

## Downloading texts

```
~/Documents/work/distributed_calculations_hw1 13:49:13
$ curl -o the_brothers_karamazov.txt https://www.gutenberg.org/cache/epub/28054/pg28054.txt
  % Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
                                 Dload  Upload   Total   Spent    Left   Speed
100 1994k  100 1994k    0     0 1234k      0  0:00:01  0:00:01 --:--:-- 1243k
```

```
~/Documents/work/distributed_calculations_hw1/inputs 13:50:43
$ cat * > combined.txt
```

```
~/Documents/work/distributed_calculations_hw1/inputs 13:50:20
$ find . -maxdepth 1 -type f -print0 | xargs -0 du -ch | grep total$

22M    total
```

## Writing program

```
11 def main():
12     base_path = get_base_path()
13     input_text_path = base_path + 'inputs/combined.txt'
14     output_path = base_path + 'results'
15
16     if os.path.exists(output_path):
17         for filename in os.listdir(output_path):
18             file_path = os.path.join(output_path, filename)
19             if os.path.isfile(file_path):
20                 os.remove(file_path)
21         os.rmdir(output_path)
22
23     sc = SparkContext('local', 'Simple WordCount on one machine')
24
25     words = sc.textFile(input_text_path).flatMap(
26         lambda text:
27             text.lower().translate(str.maketrans('', '', string.punctuation)).split()
28     )
29
30     result = words.map(lambda word: (word, 1))\
31         .reduceByKey(lambda a, b: a + b)\
32         .sortBy(lambda a: a[1], ascending=False)
33     result.saveAsTextFile(output_path)
34
35
36 def show_results():
37     path = get_base_path() + 'results/part-00000'
38     with open(path, 'r', encoding='utf-8') as f:
39         for i, line in enumerate(f, start=1):
40             if i > 10:
41                 break
42             print(line.rstrip())
43
44
45 if __name__ == "__main__":
46     main()
47     show_results()
```

## Testing

```
/Users/maxim-movshin/.local/share/virtualenvs/pythonProject2-8u_Nj7EF/bin/python /Users/maxim-movshin/Documents/work/distributed_calculations_hw1/main.py
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/06/19 13:52:58 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
/Users/maxim-movshin/.local/share/virtualenvs/pythonProject2-8u_Nj7EF/lib/python3.11/site-packages/pyspark/python/lib/pyspark.zip/pyspark/shuffle.py:65: UserWarning: Please install psutil to have better
('the', 93871)
('and', 63996)
('to', 52720)
('of', 45665)
('a', 37139)
('he', 35528)
('in', 29444)
('that', 26552)
('was', 24329)
('his', 23933)

Process finished with exit code 0
```