

# TASK: parallel grep Facebook posts Log

---

- **TASK** find all occurrences of a word, say, “apple” in the Facebook posts log file
- **In Unix solve using**
  - *grep “apple” facebook.log*

*[2007/12/25 09:15] Ate an apple*

*[2008/02/12 12:37] Thought about apples*

*[2009/01/09 19:55] MMmmm.. Apples*

- **Assume 1 Terabyte of log files a single core will take about 5 hours to grep the entire file**
- **How can you bring this down to a minute or two using 100 core machine?**
- **Using the Unix map-reduce framework:**
  - Given a mapper and a reducer function in python fill in the blanks

# Challenge

---

- **Fill in the blanks with comments and python code in the following to get this version of the parallel grep to work:**
  - Mapper
  - Reducer

## Python Notebook Mapper and reducers in Command Line

```

1  %%writefile pGrepCount.sh
2  ORIGINAL_FILE=$1
3  FIND_WORD=$2
4  BLOCK_SIZE=$3
5  CHUNK_FILE_PREFIX=$ORIGINAL_FILE.split
6  SORTED_CHUNK_FILES=$CHUNK_FILE_PREFIX*.sorted
7  usage()
8  {
9      echo Parallel grep
10     echo usage: pGrepCount filename word chunksize
11     • echo greps file file1 in $ORIGINAL_FILE and counts the number of lines
12     echo Note: file1 will be split in chunks up to $ BLOCK_SIZE chunks each
13     echo $FIND_WORD each chunk will be grepCounted in parallel
14 }
15 #Splitting $ORIGINAL_FILE INTO CHUNKS
16 split -b $BLOCK_SIZE $ORIGINAL_FILE $CHUNK_FILE_PREFIX
17 #DISTRIBUTE
18 for file in $CHUNK_FILE_PREFIX*
19 do
20     grep -i $FIND_WORD $file|wc -l >$file.intermediateCount &
21 done
22 wait
23 #MERGEING INTERMEDIATE COUNT CAN TAKE THE FIRST COLUMN AND TOTOL...
24 numOfInstances=$(cat *.intermediateCount | cut -f 1 | paste -sd+ - |bc)
25 echo "found [$numOfInstances] occurences of [$FIND_WORD] in the file [$ORIGINAL_FILE]"

```

Mapper

Reducer

Writing pGrepCount.sh

## Run the file

```
1  !chmod a+x pGrepCount.sh
```

Usage: usage: pGrepCount filename word chunksize

```
1  !./pGrepCount.sh License COPYRIGHT 4k
```

Large-Sc found [57] occurences of [COPYRIGHT] in the file [License]

## Mapper for grep

```
1 %%writefile mapper.py
2 #!/usr/bin/python
3 import sys
4 import re
5 count = 0
6 filename = sys.argv[2]
7 findword = sys.argv[1]
8 with open (filename, "r") as myfile:
9     for line in myfile.readlines():
10
11
12 print count
```

Read each line  
Increment counter if desired is present in the line

**Python Notebook  
Mapper and reducers in  
Python!!**

**1 of 2: Fill in the blanks**

Overwriting mapper.py

```
1 !chmod a+x mapper.py
```

## Reducer for grep

```
1 %%writefile reducer.py
2 #!/usr/bin/python
3 import sys
4 sum = 0
5 for countStr in sys.stdin:
6
7
8 print sum
```

Read mapper output (single number per line)  
Sum up  
Output total count

**2 of 2: Fill in the blanks**

Overwriting reducer.py

```
1 !chmod a+x reducer.py
```

## Map-Reduce framework in the Com

Python Notebook  
Mapper and reducers in  
Python!!

```
1 %%writefile pGrepCount.sh
2 ORIGINAL_FILE=$1
3 FIND_WORD=$2
4 BLOCK_SIZE=$3
5 CHUNK_FILE_PREFIX=$ORIGINAL_FILE.split
6 SORTED_CHUNK_FILES=$CHUNK_FILE_PREFIX*.sorted
7 usage()
8 {
9     echo Parallel grep
10    echo usage: pGrepCount filename word chunksize
11    echo greps file in $ORIGINAL_FILE and counts the number of lines
12    echo Note: file will be split in chunks up to $ BLOCK_SIZE chunks each
13    echo $FIND_WORD each chunk will be grepCounted in parallel
14 }
15 #Splitting $ORIGINAL_FILE INTO CHUNKS
16 split -b $BLOCK_SIZE $ORIGINAL_FILE $CHUNK_FILE_PREFIX
17 #DISTRIBUTE
18 for file in $CHUNK_FILE_PREFIX*
19 do
20     #grep -i $FIND_WORD $file/wc -l >$file.intermediateCount &
21     ./mapper.py $FIND_WORD $file >$file.intermediateCount &
22 done
23 wait
24 #MERGEING INTERMEDIATE COUNT CAN TAKE THE FIRST COLUMN AND TOTOL...
25 #numOfInstances=$(cat *.intermediateCount | cut -f 1 | paste -sd+ - |bc)
26 numOfInstances=$(cat *.intermediateCount | ./reducer.py)
27 echo "found [$numOfInstances] occurences of [$FIND_WORD] in the file [$ORIGINAL_FILE]"
```

Mapper

Reducer

Overwriting pGrepCount.sh

## Run the parallel grep using python-based mapper and reducer

```
1 !chmod a+x pGrepCount.sh
```

Usage: usage: pGrepCount filename word chunksize

```
1 !./pGrepCount.sh License COPYRIGHT 4k
```

found [57] occurences of [COPYRIGHT] in the file [License]