

MapReduce 高级程序设计



摘要

2

- 复合键值对的使用
- 用户自定义数据类型
- 用户自定义输入输出格式
- 用户自定义Partitioner和Combiner
- 迭代完成MapReduce计算
- 链式MapReduce任务
- 全局参数/数据文件的传递
- 其它处理技术



复合键值对的使用

3

□ 思路：用复合键让系统完成排序

- ▣ 问题：map计算过程结束后进行Partitioning处理时，系统自动按照map的输出键进行排序，因此，进入Reduce节点的(key, [value])对将保证是按照key进行排序的，而[value]则不保证是排好序的。为了解决这个问题，可以在Reduce过程中对[value]列表中的各个value进行本地排序。但当[value]列表数据量巨大、无法在本地内存中进行排序时，将出现问题。
- ▣ 改进方法：将value中需要排序的部分加入到key中形成复合键，这样将能利用MapReduce系统的排序功能完成排序。
- ▣ 代价：但需要实现一个新的Partitioner，保证原来同一key值的键值对最后分区到同一个Reduce节点上。



复合键值对的使用

4

□ 带频率的倒排索引示例

1: **class Mapper**

2: **procedure** Map(docid n, doc d)

3: $H \leftarrow \text{new AssociativeArray}$

4: **for all** term $t \in \text{doc } d$ **do**

5: $H\{t\} \leftarrow H\{t\} + 1$

6: **for all** term $t \in H$ **do**

7: Emit(term t , posting $\langle n, H\{t\} \rangle$)

1: **class Reducer**

2: **procedure** Reduce(term t , postings $[\langle n_1, f_1 \rangle, \langle n_2, f_2 \rangle \dots]$)

3: $P \leftarrow \text{new List}$

4: **for all** posting $\langle a, f \rangle \in \text{postings}$ $[\langle n_1, f_1 \rangle, \langle n_2, f_2 \rangle \dots]$ **do**

5: Append($P, \langle a, f \rangle$)

6: Sort(P); // 进入Reduce节点的postings不保证按照文档序号排序,因而需要对postings进行一个本地排序

7: Emit(term t ; postings P)

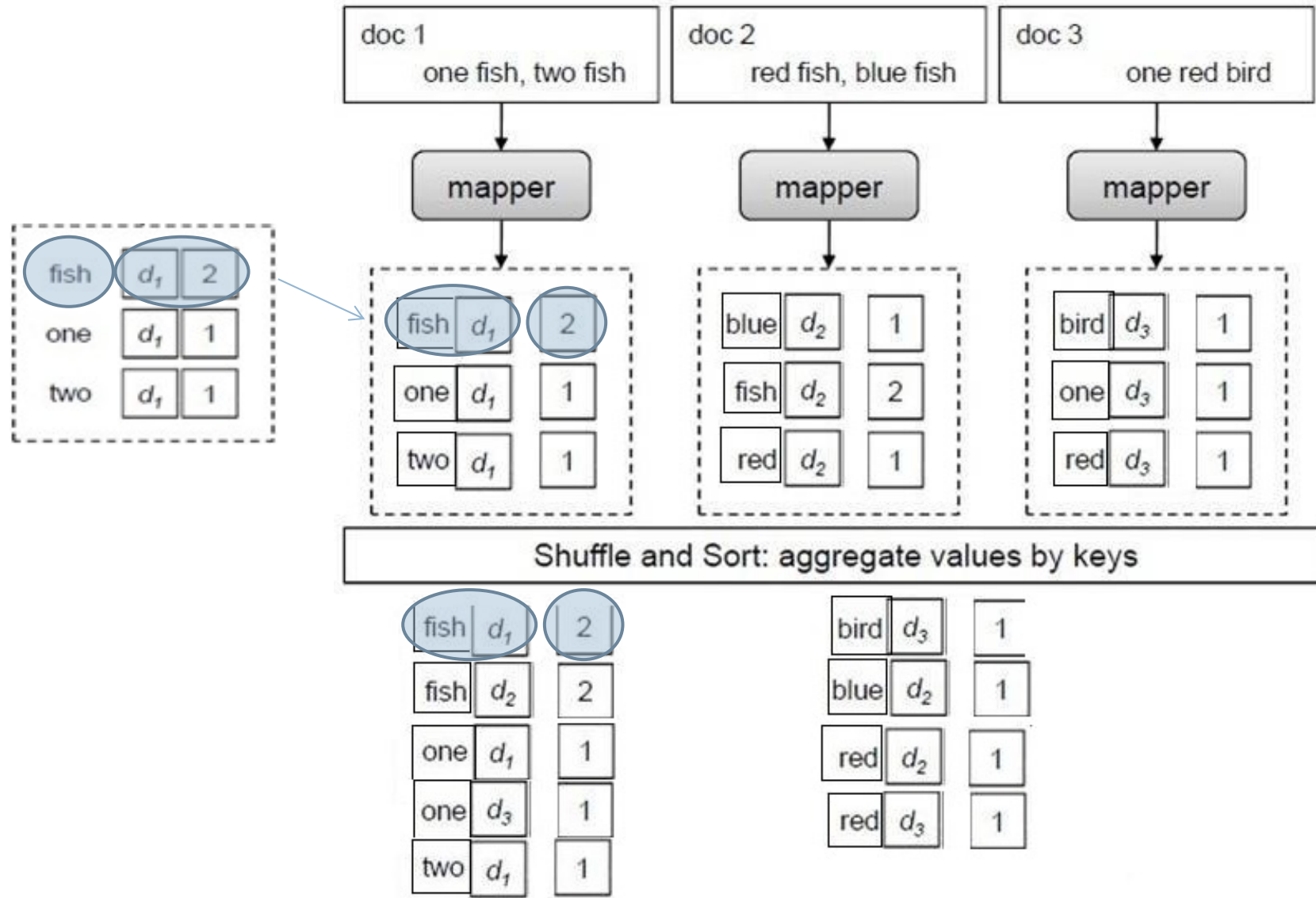


复合键值对的使用

5

□ 带频率的倒排索引示例

- ▣ 为了能利用系统自动对`docid`进行排序，解决方法是：代之以生成 (`term`, `<docid, tf>`) 键值对，`map`时将`term`和`docid`组合起来形成复合键`<term, docid>`。
- ▣ 但会引起新的问题，同一个`term`下的所有`posting`信息无法被分区到同一个`Reduce`节点，为此，需要实现一个新的`Partitioner`：从`<term, docid>`中取出`term`，以`term`作为`key`进行分区。



Customized Partitioner

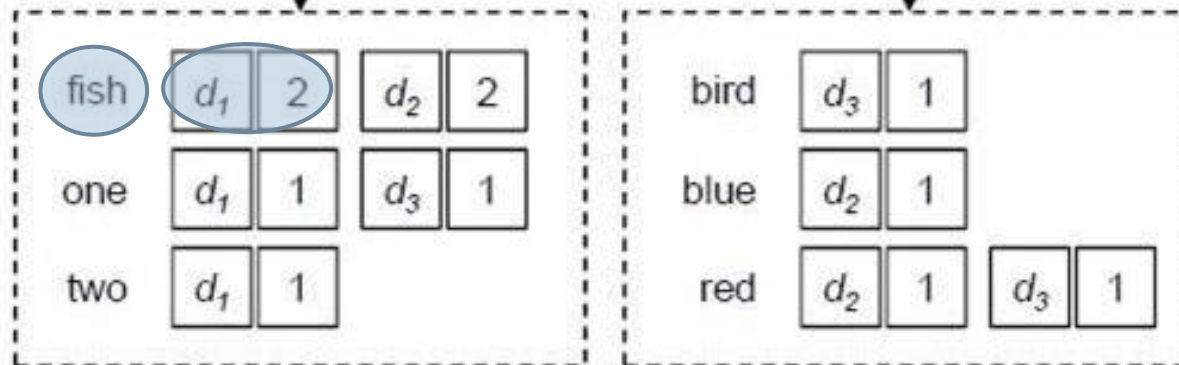
fish	d_1	2
fish	d_2	2
one	d_1	1
one	d_3	1
two	d_1	1

bird	d_3	1
blue	d_2	1
red	d_2	1
red	d_3	1

进入reduce的键值对按照(term, docid)排序

reducer

reducer





复合键值对的使用

8

- 思路：把小的键值对合并成大的键值对
 - ▣ 通常一个计算问题会产生大量的键值对，为了减少键值对传输和排序的开销，一些问题中的大量小的键值对可以被合并成一些大的键值对(**pairs**->**stripes**)。



复合键值对的使用

□ 例如：单词同现矩阵算法

- 一个Map可能会产生单词a与其它单词间的多个键值对，这些键值对可以在Map过程中合并成右侧的一个大的键值对(条):

$$\begin{array}{l} (a, b) \rightarrow 1 \\ (a, c) \rightarrow 2 \\ (a, d) \rightarrow 5 \\ (a, e) \rightarrow 3 \\ (a, f) \rightarrow 2 \end{array} \quad \longrightarrow \quad a \rightarrow \{ b: 1, c: 2, d: 5, e: 3, f: 2 \}$$

- 然后，在Reduce阶段，把每个单词a的键值对(条)进行累加:

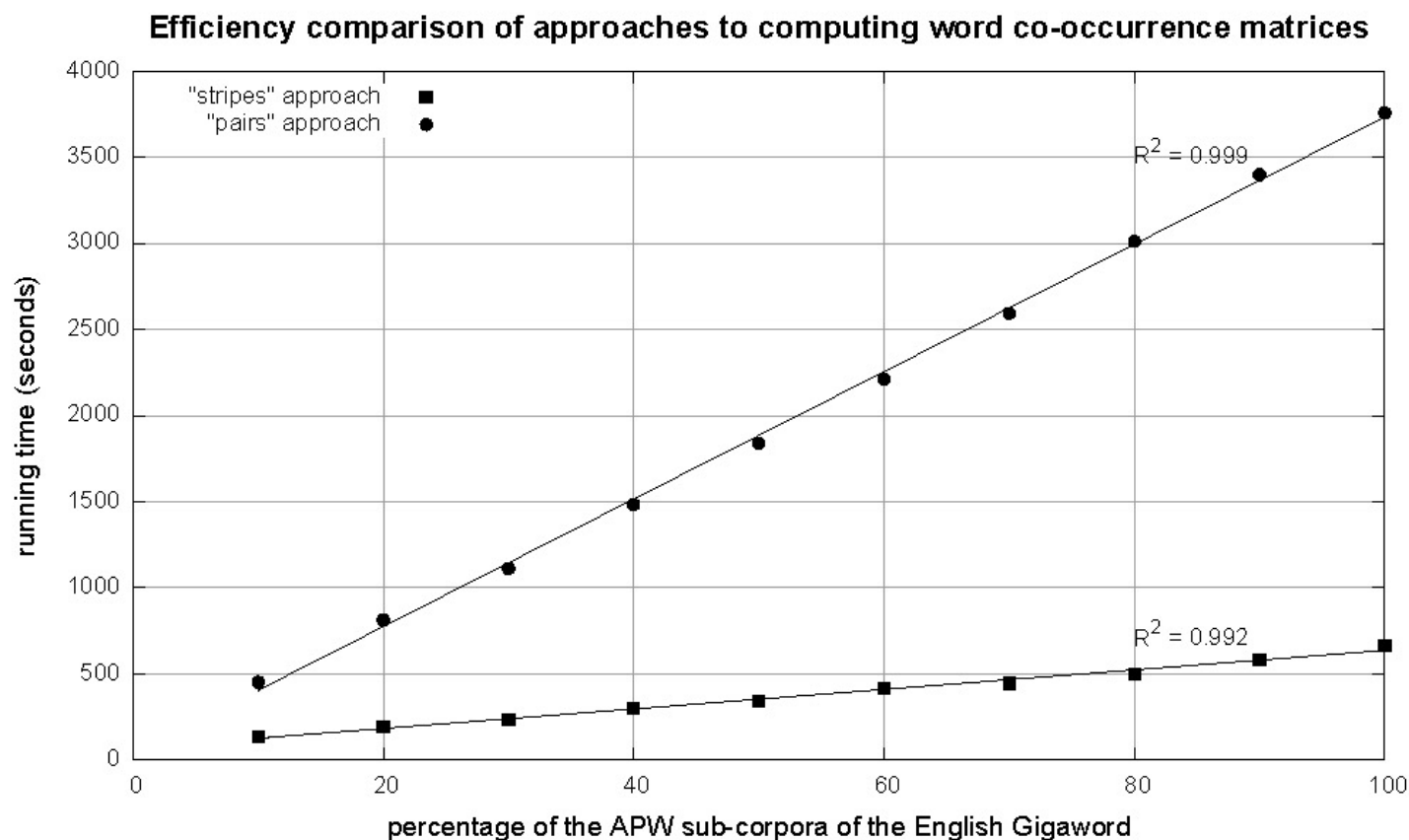
$$\begin{array}{r} a \rightarrow \{ b: 1, \quad d: 5, e: 3 \} \\ + \quad a \rightarrow \{ b: 1, c: 2, d: 2, \quad f: 2 \} \\ \hline a \rightarrow \{ b: 2, c: 2, d: 7, e: 3, f: 2 \} \end{array}$$



复合键值对的使用

10

□ 单词同现矩阵算法



Cluster size: 38 cores

Data Source: Associated Press Worldstream (APW) of the English Gigaword Corpus (v3), which contains 2.27 million documents (1.8 GB compressed, 5.7 GB uncompressed)



用户自定义数据类型

11

- Hadoop内置的数据类型，这些数据类型都实现了**WritableComparable**接口，以便进行网络传输和文件存储，以及进行大小比较。

Class	Description
BooleanWritable	Wrapper for a standard Boolean variable
ByteWritable	Wrapper for a single byte
DoubleWritable	Wrapper for a Double
FloatWritable	Wrapper for a Float
IntWritable	Wrapper for a Integer
LongWritable	Wrapper for a Long
NullWritable	Placeholder when the key or value is not needed
Text	Wrapper to store text using the UTF-8 format



用户自定义数据类型

12

- 需要实现**Writable**接口，作为**key**或者需要比较大小时则需要实现**WritableComparable**接口。

```
public class Point3D implements WritableComparable <Point3D>{
    private int x, y, z;
    public int getX() { return x; }
    public int getY() { return y; }
    public int getZ() { return z; }
    public void write(DataOutput out) throws IOException{
        out.writeFloat(x);
        out.writeFloat(y);
        out.writeFloat(z);
    }
    public void readFields(DataInput in) throws IOException{
        x = in.readFloat();
        y = in.readFloat();
        z = in.readFloat();
    }
    public int compareTo(Point3D p){
        //compares this(x, y, z) with p(x, y, z) and
        //outputs -1(小于), 0(等于), 1(大于)
    }
}
```



用户自定义数据类型

13

```
public class Edge implements WritableComparable<Edge>
{
    private String departureNode;
    private String arrivalNode;
    public String getDepartureNode() { return departureNode;}
    @Override
    public void readFields(DataInput in) throws IOException
    {
        departureNode = in.readUTF();
        arrivalNode = in.readUTF();
    }
    @Override
    public void write(DataOutput out) throws IOException
    {
        out.writeUTF(departureNode);
        out.writeUTF(arrivalNode);
    }
    @Override
    public int compareTo(Edge o)
    {
        return (departureNode.compareTo(o.departureNode)!=0)
            ?departureNode.compareTo(o.departureNode):arrivalNode.compareTo(o.arrivalNode);
    }
}
```



用户自定义输入输出格式

14

- 数据输入格式 (**InputFormat**) 用于描述**MapReduce**作业的数据输入规范。
- **MapReduce**框架依靠数据输入格式完成输入规范检查（比如输入文件目录的检查）、对数据文件进行输入分片 (**InputSplit**)，以及提供从输入分块中将数据记录逐一读出，并转换为**Map**过程的输入键值对等功能。
- **TextInputFormat**是系统缺省的数据输入格式。



用户自定义输入输出格式

15

□ Hadoop内置的文件输入格式

InputFormat:	Description:	Key:	Value:
TextInputFormat	Default format; reads lines of text files	The byte offset of the line	The line contents
KeyValueTextInputFormat	Parses lines into key-val pairs	Everything up to the first tab character	The remainder of the line
SequenceFileInputFormat	A Hadoop-specific high-performance binary format	user-defined	user-defined



用户自定义输入输出格式

16

- Hadoop内置的文件输入格式
- AutoInputFormat, CombineFileInputFormat, CompositeInputFormat, DBInputFormat, FileInputFormat, KeyValueTextInputFormat, LineDocInputFormat, MultiFileInputFormat, NLineInputFormat, SequenceFileAsBinaryInputFormat, SequenceFileAsTextInputFormat, SequenceFileInputFilter, SequenceFileInputFormat, StreamInputFormat, TextInputFormat



用户自定义输入输出格式

17

□ Hadoop内置的RecordReader

RecordReader:	InputFormat	Description:
LineRecordReader	default reader for TextInputFormat	reads lines of text files
KeyValueLineRecordReader	default reader for KeyValueTextInputFormat	parses lines into key-val pairs
SequenceFileRecordReader	default reader for SequenceFileInputFormat	User-defined methods to create keys and values



用户自定义输入输出格式

18

- **Hadoop 内置的 RecordReader**
- CombineFileRecordReader, DBInputFormat.DBRecordReader, InnerJoinRecordReader, JoinRecordReader, KeyValueLineRecordReader, LineDocRecordReader, MultiFilterRecordReader, OuterJoinRecordReader, OverrideRecordReader, SequenceFileAsBinaryInputFormat.SequenceFileAsBinaryRecordReader, SequenceFileAsTextRecordReader, SequenceFileRecordReader, StreamBaseRecordReader, StreamXmlRecordReader, WrappedRecordReader



用户自定义输入输出格式

19

□ 用户自定义InputFormat和RecordReader

简单的文档倒排索引

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class InvertedIndexMapper extends Mapper<Text, Text, Text, Text>
{
    @Override
    protected void map(Text key, Text value, Context context)
        throws IOException, InterruptedException
    {
        // default RecordReader: LineRecordReader;
        // key: line offset; value: line string
        Text word = new Text();
        FileSplit fileSplit = (FileSplit)context.getInputSplit();
        String fileName = fileSplit.getPath().getName();
        Text fileName_lineOffset = new Text(fileName+"@"+key.toString());
        StringTokenizer itr = new StringTokenizer(value.toString());
        for(; itr.hasMoreTokens(); )
        {
            word.set(itr.nextToken());
            context.write(word, fileName_lineOffset);
        }
    }
}
```

由于采用了缺省的
TextInputFormat和
LineRecordReader,
需要增加此段代码
完成特殊处理



用户自定义输入输出格式

□ 用户自定义 **InputFormat** 和 **RecordReader**

简单的文档倒排索引

可以自定义一个 **InputFormat** 和 **RecordReader** 实现同样的效果

```
public class FileNameLocInputFormat extends FileInputFormat<Text, Text>
{
    @Override
    public RecordReader<Text, Text> createRecordReader(InputSplit split,
                                                         TaskAttemptContext context)
    {
        FileNameLocRecordReader fnrr = new FileNameRecordReader();
        try
        {
            fnrr.initialize(split, context);
        }
        catch (IOException e) { e.printStackTrace(); }
        catch (InterruptedException e) { e.printStackTrace(); }
        return fnrr;
    }
}
```



用户自定义输入输出格式

21

```
public class FileNameLocRecordReader extends RecordReader<Text, Text>
{   String fileName;
    LineRecordReader lrr = new LineRecordReader();
    .....
    @override
    public Text getCurrentKey() throws IOException, InterruptedException
    {   return new Text("(" + fileName + "@" + lrr.getCurrentKey() + ")");   }
    @override
    public Text getCurrentValue() throws IOException, InterruptedException
    {   return lrr.getCurrentValue(); }
    @override
    public void initialize(InputSplit arg0, TaskAttemptContext arg1)
        throws IOException, InterruptedException
    {   lrr.initialize(arg0, arg1);
        fileName = ((FileSplit)arg0).getPath().getName();
    }
}
```



用户自定义输入输出格式

22

```
public class InvertedIndexer
{
    public static void main(String[] args)
    {
        try {
            Configuration conf = new Configuration();
            job = new Job(conf, "invert index");
            job.setJarByClass(InvertedIndexer.class);
            job.setInputFormatClass(FileNameLocInputFormat.class);
            job.setMapperClass(InvertedIndexMapper.class);
            job.setReducerClass(InvertedIndexReducer.class);
            job.setOutputKeyClass(Text.class);
            job.setOutputValueClass(Text.class);
            FileInputFormat.addInputPath(job, new Path(args[0]));
            FileOutputFormat.setOutputPath(job, new Path(args[1]));
            System.exit(job.waitForCompletion(true) ? 0 : 1);
        } catch (Exception e) {
            e.printStackTrace();
        }
    }
}
```



用户自定义输入输出格式

23

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class InvertedIndexMapper extends Mapper<Text, Text, Text, Text>
{   @Override
    protected void map(Text key, Text value, Context context)
                                throws IOException, InterruptedException
    // InputFormat: FileNameLocInputFormat
    // RecordReader: FileNameLocRecordReader
    // key: filename@lineoffset; value: line string
    {   Text word = new Text();
        StringTokenizer itr = new StringTokenizer(value.toString());
        for(; itr.hasMoreTokens(); )
        {   word.set(itr.nextToken());
            context.write(word, key);
        }
    }
}
```



用户自定义输入输出格式

24

- 数据输出格式（**OutputFormat**）用于描述**MapReduce**作业的数据输出规范。
- **MapReduce**框架依据数据输出格式完成输出规范检查（如检查输出目录是否存在）以及提供作业结果数据输出等功能。
- **TextOutputFormat**是系统缺省的数据输出格式。



用户自定义输入输出格式

25

□ Hadoop内置的OutputFormat和RecordWriter

OutputFormat:	Description
TextOutputFormat	Default; writes lines in "key \t value" form
SequenceFileOutputFormat	Writes binary files suitable for reading into subsequent MapReduce jobs
NullOutputFormat	Disregards its outputs

[DBOutputFormat](#), [FileOutputFormat](#), [FilterOutputFormat](#), [IndexUpdateOutputFormat](#), [LazyOutputFormat](#), [MapFileOutputFormat](#), [MultipleOutputFormat](#), [MultipleSequenceFileOutputFormat](#), [MultipleTextOutputFormat](#), [NullOutputFormat](#), [SequenceFileAsBinaryOutputFormat](#), [SequenceFileOutputFormat](#), [TextOutputFormat](#)



用户自定义输入输出格式

26

□ Hadoop内置的OutputFormat和RecordWriter

RecordWriter:	Description
LineRecordWriter	Default RecordWriter for TextOutputFormat writes lines in "key \t value" form

[DBOutputFormat.DBRecordWriter](#), [FilterOutputFormat.FilterRecordWriter](#),
[TextOutputFormat.LineRecordWriter](#)

与InputFormat和RecordReader类似，用户可以根据需要定制OutputFormat和RecordWriter



用户自定义输入输出格式

27

□ 划分多个输出文件集合

- 缺省情况下，**MapReduce**将产生包含一至多个文件的单个输出数据文件集合。但有时候作业可能需要输出多个文件结合。
 - 比如：在处理巨大的访问日志文件时，由于文件太大我们可能希望按每天的日期将访问日志记录输出为每天日期下的文件。在处理专利数据集时，我们希望根据不同国家，将每个国家的专利数据记录输出到不同国家的文件目录中。
 - Hadoop提供了**MultipleOutputFormat**类([org.apache.hadoop.mapred.lib.MultipleOutputFormat](#))来快速完成这一处理功能。在**Reduce**进行数据输出前，**MultipleOutputFormat**将调用一个内部方法以决定输出的文件名是什么。通常需要继承并实现**MultipleOutputFormat**的一个子类并实现其中的**generateFileNameForKeyValue()**方法以根据当前的键值对由程序产生并返回一个输出文件路径：

`protected String generateFileNameForKeyValue(K key, V value, String name)`



用户自定义输入输出格式

28

□ 划分多个输出文件集合

▣ 例如：将专利描述文件数据集按照国家进行多文件集合输出

"PATENT","GYEAR","GDATE","APPYEAR","**COUNTRY**", "POSTATE","ASSIGNEE", "ASSCODE","CLAIMS","NCLASS","CAT","SUBCAT","CMADE","CRECEIVE", "RATIOCIT","GENERAL","ORIGINAL","FWDAPLAG","BCKGTLAG","SELFCTUB", "SELFCTLB","SECDUPBD","SECDLWBD"

3070801,1963,1096,,"**BE**",",",1,,269,6,69,,1,,0,,,,,

3070802,1963,1096,,"**US**","TX",1,,2,6,63,,0,,,,,



用户自定义输入输出格式

29

□ 划分多个输出文件集合

▣ 例如：将专利描述文件数据集按照国家进行多文件集合输出

```
public static class MapClass extends Mapper<LongWritable, Text, NullWritable, Text>
{
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException
    {
        context.write (NullWritable.get(), value); }
        // NullWritable.get() 返回 singleton 单一实例
    }
    public static class SaveByCountryOutputFormat
        extends MultipleTextOutputFormat<NullWritable,Text>
    {
        protected String generateFileNameForKeyValue
            (NullWritable key, Text value, String filename) {
            String[] arr = value.toString().split(",", -1);
            String country = arr[4].substring(1,3);
            return country + "/" + filename;
        }
    }
}
```



用户自定义输入输出格式

30

□ 划分多个输出文件集合

```
public class MultiFileDemo {  
    public static void main(String[] args) throws Exception  
    {  
        Configuration conf = new Configuration();  
        Job job = new Job(conf, MultiFileDemo.class);  
        Path in = new Path(args[0]);  
        Path out = new Path(args[1]);  
        FileInputFormat.setInputPaths(job, in);  
        FileOutputFormat.setOutputPath(job, out);  
        job.setJobName("MultiFileDemo");  
        job.setMapperClass(MapClass.class);  
        job.setInputFormat(TextInputFormat.class);  
        job.setOutputFormat(SaveByCountryOutputFormat.class);  
        job.setOutputKeyClass(NullWritable.class);  
        job.setOutputValueClass(Text.class);  
        job.setNumReduceTasks(0);  
        Job.waitForCompletion(true); }  
}
```



用户自定义输入输出格式

31

□ 执行结果

```
ls output/
```

AD	BN	CS	GE	IN	LC	MT	PH	SV	VE
AE	BO	CU	GF	IQ	LI	MU	PK	SY	VG
AG	BR	CY	GH	IR	LK	MW	PL	SZ	VN
AI	BS	CZ	GL	IS	LR	MX	PT	TC	VU
AL	BY	DE	GN	IT	LT	MY	PY	TD	YE
AM	BZ	DK	GP	JM	LU	NC	RO	TH	YU
AN	CA	DO	GR	JO	LV	NF	RU	TN	ZA
AR	CC	DZ	GT	JP	LY	NG	SA	TR	ZM
AT	CD	EC	GY	KE	MA	NI	SD	TT	ZW
AU	CH	EE	HK	KG	MC	NL	SE	TW	
AW	CI	EG	HN	KN	MG	NO	SG	TZ	
AZ	CK	ES	HR	KP	MH	NZ	SI	UA	
BB	CL	ET	HT	KR	ML	OM	SK	UG	
BE	CM	FI	HU	KW	MM	PA	SM	US	
BG	CN	FO	ID	KY	MO	PE	SN	UY	
BH	CO	FR	IE	KZ	MQ	PF	SR	UZ	
BM	CR	GB	IL	LB	MR	PG	SU	VC	

```
ls output/AD
```

```
part-00003      part-00005      part-00006
```

```
head output/AD/part-00006
```

```
5765303,1998,14046,1996,"AD",",",1,12,42,5,59,11,1,0.4545,0,0,1,67.3636,,,,  
5785566,1998,14088,1996,"AD",",",1,9,441,6,69,3,0,1,,0.6667,,4.3333,,,,  
5894770,1999,14354,1997,"AD",",",1,,82,5,51,4,0,1,,0.625,,7.5,,,,
```



用户自定义输入输出格式

32

- **MultipleOutputFormat** ([org.apache.hadoop.mapred.lib.MultipleOutputFormat](#))
 - ▣ **MultipleOutputFormat**是Hadoop的**OutputFormat**的一个扩展，用于处理多个输出文件。
 - ▣ 它允许你为每个**Reducer**任务定义一个不同的**OutputFormat**。这意味着你可以为每个**Reducer**任务指定不同的输出目录和输出文件格式。
 - ▣ 这对于根据特定的数据或逻辑将数据分发到不同的输出目录非常有用，更适合在整个**MapReduce**作业级别控制多个输出文件的格式和目录。
- **MultipleOutputs** ([org.apache.hadoop.mapreduce.lib.output.MultipleOutputs](#))
 - ▣ **MultipleOutputs**是一个更高级别的API，用于在**Mapper**或**Reducer**内部根据某些条件将数据输出到多个文件或目录。
 - ▣ 它允许你在**Mapper**或**Reducer**内部为不同的输出文件指定不同的键值对，而不需要为每个**Reducer**任务创建不同的**OutputFormat**。
 - ▣ 这对于根据数据的属性或业务逻辑将数据分发到多个输出目录非常有用，而无需创建多个**Reducer**任务，更适合在**Mapper**或**Reducer**内部根据条件动态控制多个输出文件。



用户自定义输入输出格式

33

- `org.apache.hadoop.mapred.*` vs `org.apache.hadoop.mapreduce.*`
- 历史演变:
 - ▣ `mapred`包是Hadoop的早期版本中使用的包，用于实现MapReduce编程模型。
 - ▣ `mapreduce`包是Hadoop 0.20版本之后引入的，用于替代`mapred`包。这一改变是为了改进和优化MapReduce的性能以及提供更灵活的编程接口。
- 性能优化:
 - ▣ `mapreduce`包相对于`mapred`包进行了许多性能优化，包括更好的资源管理、任务调度、错误处理等方面的改进，以提高MapReduce作业的执行效率。
- API灵活性:
 - ▣ `mapreduce`包引入了一种新的API，使得编写MapReduce作业更加灵活和容易。这一API的设计更加现代化，与标准Java编程实践更加一致，因此编写和维护MapReduce作业变得更容易。



用户自定义Partitioner和Combiner

34

□ 定制Partitioner

- 程序员可以根据需要定制Partitioner来改变Map中间结果到Reduce节点的分区方式，并在Job中设置新的Partitioner

```
class NewPartitioner extends HashPartitioner<K,V>
{ // override the method
    getPartition(K key, V value, int numReduceTasks)
    { term = key.toString().split(",")[0]; //<term, docid>=>term
      super.getPartition(term, value, numReduceTasks);
    }
}
```

并在Job中设置新的Partitioner：

Job.setPartitionerClass(NewPartitioner)



定制Combiner

- 程序员可以根据需要定制**Combiner**来减少网络数据传输量，提高系统效率，并在**Job**中设置新的**Combiner**

例如，每年申请美国专利的国家数统计

Patent description data set “apat63_99.txt”

“PATENT”, “**GYEAR**”, “GDATE”, “APPYEAR”, “**COUNTRY**”, “POSTATE”, “ASSIGNEE”, “ASSCODE”, “CLAIMS”, “NCLASS”, “CAT”, “SUBCAT”, “CMADE”, “CRECEIVE”, “RATIOCIT”, “GENERAL”, “ORIGINAL”, “FWDAPLAG”, “BCKGTLAG”, “SELFCTUB”, “SELFCTLB”, “SECDUPBD”, “SECDLWBD”

3070801,1963,1096,,**BE**,"",,,1,,269,6,69,,1,,0,,,,,,,,,

3070802,1963,1096,"US","TX",1,2,6,63,0,,,,,,,,,,,,,

3070803,1963,1096,"US","IL",1,2,6,63,9,0.3704,,,,,,,,

3070804,1963,1096,, "US", "OH", 1,, 2,6,63,, 3,, 0.6667,,,,,,,,

3070805,1963,1096,"US","CA",1,2,6,63,1,0,,,,,,,,



用户自定义Partitioner和Combiner

36

□ 定制Combiner

每年申请美国专利的国家数统计

1. Map中用 $\langle \text{year}, \text{country} \rangle$ 作为key输出， $\text{Emit}(\langle \text{year}, \text{country} \rangle, 1)$

$(\langle 1963, \text{BE} \rangle, 1), (\langle 1963, \text{US} \rangle, 1), (\langle 1963, \text{US} \rangle, 1), \dots$

2. 实现一个定制的Partitioner，保证同一年份的数据划分到同一个Reduce节点

3. Reduce中对每一个 $(\langle \text{year}, \text{country} \rangle, [1, 1, 1, \dots])$ 输入，忽略后部的出现次数，仅考虑key部分： $\langle \text{year}, \text{country} \rangle$

问题：如每碰到一个 $\langle \text{year}, \text{country} \rangle$ ，即 $\text{emit}(\text{year}, 1)$ 有问题吗？

答案：有问题。因为可能会有从不同Map节点发来的同样的 $\langle \text{year}, \text{country} \rangle$ ，因此会出现对同一国家的重复计数

解决办法：在Reduce中仅计数同一年份下不同的国家个数

问题：Map结果 $(\langle \text{year}, \text{country} \rangle, [1, 1, 1, \dots])$ 数据通信量较大

解决办法：实现一个Combiner将 $[1, 1, 1, \dots]$ 合并为1



用户自定义Partitioner和Combiner

37

□ 定制Combiner

每年申请美国专利的国家数统计

public static class **NewCombiner** extends Reducer

< Text, IntWritable, Text, IntWritable >

{

public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException, InterruptedException

{ // 忽略(<year, country>, [1,1,1,...])后部大量重复的[1,1,1,...],

// 归并为<year, country>的1次出现

context.write(key, new IntWritable(1));

} // 输出key: <year, country>; value: 1

}



迭代MapReduce计算

□ 基本问题

- 一些求解计算需要用迭代方法求得逼近结果（求解计算必须是收敛性的）。当用MapReduce进行这样的问题求解时，运行一趟MapReduce过程无法完成整个求解过程，因此，需要采用迭代方法循环运行该MapReduce过程，直到达到一个逼近结果。

□ 例如：页面排序算法PageRank

- 随机浏览模型：假设一位上网者随机地浏览一些网页
 - 有可能从当前网页点击一个链接继续浏览（概率为d）；
 - 有可能随机跳转到其它N个网页中的任一个（概率为1-d）。
- 每个网页的PageRank值可以看成该网页被随机浏览的概率：

$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)}$$

$L(p_i)$ 为网页 p_i 上的超链个数



迭代MapReduce计算

39

□ 页面排序算法PageRank

$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)}$$

- 问题是在求解 $PR(p_i)$ 时，需要递归调用 $PR(p_j)$ ，而 $PR(p_j)$ 本身也是待求解的。因此，我们只能先给每个网页赋一个假定的 PR 值，如0.5。但这样求出的 $PR(p_i)$ 肯定不准确。然而，当用求出的 PR 值反复进行迭代计算时，会越来越趋近于最终的准确结果。
- 因此，需要用迭代方法循环运行MapReduce过程，直至第 n 次迭代后的结果与第 $n-1$ 次的结果小于某个指定的阈值时结束，或者通过经验控制循环固定的次数。



多趟MapReduce的处理

40

```
public class PageRankDriver {  
    private static int times = 10;  
    public static void main(String args[]) throws Exception{  
        String[] forGB = {"", args[1]+"/Data0"};  
        forGB[0] = args[0];  
        GraphBuilder.main(forGB);  
        String[] forltr = {"Data","Data"};  
        for (int i=0; i<times; i++) {  
            forltr[0] = args[1]+"/Data"+(i);  
            forltr[1] = args[1]+"/Data"+(i+1);  
            PageRankIter.main(forltr);  
        }  
        String[] forRV = {args[1]+"/Data"+times, args[1]+"/FinalRank"};  
        PageRankViewer.main(forRV);  
    }  
}
```

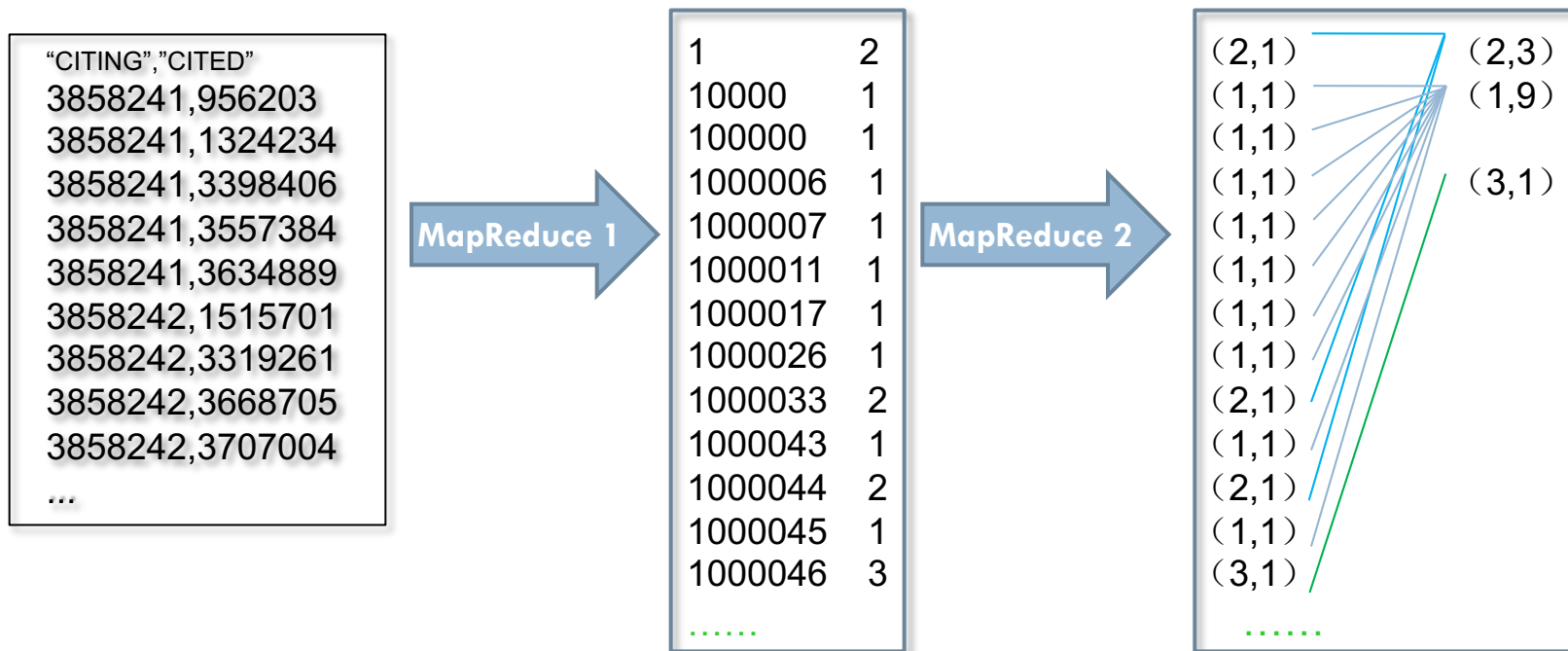



链式MapReduce任务

41

□ 基本问题

- 一些复杂任务难以用一趟MapReduce处理过程来完成，需要将其分为多趟简单些的MapReduce子任务完成。如：
- 专利文献引用直方图统计，需要先进行被引次数统计，然后在被引次数上再进行被引直方图统计





链式MapReduce任务

□ MapReduce子任务的顺序化执行

- 多个MapReduce子任务可以用手工逐一执行，但更方便的做法是将这些子任务穿起来，前面MapReduce任务的输出作为后面MapReduce的输入，自动地完成顺序化的执行，如：

mapreduce-1 → mapreduce-2 → mapreduce-3 → ...

MapReduce作业控制执行代码：

```
Configuration jobconf = new Configuration();  
job = new Job(jobconf, "invert index");  
job.setJarByClass(InvertedIndexer.class);  
.....  
FileInputFormat.addInputPath(job, new Path(args[0]));  
FileOutputFormat.setOutputPath(job, new Path(args[1]));  
job.waitForCompletion(true);
```

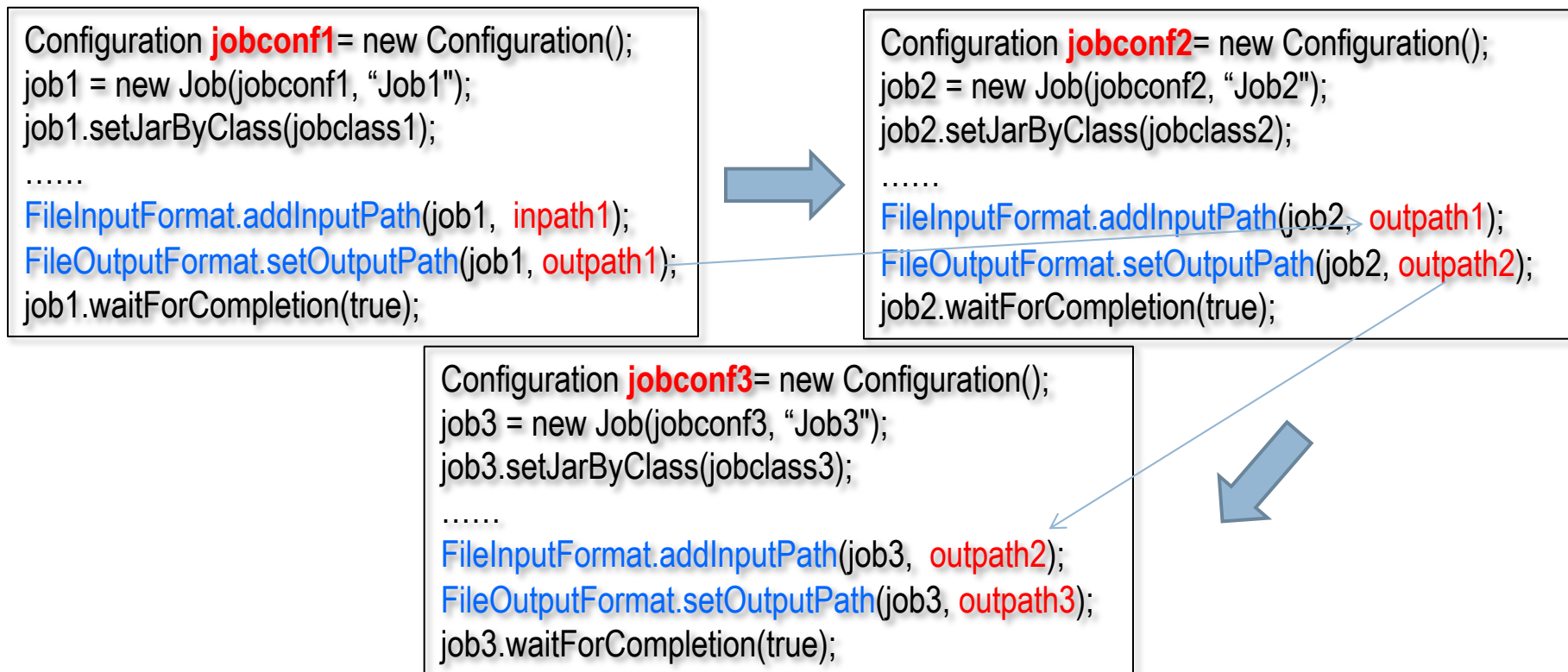


链式MapReduce任务

43

□ MapReduce子任务的顺序化执行

- 同样，链式MapReduce中的每个子任务需要穿件独立的**jobconf**，并按照前后子任务间的输入输出关系设置输入输出路径，而任务完成后所有中间过程的输出结果路径都可以删除掉。





链式MapReduce任务

44

- MapReduce前处理和后处理步骤的链式执行
 - ▣ 一个MapReduce作业可能会有一些前处理和后处理步骤，比如，文档倒排索引处理前需要一个去除Stop-word的前处理，倒排索引处理后需要一个变形词后处理步骤(making, made → make)。将这些前后处理步骤实现为单独的MapReduce任务可以达到目的，但将增加很多I/O操作，因而效率不高。
 - ▣ 一个办法是在核心的Map和Reduce过程之外，把这些前后处理步骤实现为一些辅助的Map和Reduce过程，将这些辅助的Map和Reduce过程与核心的Map和Reduce过程合并为一个过程链，从而完成整个作业。



链式MapReduce任务

45

□ MapReduce前处理和后处理步骤的链式执行

- ▣ Hadoop提供了链式Mapper(org.apache.hadoop.mapred.lib.ChainMapper)和链式Reducer(org.apache.hadoop.mapred.lib.ChainReducer)来完成这种处理
- ▣ ChainMapper和ChainReducer分别提供了addMapper方法加入一系列Mapper:
ChainMapper.addMapper (.....) ; **ChainReducer.addMapper** (.....)

public static void **addMapper**

(JobConf job,

//主作业

Class<? extends Mapper> mclass,

//待加入的map class

Class<?> inputKeyClass,

//待加入的map输入键class

Class<?> inputValueClass,

//待加入的map输入键值class

Class<?> outputKeyClass,

//待加入的map输出键class

Class<?> outputValueClass,

//待加入的map输出键值class

boolean byValue

//指示键/值是否应按值传递给链中的下一个Mapper

JobConf mapperConf

// 待加入的map的conf

) throws IOException



链式MapReduce任务

46

□ MapReduce前处理和后处理步骤的链式执行

- ▣ 设有一个完整的MapReduce作业，由Map1 , Map2 , Reduce, Map3, Map4构成。

```
Configuration conf = new Configuration();
Job job = new Job(conf);
job.setJobName("ChainJob");
job.setInputFormat(TextInputFormat.class);
job.setOutputFormat(TextOutputFormat.class);
FileInputFormat.setInputPaths(job, in);
FileOutputFormat.setOutputPath(job, out);
JobConf map1Conf = new JobConf(false);
ChainMapper.addMapper(job, Map1.class, LongWritable.class, Text.class,
    Text.class, Text.class, true, map1Conf);
JobConf map2Conf = new JobConf(false);
ChainMapper.addMapper(job, Map2.class, Text.class, Text.class, LongWritable.class,
    Text.class, true, map2Conf);
```



链式MapReduce任务

47

□ MapReduce前处理和后处理步骤的链式执行

```
JobConf reduceConf = new JobConf(false);
ChainReducer.setReducer(job, Reduce.class, LongWritable.class, Text.class,
                        Text.class, Text.class, true, reduceConf);
JobConf map3Conf = new JobConf(false);
ChainReducer.addMapper(job, Map3.class, Text.class, Text.class,
                       LongWritable.class, Text.class, true, map3Conf);
JobConf map4Conf = new JobConf(false);
ChainReducer.addMapper(job, Map4.class, LongWritable.class, Text.class,
                       LongWritable.class, Text.class, true, map4Conf);
JobClient.runJob(job);
```

`ChainReducer.setReducer()`方法必须在`ChainReducer`最开始的地方使用，其后方可加入后续的辅助处理`Mapper`；另一个需要注意的问题是，这些链式`Mapper`和`Reducer`之间传递的键值对数据类型必须保持前后一致。



全局参数/数据文件的传递

□ 全局作业参数的传递

- 为了能让用户灵活设置某些作业参数，避免作业参数在程序中的硬编码，一个 **MapReduce** 计算任务可能需要在执行时从命令行输入这些作业参数，并将这些参数传递给各个计算节点。
- 比如，对两个关系进行自然连接时程序用硬编码方式指定第一个数据列为 **join** 的主键。但为了要实现一个具有一定通用性的程序，可以任意指定一个列为 **join** 主键的话，就需要在程序运行时在命令行中指定 **join** 主键所在的数据列。然后该输入参数可以作为一个属性保存在 **Configuration** 对象中，并允许 **Map** 和 **Reduce** 节点从 **Configuration** 对象中获取和使用该属性值。



全局参数/数据文件的传递

49

□ 全局作业参数的传递

Configuration类专门提供以下用于保存和获取属性的方法：

- `public void set(String name, String value)` //设置字符串属性
- `public String get(String name)` // 读取字符串属性
- `public String get(String name, String defaultValue)` // 读取字符串属性
- `public void setBoolean(String name, boolean value)` //设置布尔属性
- `public boolean getBoolean(String name, boolean defaultValue)` //读取布尔属性
- `public void setInt(String name, int value)` //设置整数属性
- `public int getInt(String name, int defaultValue)` // 读取整数属性
- `public void setLong(String name, long value)` //设置长整数属性
- `public long getLong(String name, long defaultValue)` // 读取长整数属性
- `public void setFloat(String name, float value)` //设置浮点数属性
- `public float getFloat(String name, float defaultValue)` //读取浮点数属性
- `public void setStrings(String name, String... values)` //设置一组字符串属性
- `public String[] getStrings(String name, String... defaultValue)` //读取一组字符串属性

- 需要说明的是，**setStrings**方法将把一组字符串转换为用“,”隔开的一个长字符串，然后**getStrings**时自动再根据“,” **split**成一组字符串，因此，在该组中的每个字符串都不能包含“,”，否则会出错。



全局参数/数据文件的传递

□ 全局作业参数的传递

例：专利文献数据集Join时主键所在数据列参数的设置

```
Configuration conf = new Configuration();
```

```
Job job = new Job(conf, "naturalJoinJob");
```

```
...
```

```
// 将第三个输入参数设置为Join Key属性
```

```
job.getConfiguration().setInt("JoinKeyColIdx", Integer.parseInt(args[2]));
```

```
.....
```

```
job.waitForCompletion(true);
```



全局参数/数据文件的传递

51

□ 全局作业参数的传递

- ▣ 在mapper类的初始化方法setup()中从configuration对象中读出属性

```
public static class MapClass extends Mapper <Text, Text, Text, Text>
{
    int join_key_col_idx;
    protected void setup(Mapper.Context context)
    {
        Configuration jobconf = context.getConfiguration();
        join_key_col_idx = jobconf.getInt("JoinKeyColIdx", -1); // 无值时置为-1
    }
    protected void map(Text key, Text value, Context context)
        throws IOException, InterruptedException
    { //使用join_key_col_idx完成数据处理;
        .....
    }
}
```



全局参数/数据文件的传递

□ 全局作业参数的传递

- 同样需要时在reducer类的初始化方法setup()中从configuration对象中读出属性

```
public static class ReduceClass extends Reducer <Text, Text, Text, Text>
{
    int join_key_col_idx;
    protected void setup(Mapper.Context context)
    {
        Configuration jobconf = context.getConfiguration();
        join_key_col_idx = jobconf.getInt("JoinKeyColIdx", -1); // 无值时置为-1
    }
    protected void reduce(Text key, Text value, Context context)
        throws IOException, InterruptedException
    { //使用join_key_col_idx完成数据处理;
        .....
    }
}
```



全局参数/数据文件的传递

□ 全局数据文件的传递

- 有时候一个MapReduce作业可能会使用一些较小的并且需要复制到各个节点的数据文件。为此，可以使用DistributedCache文件传递机制，先将这些文件传送到DistributedCache中，然后各个节点从DistributedCache中将这些文件并复制到本地的文件系统中使用。具体使用时，为提供访问速度，可将这些较小的文件数据读入内存。
- Job类中：`public void addCacheFile(URI uri)`：将一个文件放到distributed cache file中
- Mapper或Reducer的context类中：
`public Path[] getLocalCacheFiles()`：获取设置在distributed cache files中的文件路径，以便能将这些文件读入到每个节点内存中



全局参数/数据文件的传递

□ 全局数据文件的传递

- ▣ 在作业Configuration时将文件存入Distributed Cache:

.....

```
Configuration conf = getConf();
```

```
Job job = Job.getInstance(conf, "word count 2.0");;
```

```
// 将命令行参数中的reamingsArgs列表里指定的文件依次放置到distributed cache file中
```

```
List<String> otherArgs = new ArrayList<String>();
```

```
for (int i = 0; i < remainingArgs.length; ++i) {
```

```
    if ("-skip".equals(remainingArgs[i])) {
```

```
        job.addCacheFile(new Path(remainingArgs[++i]).toUri());
```

```
        job.getConfiguration().setBoolean("wordcount.skip.patterns", true);
```

```
    } else {
```

```
        otherArgs.add(remainingArgs[i]);
```

```
    }
```

```
}
```



全局参数/数据文件的传递

□ 全局数据文件的传递

```
public static class MapClass extends Mapper<Text, Text, Text, Text> {  
    public void setup(Mapper.Context context) throws IOException, InterruptedException  
    {  
        conf = context.getConfiguration();  
        caseSensitive = conf.getBoolean("wordcount.case.sensitive", true);  
        if (conf.getBoolean("wordcount.skip.patterns", true)) {  
            URI[] patternsURIs = Job.getInstance(conf).getCacheFiles();  
            if (patternsURIs != null) {  
                for (URI patternsURI : patternsURIs) {  
                    Path patternsPath = new Path(patternsURI.getPath());  
                    String patternsFileName = patternsPath.getName().toString();  
                    parseSkipFile(patternsFileName);  
                }  
            }  
        }  
    }  
}
```



其它处理技术

56

□ 查询任务相关信息

- 可以通过**Configuration**对象，使用预定义的属性名称查询计算作业相关的信息。

Property	Type	Description
<code>mapred.job.id</code>	String	The job ID
<code>mapred.jar</code>	String	The jar location in job directory
<code>job.local.dir</code>	String	The job's local scratch space
<code>mapred.tip.id</code>	String	The task ID
<code>mapred.task.id</code>	String	The task attempt ID
<code>mapred.task.is.map</code>	boolean	Flag denoting whether this is a map task
<code>mapred.task.partition</code>	int	The ID of the task within the job
<code>map.input.file</code>	String	The file path that the mapper is reading from
<code>map.input.start</code>	long	The offset into the file of the start of the current mapper's input split
<code>map.input.length</code>	long	The number of bytes in the current mapper's input split
<code>mapred.work.output.dir</code>	String	The task's working (i.e., temporary) output directory



其它处理技术

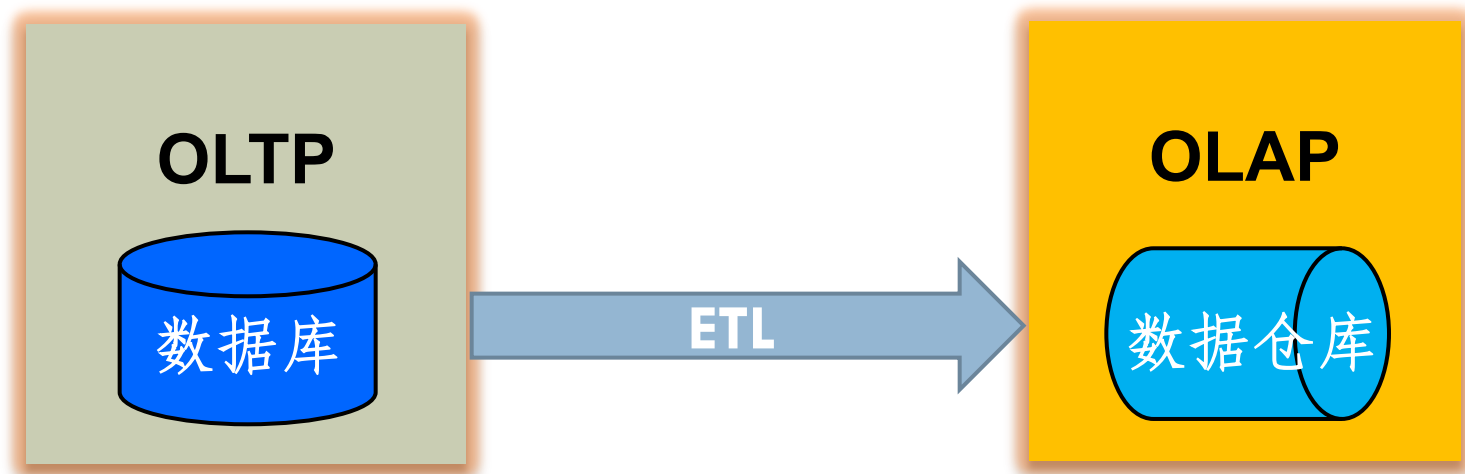
57

□ 输入输出到关系数据库

- **MapReduce**用于处理存储在**HDFS**中的大规模数据，但现实环境中有很多应用数据保存在关系数据库中，因此，**Hadoop**提供了访问关系数据库的能力以便在需要时能用**MapReduce**技术处理关系数据库中的数据。这在基于**MapReduce**进行联机数据分析处理时尤为有用。
- **OLTP (online transaction processing)**
 - 联机事务处理：主要是关系数据库应用系统中前台常规的各种事务处理
- **OLAP (online analytical processing)**
 - 联机分析处理：主要是进行基于数据仓库的后台数据分析和挖掘，提供优化的客户服务和运营决策支持
- **OLTP**与**OLAP**一般采用分离的数据库，前者数据库负责大量的常规的事务处理，后者用数据仓库应对大量的数据分析处理负载。

其它处理技术

企业数据库应用系统



Extract: 从OLTP数据库中抽取事务数据

Transform: 转换为数据仓库中的数据格式

Load: 装载到数据仓库中

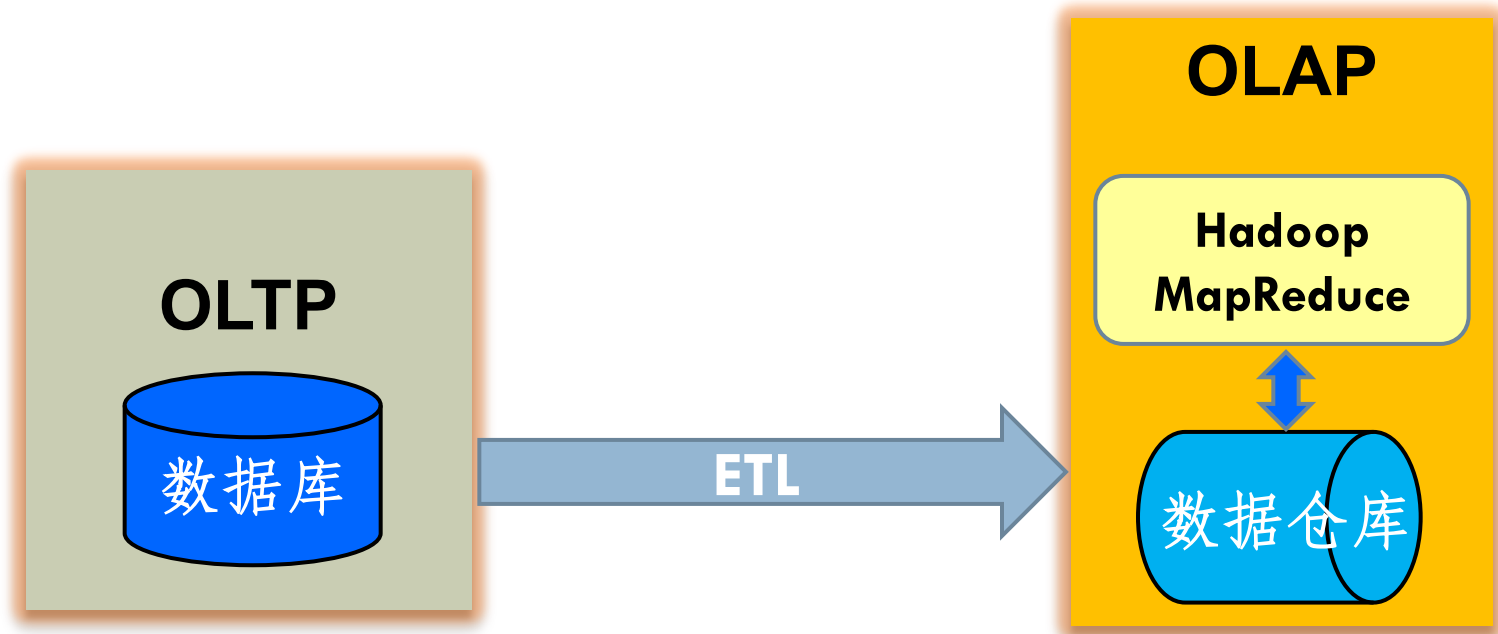
问题: **OLAP**端基于关系数据库的数据仓库解决方案, 在数据量巨大的情况下, 复杂数据分析和挖掘处理的负载很大, 速度性能跟不上



其它处理技术

59

企业数据库应用系统



解决方案：提供基于MapReduce大规模数据并行处理的OLAP！

问题：如何从MapReduce访问关系数据库？



其它处理技术

60

□ 输入输出到关系数据库

▣ 从数据库中输入数据

- Hadoop提供了相应的从关系库查询和读取数据的接口(`org.apache.hadoop.mapred.lib.db.*`)
- `DBInputFormat`: 提供从数据库读取数据的格式
- `DBWritable`: 提供读取数据记录的接口

- ▣ 虽然Hadoop允许用以上接口从数据库中直接读取数据记录作为MapReduce的输入，但处理效率不理想，因此，仅适合读取小量数据记录的计算和应用，不适合OLAP数据仓库大量数据的读取处理。
- ▣ 读取大量数据记录一个更好的解决办法是，用数据库中的Dump工具将大量待分析数据输出为文本数据文件，并上载到HDFS中进行处理。



其它处理技术

61

□ 输入输出到关系数据库

▣ 向数据库中输出计算结果

- 基于数据仓库的数据分析和挖掘输出结果的数据量一般不会太大，因而可能适合于直接向数据库写入。
Hadoop提供了相应的向关系库直接输出计算结果的编程接口
- `DBOutputFormat`：提供向数据库输出数据的格式
- `DBConfiguration`：提供数据库配置和创建连接的接口

▣ 创建数据库连接

- `DBConfiguration` 类中提供了一个静态方法创建数据库连接：
`public static void configureDB(Job job, String driverClass, String dbUrl, String userName, String passwd)`

▣ 指定写入的数据表和字段

- `DBOutputFormat`中提供了一个静态方法完成这一工作：
`public static void setOutput(Job job, String tableName, String... fieldNames)`



其它处理技术

62

□ 输入输出到关系数据库

▣ 向数据库中输出计算结果

▣ Configuration示例

```
Configuration conf = new Configuration();  
Job job = new Job(conf, JobClass.class);  
job.setOutputFormat(DBOutputFormat.class);  
DBConfiguration.configureDB(job, "com.mysql.jdbc.Driver",  
                             "jdbc:mysql://db.host.com/mydb", "username", "password")  
DBOutputFormat.setOutput(job, "Events", "event_id", "time"); // 向Events表输出event_id和time字段
```



其它处理技术

63

□ 输入输出到关系数据库

▣ 向数据库中输出计算结果

■ 实现DBWritable

■ 为了实际完成向数据库中数据写入，程序员要实现DBWritable：

```
public class EventsDBWritable implements Writable, DBWritable
{
    private int id;
    private long timestamp;
    public void write(DataOutput out) throws IOException
    {
        out.writeInt(id);  out.writeLong(timestamp);  }
    public void readFields(DataInput in) throws IOException
    {
        id = in.readInt();  timestamp = in.readLong();  }
    public void write(PreparedStatement statement) throws SQLException
    {
        statement.setInt(1, id); statement.setLong(2, timestamp);  }
    public void readFields(ResultSet resultSet) throws SQLException
    {
        id = resultSet.getInt(1);          timestamp = resultSet.getLong(2);}
    // 除非使用DBInputFormat直接从数据库输入数据,否则readFields方法不会被调用
}
```

THANK YOU



南京大學
NANJING UNIVERSITY

南京大学计算机软件研究所
Institute of Computer Software, Nanjing University