# 分散処理アプリ演習 第2回 MapReduceアプリケーションの概要

(株)NTTデータ

## 講義内容

- 1. WordCountプログラムの概要
- 2. MapReduceプログラムのコンポーネント
- 3. 演習:WordCountプログラムの実行
- 4. Streaming API
- 5. レポート課題







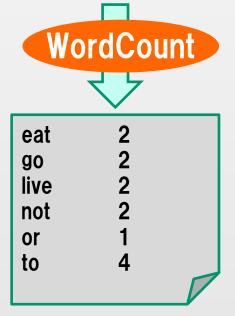
# 1. WordCountプログラムの概要

### MapReduceプログラム事例:WordCount

- 4 CHERS EDUCATION PROPERTY OF THE NGINEERS OF
- 入力のテキストファイル中にある各単語の出現回数をカウントする
  - 本講義では、入力ファイルとしてWikipedia英語版のabstract版のダンプファイル (xmlファイル)の一部とする
    - さらにこのファイルの中で、<abstract>タグで囲まれた文章のみを対象とする

<abstract>To go, or not to go. Eat to live, not live to eat.</abstract>

入力ファイル



出力ファイル

# 5 CHARERS EDUCATION OF THE NGINEERS OF THE NGI

# Map関数/Reduce関数の設計

#### 擬似コード

#### Mapper

```
map(offset, line)
  words = line.split
  for w in words:
      emit(w, 1)
```

#### Reducer

```
reduce(key, values)
    sum = 0
    for v in values:
        sum += v
    emit(key, sum)
```



# 2. MapReduceプログラムのコンポー ネント

# MapReduceプログラムのコンポーネント

AND DE TWARE WITH THE PROPERTY OF A OL HO J NOT HOLD THE PROPERTY OF A OL HO J NOT HOLD THE PROPERTY OF A OL HO J NOT HOLD THE PROPERTY OF A OL HO J NOT HOLD THE PROPERTY OF A OL HO J NOT HOLD THE PROPERTY OF A OLD THE P

- MapReduceプログラムは通常3つの部分から構成される
  - Mapper
  - Reducer
  - ドライバー

# B GINEERS EDUCATION PROPERTY OF THE PROPERTY O

### コードの全体構成

#### ■ コードの概要

```
import •••
                                                                 Importステートメント
• • •
public class WordCount {
 public static class Map extends Mapper < Long Writable, Text, Text,
  IntWritable> {
                                                                         Mapper
 public static class Reduce extends Reducer<Text, IntWritable, Text,
  IntWritable> {
                                                                         Reducer
 public static void main(String[] args) throws Exception {
                                                                         ドライバー
```

## importステートメント

#### ■ コードの内容

```
import java.io.IOException;
import java.util.*;
import java.util.regex.*;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

(赤字部分は本講義用に独自に追加)



# O GINEERS EDUCATION PROPERTY OF THE NGINEERS O

## **Mapper**

#### ■ コードの内容

入力のKey、Value、出 力のKey、Valueの型 を指定

```
public static class Map extends Mapper < Long Writable, Text, Text,
 IntWritable> {
                                                                書き込みの2つの
  private final static IntWritable one = new IntWritable(1);
  private Text word = new Text();
  public void map (LongWritable key, Text value, Context context)
    throws IOException, InterruptedException {
                                                              <abstract>タグで囲ま
       String line = value.toString();
                                                              れた文字列を抽出
      Matcher m =
         Pattern.compile("<abstract>(.*)</abstract>").matcher(line);
      line = (m.matches()) ? m.group(1) : "";
       StringTokenizer tokenizer = new StringTokenizer(line);
      while (tokenizer.hasMoreTokens()) {
                                                        文字列を単語に分解
          word.set(tokenizer.nextToken());
                                                        し、(単語、1)を出力
           context.write(word, one);
```

#### (赤字部分は本講義用に独自に追加)

# HANNERS EDUCATION PROPERTY OF THE MISTINGERS OF

#### Reducer

#### ■ コードの内容

入力のKey、Value、出 力のKey、Valueの型 を指定

```
public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {

   public void reduce(Text key, Iterable<IntWritable> values,
        Context context) throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        context.write(key, new IntWritable(sum));
        }
}
```

# 2 SINEERS EDUCATION PROPERTIES OF THARE PROPERTIES OF THE PROPERTI

### ドライバー

#### ■ コードの内容

```
public static void main(String[] args) throws Exception {
                                                           Jobに名前を付与
  Configuration conf = new Configuration();
                                                           し生成
  Job job = new Job(conf, "wordcount");
                                                           Mapフェーズで出
                                            main関数のクラス
   job.setJarByClass(WordCount.class);
                                            を指定
                                                           力されるKey、
                                                           Valueそれぞれの
   job.setOutputKeyClass(Text.class);
                                                           型情報を指定
  job.setOutputValueClass(IntWritable.class);
                                                           Mapper, Reducer
   job.setMapperClass(Map.class);
                                                           のクラスを指定
   job.setReducerClass(Reduce.class);
   job.setInputFormatClass(TextInputFormat.class);
                                                           入力と出力の
   job.setOutputFormatClass(TextOutputFormat.class);
                                                           フォーマット情報を
                                                           指定
  FileInputFormat.addInputPath(job, new Path(args[0]));
  FileOutputFormat.setOutputPath(job, new Path(args[1]));
                                                           入力ディレクトリと
                                                           出力ディレクトリを
   job.waitForCompletion(true);
                                                           指定
```



# 3. 演習:WordCountプログラムの実行

# 演習:WordCountプログラムの実行

- 1. クライアントの仮想マシンにログインする。
- 2. 第1回のHDFSの演習で行ったのと同様に、ローカルファイルシステムにある 入力データファイルをHDFS上へコピーする。(第1回に続けて行う場合は不要)

```
$ hadoop fs -put /root/hadoop_exercise/01/enwiki /enwiki
$ hadoop fs -lsr /enwiki
```

■ 以下の2つのデータファイルが存在

■ enwiki-long1.xml : 150MB程度のファイル

■ enwiki-short1.xml : 1MB程度のファイル

■ 3. WordCountのJavaプログラムファイルがローカルファイルシステムにあることを確認する。また、内容を確認する。

```
$ cd /root/hadoop_exercise/02/wordcount
$ ls -1
```

\$ less WordCount.java



# SOLUTION PROPERTY OF THE NGINEERS OF THE NGINE

# 演習:WordCountプログラムの実行(続き)

■ 4. Javaプログラムをコンパイルする。コンパイル後、Classファイルが生成された ことを確認する。

- \$ mkdir wordcount classes
- \$ javac -classpath /usr/lib/hadoop-0.20/hadoop-core.jar ¥
   -d wordcount classes WordCount.java
- \$ ls -l wordcount classes

#### ■ 以下のClassファイルが生成される

- WordCount.class:ドライバークラス
- WordCount\$Map.class: Mapクラス
- WordCount\$Reduce.class: Reduceクラス
- 5. コンパイルしたプログラムをJARファイルに入れる。

\$ jar cvf wordcount.jar -C wordcount classes .

### 演習:WordCountプログラムの実行(続き)



■ 6. JARファイルを使ってMapReduceジョブをサブミットする。入力ファイルは1つ を指定する。

```
$ hadoop jar wordcount.jar WordCount ¥
   /enwiki/enwiki-short1.xml /wordcount01
```

- hadoop jar wordcount.jar: wordcount.jarファイルを使ってMapReduceジョブをサブミット
- WordCount: WordCountクラスのmainメソッドを実行
- /enwiki /enwiki-short1.xml:HDFS上の入力ファイル
- /wordcount 01: HDFS上の出力ディレクトリ
- MapとReduceの処理の進行状況が以下のように標準出力に表示されるのを確認する(mapが完了した後、reduceが開始される)

```
yy/mm/dd hh:mm:ss INFO mapred.JobClient: map 0% reduce 0% yy/mm/dd hh:mm:ss INFO mapred.JobClient: map 100% reduce 0% yy/mm/dd hh:mm:ss INFO mapred.JobClient: map 100% reduce 33% yy/mm/dd hh:mm:ss INFO mapred.JobClient: map 100% reduce 100%
```

## 演習:WordCountプログラムの実行(続き)

■ 7. 出力ディレクトリを確認する。

```
$ hadoop fs -lsr /wordcount01
```

- \_SUCCESS:ジョブが成功して完了したことを示すファイル
- \_logs: ログディレクトリ
- part-r-00000: 出力ファイル(Reducerが1つの場合)
- 8. 出力ファイルの内容を確認する。

```
$ hadoop fs -cat /wordcount01/part-r-00000 | less
```

- 各行は、単語 [空白] 出現回数
- 単語は昇順にソートされている
- 9. 出力ファイルの内容を出現回数(2フィールド目)の降順にソートして表示する。 出現回数の最も多い単語と、その出現回数を確認する。

```
$ hadoop fs -cat /wordcount01/part-r-00000 $ | sort -k2 -n -r | less
```



# 演習:WordCountプログラムの実行(続き)



■ 10. JARファイルを使ってMapReduceジョブをサブミットする。入力ファイルは6とは別のファイルを指定する。

```
$ hadoop jar wordcount.jar WordCount ¥
    /enwiki/enwiki-long1.xml /wordcount02
```

■ 11. 出力ディレクトリを確認する。

```
$ hadoop fs -lsr /wordcount02
```

■ 12. 出力ファイルの内容を出現回数(2フィールド目)の降順にソートして表示する。出現回数の最も多い単語と、その出現回数を確認する。

# HARE EDUCATION DE LA COLUMNA D

## 演習:WordCountプログラムの実行(続き)

■ 13. JARファイルを使ってMapReduceジョブをサブミットする。入力はディレクトリ 全体(配下の全てのファイルが対象となる)を指定する。

\$ hadoop jar wordcount.jar WordCount /enwiki /wordcount03

■ 14. 出力ディレクトリを確認する。

\$ hadoop fs -lsr /wordcount03

■ 15. 出力ファイルの内容を出現回数(2フィールド目)の降順にソートして表示する。出現回数の最も多い単語と、その出現回数を確認する。

■ 出現回数が、9と12で確認した各出現回数の合計(各入力ファイルの結果の合計) となっていることを確認する



# 4. Streaming API

# **Streaming API**

- TOP SOLUTION OF THE PROPERTY O
- MapReduceを利用したいが、Java以外の言語を得意とする開発者も数多く存在している
  - Perl
  - Python
  - その他
- HadoopのStreaming APIは、Java以外の言語を使用してMapReduceプログラムを書くことができる
  - 標準入力から読み込み、標準出力に書き出すことができる言語であればよい

# Streaming APIの仕組み

- MapperとReducerのプログラムをそれぞれ記述する
  - 標準入力から入力を受け取る
  - 標準出力に出力を書き込む
- 入力フォーマットは、Key [TAB] Value
- 出力フォーマットは、Key [TAB] Value [改行]
- TAB以外の区切り文字も指定可能



23

EDUCATION PROGRAM FOR TOP SOFTWARE ENGINEERS

# Streamingジョブの実行

#### ■ Streamingジョブの実行方法

```
$ hadoop jar ¥
/usr/lib/hadoop-0.20/contrib/streaming/hadoop-streaming*.jar ¥
-input 入力ディレクトリ ¥
-output 出力ディレクトリ ¥
-file Mapperプログラムのパス ¥
-file Reducerプログラムのパス ¥
-mapper Mapperプログラムのファイル名 ¥
-reducer Reducerプログラムのファイル名
```

#### ■ 例

```
$ hadoop jar ¥
/usr/lib/hadoop-0.20/contrib/streaming/hadoop-streaming*.jar ¥
-input /enwiki ¥
-output /wordcount04 ¥
-file mapper.py ¥
-file reducer.py ¥
-mapper mapper.py ¥
-reducer reducer.py
```

### まとめ

#### 本講義で学んだ内容

- **WordCountプログラムの概要** 
  - 文章中の単語の出現回数をカウント、入力データとしてWikipediaのダンプファイルを 利用
- MapReduceプログラムのコンポーネント
  - Mapper、Reducer、ドライバー
- 演習:WordCountプログラムの実行
  - hadoop fsコマンドで入力データファイルをHDFSへコピー、javacコマンドでプログラムをコンパイル、hadoop jarコマンドでジョブをサブミット(プログラムの実行)、hadoop fsコマンドで出力ファイルを確認
- Streaming API
  - Java以外の言語を使用してMapReduceプログラムを書くことが可能





# 5. レポート課題

### レポート課題: Streaming APIを用いたWordCountプログラムの実装



- Java以外で自分の得意なプログラミング言語を用いて、本講義と同じ仕様の WordCountプログラムを実装する
  - 入力ファイルは、本講義で用いたWikipedia英語版のabstract版のダンプファイルと する
  - このファイルの中で、<abstract>タグで囲まれた文章のみを対象とする
  - 本講座の演習環境に実行環境がインストール済みのプログラミング言語を選択すること
- 提出物
  - ソースコード(コメントを含めること)
  - 実行コマンド
  - 実行結果(カウント数の多いものから降順にソートして上位の一部のみを表示したもの)
  - 未完成であっても部分点を出すので提出すること
- 提出方法:上記のファイル群をtar.gzまたはzipで1つにまとめてLMSに提出
- 提出期限:当日周知
- 質問先:当日周知