

**CS 2302 - Data Structures**  
**Fall 2018**  
**Project 4 - Option A**

**Overview**

In the previous project (lab), you implemented an algorithm to compute the similarity between two words using their word embeddings. For a word corpus size  $v$ , it took  $O(v \log v)$  to build the BST that holds all the words and their embeddings. The purpose of this assignment is to improve the running time of the algorithm by using a hash table instead of a BST.

Your hash table must use chaining to solve collisions. For comparison purposes, you must write multiple hash functions. Also, write a method to determine the average number of comparisons required to perform a successful retrieve operation, and a method to compute the table's load factor. In the previous project, we used a Balanced BST, so on average, the retrieve operation needed to perform  $\log(400k) \approx 18$  comparisons. Your hash table (to be successful) needs to perform significantly fewer operations than that (in the ideal case, you'd need exactly 1 comparison per access).

**Hint 1:** The choice of table size and hashing function has a large impact in the number of collisions and thus running time.

**Hint 2:** There are 400k words in *glove.6B.50d.txt*.

**Hint 3:** You can think of a word in the English language as a base-26 number.

**What you need to do**

**Part 1 - Due Sunday, November 11, 2018**

Implement the program described above and upload your code to GitHub.

**Part 2 - Due Wednesday, November 14, 2018**

Add your team members as collaborators to your GitHub repo. They will add you to their projects as a collaborator as well. Read their code and give them feedback. Use *pull requests* and/or the *Issues* section to do so .

**Rubric**

Criteria	Proficient	Neutral	Unsatisfactory
----------	------------	---------	----------------

<b>Correctness</b>	The code compiles, runs, and solves the problem.	The code compiles, runs, but does not solve the problem (partial implementation).	The code does not compile/run, or little progress was made.
<b>Space and Time complexity</b>	Appropriate for the problem.	Can be greatly improved.	Space and time complexity not analyzed
<b>Problem Decomposition</b>	Operations are broken down into loosely coupled, highly cohesive methods	Operations are broken down into methods, but they are not loosely coupled/highly cohesive	Most of the logic is inside a couple of big methods
<b>Style</b>	Variables and methods have meaningful/appropriate names	Only a subset of the variables and methods have meaningful/appropriate names	Few or none of the variables and methods have meaningful/appropriate names
<b>Robustness</b>	Program handles erroneous or unexpected input gracefully	Program handles some erroneous or unexpected input gracefully	Program does not handle erroneous or unexpected input gracefully
<b>Documentation</b>	Non-obvious code segments are well documented	Some non-obvious code segments are documented	Few or none non-obvious segments are documented
<b>Code Review</b>	Useful feedback was provided to team members.  Feedback received from team members was used to improve the code.	Feedback was provided to team members, but it was not very useful. Feedback received from team mates was partially used to improve the code	Little to no feedback was provided to team mates.  Received feedback was not used to improve the code.
<b>Report</b>	Covers all required material in a concise and clear way with proper grammar and spelling.	Covers a subset of the required material in a concise and clear way with proper grammar and spelling.	Does not cover enough material and/or the material is not presented in a concise and clear way with proper

			grammar and spelling.
--	--	--	-----------------------