



財團法人國家實驗研究院

國家高速網路與計算中心

NATIONAL CENTER FOR HIGH-PERFORMANCE COMPUTING

# Map Reduce 介紹



王耀聰 陳威宇

[jazz@nchc.org.tw](mailto:jazz@nchc.org.tw)

[waue@nchc.org.tw](mailto:waue@nchc.org.tw)

國家高速網路與計算中心  
(NCHC)



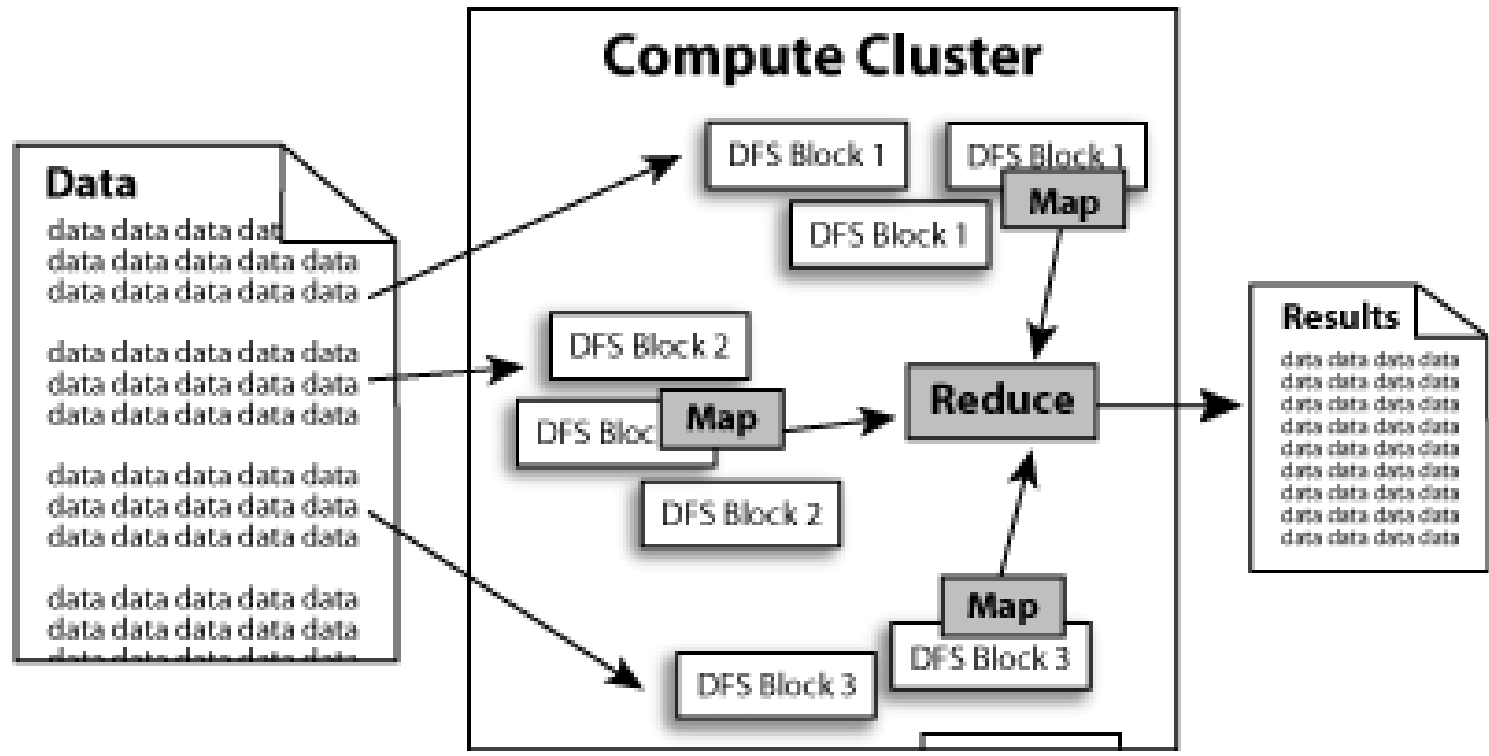
自由軟體實驗室

# Outline

- What is MapReduce ?
- Where does it fix ?
- What is its benefit ?
- How does it work ?
- Must be in Java ?

What is  
MapReduce ?

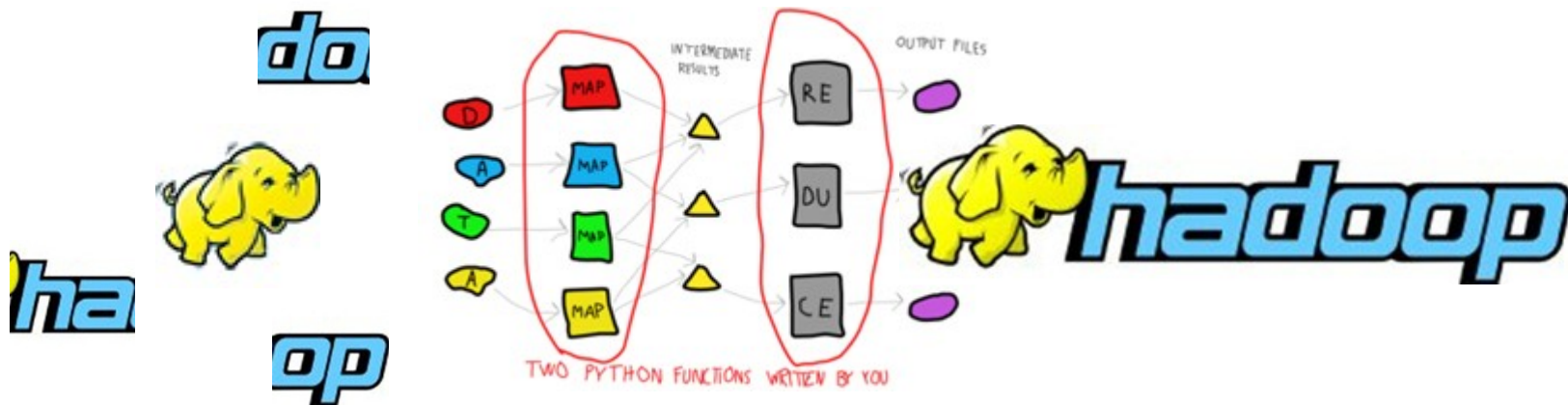
# Google 原生定義



MapReduce is a framework for computing certain kinds of **distributable problems** using a large number of computers (nodes), collectively referred to as a cluster.

What is  
MapReduce ?

# Hadoop MapReduce 定義



Hadoop Map/Reduce 是一個易於使用的軟體平台，以 MapReduce 為基礎的應用程序，能夠運作在由上千台 PC 所組成的大型叢集上，並以一種**可靠容錯**的方式**平行處理**上 Peta-Bytes 數量級的資料集。

Where does  
it fix ?

# 應用範圍

- 大規模資料集
- 可拆解
- Text tokenization
- Indexing and Search
- Data mining
- machine learning
- ...

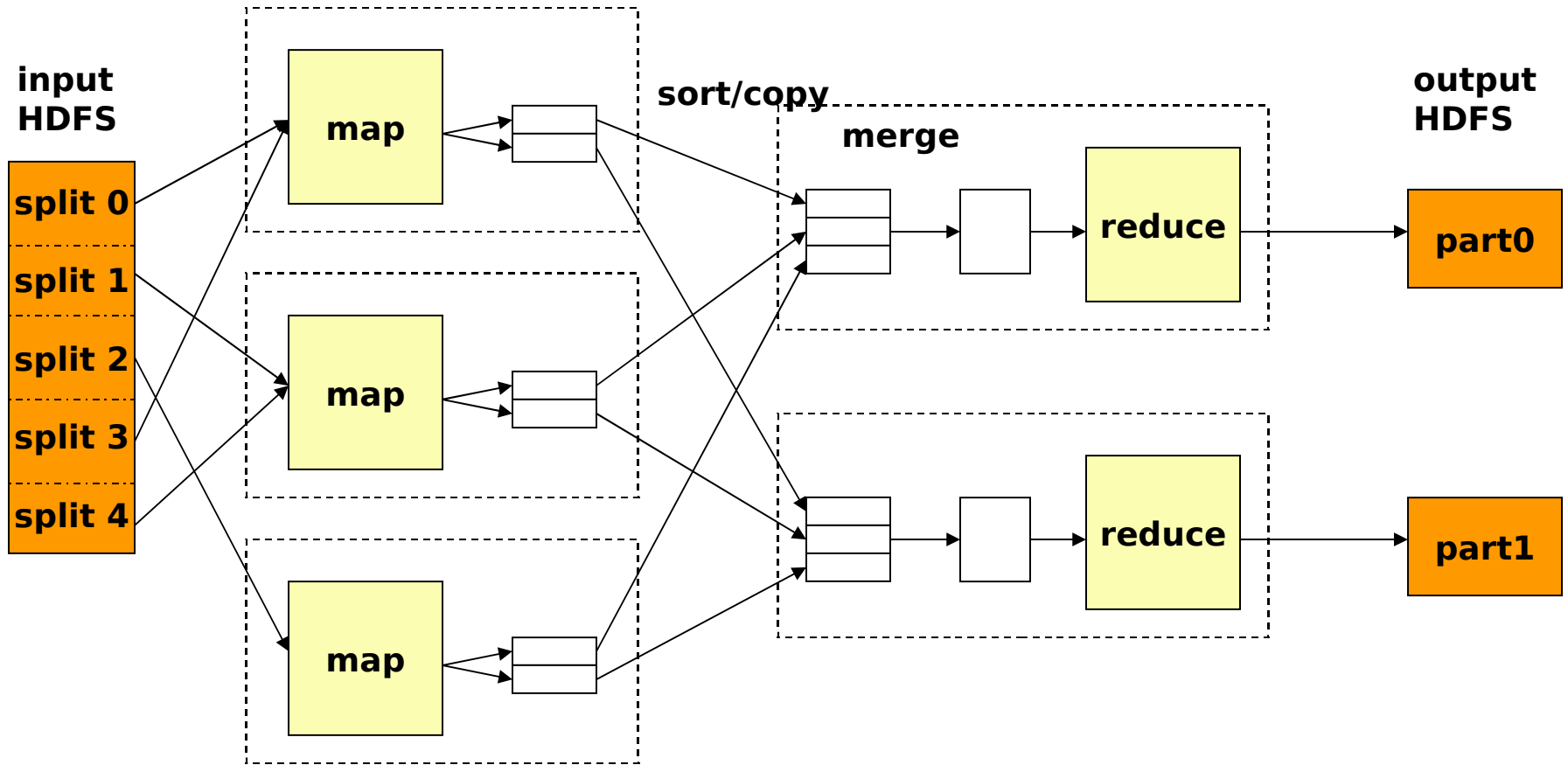


# MapReduce 由來

- Functional Programming : Map Reduce
  - map(...):
    - [ 1,2,3,4 ] - (\*2) -> [ 2,4,6,8 ]
  - reduce(...):
    - [ 1,2,3,4 ] - (sum) -> 10
  - 對應演算法中的 Divide and conquer
  - 將問題分解成很多個小問題之後，再做總和
- 首先被 Google 引用到程式設計的軟體架構內，使用在大規模數據的運算中

How does  
it work ?

# MapReduce 運作流程



JobTracker 跟 NameNode 取得需要運算的 blocks

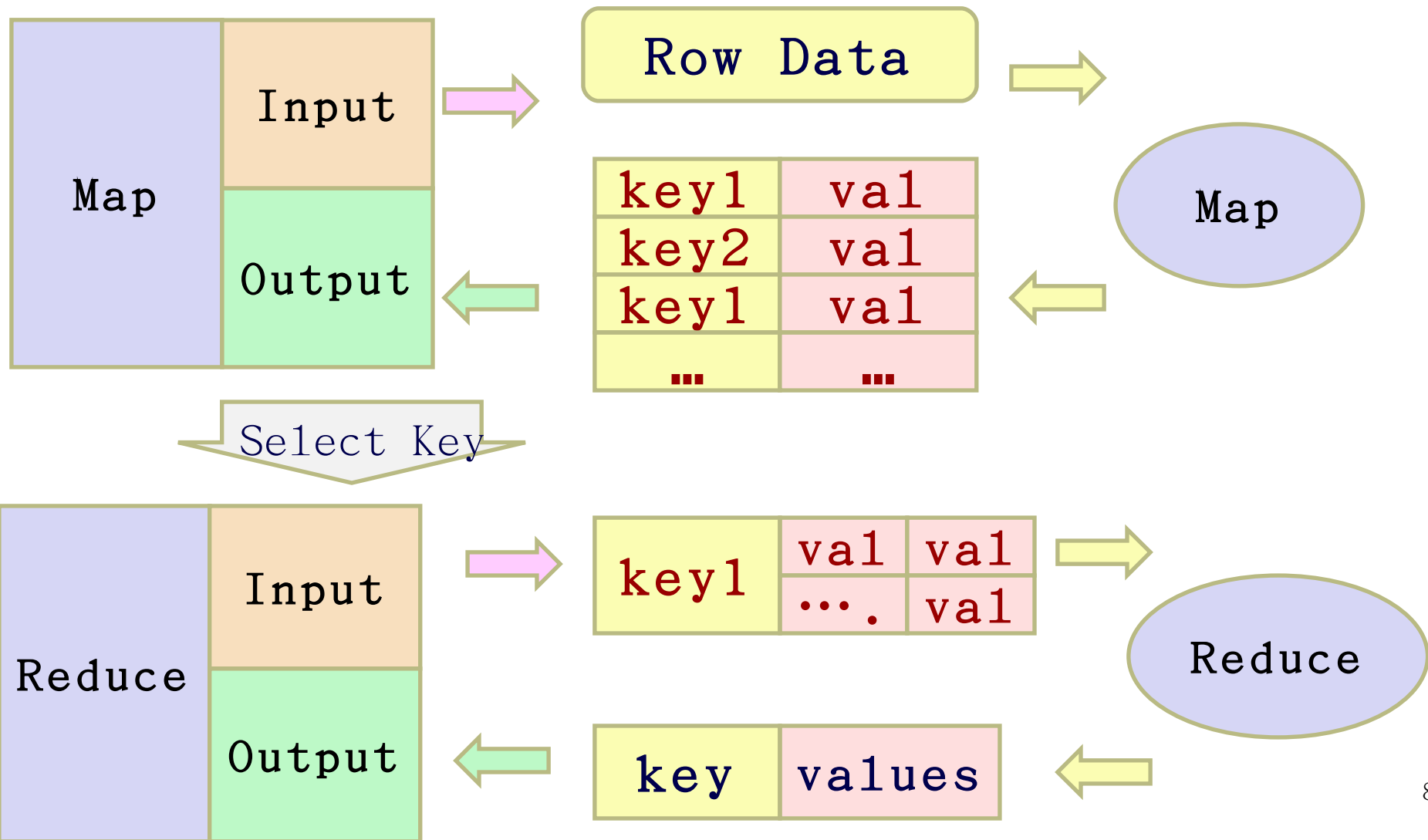
JobTracker 選數個 TaskTracker 來作 Map 運算，產生些中間檔案

JobTracker 將中間檔案整合排序後，複製到需要的 TaskTracker 去

JobTracker 派遣 TaskTracker 作 reduce

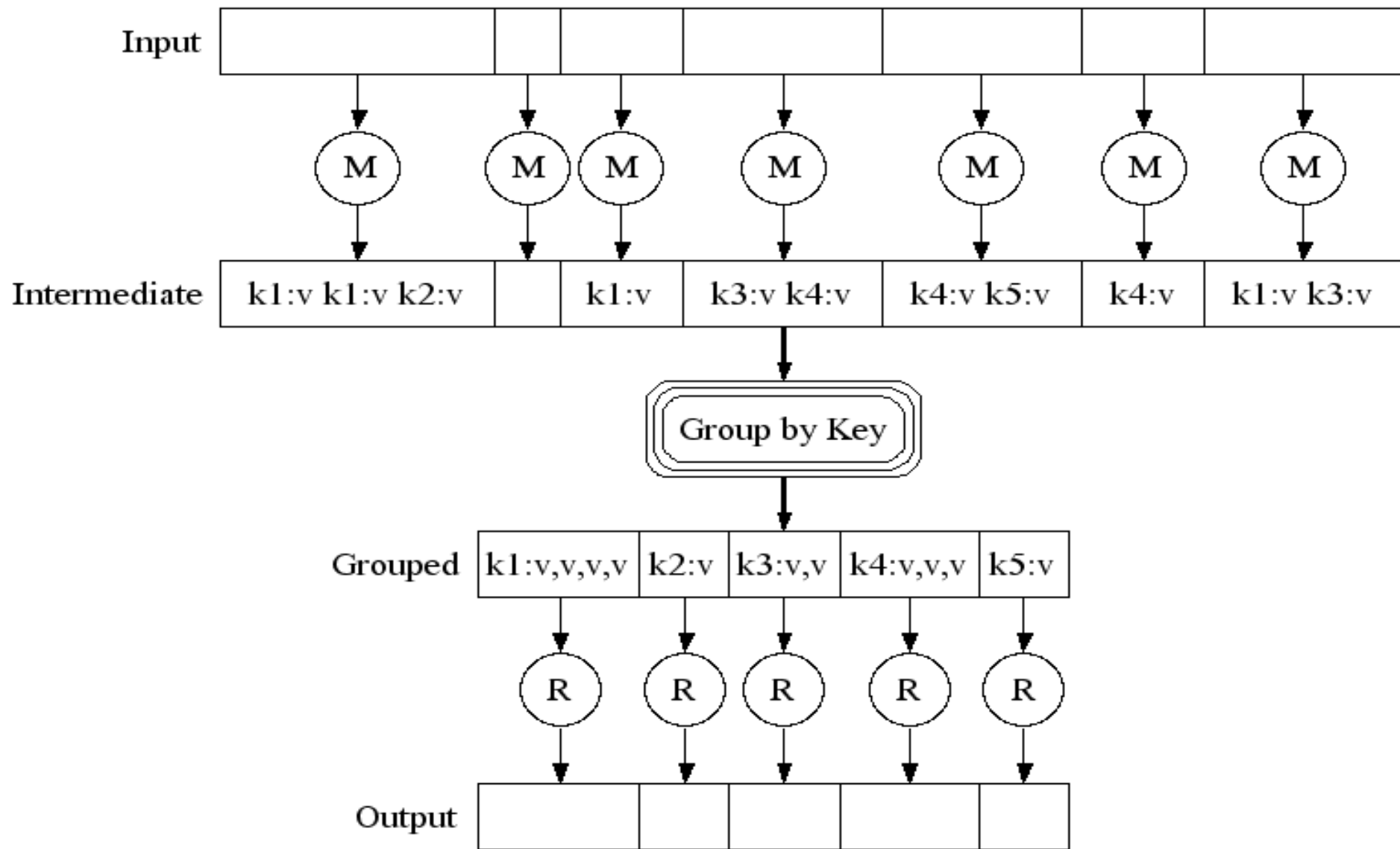
reduce 完後通知 JobTracker 與 Namenode 以產生 output

# <Key, Value> Pair

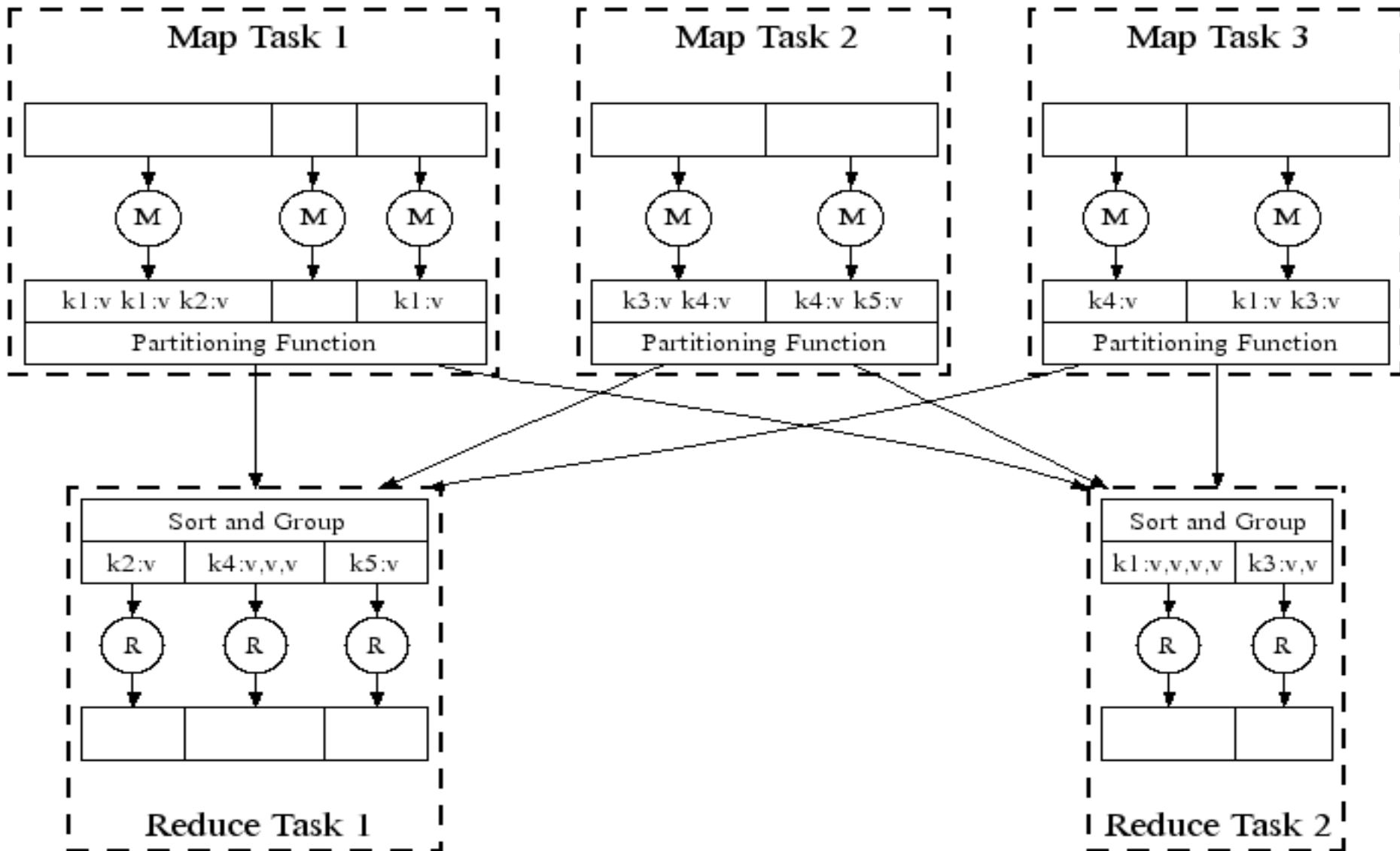




# MapReduce 圖解



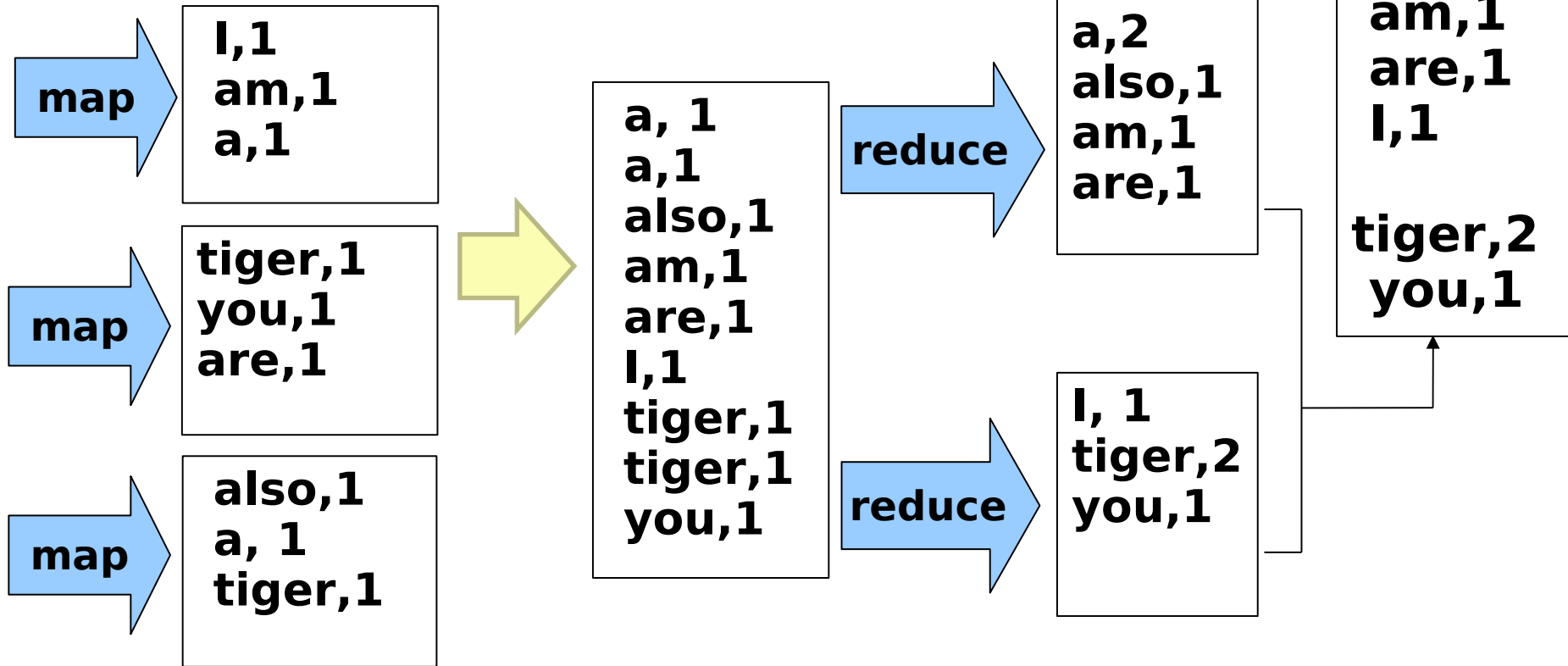
# MapReduce in Parallel



How does  
it work ?

# 範例

**I am a tiger, you are  
also a tiger**



JobTracker 先選了  
三個 Tracker 做 map

Map 結束後，hadoop 進  
行中間資料的整理與排序

JobTracker 再選兩個  
TaskTracker 作 reduce

# Options without Java

- 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