

# Parallel Autocomplete

Authors: Sambhav Anand (sa3433), Kanishk Vashisht (kv2295)

## Abstract

The task of autocomplete is returning an arbitrary number of popular words that begin with a suggested prefix. Our project was centered around parallelising the entire autocomplete workflow. We look at parallelly constructing a word frequency table and then traversing a trie to autocomplete a given word. We compare single core implementations versus multiple core implementations on small and large bodies of text.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data Gathering</b>	<b>1</b>
<b>3</b>	<b>Sequential approach</b>	<b>1</b>
3.1	Word Frequencies . . . . .	1
3.2	Word Completion . . . . .	1
<b>4</b>	<b>Parallel approach</b>	<b>1</b>
4.1	Word Frequencies . . . . .	1
4.2	Word Completion . . . . .	2
<b>5</b>	<b>Results</b>	<b>2</b>
5.1	Observations . . . . .	2
5.2	Results . . . . .	2
5.2.1	Word Frequencies . . . . .	2
5.3	Conclusion . . . . .	2
<b>6</b>	<b>Appendix</b>	<b>3</b>
6.1	Core Usage . . . . .	3
6.2	Code Listing . . . . .	4

# 1 Introduction

The problem we undertook was that given a corpus of text, could we return the k most popular words that began with a suggested prefix. Popularity is defined on the basis of frequency of the word in the original corpus.

We break down the original problem into two sub problems.

1. **Word Frequencies:** This is the step where the corpus is examined and reduced to a dictionary of words that includes the frequencies. This is an intermediate step that only needs to be run once for an infinite number of autocomplete queries.
2. **Word Completion:** In this step we get two inputs from the user - the prefix p and k, the number of words to return. We load the words as a trie and then pick the k most popular children with the prefix p and return them.

## 2 Data Gathering

We used a dataset that contains 1,097,592 Amazon reviews spanning May 1996 - July 2014 in the CDs/Vinyl section. There are a total of 200M words of text with 200K unique words in this dataset. This dataset can be found [here](#).

## 3 Sequential approach

### 3.1 Word Frequencies

We implement word frequencies using map-reduce. Our mapping function takes a filename, reads the content from it and then generates a frequency table. Our reducer function takes these frequency tables and then combines them into one bigger frequency table. We convert this frequency table to a list of (word, frequency) pairs which we then write to a file.

### 3.2 Word Completion

When its time for word completion, we load our file into a trie. At this point we ask the user for an input and return the k most common words that begin with the provided input.

## 4 Parallel approach

### 4.1 Word Frequencies

Our parallel approach is similar to our sequential approach except we run the mapping functions in parallel on different files.

## 4.2 Word Completion

We tried to parallelize our word completion by querying the trie in parallel. We find the parent node that represents the prefix, then we traverse all of its children nodes in parallel and return the k most common words from each. After this we find the k most common words in total and return those. However, we found that for even the most basic prefixes such as "pre" or "pr" the sequential approach was outperforming our rudimentary parallel approaches as there were only close to 4-5000 words to sort through and pick k from. We believed that parallelising this step would not change much so instead we focused on improving the efficiency of word frequencies to scale.

## 5 Results

### 5.1 Observations

One interesting thing we observed was that if the files were not distributed well then the cores did not provide much optimization. This is obvious since we parallelize over files. This is why for the general results that follow we made sure to divide our dataset into 55 evenly distributed files.

### 5.2 Results

#### 5.2.1 Word Frequencies

The following table shows our results for the word frequency part. We look at how things change by making the corpus larger and by adding more cores.

# Words / Cores	1 core	4 cores	8 cores
417,017,098	966s	416s	238s
208,508,549	355s	140s	136s
75,667,556	160s	49s	37s
1,899,4220	40s	16s	16s

We've also attached the output artifacts (in the appendix section) showing the usage of our cores for the 400 million word data set. As one can see all our cores are being used well and as we're adding more cores they are getting busier as well. This shows that our algorithm is using parallelization efficiently. This is also evident in the fact that we can deal with a 417m word corpus in 238s.

### 5.3 Conclusion

In conclusion, we have:

1. Come up with a full algorithm for word completion that is broken into word frequencies and word completion
2. We are able to parallelise word frequencies to a large scale ( 417m words).

# 6 Appendix

## 6.1 Core Usage

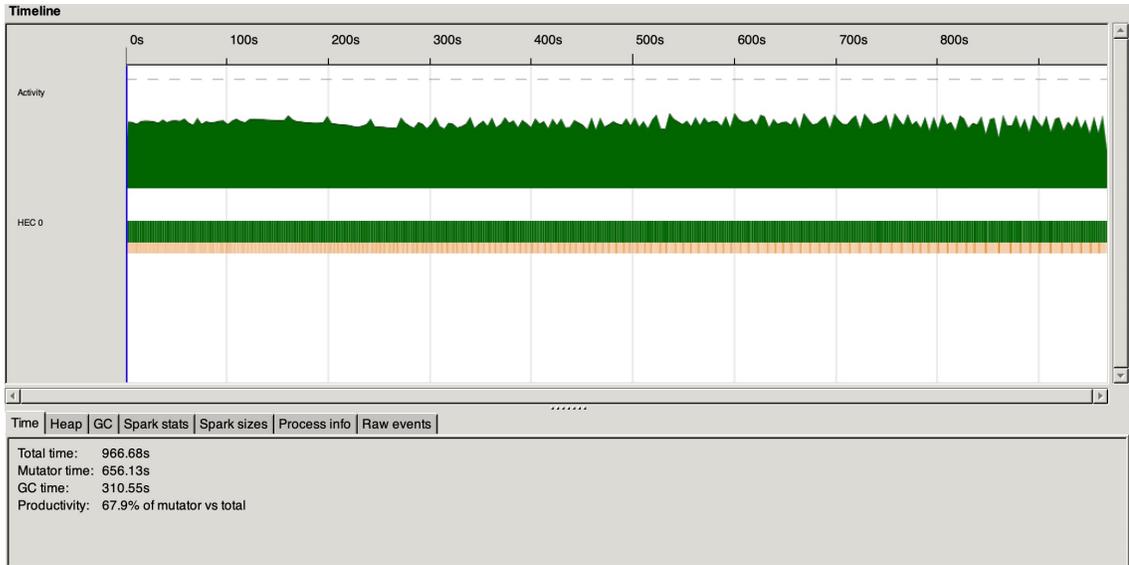


Figure 1: One Core

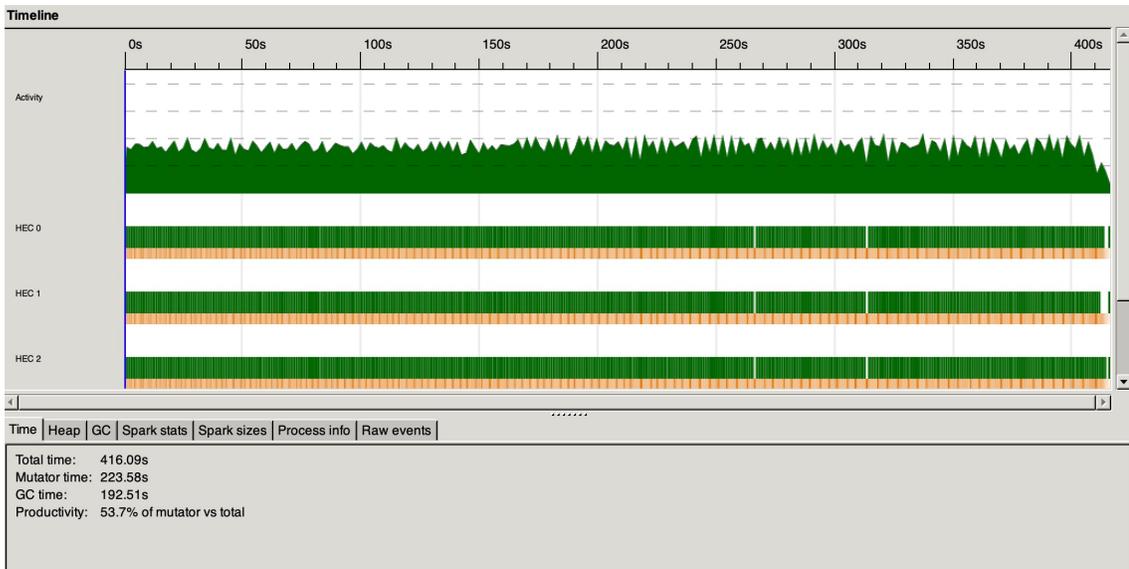


Figure 2: Four Cores

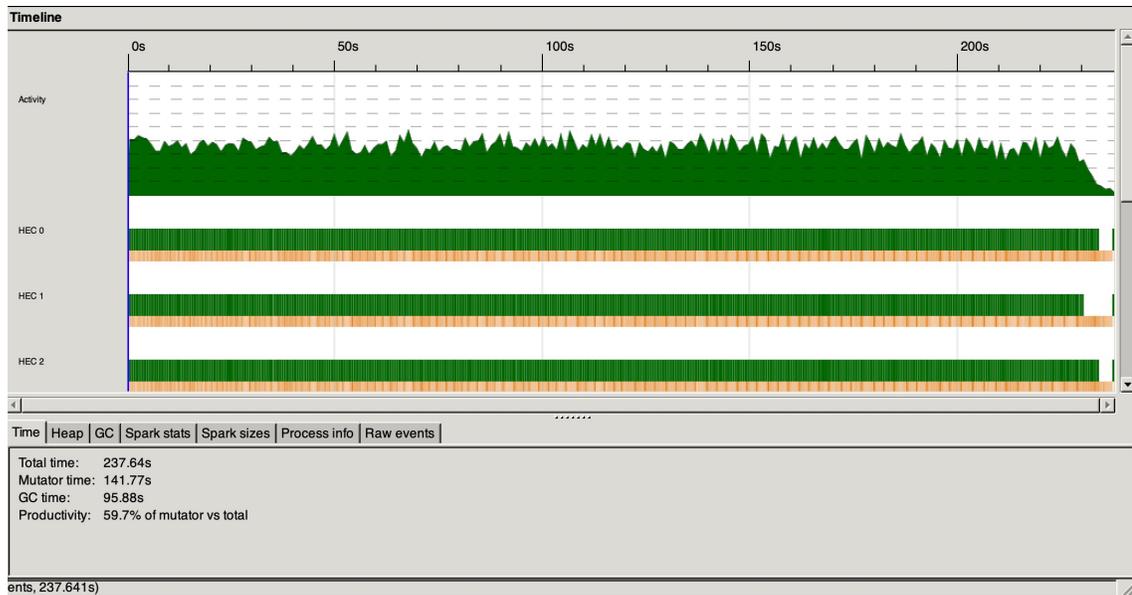


Figure 3: Eight Cores

## 6.2 Code Listing

```

1 module Main where
2
3 import Data.List (sort, group, sortBy, groupBy, isPrefixOf, isInfixOf, sortOn)
4 import Data.Text (isSuffixOf)
5 import Control.Parallel
6 import Control.Parallel.Strategies
7 import Data.Char (isAlpha, isSpace, toLower)
8 import qualified Data.ByteString.Char8 as BC
9 import qualified Data.Trie as T
10 import qualified Data.Map as M
11 import System.Directory
12 import System.IO
13
14 mapReduce :: NFData b1 => (a -> b1) -> ([b1] -> b2) -> [a] -> b2
15 mapReduce mapper reducer input = pseq mOutput rOutput
16   where mOutput = parMap (rpar 'dot' rdeepseq) mapper input
17         rOutput = reducer mOutput 'using' rseq
18
19
20 mappingFunc :: String -> M.Map String Int
21 mappingFunc document = getWordFreqMap $ (strToList) document
22
23
24 getWordFreqMap :: [String] -> M.Map String Int
25 getWordFreqMap tokens = M.fromListWith (+) (map (\x -> (x, 1)) tokens)
26
27 strToList :: String -> [String]
28 strToList str = words $ filter (\char -> isAlpha char || isSpace char) $ map toLower str
29

```

```

30 reducingFunc :: [M.Map String Int] -> [(String, Int)]
31 reducingFunc maps = M.toList $ foldl (\foldMap (word, count) -> M.insertWith (+) word count foldMap) (M.empty)
  listOfMaps
32       where listOfMaps = concat (map M.toList maps)
33
34 loadList :: FilePath -> IO (M.Map String Int)
35 loadList fname = do
36     filedata <- readFile fname
37     contents <- return (read filedata :: [(String, Int)])
38     let loadedMap = M.fromList contents
39     return loadedMap
40
41 saveList :: Show a => a -> FilePath -> IO ()
42 saveList ls fname= writeFile fname (show ls)
43
44 FilePath -> IO ()
45 buildWordList directory = do
46     files <- listDirectory directory
47     let myFiles = filter (\x -> not (elem x [".", ".."])) (map (\x -> directory ++ x) files)
48     parsedFiles <- mapM readFile myFiles
49     let freqMap = mapReduce mappingFunc reducingFunc parsedFiles
50     putStrLn $ "Found " ++ (show.length) freqMap ++ " words!"
51     saveList freqMap "word_counts.txt"
52
53 buildTrie :: M.Map String Int -> T.Trie Int
54 buildTrie loadedMap = T.fromList (map (\x -> (BC.pack (fst x), snd x::Int)) (M.toList loadedMap))
55
56
57 Ord b1 => T.Trie b1 -> IO b2
58 userInput trie = do
59     putStrLn "enter word"
60     pre <- getLine
61     putStrLn "Enter number of words you want..."
62     k <- getLine
63     let topk = take (read k::Int) (sortBy (\(-,a) (-, b) -> flip compare a b) $ T.toList (T.submap (BC.pack pre) trie))
64     putStrLn "Found the following words"
65     putStrLn " - Start - "
66     mapM_ (\(a,-) -> putStrLn $ BC.unpack a) topk
67     putStrLn " - Finish - "
68     userInput trie
69
70
71 main :: IO ()
72 main = do
73     -- buildWordList "./data/"
74     loadedMap <- loadList "data.txt"
75     let trie = buildTrie loadedMap
76     putStrLn $ "Found " ++ (show $ T.size trie) ++ " words"
77     userInput trie

```

Citation: the format for this code can be found here