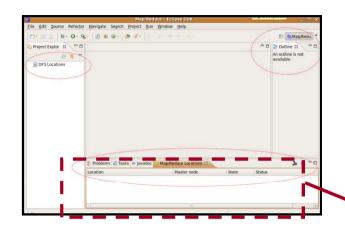


#### 8-1 選擇之前設定好所要運算的主機





#### 8.2 運算資訊出現於Eclipse 右下方

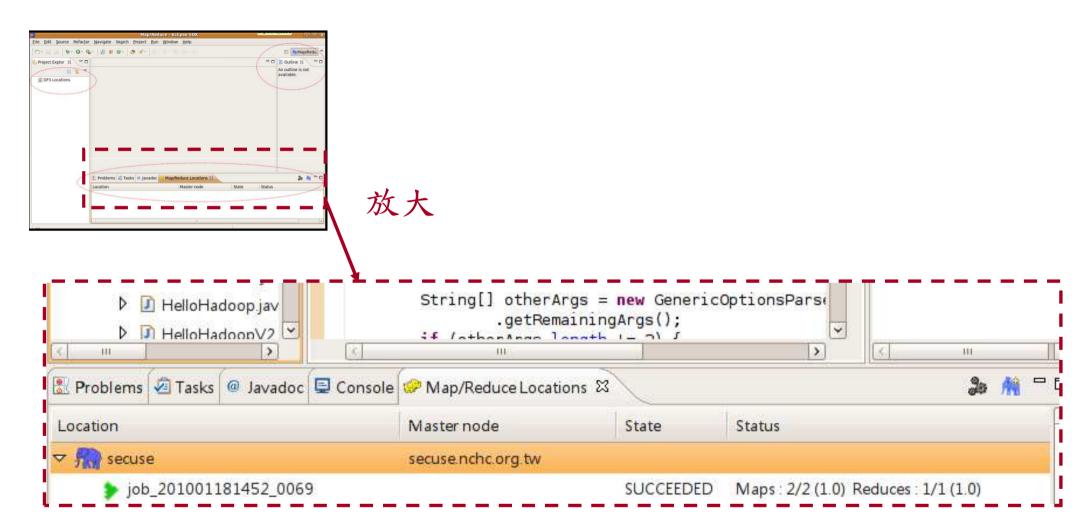


#### 的Console 視窗

放大

```
🖺 Problems 🙆 Tasks @ Javadoc 📮 Console 🛭
                                           Map/Reduce Locations
 <terminated> WordCount (1) [Java Application] /usr/lib/jvm/java-6-sun-1.6.0.16/bin/java (2010/1/20 下午 6:15:07)
 10/01/20 18:15:08 WARN conf.Configuration: DEPRECATED: hadoop-site.xml found in the classpath. Usage of hado
 10/01/20 18:15:08 WARN mapred.JobConf: The variable mapred.task.maxvmem is no longer used. Instead use mapre
 10/01/20 18:15:08 WARN mapred.JobConf: The variable mapred.task.maxvmem is no longer used. Instead use mapre
 10/01/20 18:15:08 INFO input.FileInputFormat: Total input paths to process : 2
 10/01/20 18:15:08 INFO mapred.JobClient: Running job: job 201001181452 0078
 10/01/20 18:15:09 INFO mapred.JobClient: map 0% reduce 0%
 10/01/20 18:15:16 INFO mapred.JobClient:
                                           map 100% reduce 0%
                                           map 100% reduce 100%
 10/01/20 18:15:28 INFO mapred.JobClient:
 10/01/20 18:15:30 INFO mapred.JobClient: Job complete: job 201001181452 0078
 10/01/20 18:15:30 INFO mapred.JobClient: Counters: 17
 _10/01/20 18:15:30 INFO mapred.JobClient:
                                             Job Counters
                                            Launched reduce tasks=1
 10/01/20 18:15:30 INFO mapred.JobClient:
 10/01/20 18:15:30 INFO mapred.JobClient:
                                            Launched map tasks=2
 10/01/20 18:15:30 INFO mapred.JobClient:
                                               Data-local map tasks=2
 10/01/20 18:15:30 INFO mapred.JobClient:
                                            FileSystemCounters
10/01/20 18:15:30 INFO mapred.JobClient:
                                               FILE BYTES READ=153
```

#### 8.3 剛剛運算的結果出現如下圖





#### **Conclusions**

#### ●優點

- ◆快速開發程式
- ◆易於除錯
- ◆智慧尋找函式庫
- ◆自動鍊結API
- ◆直接操控 HDFS 與 JobTracker
- **♦**...

#### ●缺點

◆Plugin 並會因Eclipse 版本而有不同的狀況

