Yiqing Huang – SEC 01 (NUID 001525629)

# Big Data System Engineering with Scala
# Spring 2022
# Assignment No. 7

**Task**

*Web Crawler*

− **Implement the 3 'to be implemented' parts in *WebCrawler.scala*.**
1. **WebCrawler.wget.getURLs(ns: Node): Seq[Try[URL]]**
2. **WebCrawler.wget(u: URL): Future[Seq[URL]] for-comprehension**
3. **WebCrawler.wget(us: Seq[URL:): Future[Seq[Either[Throwable, Seq[URL]]]]**
− **Implement the 2 'to be implemented' parts in MonadOps.scala**
1. **MonadOps.mapFuture**
2. **MonadOps.sequence**
− **[Optional]** **Suggestions on how to improve the web crawler**
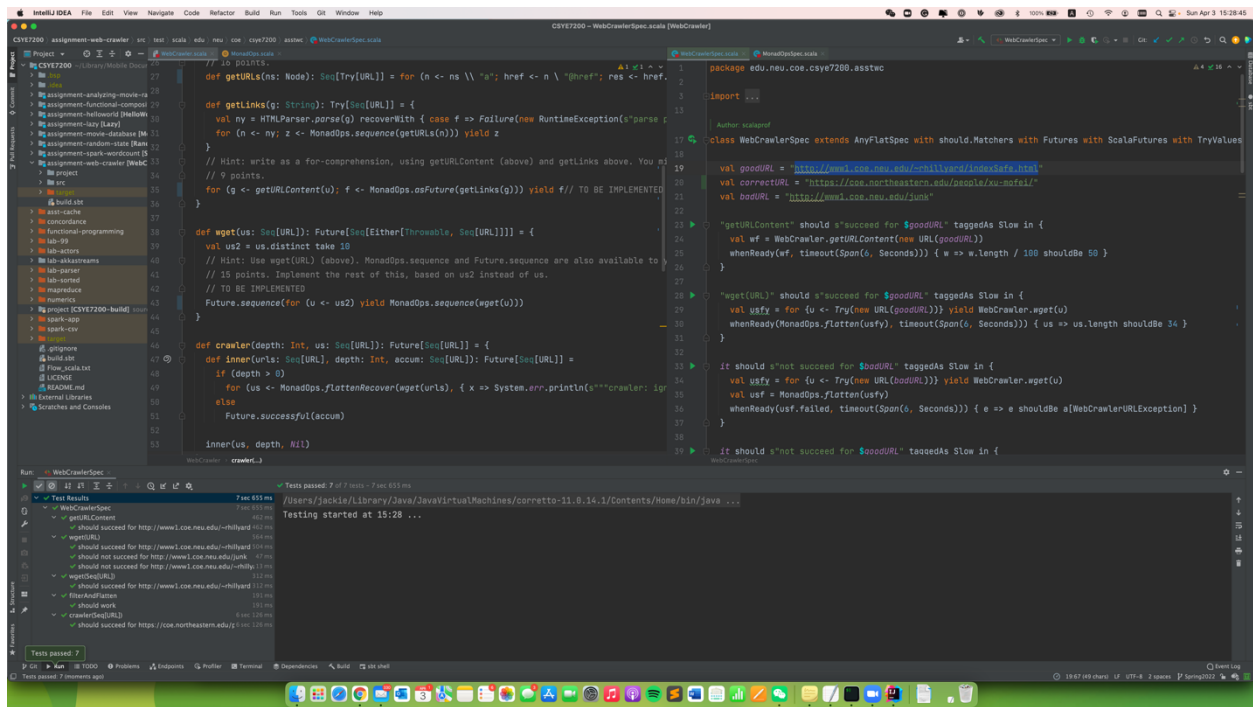
**Task**

*Suggestions*

I suggest refactoring the code by embedding an open-source toolkit called *akka* and using *Apache Kafka* as a buffer and load distributor for the crawling jobs. Firstly, *Akka* provides multiple high-performance programming models for concurrency (e.g., actor-based concurrency). Secondly, it provides an asynchronous, non-blocking HTTP client that enables a high degree of parallelism while keeping resource requirements to a minimum. Moreover, thanks to its high-level encapsulation, the whole web-crawling process can be simplified. As a result, our code can be much more readable and easier to develop and maintain.

According to a case study of the improvement of *PayPal*'s web crawler. It shows that by using the *akka* toolkit, the performance has increased tenfold for the reason that *akka* uses a much lower number of threads and achieves an incredible CPU utilization rate of around 90%. *(PayPal scales to billions of daily transactions with Scala and Akka Platform from Lightbend)*

**Solution/Unit test (screenshot)**

   **A. Unit Test**

## Project Source

https://github.com/eikcaj16/CSYE7200/tree/Spring2022/assignment-web-crawler