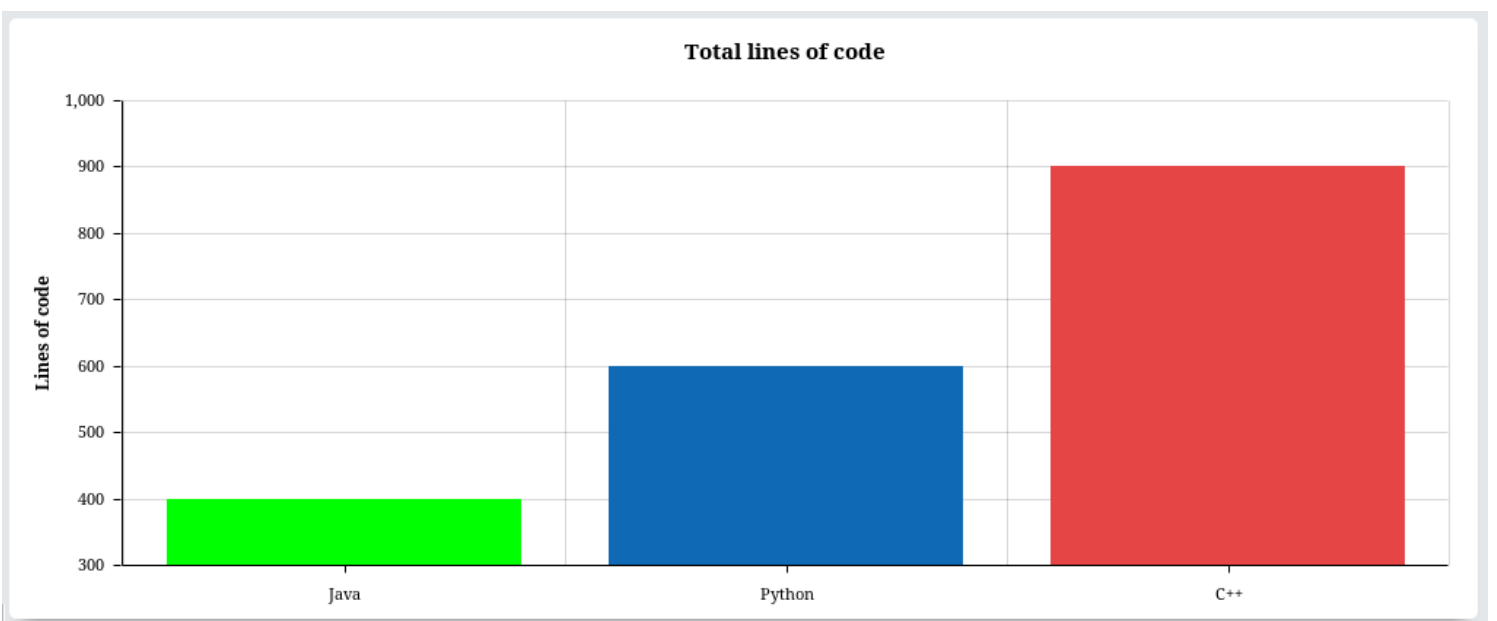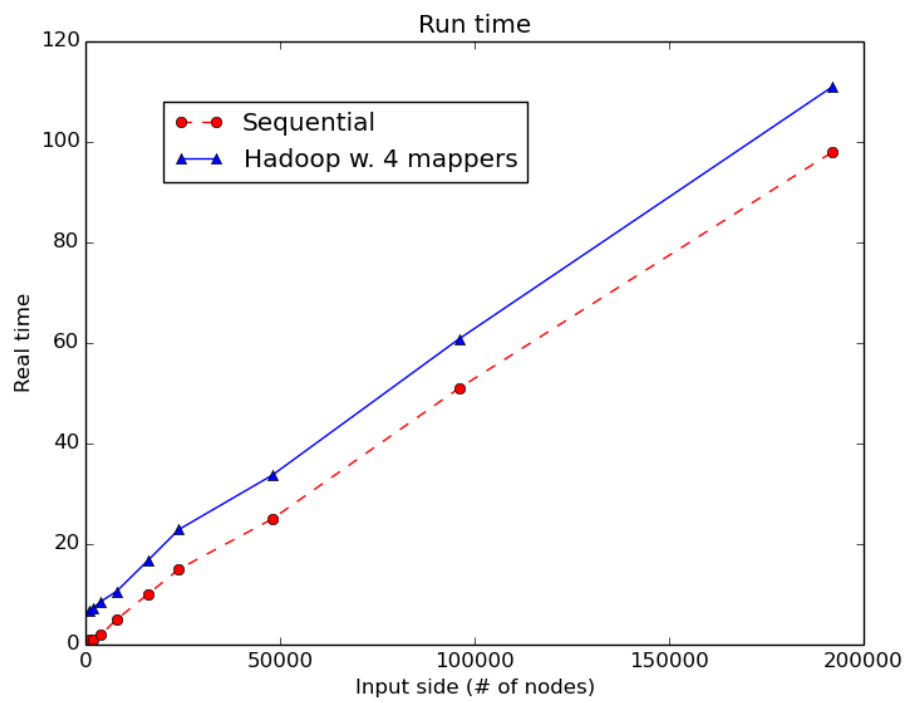# CSC494 Spring 2015 final report
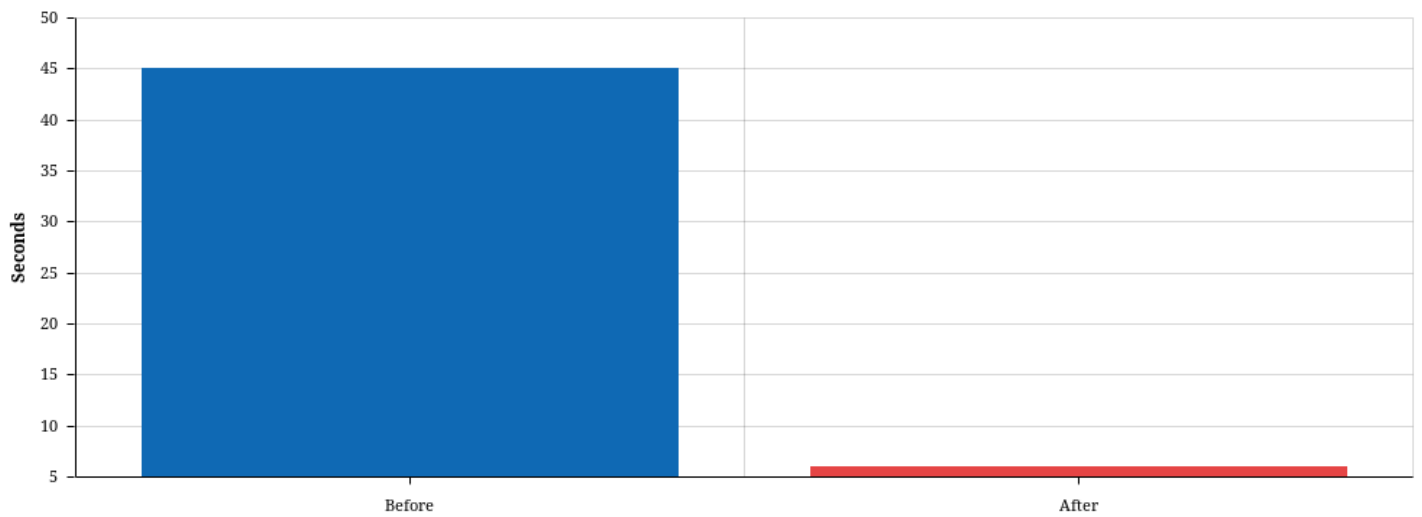
Supervisors: Renée J. Miller, Fatemeh Nargesian

Authors: Huy Bui, Mihai Nicolae

## Research highlights

1. Gained working knowledge of Apache Hadoop framework.
2. Implemented and experimented with 2 different ways to parallelize search in Hadoop. Concluded that Hadoop is too slow given the small input size.
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 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+

**Run time**



**Total lines of code**

## Search speed-up with C++ code rewrite



## Dictionary creation speed-up