# FunPar Final Project Parallel Web Scraping
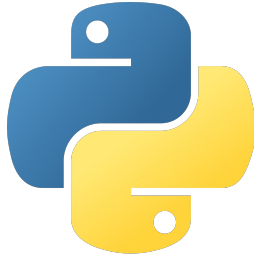
Porter Fredrickson

# The Basic Task

- Scrape lots of data from the web which would normally take long in a
  sequential fashion

- Analyze the data in some useful way

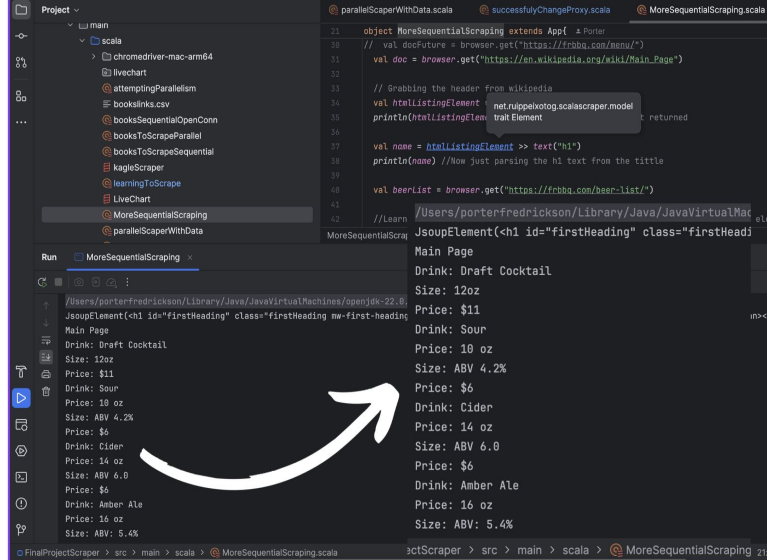- Find a way to speedup inherently slow sequential code

# Technologies used

# Phase 1: Figuring out how to Scrape

- My only experience "web scraping" was from assignment 3 in this class
- First thing I needed to do was learn exactly what web scraping is and how to do it
- Followed various scraping with Scala + JSoup tutorials

# Went on scraping journey

- After learning the basics with Jsoup and

    Scala I began testing my skills

- Connected to the website of my previous

    employer

- Scraped menu items

- Tried scraping other more secure websites

    and ran into my first blocks

# Books to scrape

# Books to Scrape Cont.

- Takes around ~315 seconds to scrape 1000's pages sequentially.

- No aggregate calculations, only storing price and title

- With Open session: ~307 seconds

- Minimal 8 second speedup

- Obviously slow

# Books to Scrape Parallel

- Scraped the same 1000 book pages
- Obtained useful data beyond title such as price, stock,
- Calculated aggregate statistics
  - Max book price = $59.99 and Min book price = $10.00
  - Average Stock = 8 books
  - Low stock count = 308 titles
  - Medium stock count = 153 titles
  - High stock count = 356 titles
- All said in done the parallel solution which did considerable amounts of work more than the sequential solution ran in **36 seconds**

```
------------------------------------------------------------------
The book with the highest price is: The Perfect Play (Play by Play #1) at $59.99, URL: http://books.toscrape.com/catalogue/the-perfect-play-play-by-play-1_352/index.htm
The book with the lowest price is: An Abundance of Katherines at $10.0, URL: http://books.toscrape.com/catalogue/an-abundance-of-katherines_362/index.html
The average price of all books is: 35.07034999999999
The average stock for books is: 8
There are 308 low stock titles
There are 153 medium stock titles
There are 356 high stock titles

------------------------------------------------------------------
Time taken: 36.1545745 seconds
```

# Proxies

- On a mission to compare my scraped data with Amazon's data
- Successfully searched book titles on amazon and obtained their price
- However, when scaling to 1000 books/requests I got blocked
- Introduced timeouts - partially successful
- Discovered the concept of using proxies (also recommended by Ajarn)
- Used a free proxy API
- Checked proxy validity in parallel
- Wrote valid proxies out to a csv file
- Checking over 1200 proxies takes 193.66 seconds
- A sequential solution should take about 75 minutes!

# Demo

# Difficulties

- Numerous difficulties and learning points
- Scala isn't the greatest to scrape in
- Obtaining large amounts of content from web pages can be very hard because it's dynamically loaded.
- High traffic sites are very difficult to scrape (Amazon)
- One solution to being blocked is proxies
- There are free proxies and paid for proxies
- Free proxies are for the most part useless
  - Most don't work
  - For those that do, major sites already block the free proxies

# Using Selenium  Python

- Determined to find more success and learn more about web scraping
- Selenium Webdriver is a software used for browser automation
- Works well with python
- Capable of loading dynamic content
- Used to scrape fivestars-thailand realty website
- Has realty listings for multiple areas in Thailand

# Another Demo

# Ultimate outcome

- Produced many different parallel web scraping solutions for this final project

- What better way to test them than to run everything at once

- Scraping 1000 website from books to scrape, all the current properties in bangkok,
  and all of the current properties in Phuket

- Doing all of this sequentially takes a gruelling ~370 seconds

- Doing all of this in parallel with Scala (futures) takes 54 seconds

- Shows a massive 6.85X speedup and a huge success

# How I was successful

Definition of Success in Proposal:

"The first component of success for this project is developing a web scraper that connects to a website (or list of websites) and obtains target data."

"The second component has a wider definition of success. This concurrent solution can be successful in several ways. First, the concurrent solution could allow me to scrape multiple sites at once at a similar rate to the first solution. Second, the concurrent solution could be successful at scraping the same data as the first solution, but much faster. Finally, the third solution could be successful in scraping so fast that it actually gets blocked by the website for making too many requests."

"Finally, the most important goals of this project are to learn about web scraping since I currently have no scaping knowledge and to learn about creating my own concurrent speedups."