

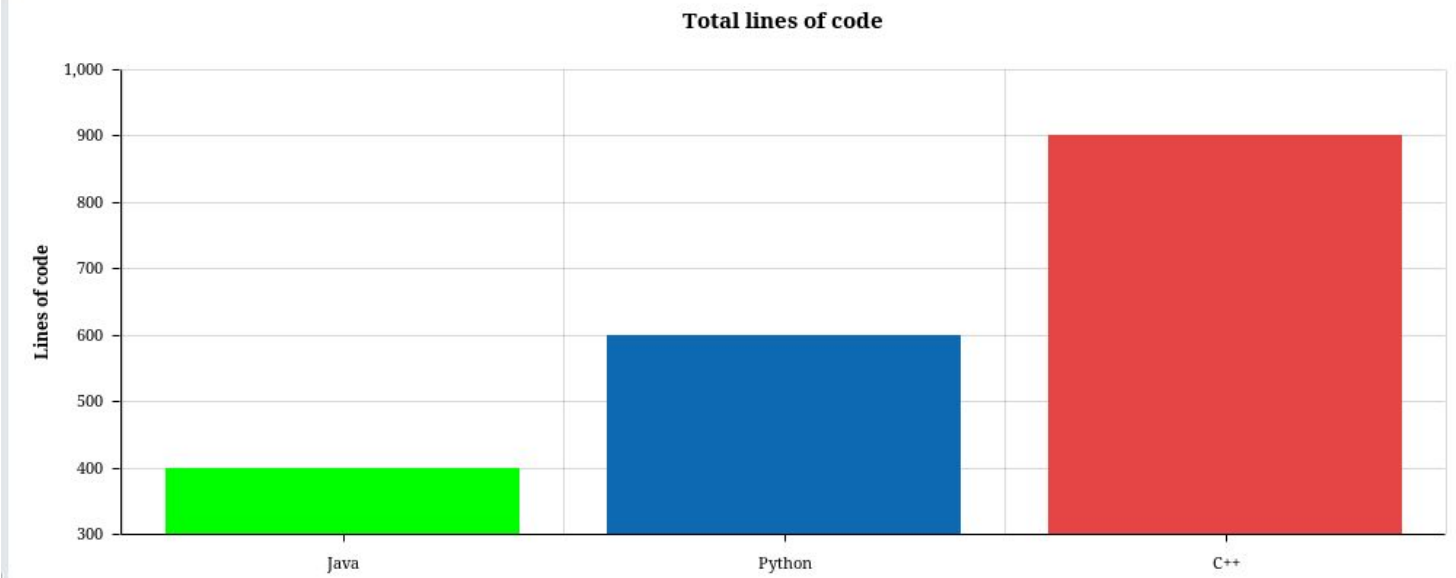
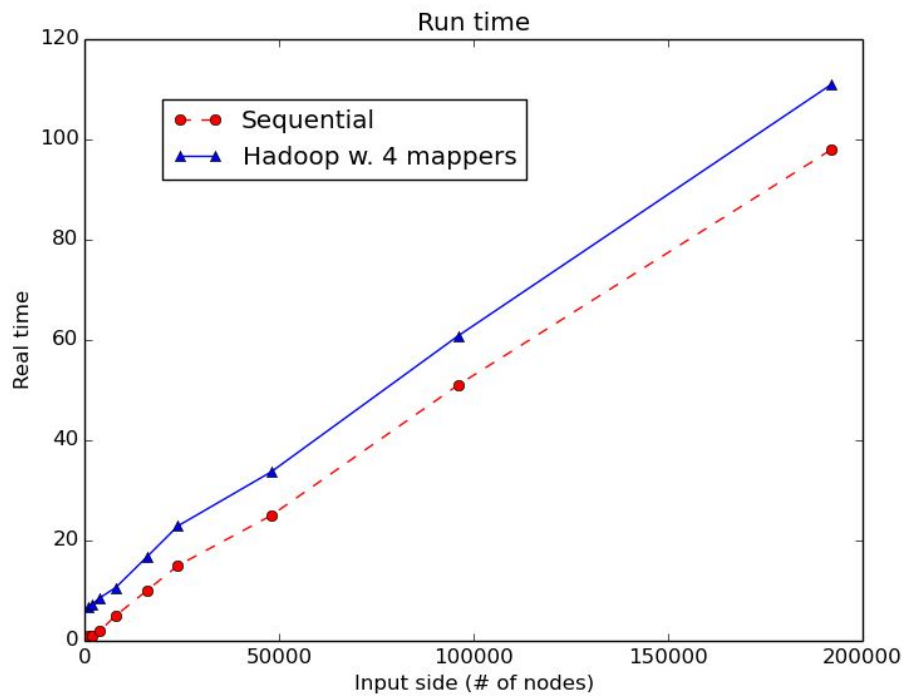
CSC494 Spring 2015

Supervisors: Renée J. Miller, Fatemeh Nargesian

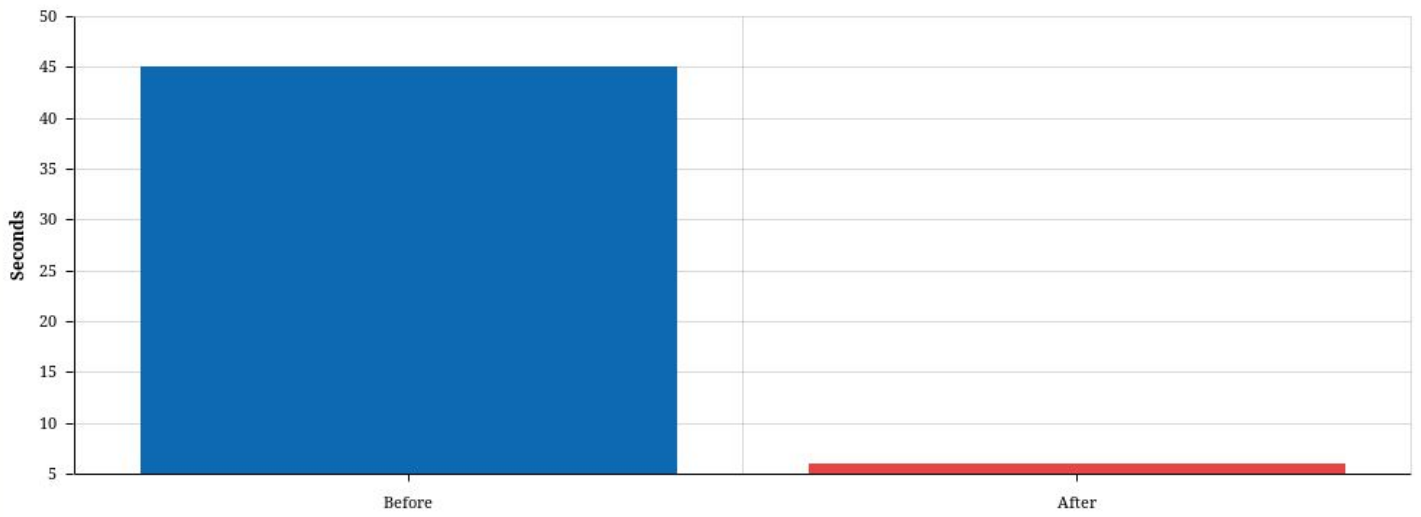
Huy Bui - g8buihuy - 996594421

Mihai Nicolae - g1mihai - 998584367

1. ✓ Gained working knowledge of the Hadoop framework.
2. ✓ Implemented and experimented with 2 different ways to parallelize search in Hadoop.
Found out that Hadoop is too slow given the small size of input.
3. ✓ Searching graph in parallel in order to achieve acceptable web application response time
 - researched different C++ graph libraries and chosen Boost Graph Library (BGL)
 - Rewrote search algorithm in C++
 - Found the bottleneck of the search algorithm: finding the initial rootQueue.
We parallelized this part and achieved 8X speed-up.
 - Removed irrelevant search results by adding logic to calculate similarity of concepts to the query.
 - Researched and applied other Boost libraries to add helpful features such as profiling.
 - Total lines of code: 900+
 - Future to-dos: serialize search DAG to disk, experiments against high quality dictionary
4. ✓ Creating a high quality dictionary
 - Experimented with more sophisticated heuristics, e.g. part of speech tagging to help remove unwanted words and stemming to reduce duplications, etc.
 - Parallelized dictionary creation. Achieved 180X speed-up.
 - This helps facilitate experiments of different heuristics to improve dictionary quality and, indirectly, search result.
 - Total lines of code: 600+



Search speed-up with C++ code rewrite



Dictionary creation speed-up

