

Map-Reduce и параллельные аналитические базы данных

Сергей Кузнецов
ИСП РАН, ЦИТФорум
kuzloc@ispras.ru



Введение (1)

- Клермонтский отчет:
 - ... сбор, интеграция и анализ данных больше не считаются расходами на ведение бизнеса; данные – это ключ к достижению эффективности и прибыльности бизнеса. В результате быстро развивается индустрия, поддерживающая анализ данных

апреля

2010 г.

данных

2010

2

Введение (2)

- К концу прошлого века аналитические средства баз данных можно было пересчитать по пальцам одной руки
- IBM DB2, Teradata, Sybase IQ, Oracle, частично Microsoft SQL Server
 - только в DB2 и Teradata поддерживалась архитектура sharing-nothing
 - только в Sybase IQ использовалось поколоночное хранение таблиц

23
апреля

2010 г.

Корпоративные
Базы
данных

2010

3

Введение (3)

- С начала 2000-х активизировалось направление Data Warehouse Appliance или Analytic Appliance
 - аббревиатура DWAA
- Основная цель - создание аппаратно-программных средств
 - существенно более дешевых, чем у предлагаемых поставщиками универсальных СУБД, но при этом
 - обеспечивающих большую производительность и масштабируемость при работе со сверхбольшими хранилищами данных

23
апреля

Корпоративные
базы
данных

2010 г.

2010

4

Введение (4)

Аналитические параллельные СУБД (1)

- Направление DWAA появилось еще в 1980-е гг.
- Соответствующие пионерские продукты были созданы в компании Britton Lee Inc., которая в 1989 г. была сначала
 - переименована в ShareBase Corporation, а затем
 - поглощена компанией Teradata, которая к этому времени тоже придерживалась подхода DWAA

Корпоративные
базы
данных

23
апреля

2010 г.

2010

5

Введение (5)

Аналитические параллельные СУБД (2)

- Аппаратно-программное решение, основанное на ассоциативной адресации элементов хранения данных, имелось у компании ICL – Content Addressable File Store
- Однако на рынке систем поддержки хранилищ данных на основе подхода DWA до сих пор осталась только Teradata

23
апреля

Корпоративные
базы
данных

2010 г.

2010

6

Введение (6)

Аналитические параллельные СУБД (3)

- Возрождение направления DWAA в начале 2000-х связано с ростом заинтересованности компаний в недорогих и эффективных решениях для поддержки хранилищ данных и их анализа
- Стали возникать софтверные стартапы, первым из которых стала компания Netezza

23
апреля

2010 г.

Корпоративные
базы
данных

2010

7

Введение (7)

Аналитические параллельные СУБД (4)

- Эффективное DWAA-решение на основе
 - программируемых вентильных матриц (Field Programmable Gate Array, FPGA) и
 - процессоров PowerPC
- Использование FPGA в контроллерах магнитных дисков позволяет осуществлять "на лету" первичную фильтрацию данных
- Применение PowerPC вместо Intel (по утверждению компании) позволяет снизить энергопотребление и расходы на охлаждение данных

23
апреля
2010 г.

Корпоративные
базы
данных
2010

8

Введение (8)

Аналитические параллельные СУБД (5)

- С тех пор появилось еще около десяти новых компаний, ориентирующихся на разработку DWAA с применением почти всегда
 - разновидности массивно-параллельной архитектуры (MPP)
 - "sharing-nothing"
- Vertica Systems
 - MPP, колоночное хранение таблиц

23
апреля
2010 г.

2010

9

Введение (9)

Аналитические параллельные СУБД (6)

- DATAlegro Inc.
 - недавно поглощена Microsoft
 - проект Madison, ставший основой SQL Server 2008 R2 Parallel Data Warehouse
 - MPP
 - основана на использовании СУБД Ingres
 - тем самым, таблицы хранятся по строкам

Корпорат

ивные

базы

данных

апреля

2010 г.

2010

10

Введение (10)

Аналитические параллельные СУБД (7)

- Greenplum
 - MPP
 - система основана на использовании СУБД PostgreSQL
 - тем самым, таблицы хранятся по строкам
- Aster Data Systems
 - MPP
 - таблицы хранятся по строкам
- Kognitio
 - MPP
 - таблицы хранятся по строкам

Корпоративные

базы

данных

23
апреля

2010

2010 г.

2010

11

Введение (11)

Аналитические параллельные СУБД (8)

- EXASOL AG
 - MPP
 - поколонное хранение таблиц
- Calpont Corporation
 - MPP, поколонное хранение таблиц
 - система (InfiniDB) внешне схожа с MySQL
- Dataupia Corporation
 - MPP

Корпоративные
базы

таблицы хранятся по строкам

апреля

данных

2010 г.

2010

12

Введение (12)

Аналитические параллельные СУБД (9)

- Infobright
 - поколоночное хранение таблиц
 - система основана на MySQL
 - ориентирована на использование многоядерных процессоров
 - массивный параллелизм не используется
- Kickfire
 - поколоночное хранение таблиц
 - используется специальная аппаратура, ускоряющая выполнение SQL-запросов
 - система создана на основе MySQL
 - не основана на массивно-параллельной архитектуре

Корпорат

ивные

базы

данных

23
апреля

2010 г.

2010

13

Введение (13)

Аналитические параллельные СУБД (10)

- Подход DWAA проникает и в продукты основных поставщиков SQL-ориентированных СУБД
- Разработка компании DATAlegro стала основой массивно-параллельного варианта Microsoft SQL Server
 - SQL Server 2008 R2 Parallel Data Warehouse
- Oracle обеспечивает специализированное массивно-параллельное хранилище табличных данных Oracle Exadata Storage Server
 - позволяющее значительно ускорить работу основной СУБД

Корпоративные
базы

23
апреля

данных

2010 г.

2010

14

Введение (14)

Аналитические параллельные СУБД (11)

- У разных решений категории DWAA имеются свои интересные технические особенности, заслуживающие более глубокого обсуждения, анализа и сравнения
- Их можно классифицировать и сравнивать по разным критериям
- Однако это не является целью доклада
- Некоторую попытку такого анализа представляет собой обзор
 - *Richard Hackathorn, Colin White. Data Warehouse Appliances: Evolution or Revolution?*
 - <http://www.beyerresearch.com/study/4639>
- Значительный рост интереса к направлению DWAA, к специализированным СУБД вдобавител к СУБД Vertica в частности вызвала статья
 - *Майкл Стоунбрейкер, Угур Кетинтемел. Один размер пригоден для всех: идея, время которой пришло и ушло*
 - http://citforum.ru/database/articles/one_size_fits_all/

апреля

2010 г.

данных

2010

15

Введение (15)

При чем здесь MapReduce? (1)

- Сосредоточимся на частном, но очень важном в настоящее время вопросе взаимоотношений технологий массивно-параллельных аналитических СУБД и MapReduce
- Контекст DWAA является естественным, поскольку большинство СУБД, созданных на основе подхода DWAA, являются массивно-параллельными без использования общих ресурсов
- Системы создавались в расчете на использование в кластерной аппаратной архитектуре, и они сравнительно легко могут быть перенесены в "облачную" среду динамически конфигурируемых кластеров

23
апреля
2010 г.

Корпоративные
базы
данных
2010

16

Введение (16)

При чем здесь MapReduce? (2)

- Появление "родной" для "облачной" среды технологии MapReduce и в особенности
 - энтузиазм по части ее использования, проявленный многими потенциальными пользователями параллельных СУБД,
- очень озаботили представителей направления DWAA
- Сначала авторитетные представители сообщества баз данных и одновременно активные сторонники подхода DWAA Майкл Стоунбрейкер и Дэвид Девитт старались убедить общественность в том, что MapReduce
 - это технология, уступающая технологии параллельных баз данных по всем статьям

23
апреля

2010 г.

Корпоративные
базы
данных

2010

17

Введение (17)

При чем здесь MapReduce? (3)

- Потом была проведена серия экспериментов, продемонстрировавшая, что при решении типичных простых аналитических задач MapReduce
 - уступает в производительности не только поколоночной СУБД Vertica, но и традиционной массивно-параллельной СУБД с хранением таблиц по строкам
 - *Майкл Стоунбрейкер и др. Сравнение подходов к крупномасштабному анализу данных*
 - http://citforum.ru/database/articles/mr_vs_dbms/
- Доводы и результаты экспериментов были весьма солидными и убедительными, и
 - вряд ли кто-нибудь из людей, знакомых с обеими технологиями, сомневается в том, что
 - MapReduce не вытеснит параллельные СУБД, и что
 - эти технологии будут благополучно сосуществовать в "облаках" и в среде кластерных архитектур вообще

23
апреля

2010 г.

Корпоративные
базы
данных

2010

18

Введение (18)

При чем здесь MapReduce? (4)

- Однако возникает другой вопрос:
 - нет ли в технологии MapReduce каких-либо положительных черт, которых не хватает параллельным СУБД?
- Можно ли каким-либо образом добавить эти черты в параллельные СУБД, сохранив их основные качества:
 - декларативный доступ на языке SQL,
 - оптимизацию запросов и т.д.
- Понятно, что у параллельных СУБД имеется масса положительных черт, которыми не обладает MapReduce, но, похоже, что добавление их к MapReduce изменило бы суть этой технологии, превратив ее в технологию параллельных СУБД

23
апреля

Корпоративные
базы
данных

2010 г.

2010

19

Введение (19)

При чем здесь MapReduce? (5)

- На эти два вопроса удалось получить положительный ответ
- В нескольких проектах, связанных с направлением DWAA, удалось воспользоваться такими преимуществами MapReduce, как
 - масштабируемость до десятков тысяч узлов,
 - отказоустойчивость,
 - дешевизна загрузки данных,
 - возможность использования явно написанного кода, который хорошо распараллеливается
- Ни в одном проекте не удалось воспользоваться сразу всеми этими преимуществами, но имеющиеся достижения позволяют
 - добавить в параллельные СУБД важные качества, которыми они до сих пор не обладали

апреля

2010 г.

корпоративные
базы
данных

2010

20

Введение (20)

При чем здесь MapReduce? (6)

- Рассмотрим три подхода к интеграции технологий MapReduce и параллельных СУБД, предложенных и реализованных специалистами
 - компаний Greenplum и Aster Data
 - университетов Yale и Brown
 - компании Vertica,
- которые можно было бы назвать:
 - MapReduce внутри параллельной СУБД
 - СУБД внутри среды MapReduce и
 - MapReduce сбоку от параллельной СУБД

23
апреля
2010 г.

Корпоративные
базы
данных
2010

21

Введение (21)

При чем здесь MapReduce? (7)

- Первый подход ориентирован на поддержку написания и выполнения хранимых на стороне сервера баз данных пользовательских функций, которые хорошо распараллеливаются в кластерной среде
 - используется преимущество MapReduce по применению явно написанного кода и его распараллеливанию
- Второй подход направлен на использование MapReduce в качестве инфраструктуры параллельной СУБД, в качестве базовых компонентов которой используются традиционные не параллельные СУБД
 - применение MapReduce позволяет добиться неограниченной масштабируемости получаемой системы и ее отказоустойчивости на уровне выполнения запросов
- При применении третьего подхода MapReduce используется для выполнения процедуры ETL над исходными данными до их загрузки в систему параллельных баз данных
 - используется преимущество MapReduce в отношении дешевой загрузки данных до их обработки

23
апреля

2010 г.

Корпоративные
базы
данных

2010

22

MapReduce: модель и реализации (1)

- Программная модель MapReduce была придумана несколько лет тому назад в компании Google, и там же была выполнена
 - первая реализация этой модели на основе распределенной файловой системы GFS (Google File System)
- Эта реализация активно используется в программных продуктах самой Google, но является
 - сугубо проприетарной и недоступна для использования вне Google

23
апреля

2010 г.

Корпорат
ивные
базы
данных

2010

23

MapReduce: модель и реализации (2)

- Свободно доступная реализация Hadoop MapReduce была выполнена в проекте Hadoop сообщества Apache
- Она основана на использовании распределенной файловой системы HDFS (Hadoop Distributed File System)
- Реальную популярность MapReduce принесла именно реализация Hadoop в силу своей доступности и открытости
- Широкое использование Hadoop MapReduce в различных исследовательских и исследовательских проектах приносит несомненную пользу этой системе,
 - стимулируя разработчиков к ее постоянному совершенствованию

23
апреля

Корпоративные
базы
данных

2010 г.

2010

24

MapReduce: модель и реализации (3)

- Однако реализация Hadoop MapReduce полностью основана на спецификациях Google, и поэтому каноническим описанием технологии была и остается статья
 - *Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified Data Processing on Large Clusters*
 - <http://labs.google.com/papers/mapreduce.html>
- Заметим, что в документации Hadoop MapReduce используется несколько отличная терминология
- В этой части доклада из уважения к первенству Google используются оригинальные термины,
 - в следующих разделах, там где будет иметься в виду конкретно реализация Hadoop MapReduce, будет использоваться терминология Hadoop

23
апреля
2010 г.

Корпоративные
базы
данных
2010

25

MapReduce: модель и реализации (3)

Общая модель программирования MapReduce (1)

- Вычисления производятся
 - над множествами входных пар "ключ-значение", и
- в результате каждого вычисления также производится
 - некоторое множество результирующих пар "ключ-значение"
- Для представления вычислений в среде MapReduce используются две основные функции:
 - *Map* и
 - *Reduce*
- Обе функции явно кодируются разработчиками приложений в среде MapReduce

Корпоративные
базы
данных

23
апреля

2010 г.

2010

26

MapReduce: модель и реализации (4)

Общая модель программирования MapReduce (2)

- Функция *Map* в цикле обрабатывает каждую пару из множества входных пар и производит множество промежуточных пар "ключ-значение"
- Среда MapReduce группирует все промежуточные значения с одним и тем же ключом *I* и передает их функции *Reduce*.
- Функция *Reduce* получает значение ключа *I* и множество значений, связанных с этим ключом
- В типичных ситуациях каждая группа обрабатывается (в цикле) таким образом, что в результате одного вызова функции образуется не более одного результирующего значения

23
апреля

2010 г.

Корпоративные
Базы
данных

2010

27

MapReduce: модель и реализации (5)

Реализация в распределенной среде (1)

- Реализации MapReduce от Google и Hadoop ориентированы на использование в кластерной распределенной среде со следующими основными характеристиками:
 - узлы среды выполнения MR-приложений представляют собой компьютеры общего назначения с операционной системой Linux;
 - используется стандартное сетевое оборудование с адаптерами, рассчитанными на скорости передачи в 100 мегабит в секунду или 1 гигабит в секунду,
 - но средняя пропускная способность существенно ниже;
 - кластер состоит из сотен или тысяч машин, так что вполне вероятны отказы отдельных узлов;
 - для хранения данных используются недорогие дисковые устройства, подключенные напрямую к отдельным машинам;
 - для управления данными, хранящимися на этих дисках, используется распределенная файловая система;
 - пользователи представляют свои задания в систему планирования; каждое задание состоит из некоторого набора задач, которые отображаются планировщиком на некоторый набор узлов кластера

23
апреля

2010 г.

Корпоративные
базы
данных

2010

28

MapReduce: модель и реализации (6)

Реализация в распределенной среде (2)

Выполнение MR-приложения (1)

- Вызовы *Map* распределяются по нескольким узлам кластера путем деления входных данных на M непересекающихся групп (split)
- Входные группы могут параллельно обрабатываться на разных машинах
- Вызовы *Reduce* распределяются путем деления пространства промежуточных ключей на R частей с использованием некоторой функции разделения – например, функции хэширования
- Число разделов R и функция разделения задаются пользователем

апреля

данных

2010 г.

2010

29

MapReduce: модель и реализации (7)

Реализация в распределенной среде (3)

Выполнение MR-приложения (2)

- Выполнение MR-программы происходит следующим образом
- Сначала среда MapReduce расщепляет входной файл на M частей,
 - размер которых может задаваться пользователем
- Затем сразу в нескольких узлах кластера запускается основная программа MapReduce
- Один из экземпляров этой программы играет специальную роль и называется *распорядителем (master)*
- Остальные экземпляры являются *исполнителями (worker)*,
 - которым распорядитель назначает работу
- Распорядитель должен назначить исполнителям для выполнения M задач *Map* и R задач *Reduce*

23
апреля

2010 г.

Корпоративные
базы
данных

2010

30

MapReduce: модель и реализации (8)

Реализация в распределенной среде (4)

Выполнение MR-приложения (3)

- Исполнитель задачи *Map*
 - читает содержимое соответствующей группы,
 - разбирает пары "ключ-значение" входных данных и
 - передает каждую пару в определенную пользователем функцию *Map*
- Промежуточные пары "ключ-значение", производимые функцией *Map*, буферизуются в основной памяти
- Периодически буферизованные пары, разделяемые на R областей на основе функции разделения, записываются в локальную дисковую память исполнителя
- Координаты этих сохраненных на диске буферизованных пар отсылаются распорядителю, который,
 - передает эти координаты исполнителям задачи *Reduce*
- i -ый *Reduce*-исполнитель снабжается координатами всех i -ых областей буферизованных пар, произведенных всеми M *Map*-исполнителями

23
апреля

Корпоративные
базы
данных

2010 г.

2010

31

MapReduce: модель и реализации (9)

Реализация в распределенной среде (5)

Выполнение MR-приложения (4)

- После получения этих координат исполнитель задачи *Reduce* с использованием механизма RPC
 - переписывает данные с локальных дисков исполнителей задачи *Map* в свою память или на локальный диск
- После переписки всех промежуточных данных выполняется их сортировка по значениям промежуточного ключа
 - для образования групп с одинаковым значением ключа
- Если объем промежуточных данных слишком велик для выполнения сортировки в основной памяти
 - используется внешняя сортировка

23
апреля

корпоративные
базы
данных

2010 г.

2010

32

MapReduce: модель и реализации (10)

Реализация в распределенной среде (5)

Выполнение MR-приложения (4)

- Далее *Reduce*-исполнитель организует цикл по отсортированным промежуточным данным и для каждого уникального значения ключа
 - вызывает пользовательскую функцию *Reduce* с передачей ей в качестве аргумента значения ключа и соответствующего множества значений
- Результирующие пары функции *Reduce* добавляются в окончательный результирующий файл данного *Reduce*-исполнителя
- После завершения всех задач *Map* и *Reduce* распорядитель активизирует программу пользователя, вызывавшую MapReduce

23
апреля

2010 г.

Корпоративные
базы
данных

2010

33

MapReduce: модель и реализации (11)

Реализация в распределенной среде (6)

Выполнение MR-приложения (5)

- После успешного завершения выполнения задания MapReduce результаты размещаются в *R* файлах распределенной файловой системы
 - имена этих результирующих файлов задаются пользователем
- Обычно не требуется объединять их в один файл, потому что часто полученные файлы используются в качестве входных
 - для запуска следующего MR-задания или
 - в каком-либо другом распределенном приложении, которое может получать входные данные из нескольких файлов

23
апреля

2010 г.

Корпоративные
базы
данных

2010

34

MapReduce: модель и реализации (12)

Реализация в распределенной среде (7)

Отказоустойчивость (1) Отказ исполнителя (1)

- Распорядитель периодически посылает каждому исполнителю контрольные сообщения
- Если некоторый исполнитель не отвечает на такое сообщение в течение некоторого установленного времени,
 - распорядитель считает его вышедшим из строя
- В этом случае все задачи *Map*, уже выполненные и еще выполнявшиеся этим исполнителем,
 - переводятся в свое исходное состояние, и
 - можно заново планировать их выполнение другими исполнителями
- Аналогично распорядитель поступает со всеми задачами *Reduce*, выполнявшимися отказавшим исполнителем к моменту отказа

Корпоративные
базы

апреля

данных

2010 г.

2010

35

MapReduce: модель и реализации (13)

Реализация в распределенной среде (8)

Отказоустойчивость (2) Отказ исполнителя (2)

- Завершившиеся задачи *Map* выполняются повторно по той причине, что
 - их результирующие пары сохранялись на локальном диске отказавшего исполнителя
 - и поэтому недоступны в других узлах
- Завершившиеся задачи *Reduce* повторно выполнять не требуется, поскольку
 - их результирующие пары сохраняются в глобальной распределенной файловой системе

23
апреля

Корпоративные
базы
данных

2010 г.

2010

36

MapReduce: модель и реализации (14)

Реализация в распределенной среде (9)

Отказоустойчивость (3) Отказ исполнителя (3)

- Если некоторая задача *Map* выполнялась исполнителем *A*, а потом выполняется исполнителем *B*, то
 - об этом факте оповещаются все исполнители, выполняющие задачи *Reduce*
- Любая задача *Reduce*, которая не успела прочитать данные, произведенные исполнителем *A*,
 - после этого будет читать данные от исполнителя *B*

23
апреля
2010 г.

Корпоративные
ИВНые
базы
данных

2010

37

MapReduce: модель и реализации (15)

Реализация в распределенной среде (10)

Отказоустойчивость (4) Отказ распорядителя (1)

- В реализациях MapReduce от Google и Hadoop какая-либо репликация распорядителя не производится
- Поскольку распорядитель выполняется только в одном узле кластера, его отказ маловероятен, и если он случается, то
 - аварийно завершается все выполнение MapReduce
- Однако отмечается, что несложно организовать периодический сброс в распределенную файловую систему всего состояния распорядителя, чтобы в случае отказа можно было
 - запустить его новый экземпляр в другом узле с данной контрольной точки

23
апреля
2010 г.

Корпоративные
базы
данных
2010

38

MapReduce: модель и реализации (16)

Реализация в распределенной среде (11)

Отказоустойчивость (5) Семантика при наличии отказов (1)

- Если обеспечиваемые пользователями функции *Map* и *Reduce* являются детерминированными
 - т.е. всегда выдают одни и те же результаты при одинаковых входных данных,
- то при их выполнении в среде распределенной реализации MapReduce при любых условиях обеспечивает тот же результат, как
 - при последовательном выполнении всей программы при отсутствии каких-либо сбоев
- Это свойство обеспечивается за счет атомарности фиксации результатов задач *Map* и *Reduce*

23
апреля

2010 г.

корпоративные
базы
данных

2010

39

MapReduce: модель и реализации (17)

Реализация в распределенной среде (12)

Отказоустойчивость (6) Семантика при наличии отказов (2)

- Каждая выполняемая задача записывает свои результаты в частные временные файлы
- Задача *Reduce* производит один такой файл, а задача *Map* – R файлов, по одной на каждую задачу *Reduce*
- По завершении задачи *Map* исполнитель посылает распорядителю сообщение, в котором указываются имена R временных файлов
- При получении такого сообщения распорядитель запоминает эти имена файлов в своих структурах данных
- Повторные сообщения о завершении одной и той же задачи *Map* игнорируются

23
апреля

Корпоративные
базы
данных

2010 г.

2010

40

MapReduce: модель и реализации (18)

Реализация в распределенной среде (13)

Отказоустойчивость (7) Семантика при наличии отказов (3)

- При завершении задачи *Reduce* ее исполнитель атомарным образом переименовывает временный файл результатов в окончательный файл
- Если одна и та же задача *Reduce* выполняется несколькими исполнителями, то
 - для одного и того же окончательного файла будет выполнено несколько операций переименования
- Если в используемой распределенной файловой системе операция переименования является атомарной, то
 - в результате в файловой системе сохранятся результаты только какого-либо одного выполнения задачи *Reduce*

23
апреля

2010 г.

Корпоративные
базы
данных

2010

41

MapReduce: модель и реализации (19)

Реализация в распределенной среде (14)

Резервные задачи (1)

- Чаще всего к увеличению общего времени выполнения задания MapReduce приводит наличие "отстающих" ("straggler")
 - узлов кластера, в которых выполнение одной из последних задач *Map* или *Reduce* занимает необычно долгое время
 - например, из-за некритичной неисправности дискового устройства
- Для смягчения проблемы "остающих" в MapReduce применяется следующий общий механизм

23
апреля
2010 г.

корпоративные
базы
данных

2010 г.

2010

42

MapReduce: модель и реализации (20)

Реализация в распределенной среде (15)

Резервные задачи (2)

- Когда задание близится к завершению, для всех еще не завершившихся задач назначаются
 - дополнительные, резервные исполнители
- Задача считается выполненной,
 - когда завершается ее первичное или резервное выполнение
- Этот механизм настраивается таким образом, чтобы потребление вычислительных ресурсов возрастало не более чем на несколько процентов
- В результате удается существенно сократить время выполнения крупных MR-заданий

апреля

данных

2010 г.

2010

43

MapReduce: модель и реализации (21)

Расширенные возможности (1)

Функция-комбинатор (1)

- В некоторых случаях в результатах задачи *Map* содержится
 - значительное число повторяющихся значений промежуточного ключа,
- а определенная пользователем задача *Reduce*
 - является коммутативной и ассоциативной
- В таких случаях пользователь может определить дополнительную функцию-комбинатор (*Combiner*),
 - выполняющую частичную агрегацию таких данных до их передачи по сети
- Функция *Combiner* выполняется на той же машине, что и задача *Map*

23

Корпоративные
базы
данных

апреля

2010 г.

2010

44

MapReduce: модель и реализации (22)

Расширенные возможности (2)

Функция-комбинатор (2)

- Обычно для реализации функции *Combiner* используется тот же самый код, что и для реализации функции *Reduce*
- Единственное различие между функциями *Combiner* и *Reduce* состоит в способе работы с их результирующими данными
- Результаты функции *Reduce* записываются в окончательный файл результатов
- Результаты же функции *Combiner* помещаются в промежуточные файлы, которые впоследствии пересылаются в задачи *Reduce*

23
апреля

2010 г.

Корпоративные
базы
данных

2010

45

MapReduce: модель и реализации (23)

Расширенные возможности (3)

Форматы входных и результирующих данных (1)

- В библиотеке MapReduce поддерживается возможность чтения входных данных в нескольких разных форматах
- Например, в режиме "text" каждая строка трактуется как пара "ключ-значение", где ключ – это смещение до данной строки от начала файла, а значение – содержимое строки
- В другом распространенном формате входные данные представляются в виде пар "ключ-значение", отсортированных по значениям ключа
- В каждой реализации формата входных данных известно, каким образом следует расщеплять данные на осмысленные части, которые обрабатываются отдельными задачами Map
 - например, данные формата "text" расщепляются только по границами строк

23
апреля

2010 г.

Корпоративные
ИВНые
базы
данных

2010

46

MapReduce: модель и реализации (24)

Расширенные возможности (4)

Форматы входных и результирующих данных (2)

- Пользователи могут добавить к реализации собственные форматы входных данных, обеспечив новую реализацию интерфейса *reader*
 - в реализации Hadoop – *RecordReader*
- *Reader* не обязательно должен читать данные из файла,
 - можно легко определить *reader*, читающий данные из базы данных или из некоторой структуры в виртуальной памяти.
- Аналогичным образом, поддерживаются возможности генерации данных в разных форматах, и имеется простая возможность определения новых форматов результирующих данных

23
апреля

2010 г.

Корпоративные
базы
данных

2010

47

MapReduce внутри параллельной СУБД (1)

- Очевидны преимущества клиент-серверных организаций СУБД: в такой архитектуре сервер баз данных поддерживает крупную базу данных, которая
 - сохраняется в одном экземпляре и
 - доступна большому числу приложений, выполняемых прямо на стороне клиентов или в промежуточных серверах приложений
- Однако даже при использовании реляционной
 - или, правильнее, SQL-ориентированной
- организации баз данных, когда от клиентов на сервер баз данных отправляются высокоуровневые декларативные запросы, в обратную сторону, от сервера к клиенту,
 - пересылаются результирующие данные, вообще говоря, произвольно большого объема

апреля

2010 г.

данных

2010

48

MapReduce внутри параллельной СУБД (2)

- Естественно, возникает вопрос: не окажется ли дешевле,
 - чем пересылать данные с сервера на клиент для их дальнейшей обработки,
 - переместить требуемую обработку данных на сервер, ближе к самим данным
- В явном виде идея перемещения вычислений на сторону сервера была высказана в статье Лоуренса Роуа и Майкла Стоунбрейкера
 - Lawrence A. Rowe, Michael R. Stonebraker. *The POSTGRES Data Model*
 - <http://www.vldb.org/conf/1987/PO83.PDF>
- Намеки на эту идею можно найти и в более ранних статьях М. Стоунбрейкера и др.
 - еще не имевших непосредственного отношения к СУБД Postgres

23
апреля

2010 г.

Корпоративные
базы
данных

2010

49

MapReduce внутри параллельной СУБД (3)

- Поддержка определяемых пользователями хранимых процедур, функций и методов, типов данных и триггеров появилась во всех развитых SQL-ориентированных СУБД
- Соответствующие языковые средства специфицированы в стандарте языка SQL
- Более того, возникла новая проблема выбора – одну и ту же функциональность приложения можно реализовать на стороне сервера, на сервере приложений и на клиенте
- Однозначных методологий и рекомендаций, способствующих простому выбору, не существует
- Например, очевидно, что если услугами одного сервера пользуется несколько приложений, то
 - перегрузка сервера хранимыми процедурами и функциями, реализующими функциональность одного приложения, может нанести ущерб эффективности других приложений

23
апреля

2010 г.

Корпоративные
базы
данных

2010

50

MapReduce внутри параллельной СУБД (4)

- Тем не менее, во всех традиционных серверных организациях СУБД возможность переноса вычислений на сторону сервера существует и не очень сложно реализуется
- Однако в параллельных СУБД дела обстоят гораздо хуже
 - в особенности, в СУБД категории sharing-nothing
- Выполнение SQL-запросов распараллеливается автоматически оптимизатором запросов
- Но оптимизатор запросов не может распараллелить определенную пользователем процедуру или функцию,
 - если она написана не на SQL,
 - а на одном из традиционных языков программирования

Корпоративные
базы

23

апреля

данных

2010 г.

2010

51

MapReduce внутри параллельной СУБД (5)

- Технически можно было бы такие процедуры и функции вообще не распараллеливать, а выполнять в каком-либо одном узле кластера
- Но тогда
 - в этом узле пришлось бы собрать все данные, требуемые для выполнения процедуры или функции, для чего потребовалась бы
 - массовая пересылка данных по сети, и
 - это свело бы на нет все преимущества параллельных СУБД,
 - производительность которых основывается именно на параллельном выполнении

апреля

2010 г.

данных

2010

52

MapReduce внутри параллельной СУБД (6)

- С другой стороны, невозможно обязать распараллеливать свои программы самих пользователей, определяющих хранимые процедуры или функции
 - например, на основе библиотеки MPI
- Во-первых, это слишком сложное занятие для разработчиков приложений баз данных, которые часто вообще не являются профессиональными программистами
- Во-вторых, при таком явном параллельном программировании требовалось бы явным же образом управлять распределением данных по узлам кластера

23
апреля

Корпоративные
базы
данных

2010 г.

2010

53

MapReduce внутри параллельной СУБД (7)

- Несмотря на эти трудности, какая-то поддержка механизма распараллеливаемых определяемых пользователями процедур и функций в параллельных аналитических СУБД все-таки требуется
 - без этого аналитики вынуждены выполнять анализ данных на клиентских рабочих станциях,
 - постоянно пересылая на них из центрального хранилища данных данные весьма большого объема
 - другого способа работы у них просто нет
- Как показывает опыт двух производственных разработок, для обеспечения возможностей серверного программирования в массивно-параллельной среде систем баз данных
 - с пользой может быть применена модель MapReduce

23
апреля
2010 г.

Корпоративные
базы
данных
2010

54

MapReduce внутри параллельной СУБД (8)

- Речь идет о параллельных аналитических СУБД
 - Greenplum Database компании Greenplum и
 - *n*Cluster компании Aster Data Systems
- Общим в подходах обеих компаний является то, что
 - модель MapReduce реализуется внутри СУБД, и
 - возможностями этих реализаций могут пользоваться разработчики аналитических приложений
- Различие состоит в том, как можно пользоваться возможностями MapReduce:
 - В Greenplum Database – наряду с SQL,
 - а в *n*Cluster – из SQL

23
апреля

Корпоративные
базы
данных

2010 г.

2010

55

MapReduce внутри параллельной СУБД (9)

Greemplum – MapReduce наравне с SQL (1)

- Сначала немного поговорим об общей философии компании Greemplum, приведшей ее, в частности, к идее поддержки технологии MapReduce наряду с технологией SQL
- По мнению идеологов Greemplum и основных архитекторов Greenplum Database
 - Джозеф Хеллерстейн и др. *МОГучие способности: новые приемы анализа больших данных*
 - http://citforum.ru/database/articles/mad_skills/
- возрастающий уровень востребованности хранилищ данных и оперативного анализа данных, возможность и целесообразность использования требуемых аппаратных средств в масштабах отдельных подразделений компаний приводят к потребности пересмотра "ортодоксального" подхода к организации хранилищ данных

23
апреля

2010 г.

Корпоративные
базы
данных

2010

56

MapReduce внутри параллельной СУБД (10)

Greenplum – MapReduce наравне с SQL (2)

MAD Skills (1)

- Предлагается и реализуется новый подход к анализу данных, который идеологи (и маркетологи!) компании связывают с аббревиатурой *MAD*
- Интересная игра слов, которую трудно выразить на русском языке
- *Mad* применительно к технологии означает, что эта технология слегка безумна и уж во всяком случае не ортодоксальна
- С другой стороны, *mad skills* означает блестящие способности, а значит, предлагаемая технология, по мнению ее творцов, обладает новыми полезнейшими качествами
- Но в Greenplum *MAD* – это еще и аббревиатура от
 - *magnetic*,
 - *agile* и
 - *deep*

Корпорат
ивные
базы

данных

23
апреля

2010 г.

2010

57

MapReduce внутри параллельной СУБД (11)

Greentplum – MapReduce наравне с SQL (3)

MAD Skills (2)

- *Magnetic (магнетичность)* применительно к хранилищу данных означает, что оно должно быть "притягательным" по отношению к новым источникам данных, появляющимся в организации
- Данные из новых источников должны легко и просто включаться в хранилище данных с пользой для аналитиков
- При использовании традиционного ("ортодоксального") подхода к организации хранилища данных, для подключения нового источника данных требуется разработка и применение соответствующей процедуры ETL, а возможно, и изменение схемы хранилища данных, — в результате чего подключение нового источника данных часто затягивается на месяцы, а иногда и вовсе кончается ничем

апреля

2010 г.

данных

2010

58

MapReduce внутри параллельной СУБД (12)

Greentplum – MapReduce наравне с SQL (4)

MAD Skills (3)

- *Agile (гибкость)* – это предоставляемая аналитикам возможность простым образом и в быстром темпе
 - воспринимать,
 - классифицировать,
 - производить и
 - перерабатывать данные
- Для этого требуется база данных, логическая и физическая структура и содержание которой могут постоянно и быстро изменяться
- В отличие от этого, традиционным хранилищам данных свойственна жесткость, связанная с потребностью долгосрочного тщательного проектирования и планирования

Корпорат

ивные

23

базы

апреля

данных

2010 г.

2010

59

MapReduce внутри параллельной СУБД (13)

Greentplum – MapReduce наравне с SQL (5)

MAD Skills (4)

- *Deep (основательность)* означает, что аналитикам должны предоставляться средства выполнения произвольно сложных статистических алгоритмов над всеми данными, находящимися в хранилище данных
 - без потребности во взятии образцов или выборок
- Хранилище данных должно служить
 - как основательным репозиторием данных,
 - так и средой, поддерживающей выполнение сложных алгоритмов

23
апреля

Корпоративные
базы
данных

2010 г.

2010

60

MapReduce внутри параллельной СУБД (14)

Greenplum – MapReduce наравне с SQL (6)

MAD Skills (5)

- Более подробно рассмотрим один аспект MAD-аналитики, который привел к реализации системы с поддержкой интерфейсов и SQL, и MapReduce
- Как считают разработчики Greenplum Database хозяевами будущего мира анализа данных должны стать аналитики
- Фактически, на это направлены все аспекты MAD-аналитики
- Но, в частности, это означает всяческую поддержку написания и использования в среде хранилища данных разнообразных аналитических алгоритмов

Корпорат

ивные

базы

данных

23
апреля

2010 г.

2010

61

MapReduce внутри параллельной СУБД (15)

Greenplum – MapReduce наравне с SQL (7)

MAD Skills (6)

- Параллельная СУБД Greenplum Database делалась на основе СУБД PostgreSQL, являющейся законной наследницей Postgres
- Помимо своих прочих достоинств, Postgres была первой *расширяемой* СУБД
- Пользователи Postgres могли определять собственные
 - процедуры и функции,
 - типы данных и даже
 - методы доступа к структурам внешней памяти
- Эти возможности расширений системы были переняты и развиты в PostgreSQL
- Наряду с традиционным в Postgres языком C, для программирования серверных расширений в PostgreSQL можно использовать, в частности, популярные скриптовые языки
 - Perl и
 - Python

Корпоративные
базы

апреля

данных

2010 г.

2010

62

MapReduce внутри параллельной СУБД (16)

Greenplum – MapReduce наравне с SQL (8)

MAD Skills (7)

- В Greenplum Database на основе этих возможностей расширений системы обеспечена расширенная среда, позволяющая на уровне языка SQL оперировать такими математическими объектами, как векторы, функции и функционалы
- Пользователи могут определять собственные статистические алгоритмы и в полуавтоматическом режиме распараллеливать их выполнение по данным в массивно-параллельной среде
 - что часто является очень нетривиальной задачей
- Однако в любом случае при использовании такого подхода к анализу данных пользователям-аналитикам приходится иметь дело с декларативным языком SQL, а как считают идеологи Greenplum,
 - для многих аналитиков и статистиков SQL-программирование является обременительным и неудобным

23
апреля

2010 г.

Корпоративные
базы
данных

2010

63

MapReduce внутри параллельной СУБД (17)

Greenplum – MapReduce наравне с SQL (9)

MAD Skills (8)

- В качестве альтернативы аналитическому SQL-программированию в Greenplum Database обеспечивается полноправная реализация MapReduce, в которой
 - предоставляется доступ ко всем данным, сохраняемым в хранилище данных
- При использовании MapReduce аналитики пишут собственный понятный для них процедурный код
 - можно использовать те же Perl и Python
- и понимают, как будет выполняться их алгоритм в массивно-параллельной среде,
 - поскольку это выполнение опирается на простую модель MapReduce

23
апреля

2010 г.

Корпоративные
базы
данных

2010

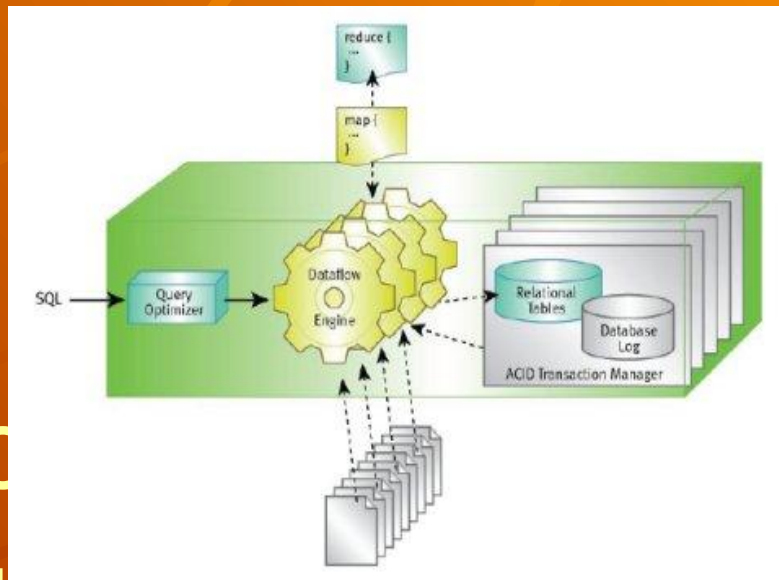
64

MapReduce внутри параллельной СУБД (18)

Greenplum – MapReduce наравне с SQL (10)

Реализация MapReduce в Greenplum Database (9)

- Ядром системы является процессор потоков данных (Dataflow Engine)
- Замена соответствующего компонента ядра PostgreSQL для обеспечения
 - массивно-параллельного выполнения запросов и
 - базовых функциональных возможностей, требуемых для поддержки модели MapReduce
- В результате SQL-ориентированная СУБД и MapReduce работают с общим ядром, поддерживающим массивно-параллельную обработку данных,
 - и механизмы SQL и MapReduce обладают интероперабельностью



23 апреля 2010 г. Корпоративные базы данных

2010 г.

2010

65

MapReduce внутри параллельной СУБД (18)

Greenplum – MapReduce наравне с SQL (10)

Реализация MapReduce в Greenplum Database (9)

- Функции *Map* и *Reduce* в среде Greenplum Database можно программировать на популярных скриптовых языках Python и Perl
- Можно использовать развитые программные средства с открытыми кодами, содержащиеся в репозиториях
 - Python Package Index (PyPI) и
 - Comprehensive Perl Archive Network (CPAN)
- В составе этих репозиториях находятся
 - средства анализа неструктурированного текста,
 - статистические инструментальные средства,
 - анализаторы форматов HTML и XML и
 - многие другие программные средства, потенциально полезные аналитикам

23
апреля

Корпоративные
базы
данных

2010 г.

2010

66

MapReduce внутри параллельной СУБД (19)

Greenplum – MapReduce наравне с SQL (11)

Реализация MapReduce в Greenplum Database (10)

- В среде Greenplum Database приложениям MapReduce обеспечивается доступ к данным,
 - хранящимся в файлах,
 - предоставляемым Web-сайтами и
 - даже генерируемым командами операционной системы
- Доступ к таким данным не влечет накладных расходов, ассоциируемых с использованием СУБД:
 - блокировок,
 - журнализации,
 - фиксации транзакций и т.д.
- С другой стороны, эффективный доступ к данным, хранимым в базе данных, поддерживается за счет выполнения MR-программ в ядре Greenplum Database
- Это позволяет избежать расходов на пересылку данных

Корпорат

ивные

базы

апреля

данных

2010 г.

2010

67

MapReduce внутри параллельной СУБД (20)

Greenplum – MapReduce наравне с SQL (12)

Реализация MapReduce в Greenplum Database (11)

- Архитектура Greenplum Database с равноправной поддержкой SQL и MapReduce позволяет смешивать стили программирования,
 - делать MR-программы видимыми для SQL-запросов и наоборот
- Например, можно выполнять MR-программы над таблицами базы данных
 - Для этого всего лишь требуется указать MapReduce, что входные данные программы должны браться из таблицы
- Поскольку таблицы баз данных Greenplum Database хранятся разделенными между несколькими узлами кластера, первая фаза MAP выполняется внутри ядра СУБД прямо над этими разделами

23
апреля
2010 г.

Корпоративные
базы данных
2010

68

MapReduce внутри параллельной СУБД (21)

Greenplum – MapReduce наравне с SQL (13)

Реализация MapReduce в Greenplum Database (12)

- Как и в автономных реализациях MapReduce, результаты выполнения MR-программ могут сохраняться в файловой системе
- Но настолько же просто сохранить результирующие данные в базе данных с обеспечением транзакционной долговечности хранения этих данных
- В дальнейшем эти данные могут анализироваться, например, с применением SQL-запросов
- Запись результирующих данных в таблицы происходит параллельным образом и не вызывает лишних накладных расходов

Корпоративные

базы

23

апреля

2010 г.

2010

69

MapReduce внутри параллельной СУБД (22)

Aster Data – MapReduce как основа UDF (1)

- У компании Aster Data имеется свой слоган *Big Data, Fast Insight*, который, по сути, означает то же самое
 - превращение массивно-параллельного хранилища данных в аналитическую платформу
- И для этого тоже используется технология MapReduce, встроенная в СУБД
 - Эрик Фридман и др. *SQL/MapReduce: практический подход к поддержке самоописываемых, полиморфных и параллелизуемых функций, определяемых пользователями*
 - http://citforum.ru/database/articles/asterdata_sql_mr/
- Однако, в отличие от Greenplum, эта технология применяется не для обеспечения альтернативного внешнего способа обработки данных, а
 - для реализации нового механизма хорошо распараллеливаемых (по модели MapReduce), самоописываемых и
 - полиморфных табличных функций, определяемых пользователями и вызываемых из операторов выборки SQL

23

апреля

2010 г.

Корпорат

ивные

базы

данных

2010

70

MapReduce внутри параллельной СУБД (23)

Aster Data – MapReduce как основа UDF (2)

Предпосылки и преимущества SQL/MapReduce (1)

- По мнению основных разработчиков СУБД nCluster, декларативный язык SQL во многом ограничивает использование аналитических СУБД
- С одной стороны, несмотря на постоянное наращивание аналитических возможностей этого языка, для многих аналитиков их оказывается недостаточно
- С другой стороны, эти возможности постепенно становятся такими сложными и непонятными, что
 - зачастую становится проще написать процедурный код, решающий частную аналитическую задачу
- Наконец, оптимизаторы запросов SQL-ориентированных СУБД постоянно отстают от развития языка, и планы сложных аналитических запросов могут быть весьма далеки от оптимальных,
 - что приводит к их недопустимо долгому выполнению, а иногда и аварийному завершению

23
апреля

Корпоративные
базы
данных

2010 г.

2010

71

MapReduce внутри параллельной СУБД (24)

Aster Data – MapReduce как основа UDF (3)

Предпосылки и преимущества SQL/MapReduce (2)

- Эти проблемы частично решаются за счет поддержки в SQL-ориентированных СУБД механизма UDF
- Такие функции позволяют пользователям решать внутри сервера баз данных свои прикладные задачи
 - путем написания соответствующего процедурного кода
- Однако традиционные механизмы UDF разрабатывались в расчете на "одноузловые" СУБД, и
 - по умолчанию предполагается чисто последовательное выполнение UDF
- Автоматическое распараллеливание последовательного кода в массивно-параллельной среде с разделением данных является
 - сложной нерешенной проблемой

апреля

2010 г.

данных

2010

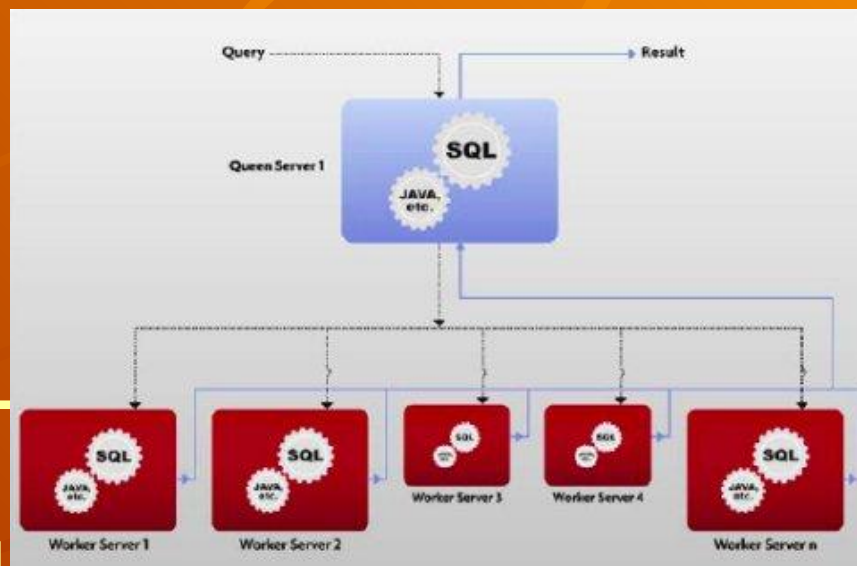
72

MapReduce внутри параллельной СУБД (25)

Aster Data – MapReduce как основа UDF (4)

Предпосылки и преимущества SQL/MapReduce (3)

- В Aster Data для обеспечения механизма естественно распараллеливаемых UDF разработана инфраструктура SQL/MapReduce,



– поддерживаемая внутри SQL-ориентированной массивно-параллельной СУБД nCluster

апреля

данных

2010 г.

2010

73

MapReduce внутри параллельной СУБД (26)

Aster Data – MapReduce как основа UDF (5)

Предпосылки и преимущества SQL/MapReduce (4)

- Организация среды SQL/MapReduce обеспечивает следующие возможности:
 - можно эффективно выполнять в "реляционном" стиле операции над таблицами с использованием SQL, а "нереляционные" задачи и оптимизации – возлагать на явно программируемые процедурные функции;
 - поскольку функции выполняются над согласованными данными из таблиц базы данных, обеспечивается согласованность вычислений;
 - оценочный (cost-based) оптимизатор может принимать решения о способе выполнения SQL-запросов, содержащих вызовы SQL/MapReduce-функций, на основе достоверной статистики данных;
 - пользователи nCluster могут формулировать SQL-запросы с использованием высокоуровневых средств анализа данных, воплощенных в SQL/MapReduce-функциях

23
апреля

2010 г.

Корпоративные
базы
данных

2010

74

MapReduce внутри параллельной СУБД (27)

Aster Data – MapReduce как основа UDF (6)

Предпосылки и преимущества SQL/MapReduce (5)

- SQL/MapReduce-функции можно программировать как на традиционных языках программирования (Java, C#, C++), так и скриптовых языках (Python, Ruby)
- Независимо от используемого языка программирования, эти табличные функции являются *самоописываемыми* и *полиморфными*
- Одна и та же функция может принимать на вход таблицы с разными схемами
 - функция настраивается на конкретную схему входной таблицы на этапе формирования плана запроса, содержащего ее вызов
- и выдавать таблицы также с разными схемами
 - функция сама сообщает планировщику запроса схему своего результата на этапе формирования плана запроса
- Это свойство SQL/MapReduce-функций упрощает процедуру их регистрации в системе и способствует повторному использованию кода

2010 г.

2010

75

MapReduce внутри параллельной СУБД (28)

Aster Data – MapReduce как основа UDF (7)

Предпосылки и преимущества SQL/MapReduce (6)

- Синтаксические особенности определения SQL/MapReduce-функций и их семантика делают эти программные объекты естественным образом параллелизуемыми по данным:
 - во время выполнения для каждой функции образуются ее экземпляры, параллельно выполняемые в узлах, которые содержат требуемые данные
- Вызовы функций подобны подзапросам SQL, что обеспечивает возможность композиции функций, при которой при вызове функции вместо спецификации входной таблицы можно задавать вызов другой функции
- Наконец, внешняя эквивалентность вызова SQL/MapReduce-функции подзапросу позволяет
 - применять при формировании плана SQL-запроса с вызовами таких функций обычную оценочную оптимизацию на основе статистики, а также
 - динамически изменять порядок выполнения функций и вычисления настоящих SQL-подзапросов

23
апреля

2010 г.

Корпоративные
базы
данных

2010

76

MapReduce внутри параллельной СУБД (29)

Aster Data – MapReduce как основа UDF (8)

Реализация SQL/MapReduce (1) Синтаксис вызова (1)

- Вызов SQL/MapReduce-функции может присутствовать только в качестве элемента списка ссылок на таблицы раздела FROM SQL-запроса
- В разделе ON, который является единственным обязательным разделом вызова, указывается любой допустимый SQL/MapReduce-запрос
 - SQL-запрос, вызов SQL/MapReduce-функции или просто имя таблицы
- Во время формирования плана запроса, содержащего вызов SQL/MapReduce-функции, схемой входной таблицы этого вызова считается схема результата запроса, указанного в разделе ON

```
SELECT ...  
FROM functionname(  
  ON table-or-query  
  [PARTITION BY expr, ...]  
  [ORDER BY expr, ...]  
  [clausename(arg, ...) ...]  
)
```

23
апреля

2010 г.

Корпоративные
базы
данных

2010

77

MapReduce внутри параллельной СУБД (30)

Aster Data – MapReduce как основа UDF (9)

Реализация SQL/MapReduce (2) Синтаксис вызова (2)

- Раздел PARTITION BY указывается только в вызовах SQL/MapReduce-функций над разделами

```
SELECT ...  
FROM functionname(  
  ON table-or-query  
  [PARTITION BY expr, ...]  
  [ORDER BY expr, ...]  
  [clauseName(arg, ...) ...]  
)
```

– аналоге функции *Reduce* исходной модели MapReduce

- В этом случае в разделе PARTITION BY указывается список выражений, на основе значений которых производится разделение таблицы, специфицированной в разделе ON

23
апреля

Корпоративные
базы
данных

2010 г.

2010

78

MapReduce внутри параллельной СУБД (31)

Aster Data – MapReduce как основа UDF (10)

Реализация SQL/MapReduce (3) Синтаксис вызова (3)

- При наличии раздела PARTITION BY в вызове может содержаться и раздел ORDER BY, указывающий на потребность в сортировке входных данных до реального вызова функции

```
SELECT ...  
FROM functionname(  
  ON table-or-query  
  [PARTITION BY expr, ...]  
  [ORDER BY expr, ...]  
  [clauseName(arg, ...) ...]  
)
```

- Наконец, вслед за разделом ORDER BY можно указать произвольное число дополнительных разделов со специальными аргументами
 - имена этих разделов и значения аргументов передаются в SQL/MapReduce-функцию при ее инициализации

MapReduce внутри параллельной СУБД (32)

Aster Data – MapReduce как основа UDF (11)

Реализация SQL/MapReduce (4) Модель выполнения (1)

- В среде SQL/MapReduce используется модель выполнения функций, являющаяся обобщением модели MapReduce
- Функция SQL/MapReduce может быть
 - либо функцией над строками (*row function*),
 - либо функцией над разделами (*partition function*)
- Функции первого типа являются аналогами функций *Map* классической модели MapReduce, а функции второго типа – аналогами функций *Reduce*
- Поскольку, как отмечалось ранее, в разделе ON вызова SQL/MapReduce-функции может содержаться вызов другой SQL/MapReduce-функции,
 - в среде SQL/MapReduce допускается любое число и любой порядок вызовов функций *Map* и *Reduce*,
 - а не только жесткая последовательность *Map-Reduce*, допускаемая классической моделью

23
апреля

Корпоративные
базы
данных

2010 г.

2010

80

MapReduce внутри параллельной СУБД (33)

Aster Data – MapReduce как основа UDF (12)

Реализация SQL/MapReduce (5) Модель выполнения (2)

- При выполнении функции над строками каждая строка входной таблицы обрабатывается ровно одним экземпляром этой функции
- С точки зрения семантики каждая строка обрабатывается независимо, поэтому входная таблица может разделяться по экземплярам функции произвольным образом,
 - что обеспечивает возможности параллелизма и масштабирования
- Для каждой строки входной таблицы функция над строками может не производить ни одной строки, а может произвести несколько строк

23
апреля

2010 г.

Корпоративные
базы
данных

2010

81

MapReduce внутри параллельной СУБД (34)

Aster Data – MapReduce как основа UDF (13)

Реализация SQL/MapReduce (6) Модель выполнения (3)

- При выполнении функции над разделами каждая группа строк, образованная на основе спецификации раздела PARTITION BY вызова функции, обрабатывается ровно одним экземпляром этой функции, и этот экземпляр получает все группу целиком
- Если в вызове функции содержится раздел ORDER BY, то экземпляры функции получают разделы в уже упорядоченном виде
- С точки зрения семантики каждый раздел обрабатывается независимо,
 - что обеспечивает возможности параллелизма на уровне разделов
- Для каждого входного раздела функция над строками может не производить ни одной строки, а может произвести несколько строк

23
апреля

2010 г.

Корпоративные
базы
данных

2010

82

MapReduce внутри параллельной СУБД (35)

Aster Data – MapReduce как основа UDF (14)

Реализация SQL/MapReduce (7) Особенности реализации (1)

- Для реализации SQL/MapReduce-функций можно использовать разные языки, но все они являются объектно-ориентированными
- Каждая SQL/MapReduce-функция реализуется в виде отдельного класса, и при выработке плана выполнения SQL-запроса, содержащего вызовы таких функций, для каждого вызова образуется объект соответствующего класса с обращением к его методу-конструктору – инициализатору функции
- Это обеспечивает настройку функции и получение требуемого описания ее результирующей таблицы

23
апреля
2010 г.

корпоративные
базы
данных
2010

83

MapReduce внутри параллельной СУБД (36)

Aster Data – MapReduce как основа UDF (15)

Реализация SQL/MapReduce (8) Особенности реализации (2)

- Взаимодействие оптимизатора запросов с инициализатором функции производится через специальный объект, называемый *контрактом времени выполнения (Runtime Contract)*
- Анализируя вызов функции, оптимизатор выявляет
 - имена и типы данных столбцов входной таблицы, а также
 - имена и значения разделов дополнительных параметров
- и соответствующим образом заполняет некоторые поля объекта-контракта, который затем передается инициализатору функции

Корпорат

ивные

базы

данных

23
апреля

2010 г.

2010

84

MapReduce внутри параллельной СУБД (37)

Aster Data – MapReduce как основа UDF (16)

Реализация SQL/MapReduce (9) Особенности реализации (3)

- Инициализатор завершает подготовку контракта путем заполнения его дополнительных полей, содержащих, в частности,
 - информацию о схеме результирующей таблицы,
- и обращается к методу `complete` объекта-контракта
- На основе готового контракта продолжается выработка плана выполнения запроса, и
 - этот контракт соблюдается при последующем выполнении SQL/MapReduce-функции всеми ее экземплярами

Корпорат

ивные

базы

23
апреля

данных

2010 г.

2010

85

MapReduce внутри параллельной СУБД (38)

Aster Data – MapReduce как основа UDF (17)

Реализация SQL/MapReduce (10) Особенности реализации (4)

- Наиболее важными методами интерфейсов классов для функций над строками и разделами являются методы OperateOnSomeRows и OperateOnPartition
- При обращении к этим методам
 - реальном выполнении соответствующей функции
- в качестве аргументов передаются
 - итератор над строками, для обработки которых вызывается функция, и
 - объект *emitter*, с помощью вызовов которого возвращаются результирующие строки

Корпоративные
базы

апреля

данных

2010 г.

2010

86

MapReduce внутри параллельной СУБД (39)

Aster Data – MapReduce как основа UDF (18)

Реализация SQL/MapReduce (11) Особенности реализации (5)

- Чтобы можно было начать использовать некоторую SQL/MapReduce-функцию, ее нужно установить
- Для этого используется общий механизм установки файлов, реализованный в nCluster
- Этот механизм реплицирует файл во всех рабочих узлах системы
- Далее проверяется, что этот файл содержит SQL/MapReduce-функцию, а также выясняются ее статические свойства:
 - является ли она функцией на строках или же над разделами,
 - содержит ли она вызовы комбинатора и т.д.

23
апреля

2010 г.

Корпоративные
базы
данных

2010

87

MapReduce внутри параллельной СУБД (40)

- Таким образом в этих двух системах обеспечивается возможность развитого анализа данных поблизости от самих данных
- Разработчики серверных аналитических приложений несколько ограничиваются моделью MapReduce
 - в большей степени в Greenplum Database, в меньшей – в nCluster,
- но зато пользовательский процедурный код хорошо распараллеливается по данным в массивно-параллельной среде

23
апреля

2010 г.

Корпоративные
базы
данных

2010

88

Параллельная СУБД на основе MapReduce (1)

- В статье Стоунбрейкера и др., посвященной сравнению эффективности технологий MapReduce и массивно-параллельных СУБД при решении аналитических задач, утверждалось,
 - что развитость и зрелость технологии параллельных СУБД категории *sharing-nothing* позволяет им обходиться стоузловыми кластерами для поддержки самых крупных современных аналитических баз данных петабайтного масштаба
- Вместе с тем, особые качества масштабируемости и отказоустойчивости технологии MapReduce проявляются при использовании кластеров **тысячами узлов**
- Из этого делался вывод, что в обозримом будущем эти качества не являются **настолько** необходимыми для параллельных СУБД

Корпоративные
базы

апреля

данных

2010 г.

2010

89

Параллельная СУБД на основе MapReduce (2)

- Однако спустя всего несколько месяцев появилась статья, в которой звучат уже совсем другие мотивы
 - Ави Зильбершац и др. *HadoopDB: архитектурный гибрид технологий MapReduce и СУБД для аналитических рабочих нагрузок*
 - <http://citforum.ru/database/articles/hadoopdb/>
- В ней говорится, что в связи с ростом объема данных, которые требуется анализировать, возрастает и число приложений,
 - для поддержки которых нужны кластеры с числом узлов, больше ста
- В то же время, имеющиеся в настоящее время параллельные СУБД не масштабируются должным образом до сотен узлов

23

базы

апреля

данных

2010 г.

2010

90

Параллельная СУБД на основе MapReduce (3)

- Это объясняется следующими причинами:
 - При возрастании числа узлов кластера возрастает вероятность отказов отдельных узлов, а массивно-параллельные СУБД проектировались в расчете на редкие отказы
 - Современные параллельные СУБД рассчитаны на однородную аппаратную среду
 - все узлы кластера обладают одной и той же производительностью
 - а при значительном масштабировании полной однородности среды добиться почти невозможно.
 - До последнего времени имелось очень небольшое число систем аналитических баз данных, для достижения требуемой производительности которых требовались кластеры с более чем несколькими десятками узлов
 - Поэтому существующие параллельные СУБД просто не тестировались в более масштабной среде, и при их дальнейшем масштабировании могут возникнуть непредвиденные технические проблемы

23
апреля
2010 г.

2010

91

Параллельная СУБД на основе MapReduce (4)

- Требуемые характеристики масштабируемости и отказоустойчивости может обеспечить технология MapReduce, поскольку
 - она с самого начала разрабатывалась с расчетом на масштабирование до тысяч узлов,
 - и ее реализация от Google эффективно используется для поддержки внутренних операций этой компании
- Несмотря на то, что изначально технология MapReduce ориентировалась на обработку неструктурированных текстовых данных,
 - известны показательные примеры ее использования и для обработки огромных объемов структурированных данных

23
апреля

2010 г.

Корпоративные
базы
данных

2010

92

Параллельная СУБД на основе MapReduce (5)

- Однако объективно при обработке структурированных данных MapReduce не может конкурировать с параллельными СУБД по производительности,
 - что объясняется отсутствием схемы у обрабатываемых данных, индексов, оптимизации запросов и т.д.
- В результате при выполнении многих типичных аналитических запросов MapReduce показывает производительность,
 - более чем на порядок уступающую производительности параллельных СУБД

23
апреля

Корпоративные
базы
данных

2010 г.

2010

93

Параллельная СУБД на основе MapReduce (6)

- В проекте HadoopDB специалисты из университетов Yale и Brown предпринимают попытку создать гибридную систему управления данными,
 - сочетающую преимущества технологий и MapReduce, и параллельных СУБД
- MapReduce обеспечивает коммуникационную инфраструктуру, объединяющую произвольное число узлов,
 - в которых выполняются экземпляры традиционной СУБД
- Запросы формулируются на языке SQL, транслируются в среду MapReduce, и
 - значительная часть работы передается в экземпляры СУБД
- Наличие MapReduce обеспечивает
 - масштабируемость и отказоустойчивость,
- а использование в узлах кластера СУБД позволяет добиться
 - высокой производительности

23
апреля

2010 г.

Корпоративные
базы
данных

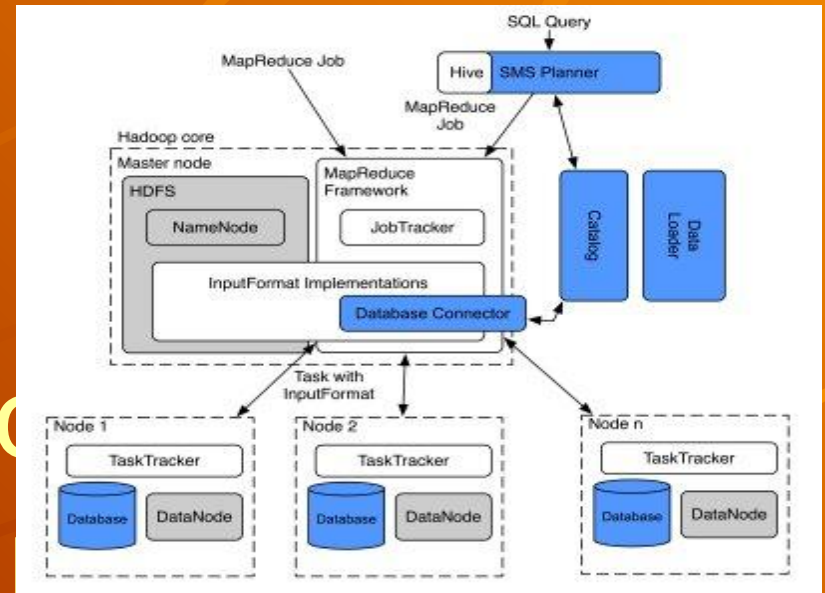
2010

94

Параллельная СУБД на основе MapReduce (7)

Общая организация HadoopDB (1)

- Основой системы является Hadoop MapReduce
- К ней добавлены компоненты компиляции поступающих в систему SQL-запросов, загрузки и каталогизирования данных, связи с СУБД и самих СУБД
- При реализации всех компонентов системы максимально использовались пригодные для этого программные средства с открытыми исходными текстами



апреля

данных

2010 г.

2010

95

Параллельная СУБД на основе MapReduce (8)

Общая организация HadoopDB (2)

Hadoop MapReduce (1)

- Hadoop MapReduce опирается на распределенную файловую систему HDFS (Hadoop Distributed File System)
- Файлы HDFS имеют блочную структуру, и блоки одного файла распределяются по узлам данных (DataNode)
- Файловая система работает под централизованным управлением выделенного узла имен (NameNode), в котором поддерживаются метаданные о файлах
 - в том числе, об их размерах, о размещении блоков и их реплик и т.д.

23
апреля

2010 г.

Корпоративные
базы
данных

2010

96

Параллельная СУБД на основе MapReduce (9)

Общая организация HadoopDB (3)

Hadoop MapReduce (2)

- В самой среде Hadoop MapReduce поддерживаются один узел-распорядитель
 - в Hadoop он называется JobTracker
- и много узлов-исполнителей
 - здесь TaskTracker
- В узле JobTracker планируется выполнение MR-заданий, а также отслеживаются данные о загрузке узлов TaskTracker и доступных ресурсах
- Каждое задание разбивается на задачи *Map* и *Reduce*, которые назначаются узлом JobTracker узлам TaskTracker

с учетом требований локальности данных и балансировки нагрузки

23
апреля
2010 г.

Корпоративные
базы
данных
2010

97

Параллельная СУБД на основе MapReduce (10)

Общая организация HadoopDB (4)

Hadoop MapReduce (3)

- Требование локальности удовлетворяется за счет того, что JobTracker пытается назначать каждую задачу *Map* тому узлу TaskTracker, – для которого данные, обрабатываемые этой задачей, являются локальными
- Балансировка нагрузки достигается путем назначения задач всем доступным узлам TaskTracker
- Узлы TaskTracker периодически посылают в узел JobTracker контрольные сообщения с информацией о своем состоянии

23
апреля

2010 г.

Корпоративные
базы
данных

2010

98

Параллельная СУБД на основе MapReduce (11)

Общая организация HadoopDB (5)

Собственные компоненты HadoopDB (1) Коннектор БД (1)

- Коннектор баз данных обеспечивает интерфейс между TaskTracker и независимыми СУБД, располагаемыми в узлах кластера
- Этот компонент расширяет класс InputFormat и является частью соответствующей библиотеки
- От каждого MR-задания в коннектор поступают SQL-запрос, а также параметры подключения к системе баз данных
 - указание драйвера JDBC,
 - размер структуры выборки данных и т.д.

23

Корпоративные
базы

апреля

данных

2010 г.

2010

99

Параллельная СУБД на основе MapReduce (12)

Общая организация HadoopDB (6)

Собственные компоненты HadoopDB (2) Коннектор БД (2)

- Теоретически коннектор обеспечивает подключение к любой JDBC-совместимой СУБД
- Однако в других компонентах HadoopDB приходится учитывать специфику конкретных СУБД, поскольку для них требуется по-разному оптимизировать запросы
- В исходных экспериментах использовалась реализация коннектора для PostgreSQL, а в позже уже упоминалась некоторая поколочная система
- В любом случае, для среды HadoopDB эта реализация обеспечивает естественное и прозрачное использование баз данных в качестве источника входных данных

23
апреля

Корпоративные
базы
данных

2010 г.

2010

100

Параллельная СУБД на основе MapReduce (13)

Общая организация HadoopDB (7)

Собственные компоненты HadoopDB (3) Каталог (1)

- В каталоге поддерживаются метаданные двух сортов:
 - параметры подключения к базе данных (ее месторасположение, класс JDBC-драйвера, учетные данные) и
 - описание наборов данных, содержащихся в кластере, расположение реплик и т.д.
- Каталог сохраняется в формате XML в HDFS
- К нему обращаются JobTracker и TaskTracker для выборки данных, требуемых для планирования задач и обработки данных

апреля

данных

2010 г.

2010

101

Параллельная СУБД на основе MapReduce (14)

Общая организация HadoopDB (8)

Собственные компоненты HadoopDB (4) Загрузчик данных (1)

- Обязанностями загрузчика данных являются:
 - глобальное разделение данных по заданному ключу при их загрузке из HDFS;
 - разбиение данных, хранимых в одном узле, на несколько более мелких разделов (*чанков*);
 - массовая загрузка данных в базу данных каждого узла с использованием чанков

апреля

2010 г.

данных

2010

102

Параллельная СУБД на основе MapReduce (15)

Общая организация HadoopDB (9)

Собственные компоненты HadoopDB (5) Загрузчик данных (2)

- Загрузчик данных состоит из компонентов GlobalHasher и LocalHasher
- GlobalHasher запускает в Hadoop MapReduce специальное задание, в котором
 - читаются файлы данных HDFS и
 - производится их разделение на столько частей, сколько имеется рабочих узлов в кластере
- Сортировка данных не производится
- Затем LocalHasher в каждом узле копирует соответствующий раздел из HDFS в свою файловую систему,
 - разделяя его на чанки в соответствии с установленным в системе максимальным размером чанка

23
апреля

Корпоративные
базы
данных

2010 г.

2010

103

Параллельная СУБД на основе MapReduce (16)

Общая организация HadoopDB (10)

Собственные компоненты HadoopDB (6) Загрузчик данных (3)

- В GlobalHasher и LocalHasher используются разные хэш-функции, – обеспечивающие примерно одинаковые размеры всех чанков
- Эти хэш-функции отличаются от хэш-функции, используемой в Hadoop MapReduce для разделения данных по умолчанию
- Это способствует улучшению балансировки нагрузки

23
апреля
2010 г.

Корпоративные
базы
данных
2010

104

Параллельная СУБД на основе MapReduce (17)

Общая организация HadoopDB (11)

Собственные компоненты HadoopDB (7) Планирование SQL (1)

- Внешний интерфейс HadoopDB позволяет выполнять SQL-запросы
- Компиляцию и подготовку планов выполнения SQL-запросов производит планировщик SMS (SMS Planner),
 - являющийся расширением планировщика Hive
- Планировщик Hive преобразует запросы, представленные на языке HiveQL (вариант SQL) в задания MapReduce,
 - которые выполняются над таблицами, хранимыми в виде файлов HDFS

Корпоративные
базы

апреля

данных

2010 г.

2010

105

Параллельная СУБД на основе MapReduce (18)

Общая организация HadoopDB (12)

Собственные компоненты HadoopDB (8) Планирование SQL (2)

- Эти задания представляются в виде ориентированных ациклических графов (DAG) реляционных операций
 - фильтрации (ограничения),
 - выборки (проекции),
 - соединения,
 - агрегации,
- каждая из которых выполняется в конвейере:
 - после обработки каждого очередного кортежа результат каждой операции направляется на вход следующей операции

Корпорат

ивные

базы

данных

апреля

2010 г.

2010

106

Параллельная СУБД на основе MapReduce (19)

Общая организация HadoopDB (13)

Собственные компоненты HadoopDB (9) Планирование SQL (3)

- Операции соединения, как правило, выполняются в задаче *Reduce* MR-задания, соответствующего SQL-запросу
- Это связано с тем, что каждая обрабатываемая таблица сохраняется в отдельном файле HDFS, и
 - невозможно предполагать совместного размещения соединяемых разделов таблиц в одном узле кластера
- Для HadoopDB это не всегда так, поскольку соединяемые таблицы
 - могут разделяться по атрибуту соединения, и тогда операцию соединения можно вытолкнуть на уровень СУБД

Корпоративные
базы
данных

23
апреля

2010 г.

2010

107

Параллельная СУБД на основе MapReduce (20)

Общая организация HadoopDB (14)

Собственные компоненты HadoopDB (10) Планирование SQL (4)

- Пример
- Пусть задан запрос:

```
SELECT YEAR(saleDate), SUM(revenue)
FROM sales
GROUP BY YEAR(saleDate);
```
- В Hive этот запрос обрабатывается следующим образом:
 - Производится синтаксический разбор запроса, и образуется его абстрактное синтаксическое дерево.
 - Далее работает семантический анализатор, который
 - выбирает из внутреннего каталога Hive *MetaStore* информацию о схеме таблицы *sales*, а также
 - инициализирует структуры данных, требуемые для сканирования этой таблицы и выборки нужных полей

23
апреля

2010 г.

Корпорат
ивные
базы
данных

2010

108

Параллельная СУБД на основе MapReduce (21)

Общая организация HadoopDB (15)

Собственные компоненты HadoopDB (11) Планирование SQL (5)

- Затем генератор логических планов запросов производит план запроса – DAG реляционных операций
- Вслед за этим оптимизатор перестраивает этот план запроса, проталкивая, например, операции фильтрации ближе к операциям сканирования таблиц
 - Основной функцией оптимизатора является разбиение плана на фазы *Map* и *Reduce*
 - В частности, перед операциями соединения и группировки добавляется операция перераспределения данных (Reduce Sink)
 - Эти операции отделяют фазу *Map* от фазы *Reduce*
 - Оценочная (cost-based) оптимизация не используется, и поэтому получаемые планы не всегда эффективны

Корпоративные
базы

23
апреля

данных

2010 г.

2010

109

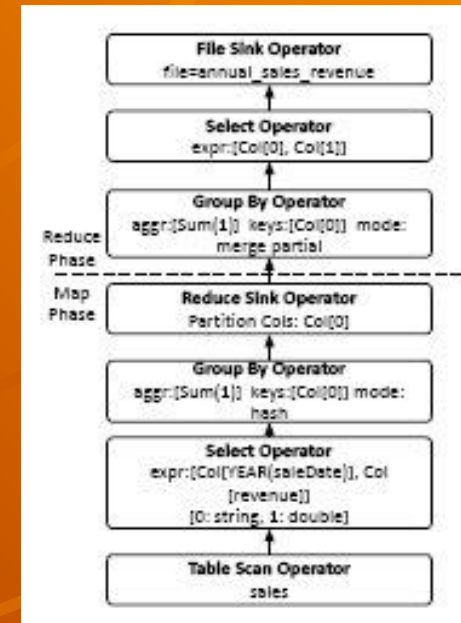
Параллельная СУБД на основе MapReduce (22)

Общая организация HadoopDB (16)

Собственные компоненты HadoopDB (12) Планирование SQL (6)

- Генератор физических планов выполнения запросов преобразует логический план в физический, допускающий выполнение в виде одного или нескольких MR-заданий

- Первая (и каждая аналогичная) операция Reduce Sink помечает переход от фазы Map к фазе Reduce некоторого задания MapReduce, а остальные операции Reduce Sink помечают начало следующего задания MapReduce



```
SELECT YEAR(saleDate),  
SUM(revenue)  
FROM sales  
GROUP BY YEAR(saleDate);
```

23
апреля

2010 г.

2010

110

Параллельная СУБД на основе MapReduce (23)

Общая организация HadoopDB (17)

Собственные компоненты HadoopDB (13) Планирование SQL (7)

– Полученный DAG сериализуется в формате XML

- Задания инициируются драйвером Hive, который

- руководствуется планом в формате SQL и

- создает все необходимые объекты, сканирующие данные в таблицах HDFS и покортежно обрабатывающие данные

Корпоративные
базы

23

апреля

данных

2010 г.

2010

111

Параллельная СУБД на основе MapReduce (24)

Общая организация HadoopDB (18)

Собственные компоненты HadoopDB (14) Планирование SQL (8)

- В планировщике SMS функциональность планировщика Hive расширяется следующим образом
- Во-первых, до обработки каждого запроса модифицируется MetaStore, куда помещается информация о таблицах базы данных
 - Для этого используется каталог HadoopDB
- Далее, после генерации физического плана запроса и до выполнения MR-заданий выполняются два прохода по физическому плану

2010 г.

2010

112

Параллельная СУБД на основе MapReduce (25)

Общая организация HadoopDB (19)

Собственные компоненты HadoopDB (15) Планирование SQL (9)

- На первом проходе
 - устанавливается, какие столбцы таблиц действительно обрабатываются запросом,
 - и определяются ключи разделения, используемые в операциях Reduce Sink
- На втором проходе DAG запроса обходится снизу-вверх от операций сканирования таблиц до формирования результата или первой операции Reduce Sink
- Все операции этой части DAG преобразуются в один или несколько SQL-запросов, которые проталкиваются на уровень СУБД
- Для повторного создания кода SQL используется специальный основанный на правилах генератор

23
апреля

Корпоративные
базы
данных

2010 г.

2010

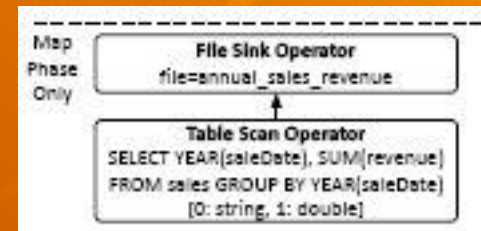
113

Параллельная СУБД на основе MapReduce (26)

Общая организация HadoopDB (20)

Собственные компоненты HadoopDB (16) Планирование SQL (10)

- Такой план производится в том случае, если таблица sales является разделенной по YEAR(saleDate)



```
SELECT YEAR(saleDate),  
SUM(revenue)  
FROM sales  
GROUP BY YEAR(saleDate);
```

- В этом случае вся логика выполнения запроса выталкивается в СУБД
- Задача Map всего лишь записывает результаты запроса в файл HDFS

23
апреля

Корпоративные
базы
данных

2010 г.

2010

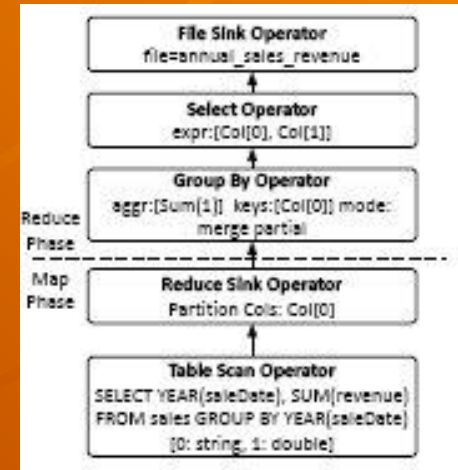
114

Параллельная СУБД на основе MapReduce (27)

Общая организация HadoopDB (21)

Собственные компоненты HadoopDB (17) Планирование SQL (11)

- В противном случае генерируется такой план
- При выполнении запроса по этому плану на уровне базы данных производится частичная агрегация данных,



```
SELECT YEAR(saleDate),  
SUM(revenue)  
FROM sales  
GROUP BY YEAR(saleDate);
```

23 апреля 2010 г. Корпоративные базы данных

– а для окончательной агрегации требуется выполнение задачи *Reduce*, производящей слияние частичных результатов группировки, которые получены в каждом узле на фазе задачи *Map*

2010 г.

2010

115

Параллельная СУБД на основе MapReduce (28)

Характеристики HadoopDB (1)

- Описан ряд экспериментов, показывающих, что гибридное использование технологий MapReduce и баз данных в реализации HadoopDB позволяет добиться от этой системы
 - производительности, сопоставимой с производительностью параллельных СУБД, и
 - устойчивости к отказам и падению производительности узлов, свойственной MapReduce

• Кратко отметим основные результаты

23 апреля
2010 г.

Корпоративные базы данных
2010

116

Параллельная СУБД на основе MapReduce (29)

Характеристики HadoopDB (2)

Производительность и масштабируемость (1)

- В большинстве экспериментов
 - параллельные СУБД существенно превосходят HadoopDB по производительности,
 - а HadoopDB оказывается значительно (иногда на порядок) производительнее связки Hive и Hadoop MapReduce
- В экспериментах использовались покладочная параллельная СУБД Vertica и некоторая коммерческая параллельная СУБД-Х с хранением таблиц по строкам
- Наибольшую производительность, естественно, продемонстрировала Vertica
 - но в ряде случаев HadoopDB уступала ей значительно меньше, чем на десятичный порядок

23
апреля

2010 г.

Корпоративные
базы
данных

2010

117

Параллельная СУБД на основе MapReduce (30)

Характеристики HadoopDB (3)

Производительность и масштабируемость (2)

- Значительное отставание HadoopDB от параллельных СУБД объясняется тем, что в HadoopDB использовалась PostgreSQL, в которой отсутствует возможность хранения таблиц по столбцам
 - как уже отмечалось, в HadoopDB уже используется поколоночная СУБД
- Кроме того, в экспериментах с HadoopDB не использовалось сжатие данных
- Наконец, в HadoopDB возникали значительные накладные расходы на взаимодействие Hadoop MapReduce и PostgreSQL, которые потенциально можно снизить

23
апреля
2010 г.

Корпоративные
базы
данных
2010

118

Параллельная СУБД на основе MapReduce (31)

Характеристики HadoopDB (4)

Производительность и масштабируемость (3)

- Так что в целом производительность HadoopDB не должна критически отставать от производительности параллельных СУБД
- Время загрузки данных в HadoopDB в десять раз больше соответствующего времени для Hadoop MapReduce
- Однако это окупается десятикратным выигрышем в производительности при выполнении некоторых запросов

23
апреля

2010 г.

Корпоративные
базы
данных

2010

119

Параллельная СУБД на основе MapReduce (32)

Характеристики HadoopDB (5)

Производительность и масштабируемость (4)

- При возрастании числа узлов в кластере при одновременном увеличении объема данных HadoopDB (как и Hadoop) масштабируется почти линейно
- Но в этом диапазоне не хуже масштабируется и Vertica
 - с СУБД-Х дела обстоят несколько хуже,
- а эксперименты на кластерах большего размера не производились
- Так что объективных данных в этом отношении пока нет

23
апреля

2010 г.

Корпоративные
ИВН
базы
данных

2010

120

Параллельная СУБД на основе MapReduce (33)

Характеристики HadoopDB (6)

Устойчивость к отказам и неоднородности среды (1)

- В экспериментах с отказоустойчивостью и падением производительности некоторого узла сравнивались HadoopDB, Hadoop MapReduce с Hive и Vertica
- В первом случае работа одного из узлов кластера искусственным образом прекращалась после выполнения 50% обработки запроса
- Во втором случае работа одного узла замедлялась за счет выполнения фонового задания с большим объемом ввода-вывода с тем же диском, на котором сохранялись файлы соответствующей системы

23
апреля

2010 г.

Корпоративные
базы
данных

2010

121

Параллельная СУБД на основе MapReduce (34)

Характеристики HadoopDB (7)

Устойчивость к отказам и неоднородности среды (2)

- При продолжении работы после отказа узла СУБД Vertica приходилось выполнять запрос заново с использованием реплик данных, и
 - время выполнения запроса возрастало почти вдвое
- В HadoopDB и Hadoop MapReduce с Hive время выполнения увеличивалось примерно на 15-20% за счет того, что
 - задачи, выполнявшиеся на отказавшем узле, перераспределялись между оставшимися узлами
- При этом относительная производительность HadoopDB оказывается несколько выше, чем у Hadoop MapReduce с Hive, поскольку
 - в первом случае обработка запроса проталкивалась на узлы, содержащие реплики баз данных,
 - а во втором приходилось копировать данные, не являющиеся локальными для обрабатывающего узла

23
апреля

2010 г.

Корпоративные
базы
данных

2010

122

Параллельная СУБД на основе MapReduce (35)

Характеристики HadoopDB (8)

Устойчивость к отказам и неоднородности среды (3)

- При замедлении работы одного из узлов производительность Vertica определялась скоростью этого узла,
 - и в экспериментах время выполнения запроса увеличивалось на 170%
- При использовании HadoopDB и Hadoop MapReduce с Hive время выполнения запроса увеличивалось
 - всего на 30% за счет образования резервных избыточных задач в недозагруженных узлах

апреля

данных

2010 г.

2010

123

Параллельная СУБД на основе MapReduce (36)

- Проект HadoopDB представляется очень интересным и перспективным
- В отличие от других систем, HadoopDB – это проект с открытыми исходными текстами, так что потенциально участие в этой работе доступно для всех желающих
- Помимо прочего, проект HadoopDB открывает путь к созданию высокопроизводительных, масштабируемых и отказоустойчивых параллельных СУБД на основе имеющихся программных средств с открытыми кодами

апреля

2010 г.

данных

2010

124

ETL с использованием MapReduce (1)

- Различных средств ETL (Extract-Transform-Load) существует великое множество, и они применяются во многих компаниях, использующих хранилища данных
 - трудно даже пытаться как-то их классифицировать и/или сравнивать
- Но, в любом случае, важность этих средств трудно переоценить, поскольку в хранилище данных по определению поступают данные из самых разнообразных источников:
 - транзакционных баз данных,
 - сообщений, участвующих в организации бизнес-процессов,
 - электронной почты,
 - журналов Web-серверов и т.д.
- Все эти данные нужно очистить, привести к единому формату, согласовать и загрузить в хранилище данных для последующего анализа

апреля

2010 г.

данных

2010

125

ETL с использованием MapReduce (2)

- Можно согласиться, что при использовании ортодоксального подхода к организации хранилищ данных подключение нового источника к хранилищу данных может занять недопустимо много времени во многом
 - из-за потребности в создании соответствующей процедуры ETL
- Можно согласиться и с тем, что для аналитиков
 - гораздо важнее получить новые данные, чем быть вынужденными ждать неопределенное время их в согласованной форме
- Но совершенно очевидно, что если данные в хранилище данных не очищать никогда, то со временем в них не разберется никакой, даже самый передовой аналитик

Корпорат

ивные

базы

данных

23
апреля

2010 г.

2010

126

ETL с использованием MapReduce (3)

- Итак, что мы имеем
- Число источников данных, пригодных для анализа в составе хранилища данных, все время растет
- Их разнородность тоже все время возрастает
- Все меньший процент составляют структурированные базы данных,
 - данные поступают из частично структурированных файлов и совсем неструктурированных текстовых документов
- Для каждой разновидности источников данных нужна своя разновидность процедуры ETL, и
 - по причине роста объемов исходных данных для обеспечения умеренного времени их загрузки в хранилище данных эти процедуры должны выполняться в массивно-параллельной среде
- И в этом может помочь технология MapReduce

апреля

2010 г.

Корпоративные
базы
данных

2010

127

ETL с использованием MapReduce (4)

MapReduce и ETL (1)

- Для канонического способа использования технологии MapReduce характерно применение следующих операций:
 - чтение журнальных данных из нескольких разных файлов-журналов;
 - разбор и очистка журнальных данных;
 - преобразования этих данных, в том числе их частичная агрегация;
 - принятие решения о схеме результирующих данных;
 - загрузка данных в хранилище данных или другую систему хранения

23
апреля
2010 г.

Корпоративные
базы
данных
2010

128

ETL с использованием MapReduce (5)

MapReduce и ETL (2)

- В точности такие же шаги выполняются в системах ETL при извлечении, преобразовании и загрузке данных
- По сути дела, MapReduce производит из исходных "сырых" данных некоторую полезную информацию
 - которую потребляет другая система хранения
- В некотором смысле можно считать любую реализацию MapReduce
 - параллельной инфраструктурой выполнения процедур ETL

23
апреля
2010 г.

Кортрат
ивные
базы
данных
2010

129

ETL с использованием MapReduce (6)

MapReduce и ETL (3)

- Имелись попытки реализации процедур ETL внутри сервера баз данных средствами языка SQL
- Разработчики параллельных СУБД с поддержкой MapReduce Greenplum Database и nCluster компании Aster Data тоже намекают, что их встроенный MapReduce можно использовать и для поддержки ETL
- Но исторически системы ETL промышленного уровня существуют отдельно от СУБД
- Обычно СУБД не пытаются выполнять ETL, а системы ETL не поддерживают функции СУБД

23
апреля

2010 г.

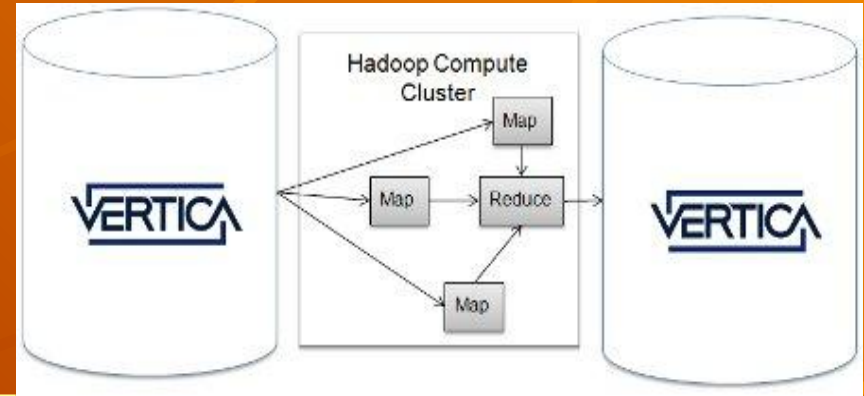
Корпоративные
базы
данных

2010

130

ETL с использованием MapReduce (7) Hadoop и Vertica (1)

- Если иметь в виду поддержку именно ETL, то наиболее грамотное



Корпорат

– и самое простое

ивные

- решение по интеграции технологий MapReduce и параллельных баз данных применяется в Vertica

23

базы

апреля

данных

2010 г.

2010

131

ETL с использованием MapReduce (8) Hadoop и Vertica (2)

- В Vertica реализован свой вариант интерфейса DBInputFormat компании Cloudera для Hadoop MapReduce, позволяющий разработчикам MapReduce
 - выбирать данные из баз данных Vertica и
 - направлять результирующие данные в эти базы данных
- При этом подходе технологии MapReduce и параллельных баз данных тесно не интегрируются, но каждая из них может использовать возможности другой

технологии
атреля

Корпоративные
базы
данных

2010 г.

2010

132

ETL с использованием MapReduce (9)

- Скорее всего, мы еще многое услышим о системах ETL, основанных на использовании технологии MapReduce, и, скорее всего,

Корпорат

– предводителем этого направления

будет Vertica

ивные

базы

апреля

данных

2010 г.

2010

133

Заключение (1)

- Еще пару лет назад было непонятно, каким образом можно с пользой применять возникающие "облачные" среды для высокоуровневого управления данными
- Многие люди считали, что в "облаках" системы управления базами данных будут просто вытеснены технологией MapReduce
- Это вызывало естественное недовольство сообщества баз данных, авторитетные представители которого старались доказать, что попытка заменить СУБД какой-либо реализацией MapReduce – если не безнравственно, то, по крайней мере, неэффективно

23
апреля
2010 г.

Корпоративные
базы
данных
2010

134

Заключение (2)

- Однако вскоре стало понятно, что технология MapReduce может быть полезна для самих параллельных СУБД
- Во многом становлению и реализации этой идеи способствовали компании-стартапы, выводящие на рынок новые аналитические массивно-параллельные СУБД и добивающиеся конкурентных преимуществ
- Свою лепту вносили и продолжают вносить и университетские исследовательские коллективы, находящиеся в тесном сотрудничестве с этими начинающими компаниями

23
апреля

Корпоративные
базы
данных

2010 г.

2010

135

Заключение (3)

- На сегодняшний день уже понятно, что технология MapReduce может
 - эффективно применяться внутри параллельной аналитической СУБД,
 - служить инфраструктурой отказоустойчивой параллельной СУБД,
 - а также сохранять свою автономность в симбиотическом союзе с параллельной СУБД
- Все это не только мешает развитию технологии параллельных СУБД, а наоборот, способствует ее совершенствованию и распространению

23
апреля

Корпоративные
базы
данных

2010 г.

2010

136

Заключение (4)

- Интересные работы ведутся и в направлении использования "облачных" сред для создания нового поколения транзакционных средств управления данными
- Но это уже совсем другая история

23

базы

апреля

данных

2010 г.

2010

137

Заключение (5)

- Полный текст:
- *Сергей Кузнецов. MapReduce: внутри, снаружи или сбоку от параллельных СУБД?*

- http://citforum.ru/database/articles/dw_appliance_and_mr/

Корпорат

ИВНЫЕ

базы

данных

23

апреля

2010 г.

2010

138

Спасибо за терпение!



23

апреля

2010 г.

2010

139