# CSYE 7200 Big Data System Engineering using Scala
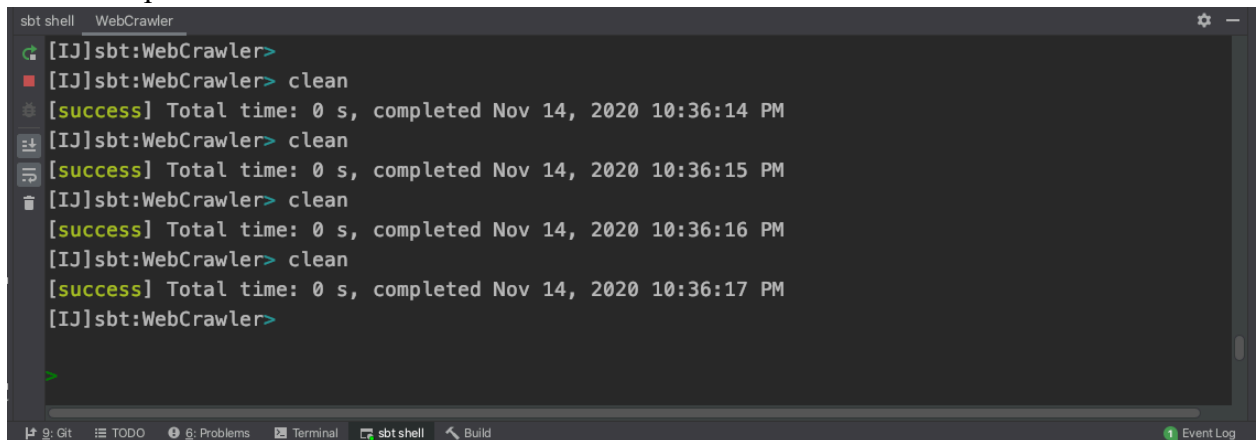
001306707
Ting-Kai Liu

A. WebCrawler - Test result:



B. sbt clean up:



C. Suggestion to improve the web crawler:
1. Scraping the data only if the content has been updated.
2. Scraping only partial of the content for a given URL.
3. Store the scraping result.