

Parallel and Distributed

Computing - P1

紧密耦合

的并行计算 (tightly-coupled parallel computing)

shared-memory symmetric multiprocessor (SMP) system

(统一内存访问) (uniform memory access)

memory connection on single bus

Cache consistency

非统一内存访问 (non-uniform memory access, NUMA)

single memory space (local / non-local memory access)

memory bus and switch architecture

松散耦合的分布式架构 (loosely-coupled distributed architecture)

use multi-SMP system connected by high-speed network

major challenge: software to utilize this architecture

典型的数据中心组织结构 (typical datacenter Organization)

服务器 (server): single SMP computer with network connection and storage

机架 (rack): physical stack of servers, with network switch

通道 (aisle): set of physically adjacent racks, with network switch

高可靠性硬件 (high-reliability hardware)

every computer/switch/etc. has redundant components (冗余备份)

high availability of each components (高度可用)

high initial cost, high operating cost (高配置/运行成本)

rare failure of computer server

商用硬件 (commodity hardware)

off-the-shelf component (现成组件)

low initial cost, no idle redundant components

typical unit of failure: sever

more frequent, software is designed to deal with

其他设计因素 (design consideration)

individual server processing and I/O speed, sever density

power, cooling, internal network latency, external network connectivity

Parallel and Distributed Computing - P2

弹性软件 (resilient software), 设计用于商用硬件系统
designed to cope with loss of servers without restarting massive computations
restart only the portion of the computation running on the lost server
文件必须驻留在多个服务器上并提供无缝切换

分布式算法 (distributed algorithm)

same algorithm suitable for SMP, not suitable for distributed solution

典型的数据中心软件组织 (example datacenter software organization)

server: Linux operating system

Hadoop scheduler Yarn: on one or more nodes

each compute server: Hadoop job scheduling node software

Distributed file system (HDFS) master /server : on one or more nodes

Distributed database server (HBase, Cassandra): on one or more nodes

User → upload data into the file and database system

→ submit job through Yarn

计算集群 (compute cluster), 在典型的计算集群中有两个部分

分布式文件系统 (distributed file system)

在 Hadoop 中为 HDFS (Hadoop Distributed File System)

files may be accessed from any node in the cluster

need special procedure, not accessible from normal Linux interface

not a shared file system, 如网络文件系统 (network file system, NFS)

工作分配者 (job scheduler)

在 Hadoop 中为 YARN

building a program for YARN (写一个基本的 Java 程序)

必须是 .jar 文件, 且包含程序运行需要的所有自定义 java 类 (user-defined java class)

使用 Maven (mvn) 或 ant build tool

需要设置环境变量 CLASSPATH, 可通过 env 命令查看

hadoop jar <parameter> run a Hadoop MapReduce job

spark-submit --master yarn --deploy-mode cluster <parameter>

: run a Spark job in the cluster

Parallel and Distributed

Computing - P3

YARN 的文件处理

jobs running on the cluster may not have access to user Linux FS
data for jobs under YARN should be stored in HDFS
output under YARN should be generated in HDFS
can copy the result to Linux FS after job completion

MapReduce : large scale data processing

automatic parallelization and distribution

fault-tolerance

I/O scheduling

status and monitoring

- ① start up copies of user programs on cluster

- ② one special is master

- ③ master assign map/reduce to idle worker

- ④ read content of corresponding input split

- ⑤ intermediate data write to local dist

- ⑥ remote read the data and do reduce task

- ⑦ reduce partition append to final output file



fine granularity task : many more map tasks than machines in the cluster

minimize time for fault recovery

can pipeline shuffling with map execution

better dynamic load balancing

Processor | Time —>

User program

MapReduce()

assign tasks to workers...

master

worker 1

Map 1 Map 3

worker 2

Map 2

Map 1 Map 3

worker 3

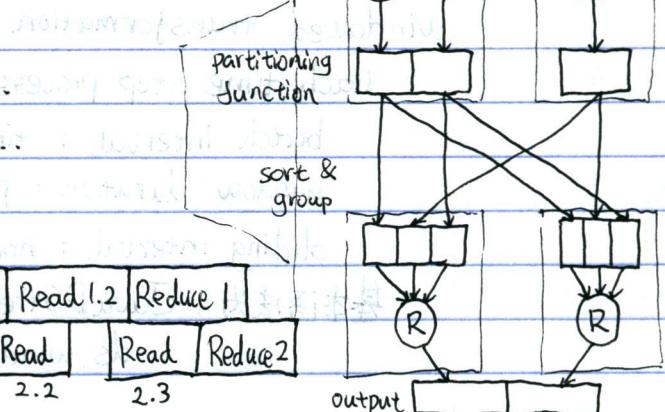
Map 1 Map 3

Read 1.1 Read 1.3 Read 1.2 Reduce 1

worker 4

Map 2

Read 2.1 Read 2.3 Reduce 2



Parallel and Distributed

Computing - P4

MapReduce fault tolerance : handle via re-execution
on worker failure (master failure unlikely)
detect failure via periodic heartbeats
re-execute completed and in-progress map tasks
re-execute in-progress reduce tasks
task competition committed through master

refinement and redundant execution slow workers significantly lengthen completion time
bad disks with soft errors transfer data slowly
processor caches disabled (weird)
solution : spawn backup copies of tasks near end of phase
dramatically shorten job completion time

locality optimization : master scheduling policy

ask GFS (global file system) for locations of replicas of input file blocks
map tasks typically split into GFS block size
scheduled so GFS input block replica on same machine or same rack
thousands of machines read input at local disk speed

without this, rack switch limit read rate

skipping bad records : Map/Reduce function fail for particular input

best solution is to debug & fix, not always possible
on seg fault, send UDP (user datagram protocol) packet to master
from signal handler include sequence num of record being processed
if master sees two failure on same record, next worker skip record
can work around bugs in third-party libraries

Sorting guarantees within each reduce partition

UDP compression of intermediate data
combiner : useful for saving network bandwidth
local execution for debugging/testing
user-defined counters

Parallel and Distributed Computing - Part 5

Common RDD types

common RDD (resilient distributed dataset) types:

注意：这里的 RDD type t 以 Java 中实现为准，且大小写敏感

JavaRDD<T>: an RDD made up of items of class T

so, JavaRDD<String> for data from a text file

假定文件中的每行都是一个文本项

JavaPairRDD<K, V>: an RDD made up of Key, Value pairs

key is class K, value is class V

注意：Java 中的 key, value pair 表示为 Tuple2 类型

即有 Tuple2<K, V>, 并且提供方法读取 key 或 value

so K key = tuple2._1();

V value = tuple2._2();

transformation require specification of function to perform

Java function is always part of a class

pass a class object which has at least a method named "call"

as a function to Spark transformation

Spark defines Java interfaces for commonly-used transformation class

simple map : Function, (org.apache.spark.api.java.function.Function)

函数类型为 Function<T, R> :: T → R

对应的 map 函数效果为

map :: JavaRDD<T> → Function<T, R> → JavaRDD<R>

其语法结构为：对于 JavaRDD<T> t

JavaRDD<R> r = t.mapC

】定义了新的 JavaRDD<R> 对象 r

new Function<T, R>() {

】传入 map 函数的 Function 实例

public R call(T arg) {

】注意扩号为空是因为其构造函数

不需要传入参数

return ret; } }] Function 实例中的 call 方法定义

return ret; }] 返回值类型为 R,

}

传入参数为类型 T 的 arg

);

return ret; }] 返回值 ret 也应为类型 R

Parallel and Distributed Computing - P6

flatMap

generate multiple outputs from each input

函数类型为 FlatMapFunction<T, R>

且有 FlatMapFunction<T, R> :: T -> Iterator<R>

flatMap 效果为 flatMap :: JavaRDD<T> -> FlatMapFunction<T, R> -> JavaRDD<R>

基本语法为，对于 JavaRDD<T> t

JavaRDD<R> r = t.flatMap(...)

```
new FlatMapFunction<T, R>() { ... }
```

```
public Iterator<R> call(T arg) { ... }
```

```
ArrayList<R> ret; ...
```

```
return ret.iterator(); } } );
```

mapToPair

generate key, value pair from each input

函数类型为 PairFunction<T, K, V>

且有 PairFunction<T, K, V> :: T -> Tuple2<K, V>

于是有 mapToPair :: JavaRDD<T> -> PairFunction<T, K, V> -> JavaPairRDD<K, V>

基本语法为，对于 JavaRDD<T> t

JavaPairRDD<K, V> p = t.mapToPair(...)

```
new PairFunction<T, K, V>() { ... }
```

```
public Tuple2<K, V> call(T arg) { ... }
```

```
Tuple2<K, V> ret; ...
```

```
return ret; } } );
```

flatMapToPair

generate multiple pairs from each input

函数类型为 PairFlatMapFunction<T, K, V>

且有 PairFlatMapFunction<T, K, V> :: T -> Iterator<Tuple2<K, V>>

于是有 flatMapToPair :: JavaRDD<T> -> PairFlatMapFunction<T, K, V> -> JavaPairRDD<K, V>

基本语法为，对于 JavaRDD<T> t

JavaPairRDD<K, V> p = t.flatMapToPair(...)

```
new PairFlatMapFunction<T, K, V>() { ... }
```

```
public Iterator<Tuple2<K, V>> call(T arg) { ... }
```

```
ArrayList<Tuple2<K, V>> ret; ...
```

```
return ret.iterator(); } } );
```

Parallel and Distributed Computing - P7

reduceByKey combine all source items which have the same key

函数类型为 $\text{Function2} < V, V, V > :: V \rightarrow V \rightarrow V$

且函数必须满足结合律 (associative) 和交换律 (commutative)

于是有 $\text{reduceByKey} :: \text{JavaPairRDD} < K, V > \rightarrow \text{Function2} < V, V, V > \rightarrow \text{JavaPairRDD} < K, V >$

基本语法为，对于 $\text{JavaPairRDD} < K, V > P1$

$\text{JavaPairRDD} < K, V > P2 = P1.\text{reduceByKey}($

$\text{Function2} < V, V, V >()$ {

$\text{public } V \text{ call } (V x, V y) \{ \dots \}$

$V \text{ ret}; \dots$

$\text{return ret;} \} \});$

serializable (可序列化的)，在 Java 中，对于 $\text{JavaRDD} < T >$ 和 $\text{JavaPairRDD} < K, V >$

其中的数据类型必须是 serializable

class object can be stored in a byte representation

from which it can be later restored elsewhere

necessary : data will be stored in file

and/or transmitted over network

class contain only serializable data members

may be made serializable with no special handling other than declaration

links, recursion, nested class and reference

may require special handling

`implements Serializable` allow definition of method

to serialize and de-serialize class object

基本语法结构为 `import java.io.Serializable;`

`class ClassName implements Serializable {`

`int d;`

`double f;`

`String s; }`

only data members matter

function method don't affect serialization

groupByKey not need to combine values for given key and just want to group

result in an iterator over values for each key , no user function needed

现有 $\text{groupByKey} :: \text{JavaPairRDD} < K, V > \rightarrow \text{JavaPairRDD} < K, \text{Iterable} < V > \rangle$

基本语法为，对于 $\text{JavaPairRDD} < K, V > P1$

$\text{JavaPairRDD} < K, \text{Iterable} < V > \rangle P2 = P1.\text{groupByKey}();$

Parallel and Distributed Computing - P8

partition (划分), RDD in Spark is parallelized in partitions

并行分区，每个分区可能在不同的工人上处理

分区数由 default number of partitions 可能取决于可用的工人、输入碎片等。

可以覆盖默认值为指定值

map 操作通常在同一个分区中完成

group 或 aggregate 操作经常导致不同的分区数

CombineByKey, general-purpose transformation by key

用于实现其他操作，如 reduceByKey, groupByKey

可能返回与源不同类型的值

调用者提供三个函数：

createCombiner : initialize a combiner with the first value

即对于处理的第一个 (key, value)，用 value 初始化 combiner

于是有 $\text{createCombiner} : V \rightarrow W$

mergeValue : add single value to an existing combiner

即描述如何将一个 (key, value) 中的 value 加入一个 combiner

于是有 $\text{mergeValue} : W \rightarrow V \rightarrow W$

mergeCombiners : combine two separate combiners into single combiner

即描述如何合并两个现有的 combiner

于是有 $\text{mergeCombiners} : W \rightarrow W \rightarrow W$

则有 $\text{combineByKey} : \text{JavaPairRDD}(K, V) \rightarrow (V \rightarrow W) \rightarrow (W \rightarrow V \rightarrow W)$

$\rightarrow (W \rightarrow W \rightarrow W) \rightarrow \text{JavaPairRDD}(K, W)$

即对 $\text{JavaPairRDD}(K, \text{int})$ p1, 实现 $\text{groupByKey} \circ \text{reduceByKey}(\text{sum})$

$\text{JavaPairRDD}(K, \text{ArrayList<int>} P1) \rightarrow \text{JavaPairRDD}(K, \text{int} P2)$

$= P1. \text{combineByKey}($

$= P1. \text{combineByKey}($

$\text{value} \rightarrow \{ \text{list} \}$, $\text{list} \rightarrow \{ \text{value} \}$, $\text{value} \rightarrow \{ \text{list} \}$

$\text{ArrayList<int>} \text{combiner} = \text{new ArrayList<int>}();$, $\text{return value}; \}$

$\text{combiner.add(value);}$

$(\text{sum}, \text{value}) \rightarrow \{$

$\text{return combiner}; \}$

$\text{return (sum + value);}$

$\text{list}, \text{value} \rightarrow \{ \text{list.add(value);}$

$(\text{sum1}, \text{sum2}) \rightarrow \{ \text{sum1} + \text{sum2}; \}$

$\text{return list}; \}$

$\}$

$(\text{list1}, \text{list2}) \rightarrow \{ \text{list1.addAll(list2);}$

$\}$

$\text{return list1}; \}$

$\}$

Parallel and Distributed Computing - Pg 9

Computing - Pg

data flow for combineByKey, within each source partition:

within each source partition:

combiner created for each key the first time the key occurs

subsequent occurrence of record with the same key add to combiner

after source partition processed:

hash of key determine partition in destination RDD

within each destination partition:

merge combiners if more than one occurrence of (key, combiner) pair

Spark object file format: an easy way to store user-defined serializable Java class
use the objectFile method of the JavaSparkContext

to create an RDD from an object file

use the saveAsObjectFile method of an RDD

to save the RDD to an object file

the class definition data members must not have changed between saving and reading the object file

must not have changed between saving and reading the object file

Java serialization less efficient than Hadoop sequence file

but much easier to use

Hadoop writable object in Spark, Hadoop define interface Writable

a custom serialization interface defined in Hadoop

can implement in own class

require implementing the method to serialize and de-serialize object

useful classes provided ending in "Writable", e.g. IntWritable, FloatWritable

non-transparent: methods need to get / set value

Spark can read/write Sequence File composed of object

reading: simple, using sequenceFile of JavaSparkContext

writing: complicated, require Spark interface to Hadoop routines

lambda expression in Java, an easy way to implement anonymous class

基本语法: arg1 | arg2, ... > 等价于 class lambdaFunc {

> -> { ... } > -> public U call(T arg){ ... }

return ret; } > -> public U ret; > -> ret;

so func :: T -> U

return ret; }

Parallel and Distributed Computing - P10

function specified for spark operation (map, reduceByKey, etc)

can be written to access variable at the outer class level

work fine when running locally

doesn't work correctly when running in distributed mode

outer class variable only get initialized in the master

worker only gets a serialized version of the function

without outer class variable

lambda run in a closure which inherits function local variable

function object : pass needed variable to the constructor

store inside the function object class as class variable

Pregel algorithm designed to solve iterative, parallel, synchronous graph algorithm

input directed graph

each vertex and edge can have associated data

user program run at each vertex

initialization → compute (superstep processing) → output

message passing model : vertex can send message

typically (not necessarily) along edge

message type specified by programmer

message sent in superstep n are received in superstep n+1

termination condition : processing terminate when all vertices 'vote' to quit

when vertex has no immediate action to take without further input processing continue if any message sent by any vertices

define combiner : multiple incoming messages can be combined

aggregator : combine specified value from each vertex at end of superstep

reduction using programmer defined code

the aggregated result is made available to all vertices at next superstep

Parallel and Distributed Computing - P11

Pregel

topology change : vertex can add/delete vertex/edge during superstep
become effective next superstep
removal first, then vertex addition, then edge addition
conflict handled by user-supplied handler

input / output : graph input/output class specified by user

can be predefined or user-defined class

any arbitrary format can be processed

data can be read from/ write to (GFS(HDFS) or Bigtable (HBase))

master task coordinate all workers

computation run on set of partitions (number specified by user)

each worker handle one or more partition

worker read part of the input in parallel

input contain vertices and edges

default assignment of vertices to partitions by hashing vertex ID

user can supply assignment to conform to graph topology

more efficient to put nearby vertices on same or nearby worker task

vertices not in worker's own partition are send by message to appropriate worker

superstep execution : at each superstep, the worker

call the compute method of the vertex for each of its vertices

passing in any arriving messages

send outgoing messages asynchronously to destination vertex worker

for delivery next superstep

tell master when all vertices have completed superstep

tell master the number of vertices still active

algorithm terminate when all vertices have voted to halt

and all vertex input queues are empty

upon termination, the output method of all vertices and edges are called to write the output

system get and release queue data no obvious id from assignment