



Распределенный кеш (Distributed Cache)

Драль Алексей, study@bigdatateam.org

CEO at BigData Team, <https://bigdatateam.org>

<https://www.facebook.com/bigdatateam>



Фильтрация по словарю



`<article_id> <tab> <article_content>`

↑
key

↑
value



...
James 2284
Thomas 1941
...





mapper.py

```
import re
import sys

def read_vocabulary(file_path):
    return set(line.strip() for line in open(file_path))

vocabulary = read_vocabulary("vocabulary.txt")

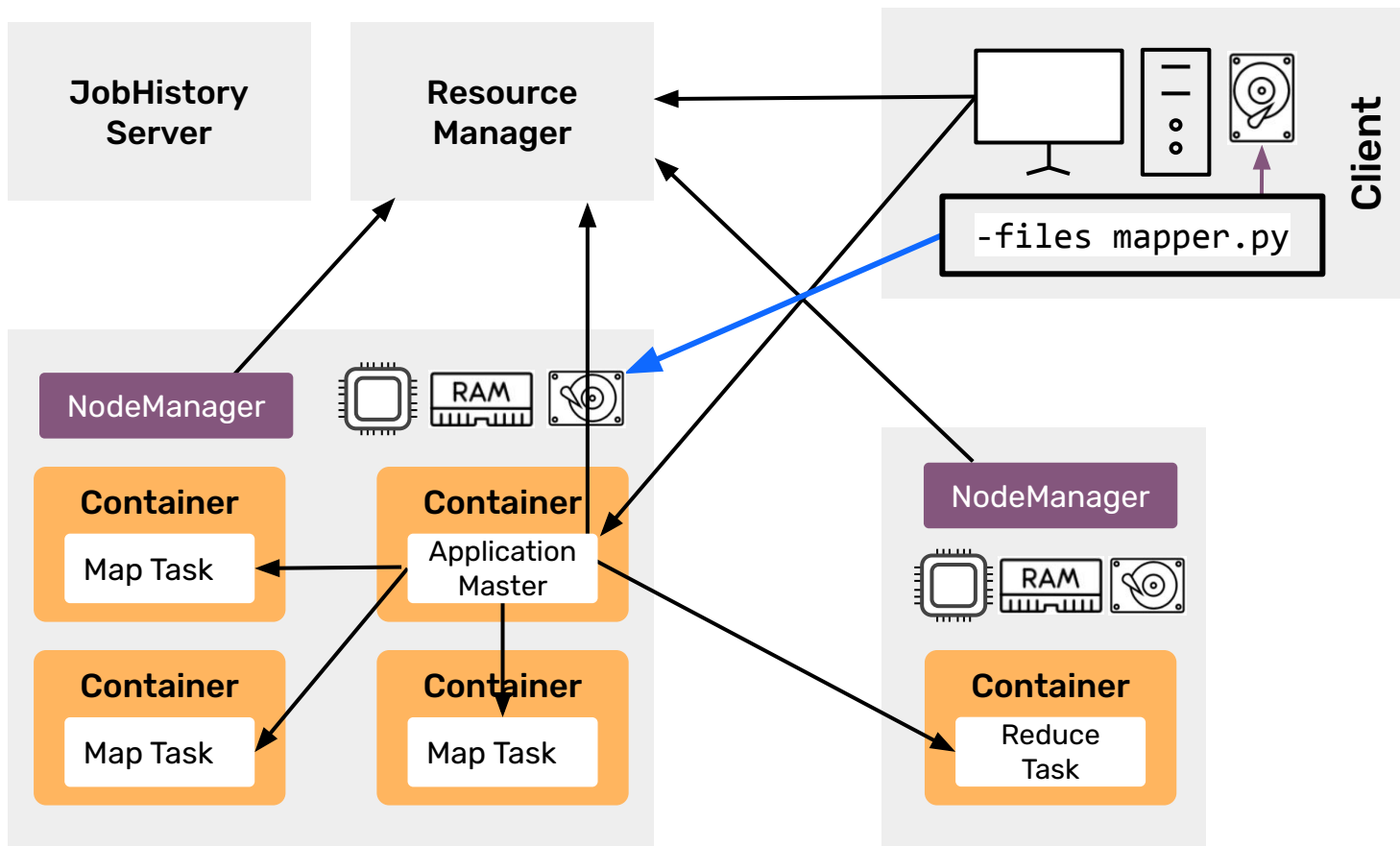
for line in sys.stdin:
    article_id, content = line.split("\t", 1)
    words = re.split("\W+", content)
    for word in words:
        if word in vocabulary:
            print(word, 1, sep="\t")
```

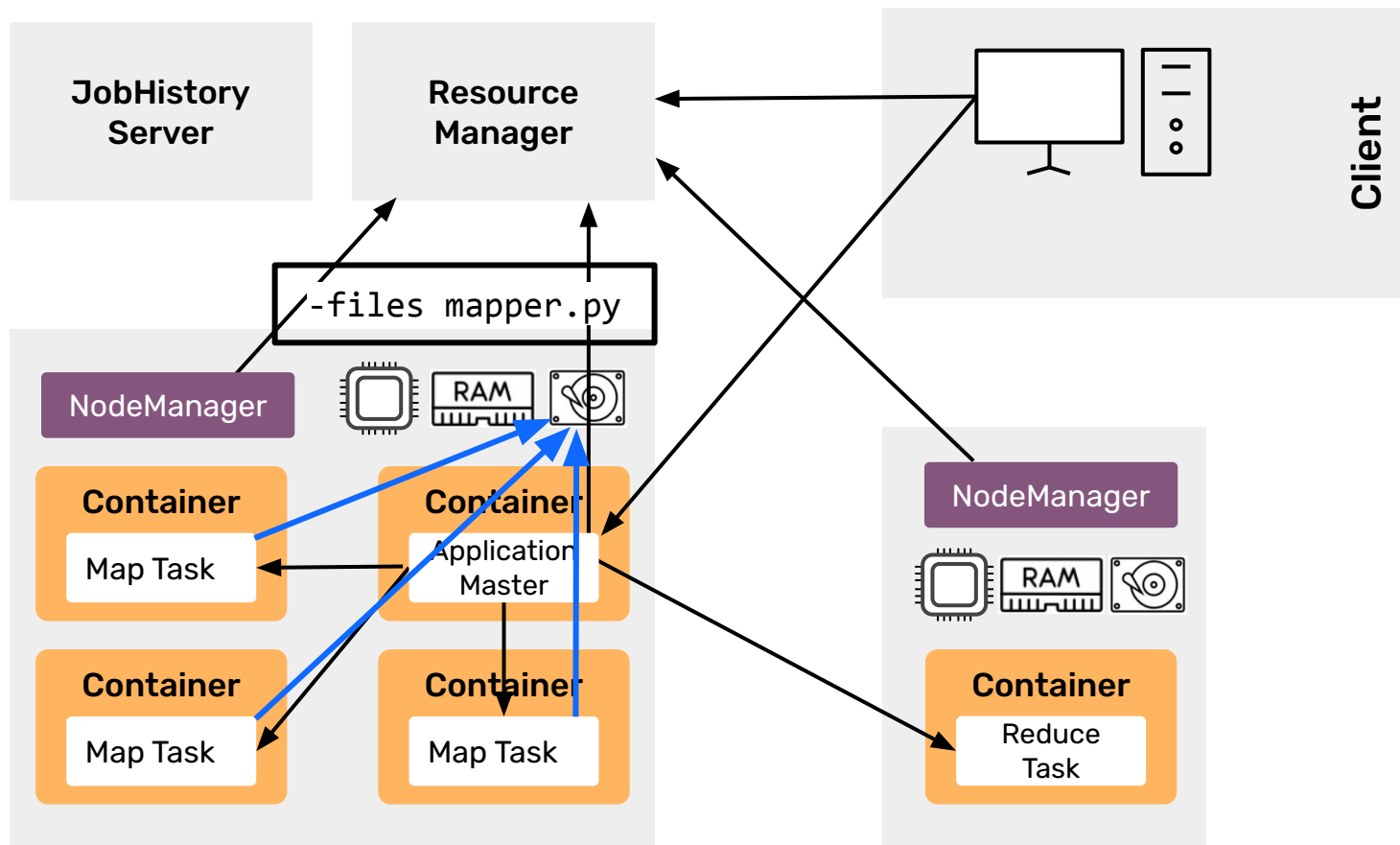


```
yarn jar $HADOOP_STREAMING_JAR \  
-files mapper.py, reducer.py, vocabulary.txt \  
-mapper "python3 mapper.py" \  
-reducer "python3 reducer.py" \  
-input /data/wiki/en_articles_part \  
-output word_count
```



Распространение данных



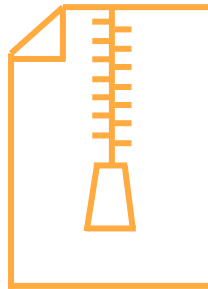




Варианты Distributed Cache



-files



-archives



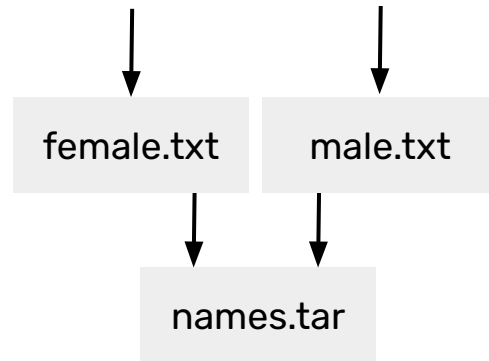
-libjars



Фильтрация по именам

TOP 10 NAMES		
GIRLS	Figures for 2015	BOYS
1 Olivia		1 Muhammad
2 Sophia		2 Oliver
3 Lily		3 Jack
4 Emily		4 Noah
5 Amelia		5 Jacob
6 Chloe		6 Harry
7 Isabelle		7 Charlie
8 Sophie		8 Ethan
9 Ella		9 James
10 Isabella		10 Thomas

Source: BabyCentre

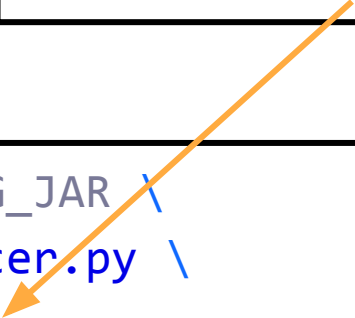


```
$ tar -cf names.tar male.txt female.txt
```




```
$ tar -cf names.tar male.txt female.txt
```

```
yarn jar $HADOOP_STREAMING_JAR \  
-files mapper.py, reducer.py \  
-archives names.tar \  
-mapper "python3 mapper.py" \  
-reducer "python3 reducer.py" \  
-input /data/wiki/en_articles_part \  
-output word_count
```





mapper.py

```
import re
import sys

def read_vocabulary(file_path):
    return set(line.strip() for line in open(file_path))

male_names = read_vocabulary("names.tar/male.txt")
female_names = read_vocabulary("names.tar/female.txt")

for line in sys.stdin:
    article_id, content = line.split("\t", 1)
    words = re.split("\W+", content)
    for word in words:
        if word in male_names or word in female_names:
            print(word, 1, sep="\t")
```



Ожидание vs Реальность?

TOP 10 NAMES

GIRLS

Figures for 2015

BOYS

- 1 Olivia
- 2 Sophia
- 3 Lily
- 4 Emily
- 5 Amelia
- 6 Chloe
- 7 Isabelle
- 8 Sophie
- 9 Ella
- 10 Isabella



- 1 Muhammad
- 2 Oliver
- 3 Jack
- 4 Noah
- 5 Jacob
- 6 Harry
- 7 Charlie
- 8 Ethan
- 9 James
- 10 Thomas

Source: BabyCentre

```
$ hdfs dfs -text word_count/* | sort -nrk2,2
```

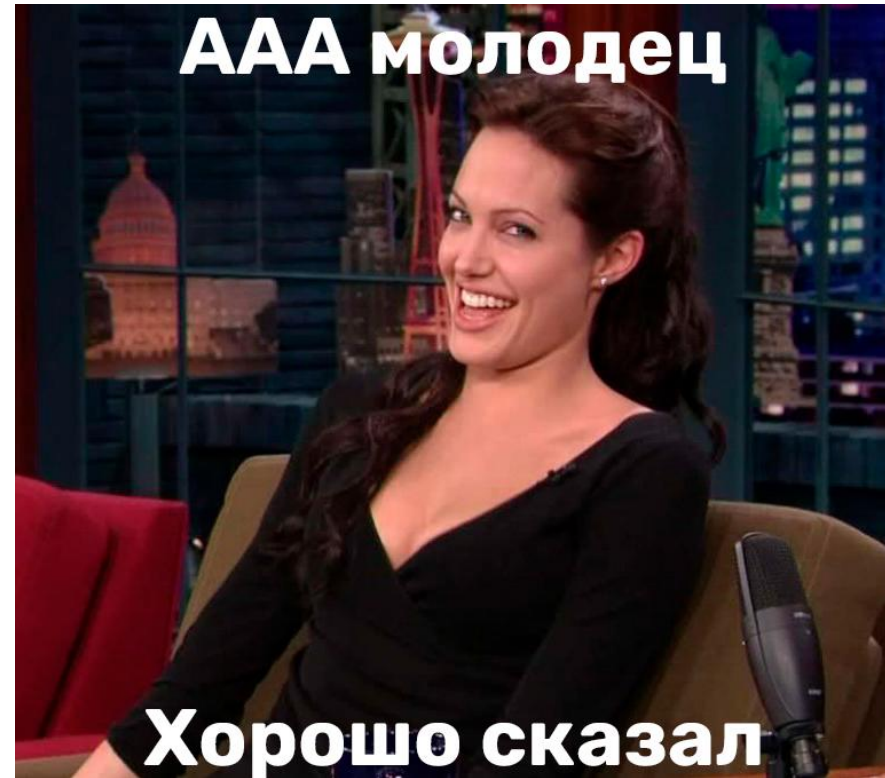
James	2284
Thomas	1941
Jack	786
Harry	504
Muhammad	444
Oliver	250
Charlie	250
Jacob	234
Emily	128
Isabella	99

Sophia	92
Noah	80
Sophie	64
Lily	31
Olivia	27
Ethan	27
Ella	25
Amelia	25
Isabelle	18
Chloe	9



Теперь вы:

- ▶ Знаете форматы использования Distributed Cache и как он работает под капотом





Больше контейнеров

Шире Shared Cache

