

# يمكن الحوسبة، ITNT404

المحاضر: د. عمر أبو سعدة البريد الإلكتروني:

[omar.abusaeeda@uot.edu.ly](mailto:omar.abusaeeda@uot.edu.ly)

Hadoop كأداة MapReduce

# ماهو تقليل الخريطة؟

يمكن تعريف Map-Reduce على أنه **جزئين**:

- **عارضة** لكتابة البرامج التي يمكن إجراؤها بسهولة لمعالجة البيانات بالتوازي.

- **إطار العمل** الذي يقوم بتشغيل هذه البرامج بالتوازي، ويتعامل مع التفاصيل تلقائياً  
تقسيم العمل والتوزيع والتزامن والتسامح مع الخطأ.

ال **نموذج** و **النطاق** العمل معاً لصنع برامج قابلة للتطوير،  
موزعة ومتسامحة مع الخطأ.

# أباتشي MapReduce

-العالم كله أصبح رقمياً

-إطار برمجي للمعالجة الموزعة لمجموعات البيانات الكبيرة

-يعتني الإطار **جدولة المهام، يراقب لهم وإعادة تنفيذ أي**

**المهام الفاشلة.**

-يقوم بتقسيم مجموعة بيانات الإدخال إلى أجزاء مستقلة تتم معالجتها بطريقة متوازية تماماً.

-يقوم إطار عمل MapReduce بفرز مخرجات الخرائط، والتي يتم إدخالها بعد ذلك إلى الملف

تقليل المهام. عادةً، يتم تخزين كل من مدخلات ومخرجات المهمة في ملف

نظام.

# نموذج البرمجة MapReduce

- يتكون MapReduce من مرحلتين والابتكار الرئيسي الخاص به هو:

- القدرة على إجراء استعلام على مجموعة بيانات وتقسيمها وتشغيلها بالتوازي على العديد من العقد.

- يحل مشكلة البيانات الكبيرة جداً (البيانات الكبيرة) التي لا يمكن احتواؤها على جهاز واحد

- الحوسبة الموزعة على العديد من الخوادم

- نموذج معالجة الدفعات

- **مرحلة الخريطة** تتم معالجة بيانات الإدخال عنصراً تلو الآخر وتحويلها إلى وسيط

مجموعة البيانات.

- **تقليل المرحلة**، يتم تقليل هذه النتائج الوسيطة إلى مجموعة بيانات مختصرة، وهي

النتيجة النهائية المرغوبة.

# توزيع البيانات

- في مجموعة MapReduce، يتم توزيع البيانات على جميع العقد في المجموعة أثناء تحميلها في.

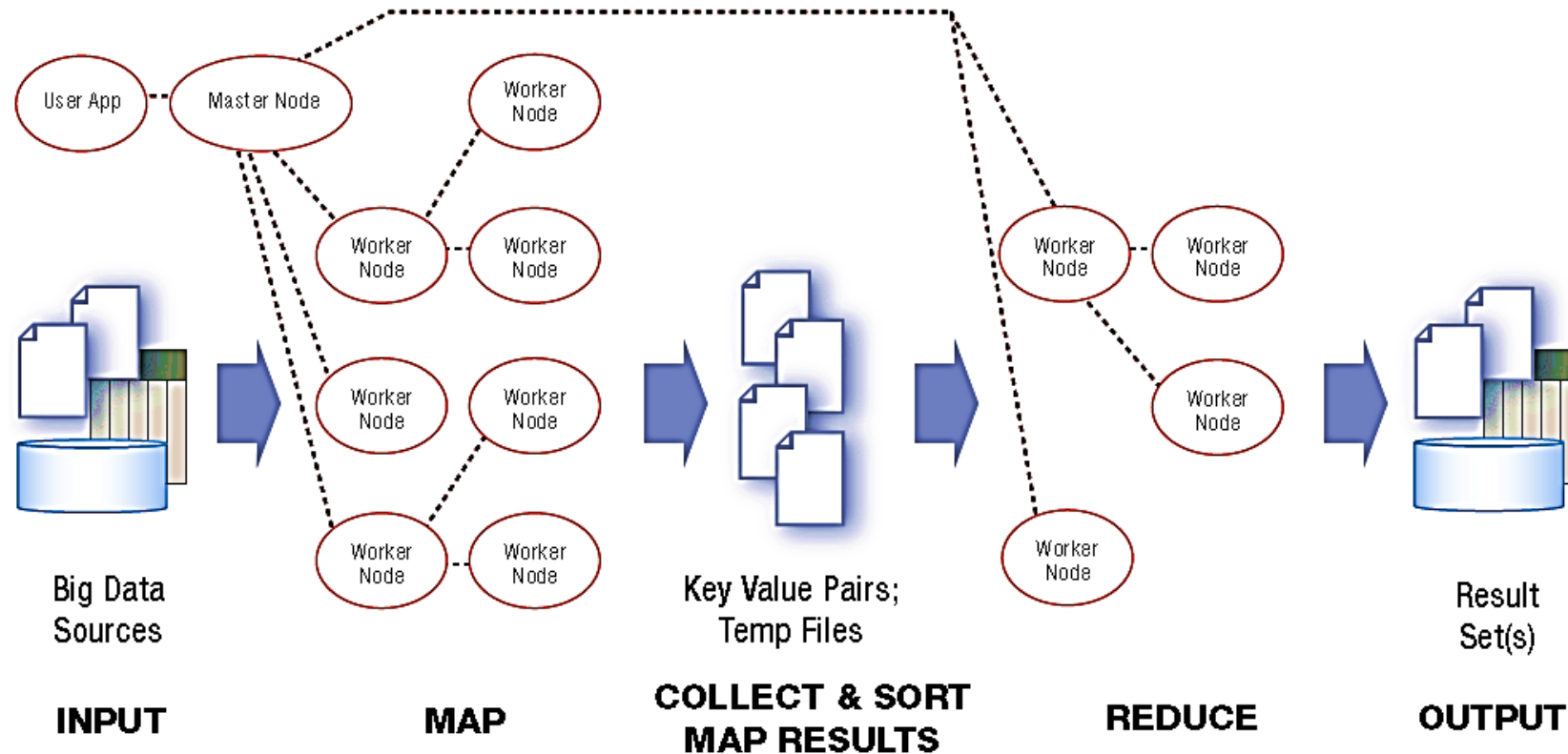
- تقوم أنظمة الملفات الموزعة الأساسية (على سبيل المثال، GFS) بتقسيم ملفات البيانات الكبيرة إلى أجزاء تدار من قبل العقد المختلفة في الكتلة

- على الرغم من أن أجزاء الملف يتم توزيعها عبر عدة أجهزة، إلا أنها تتشكل مساحة اسم واحدة

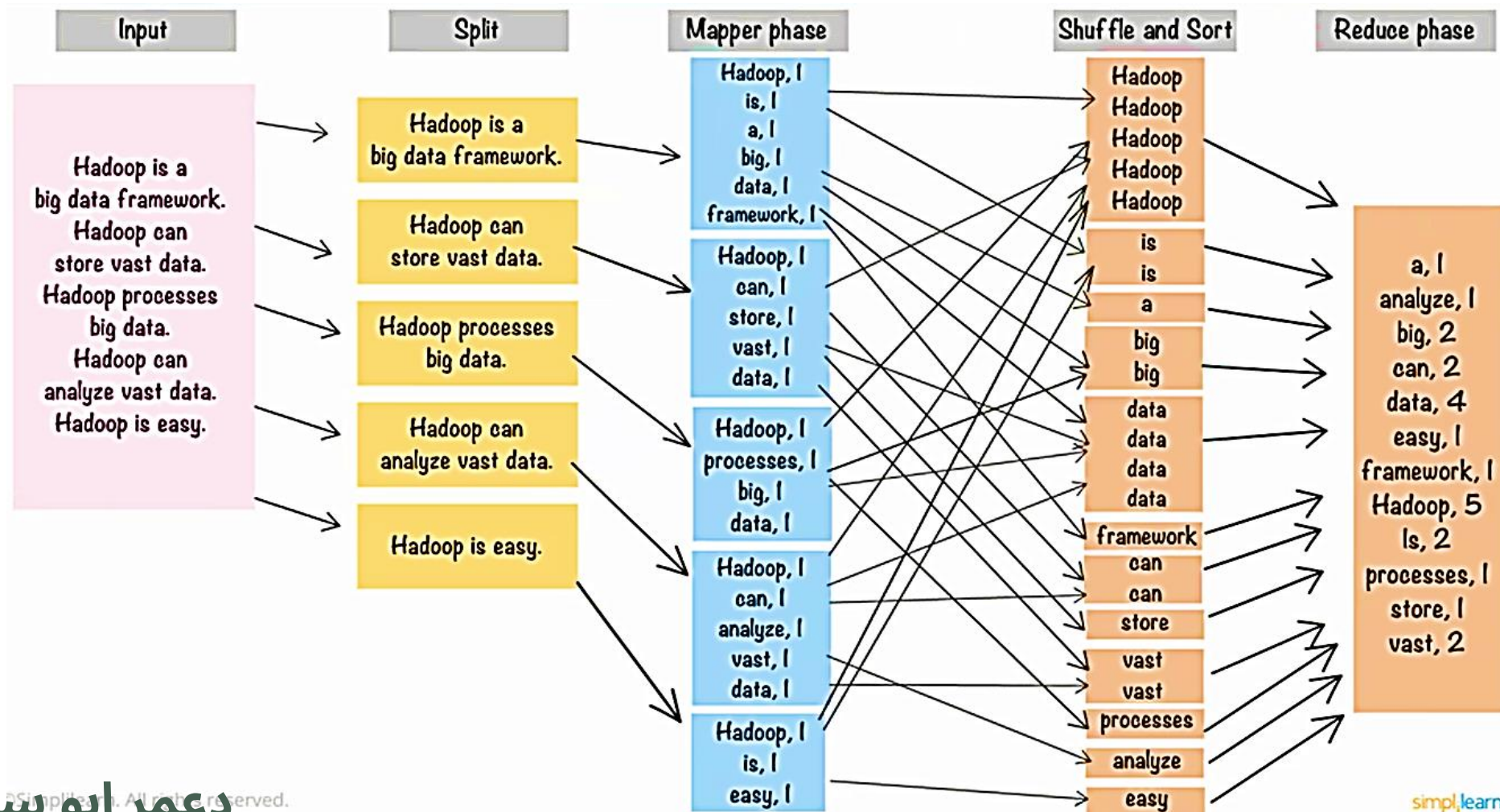


# نظرة عامة على مستوى عالٍ جداً: MapReduce

- معالجة البيانات بطريقة موجهة نحو الدفعات وقد تستغرق معالجتها دقائق أو ساعات (عادة).



# نظرة عامة على مستوى عالٍ جداً MapReduce



# نظرة عامة على مستوى عالٍ جداً: MapReduce

## - تحميل البيانات

- تسمى هذه العملية بشكل صحيح **يستخرج، تحول، حمولة (إيتل)** في مصطلحات تخزين البيانات.
- يجب استخراج البيانات من مصدرها، وتنظيمها لجعلها جاهزة للمعالجة، وتحميلها في ملف طبقة تخزين لـ MapReduce للعمل عليها.

## - MapReduce

- ستقوم هذه المرحلة باستعادة البيانات من التخزين،
- معالجتها (رسم خريطة، جمع نتائج الخريطة وفرزها، تقليلها)
- وإرجاع النتائج إلى المخزن.

## - استخراج النتيجة

- بمجرد اكتمال المعالجة، لكي تكون النتيجة مفيدة، يجب استرجاعها من المخزن و

قدم.

ITNT404



# نظرة شاملة: MapReduce

- في MapReduce، تتم معالجة القطع بشكل منفصل عن طريق مهام تسمى

مصمموا الخرائط

- تتم الإشارة إلى مخرجات مصممي الخرائط على أنها مخرجات وسيطة

ويتم إحضارها إلى مجموعة ثانية من المهام تسمى المخفضات مراحل المرحلة

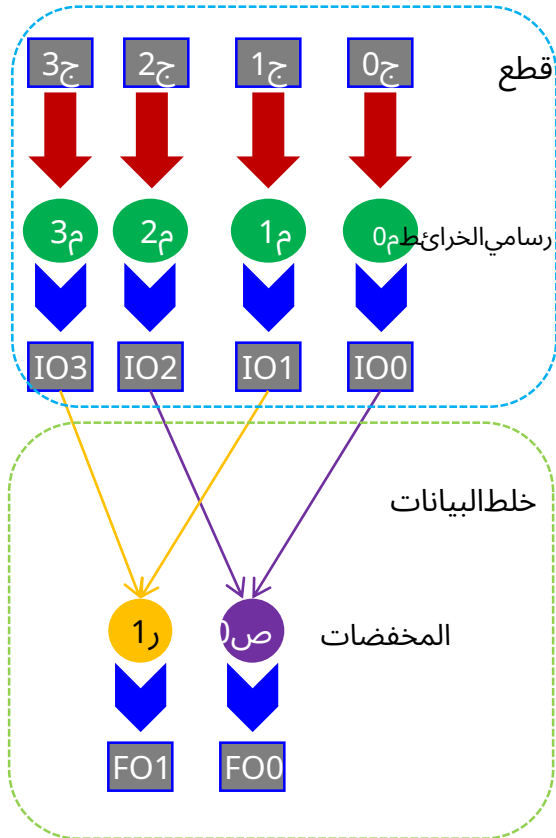
- إن عملية جمع عمليات الإدخال والإخراج في مجموعة من المخفضات معروفة

مثل عملية الخلط مرحلة الخريطة

- تنتج المخفضات المخرجات النهائية (FOs)

- بشكل عام، يقوم MapReduce بتقسيم تدفق البيانات إلى مرحلتين، مرحلة الخريطة

وتقليل المرحلة



# نظرة عامة على مستوى عالٍ جداً MapReduce:

- MapReduce Programming Model

- Data type: key-value *records*

- Map function:

$$(K_{in}, V_{in}) \rightarrow \text{list}(K_{inter}, V_{inter})$$

- Reduce function:

$$(K_{inter}, \text{list}(V_{inter})) \rightarrow \text{list}(K_{out}, V_{out})$$

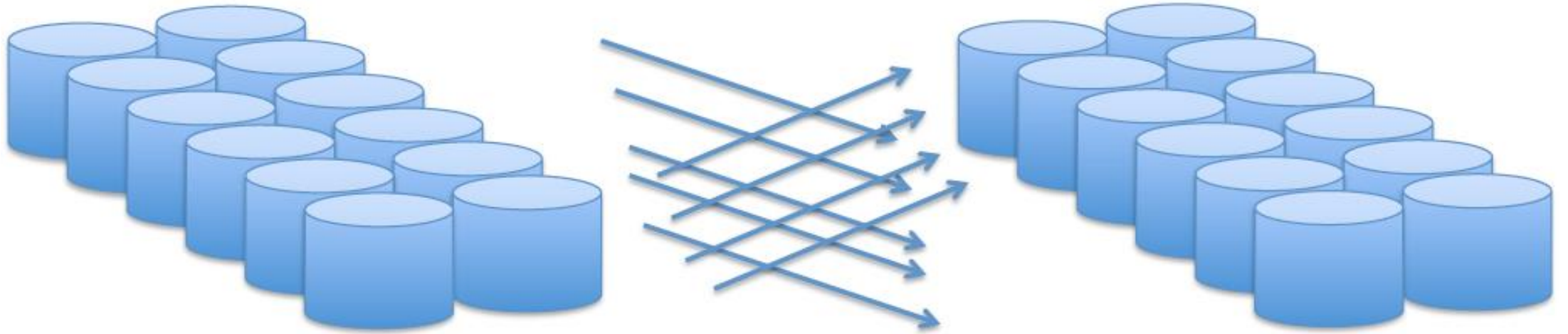
# فوائده نموذج MapReduce

- من خلال توفير نموذج برمجة متوازي للبيانات، يمكن لـ MapReduce التحكم في المهمة التنفيذ بطرق مفيدة:

- التقسيم التلقائي للوظيفة إلى مهام
  - الوضع التلقائي للحساب بالقرب من البيانات
  - موازنة الحمل التلقائي
  - التعافي من الفشل والمتعثرين
- يركز المستخدم على التطبيق، وليس على تعقيدات الحوسبة الموزعة (Implicit تماثل)

# مثال: عدد الكلمات

-النمط الأساسي: سلاسل



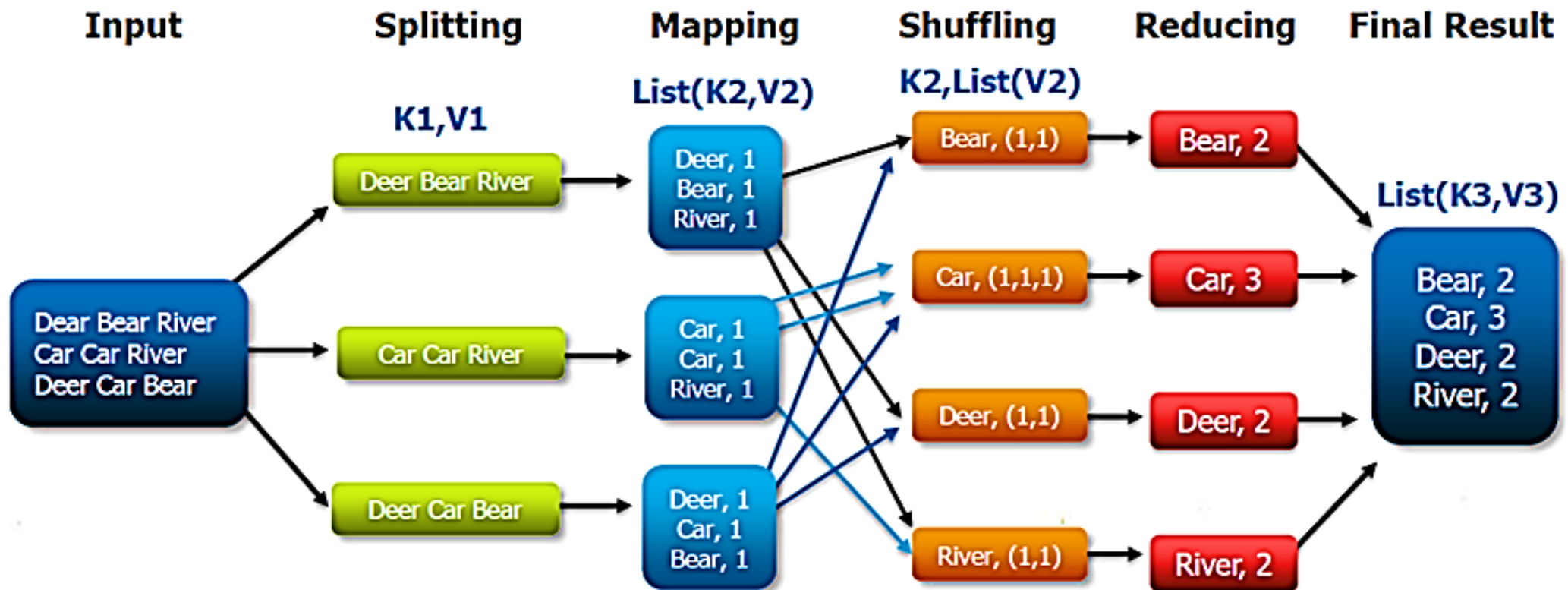
1. استخراج الكلمات من  
صفحات الويب بالتوازي.

2. تجزئة الكلمات وفرزها.

3. العد بالتوازي.

# عددا الكلمات الشائعة

## The Overall MapReduce Word Count Process



# عدداالكلمات الشائئة

خريطة باطلة (السلسلة  $i$ ، خط السلسلة):  
للكلمة في السطر:  
كلمة الطباعة، 1

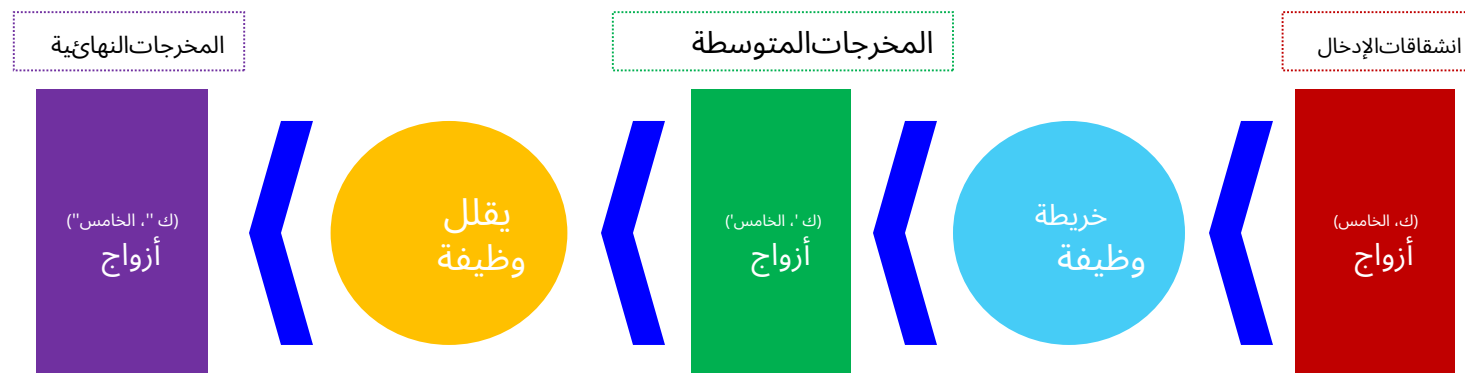
عدداالكلمات - وظيفة الخريطة

تقليل الفراغ (كلمة السلسلة، قائمة الأعداد الجزئية):  
المجموع = 0  
لـ  $c$  في التهم الجزئية:  
المجموع = + ج  
طباعة الكلمة، المجموع

عدداالكلمات - تقليل الوظيفة

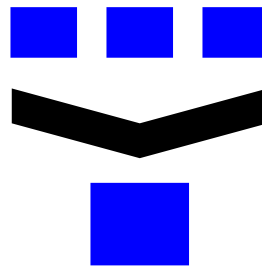
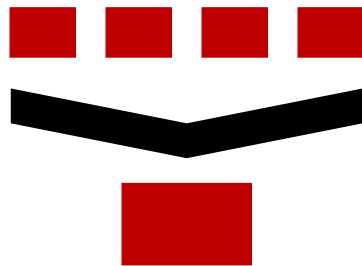
# المفاتيح والقيم

- يجب على المبرمج في MapReduce تحديد وظيفتين، وظيفة الخريطة و تقليل الوظيفة التي تنفذ Mapper و Reducer في برنامج MapReduce.
- في MapReduce، يتم تنظيم عناصر البيانات دائماً كأزواج ذات قيمة رئيسية (على سبيل المثال،  $(K, V)$ ).
- تقوم وظائف الخريطة والتقليل باستقبال وإصدار أزواج  $(K, V)$ .



# أقسام

- في MapReduce، لا يتم عادةً تقليل قيم المخرجات المتوسطة معاً
- يتم تقديم جميع القيم التي لها نفس المفتاح إلى مخفض واحد معاً
- وبشكل أكثر تحديداً، يتم تعيين مجموعة فرعية مختلفة من مساحة المفتاح المتوسطة لكل مخفض
- تُعرف هذه المجموعات الفرعية بالأقسام

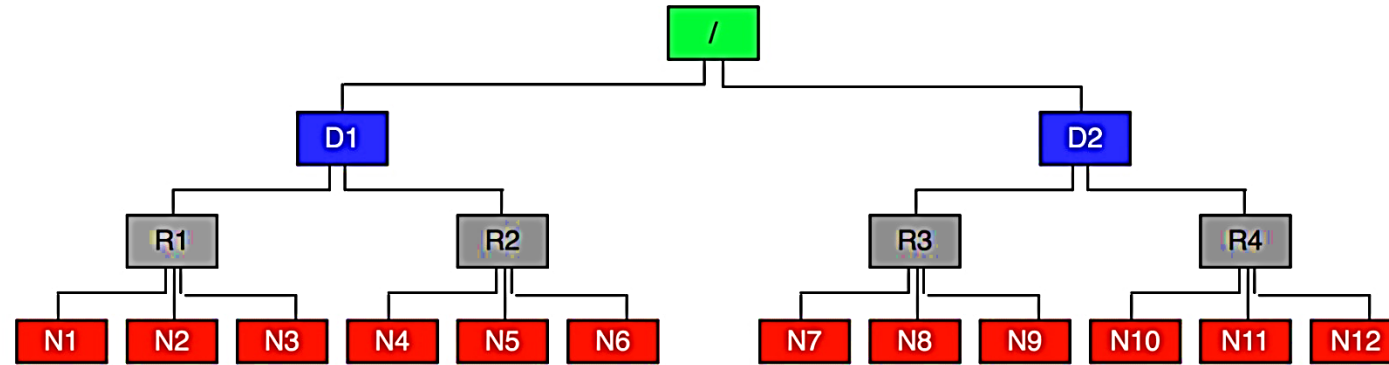


تمثل الألوان المختلفة مفاتيح  
مختلفة (من المحتمل) من  
مصممي الخرائط المختلفين

الأقسام هي المدخلات إلى المخفضات



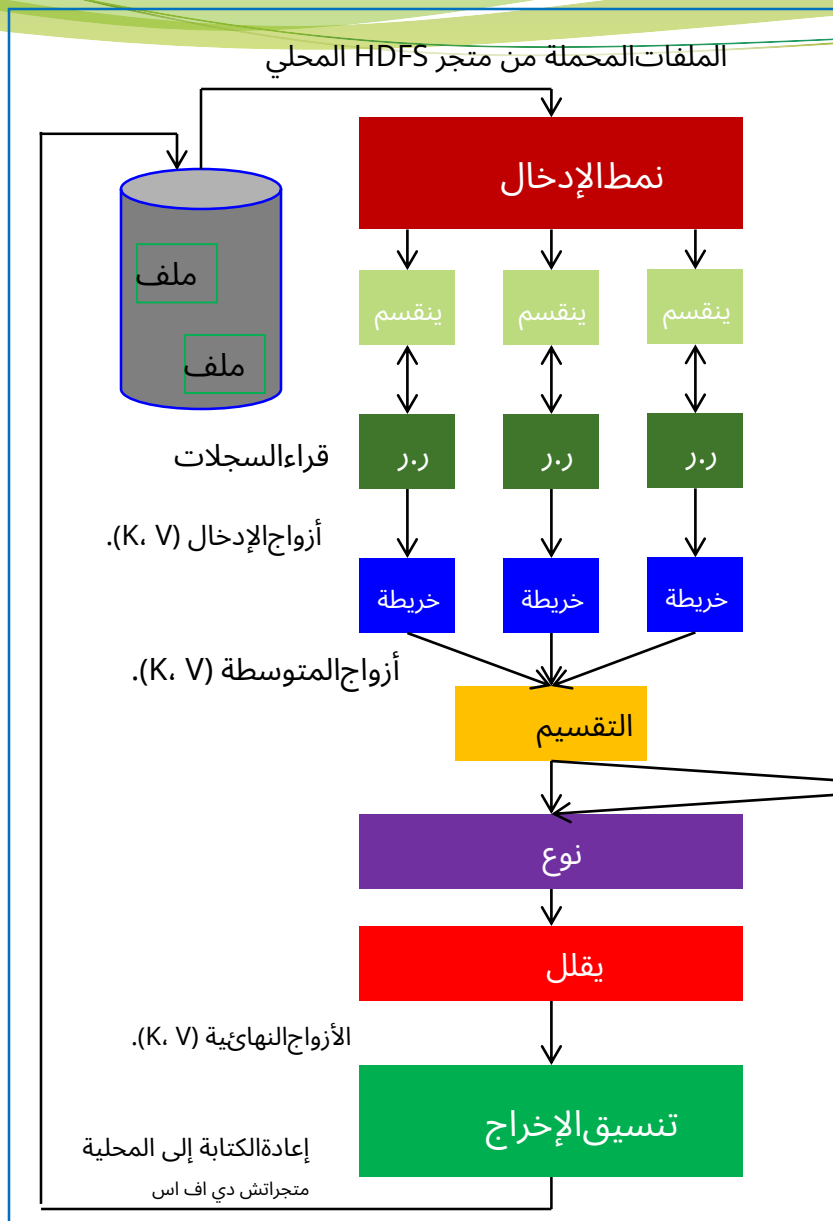
# طوبولوجيا الشبكة في MapReduce



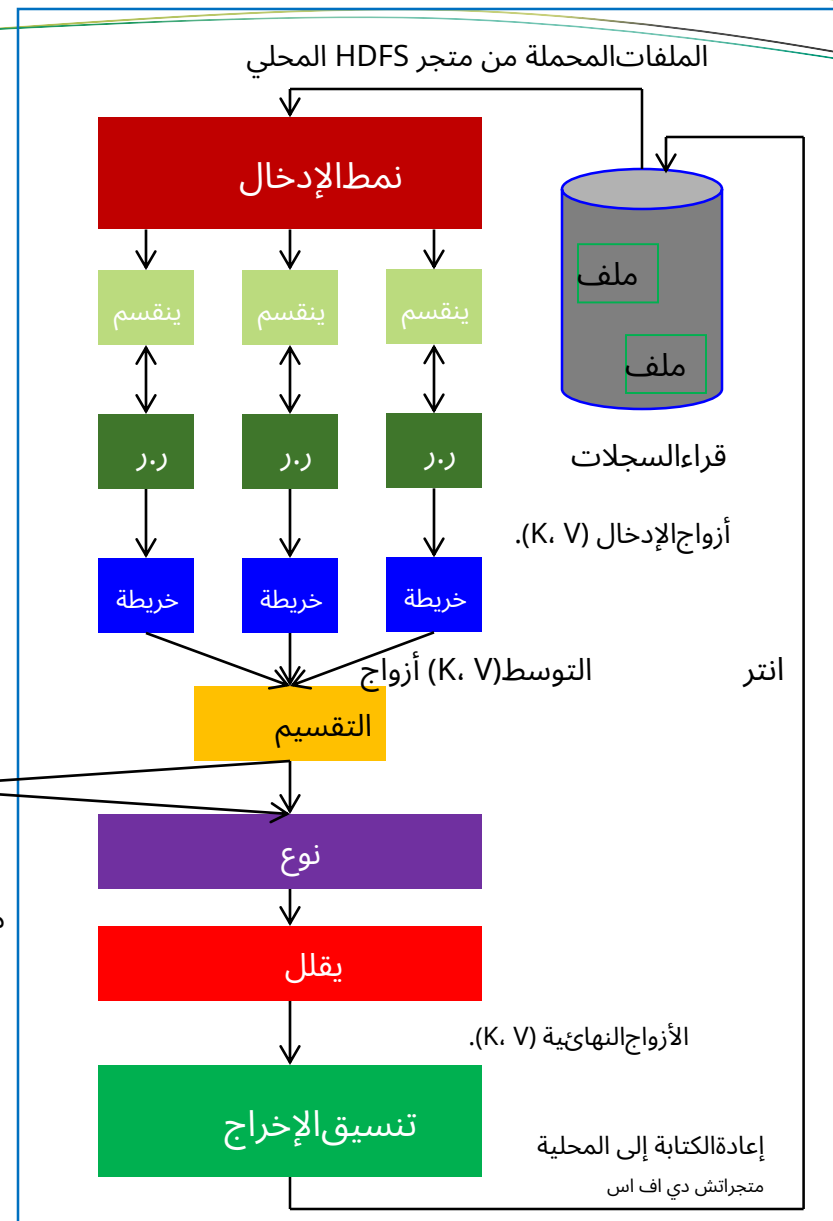
- يفترض MapReduce طوبولوجيا شبكة على شكل شجرة.
- تنتشر العقد على رفوف مختلفة متضمنة في مركز بيانات واحد أو أكثر.
- النقطة البارزة هي أن عرض النطاق الترددي بين عقدتين يعتمد على مواقعهما النسبية في طوبولوجيا الشبكة.
- على سبيل المثال، سيكون للعقد الموجودة على نفس الحامل نطاق ترددي أعلى فيما بينها على عكس العقد التي هي خارج الرف.

# نظرة فاحصة: Hadoop MapReduce

العقدة 2



العقدة 1

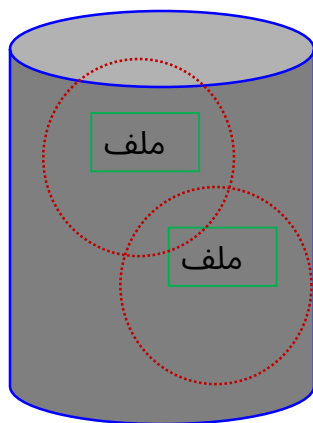


خط  
عملية

متوسط  
أزواج (K, V)  
متبادلة بواسطة  
جميع العقد

# ملفات الإدخال

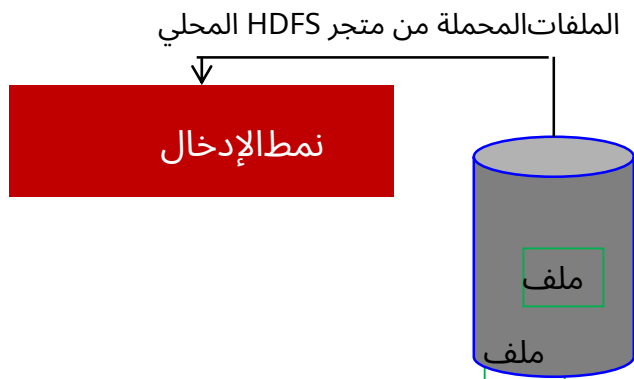
- **ملفات الإدخال** حيث البيانات ل **MapReduce** يتم تخزين المهمة في البداية
- توجد ملفات الإدخال عادةً في نظام ملفات موزع (مثل HDFS)
- تنسيق ملفات الإدخال تعسفي



- ملفات السجل المستندة إلى الخط
- الملفات الثنائية
- سجلات إدخال متعددة الأسطر
- أو أي شيء آخر تماماً

# نمط الإدخال

- تنسيق ملفات الإدخال تعسفي
- يتم تحديد كيفية تقسيم ملفات الإدخال وقراءتها بواسطة Input\_Format
- هي فئة تقوم بما يلي Input\_Format:
  - تحديد الملفات التي يجب استخدامها للإدخال
  - يحدد InputSplits التي تكسر الملف
- يوفر مصنعاً لكائنات Record\_Reader التي تقرأ الملف



# نوع تنسيق الإدخال

- يتم توفير العديد من تنسيقات الإدخال مع Hadoop:

نمط الإدخال	وصف	مفتاح	قيمة
تنسيق إدخال النص	التنسيق الافتراضي؛ يقرأ أسطر الملفات النصية	إزاحة البايت محتويات السطر	
KeyValueInputFormat	يوزع الخطوط إلى أزواج (K, V).	كل شيء إلى حرف علامة التبويب الأول	ما تبقى من الخط
SequenceFileInputFormat	خاص بـ Hadoop أداء عالي تنسيق ثنائي	تعريف المستخدم	تعريف المستخدم

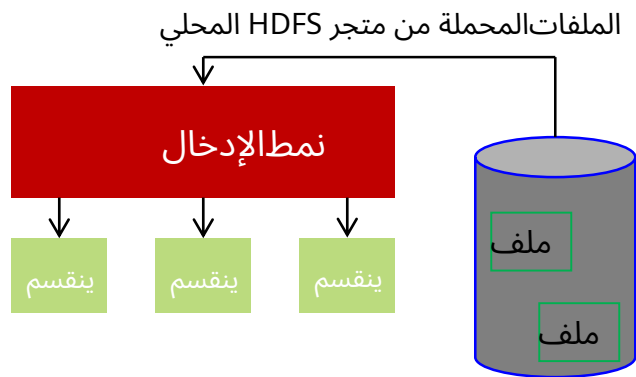
## انشقاقات الإدخال

-ان تقسيم المدخلات يصف وحدة العمل التي تشتمل على مهمة خريطة واحدة في برنامج MapReduce.

-افتراضياً، يقوم Input\_Format بتقسيم الملف إلى أجزاء يصل حجمها إلى 64 ميجابايت.  
-من خلال تقسيم الملف إلى أقسام، نسمح بتنفيذ العديد من مهام الخريطة  
تعمل على ملف واحد بالتوازي.

-إذا كان الملف كبيراً جداً، فقد يؤدي ذلك إلى تحسين الأداء بشكل ملحوظ  
من خلال التوازي.

-تتوافق كل مهمة خريطة مع أعزب تقسيم المدخلات



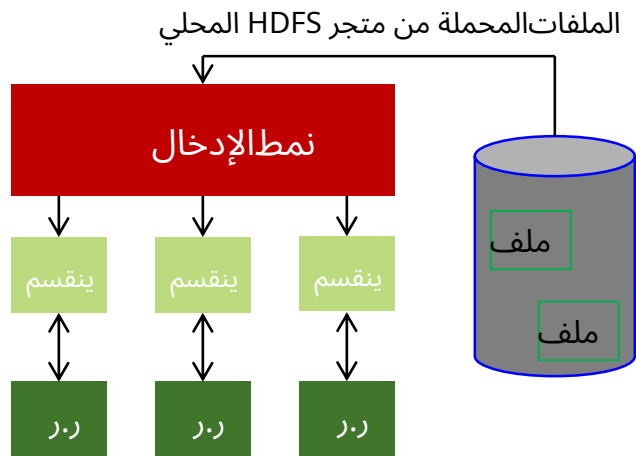
# Record\_Reader

-يحدد تقسيم الإدخال شريحة من العمل ولكنه لا يصف كيفية القيام بذلك للوصول إليه.

-ال **Record Reader** يقوم الفصل فعلياً بتحميل البيانات من مصدرها و يحولها إلى أزواج  $(K, V)$  مناسبة للقراءة بواسطة مصممي الخرائط.

-ال **Record Reader** يتم استدعاؤه بشكل متكرر على الإدخال حتى يتم استهلاك الانقسام بأكمله.

-كل دعوة من **Record Reader** يؤدي إلى استدعاء آخر لوظيفة الخريطة التي يحددها المبرمج.



# مخطط والمخفض

-ال مصمم الخرائط ينفذ العمل المحدد من قبل المستخدم للمرحلة

الأولى من برنامج MapReduce

-يتم إنشاء مثيل جديد لـ Mapper لكل تقسيم

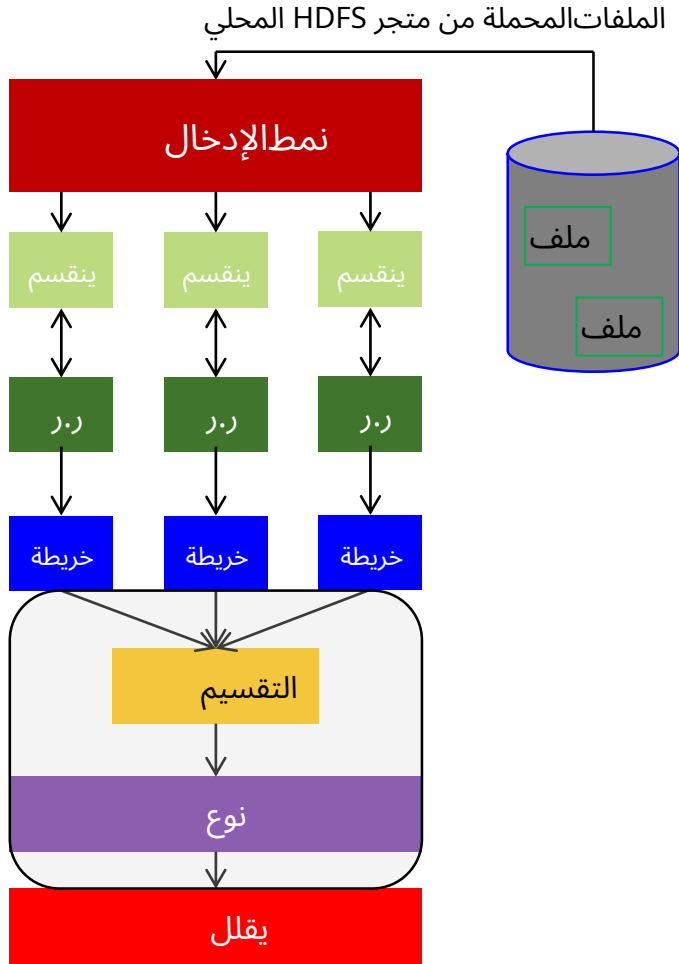
-ال المخفض ينفذ العمل المحدد من قبل المستخدم للثانية

مرحلة برنامج MapReduce

-يتم إنشاء مثيل جديد من المخفض لكل قسم

-لكل مفتاح في القسم المخصص للمخفض، فإن

يتم استدعاء المخفض مرة واحدة.





# جدولة المهام في MapReduce

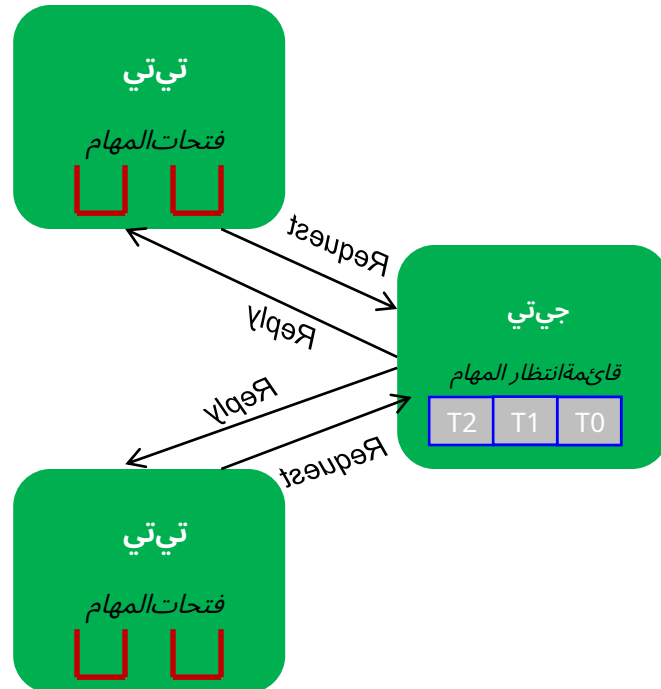
- يعتمد MapReduce بنية السيد والعبد

- يشار إلى العقدة الرئيسية في MapReduce باسم تعقب الوظيفة (JT)

- تتم الإشارة إلى كل عقدة تابعة في MapReduce باسم تعقب المهام (TT)

- يعتمد MapReduce جدولة السحب استراتيجية بدلا من دفعة واحدة

- على سبيل المثال، لا يقوم JT بدفع الخريطة وتقليل المهام إلى TTs بل يقوم TTs بسحبها عن طريق تقديم الطلبات ذات الصلة



# Job\_Tracker

-Hadoop في MapReduce هي الخدمة الروحية لإرسال وتتبع وظائف Job\_Tracker

-يقوم Job\_Tracker بتنفيذ الإجراءات التالية في Hadoop:

- يقبل وظائف MapReduce من تطبيقات العميل
- يتحدث إلى Name\_Node لتحديد موقع البيانات
- يحدد موقع عقدة Task\_Tracker المتاحة
- يرسل العمل إلى عقدة Task\_Tracker المختارة

# Task\_Tracker

-تقبل عقدة Task\_Tracker الخريطة أو تقليل العمليات أو تبديلها عشوائياً من Job\_Tracker.

-تم تكوينه بمجموعة من الفتحات، والتي تشير إلى عدد المهام التي يمكنه قبولها.

-يبحث Job\_Tracker عن الفتحة المجانية لتعيين وظيفة.

-يقوم Task\_Tracker بإعلام Job\_Tracker بحالة نجاح المهمة.

-يرسل Task\_Tracker أيضاً إشارات نبضات القلب إلى متتبع الوظائف للتأكد من توفرها، كما يقوم أيضاً بالإبلاغ عن الرقم. من الفتحات المجانية المتاحة معها.

# خريطة وتقليل جدولة المهام

- كل TT يرسل أرسالة نبضات القلب بشكل دوري إلى IT يتضمن طلباً لخريطة أو مهمة تقليل للتشغيل.

أنا. جدولة مهمة الخريطة:

- تلبية IT طلبات مهام الخريطة من خلال محاولة جدولة رسامي الخرائط في المنطقة من انقسامات المدخلات الخاصة بهم (أي أنها تعتبر المنطقة).

ثانياً. تقليل جدولة المهام:

- ومع ذلك، تقوم IT ببساطة بتعيين مهمة التخفيض التالية التي لم يتم تشغيلها بعد إلى TT الطالبة بغض النظر عن موقع شبكة TT وتأثيرها الضمني على وقت التبديل العشوائي للمخفض (أي أنه لا يأخذ في الاعتبار المنطقة المحلية).

# جدولة الوظائف في MapReduce

- في MapReduce، يتم تمثيل التطبيق على أنه عمل.
- تشمل الوظيفة خريطة متعددة وتقلل من المهام.
- يأتي MapReduce في Hadoop مع مجموعة مختارة من أدوات الجدولة:
  - الافتراضي هو جدولة FIFO الذي يقوم بجدولة الوظائف حسب ترتيب التقديم
  - يوجد أيضاً برنامج جدولة متعدد المستخدمين يسمى جدولة عادلة والتي تهدف إلى إعطاء كل استخدام حصة عادلة من سعة المجموعة مع مرور الوقت.

# التسامح مع الخطأ في Hadoop

-يمكن لـ MapReduce توجيه المهام نحو إكمال ناجح حتى عند تشغيل المهام على مجموعة كبيرة حيث تزداد احتمالية الفشل

-الطريقة الأساسية التي يحقق بها MapReduce التسامح مع الخطأ هي من خلال إعادة تشغيل المهام

-إذا فشل TT في الاتصال بـ JT لفترة من الوقت (افتراضياً، دقيقة واحدة).

المعنية قد تحطمت TT أن JT سوف تفترض ، Hadoop)

-إذا كانت المهمة لا تزال في مرحلة الخريطة، يطلب JT من TT آخر القيام بذلك إعادة تنفيذ كافة مصممي الخرائط ذلك

ركض سابقا في TT الفاشلة

-إذا كانت المهمة في مرحلة التخفيض، فإن JT يطلب TT آخر إعادة تنفيذ كافة المخفضات ذلك

كانت قيد التقدم في TT الفاشلة

# مالذي يجعل MapReduce فريداً؟

-يتميز MapReduce بما يلي:

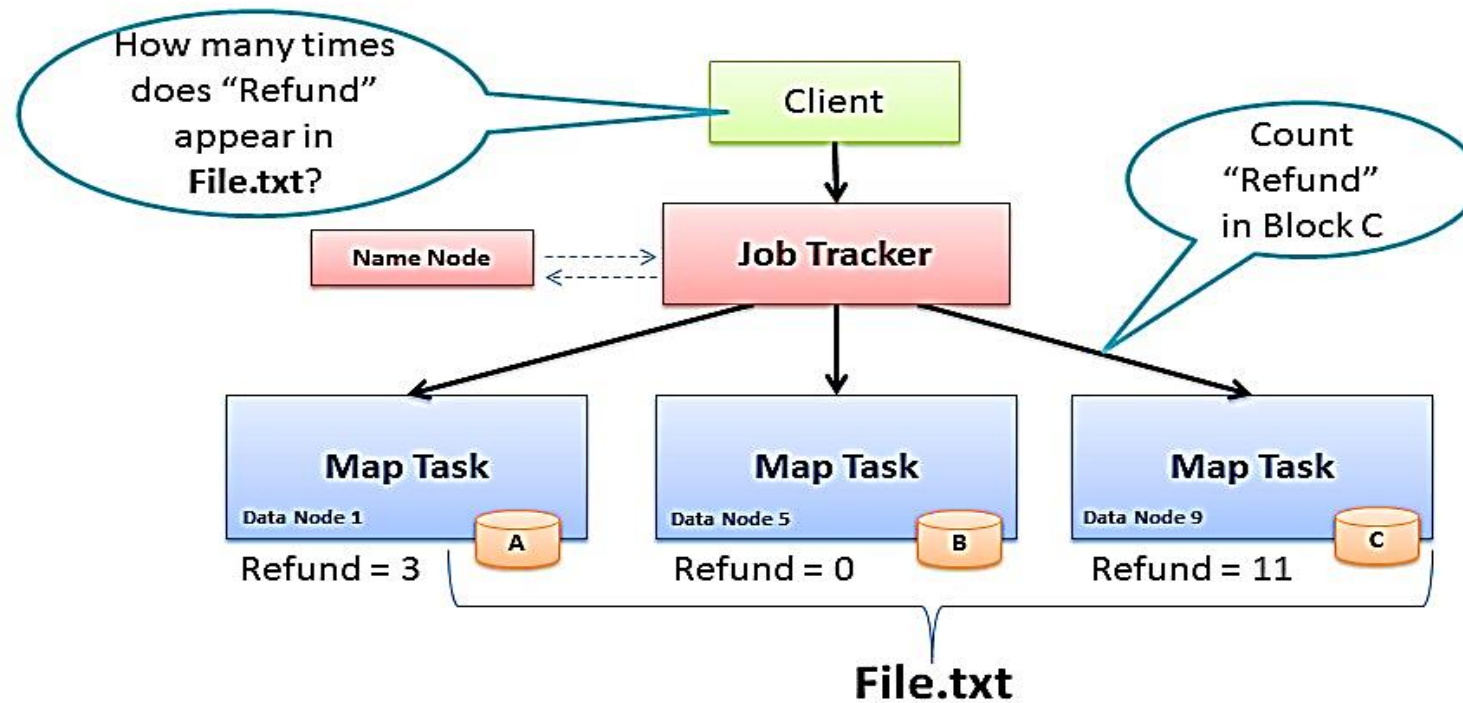
1. نموذج البرمجة المبسط الذي يسمح للمستخدم بكتابة واختبار الأنظمة الموزعة بسرعة

2. التوزيع الفعال والتلقائي للبيانات وعبء العمل عبر الأجهزة

3. منحى قابلية التوسع المسطح. على وجه التحديد، بعد كتابة برنامج MapReduce وتشغيله على 10

عقد، لا يلزم سوى القليل جداً من العمل، إن وجد، للقيام بذلك

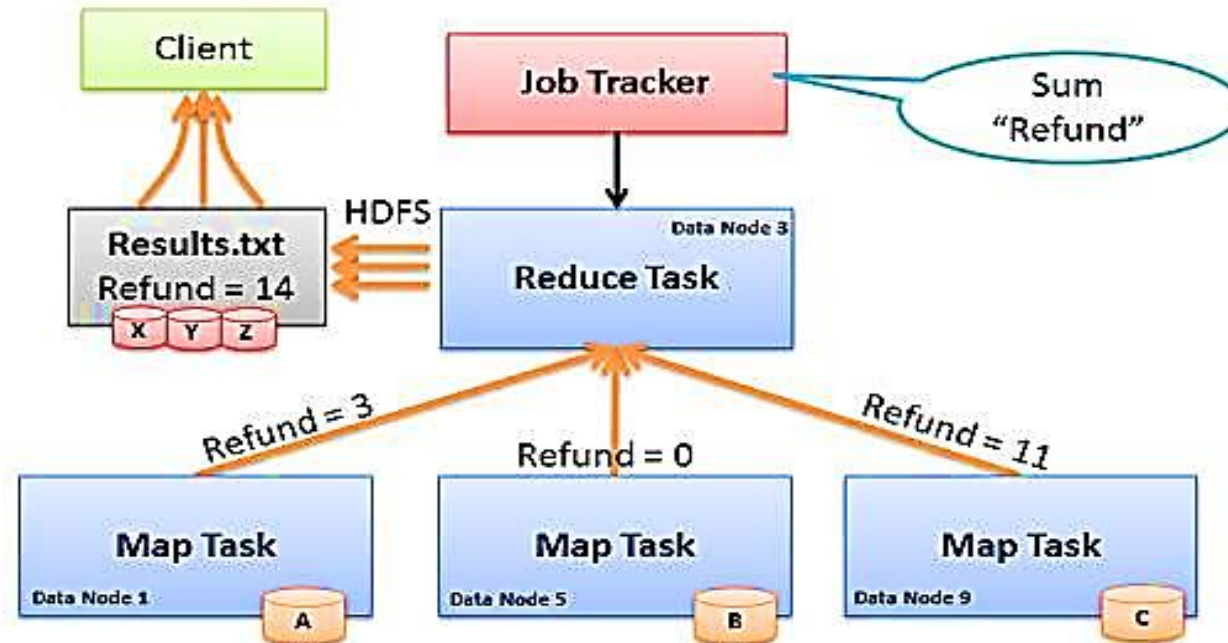
نفس البرنامج يعمل على 1000 عقدة



- **Map:** "Run this computation on your local data"
- Job Tracker delivers Java code to Nodes with local data



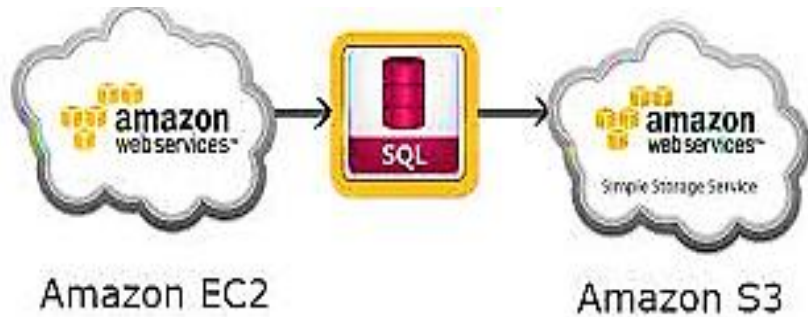
# تقليل المهمة



- **Reduce:** "Run this computation across Map results"
- Map Tasks send output data to Reducer over the network
- Reduce Task data output written to and read from HDFS

# AWS: السحابة MapReduce

- يوفر واجهة قائمة على الويب وأدوات سطر أوامر لتشغيل وظائف Hadoop على Amazon **EC2** (سحابة الحوسبة المرنة من أمازون)



- البيانات المخزنة في الأمازون **S3** (خدمة التخزين البسيطة من أمازون)

- يراقب العمل ويغلق الآلات بعد الاستخدام

- رسوم إضافية صغيرة على رأس **EC2** التسعير

- ان **EC2** المثال يشبه جهاز كمبيوتر بعيد يعمل بنظام Windows أو Linux ويمكنك ذلك

قم بتثبيت أي برنامج تريده، بما في ذلك خادم ويب يقوم بتشغيل كود PHP و خادم قاعدة البيانات.

- أمازون **S3** هي مجرد خدمة تخزين، تستخدم عادةً لتخزين الملفات الثنائية الكبيرة.

# أمازون مطا MapReduce

[Sign Up](#)[My Account / Console](#)[English](#)[AWS Products & Solutions](#)[AWS Product Information](#)[Developers](#)[Support](#)

## Amazon EMR

- [Amazon EMR Overview](#)
- [FAQs](#)
- [Pricing](#)

## Developer Resources

- [EMR Training](#)
- [AWS Management Console](#)
- [Documentation](#)
- [Release Notes](#)
- [Sample Code & Libraries](#)
- [Developer Tools](#)
- [Articles & Tutorials](#)
- [Community Forum](#)

## Amazon Elastic MapReduce (Amazon EMR)

Amazon Elastic MapReduce (Amazon EMR) is a web service that enables businesses, researchers, data analysts, and developers to easily and cost-effectively process vast amounts of data. It utilizes a hosted Hadoop framework running on the web-scale infrastructure of Amazon Elastic Compute Cloud (Amazon EC2) and Amazon Simple Storage Service (Amazon S3).

Using Amazon Elastic MapReduce, you can instantly provision as much or as little capacity as you like to perform data-intensive tasks for applications such as web indexing, data mining, log file analysis, data warehousing, machine learning, financial analysis, scientific simulation, and bioinformatics research. Amazon Elastic MapReduce lets you focus on crunching or analyzing your data without having to worry about time-consuming set-up, management or tuning of Hadoop clusters or the compute capacity upon which they sit.

New to EMR? Check out these resources:

- [Training](#)
- [Documentation](#)
- [FAQs](#)
- [EMR Forum](#)

Easy to sign up,  
pay only for what you  
use

[Sign Up Now](#)

## Featured Quote

razorfish.

"With Amazon Elastic MapReduce, there was no upfront investment in hardware, no hardware procurement delay, and no need to hire additional operations staff. Because of the flexibility of the platform, our first new online advertising campaign experienced a 500% increase in return on ad spend from a similar campaign a year before." Read the full [case study](#).

# مرونة MapReduce سير العمل

# مرونة MapReduce سير العمل

# مرونة MapReduce سير العمل

# مرونة MapReduce سير العمل

# المكون المركزي لأمازون EMR

-المكون المركزي للأمازون السجلات الطبية الإلكترونية هل تَجَمَع.

-مجموعة عبارة عن مجموعة من Amazon Compute Cloud (Amazon Elastic EC2) الحالات.

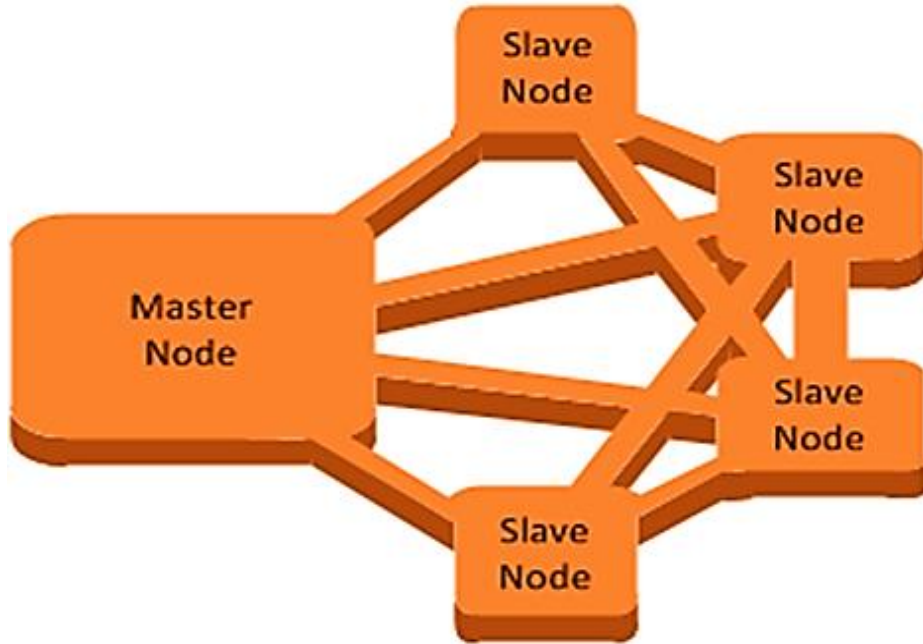
-كل مثل في الكتلة يسمى الأنود.

-كل عقدة له دور داخل المجموعة، يشار إليه بنوع العقدة. يقوم EMR

Amazon أيضاً بتثبيتات مختلفة

مكونات البرنامج على كل نوع عقدة، مع إعطاء كل منها

قمة عقد دور في تطبيق موزع مثل Apache Hadoop.



Cluster



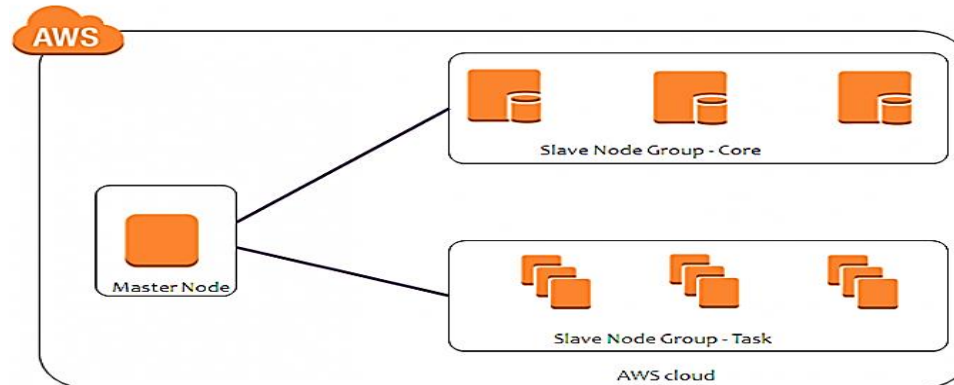
# أنواع العقد في Amazon EMR هي كما يلي:

**-العقدة الرئيسية:**عقدة تدير المجموعة عن طريق تشغيل مكونات برمجية لتنسيق توزيع البيانات والمهام بين العقد الأخرى- والتي يشار إليها مجتمعة باسم العقد التابعة -

للمعالجة.تتبع العقدة الرئيسية حالة المهام وتراقب صحة المجموعة.

**-العقدة الأساسية:**عقدة تابعة تحتوي على مكونات برمجية تقوم بتشغيل المهام وتخزين البيانات في Hadoop نظام الملفات الموزعة (HDFS) على مجموعتك.

**-عقدة المهمة:**عقدة تابعة تحتوي على مكونات برمجية تقوم بتشغيل المهام فقط. العقد المهمة هي اختياري.



... شكراً لك ...

