# VARSHA PATTIPATI, USC ID: 8244565758

## Information Integration on the Web(CSCI 548) – Homework1

## TASK1:

The website I chose to crawl is ebay.

**WEBSITE:** www.ebay.com

Description: ebay is essentially a web platform for people to buy and sell goods. It is a place where sellers advertise their products( both old and new) and buyers can buy and also bid for products they want to purchase.

## TASK2:

The query I chose is as follows:

Extract webpages with the following filter

Samsung Cell Phones and Smartphones

Price range:

- Minimum: $0
- Maximum: $400

**Sample webpages:**

a.) http://www.ebay.com/itm/Samsung-Galaxy-S-4-SGH-M919-BRAND-NEW-T-Mobile-/252074961574?hash=item3ab0d6aaa6 (Figure 1)

b.) http://www.ebay.com/itm/Samsung-Galaxy-S5-SM-G900A-16GB-AT-T-Unlocked-Smartphone-Black-White-/311370709654?var=&hash=item487f23fa96 (Figure 2)

## TASK3:

**QUESTION 1: Describe the crawler you used to get the results. If you wrote your own crawler, explain which technologies you used to develop it.**

**Answer:** I had previously used Nutch for crawling websites in some projects I did last semester and hence very familiar with how it works. But I preferred creating one on my own in Java so that I could I could learn the language. I used jsoup parser library and imported the JAR file into my project.
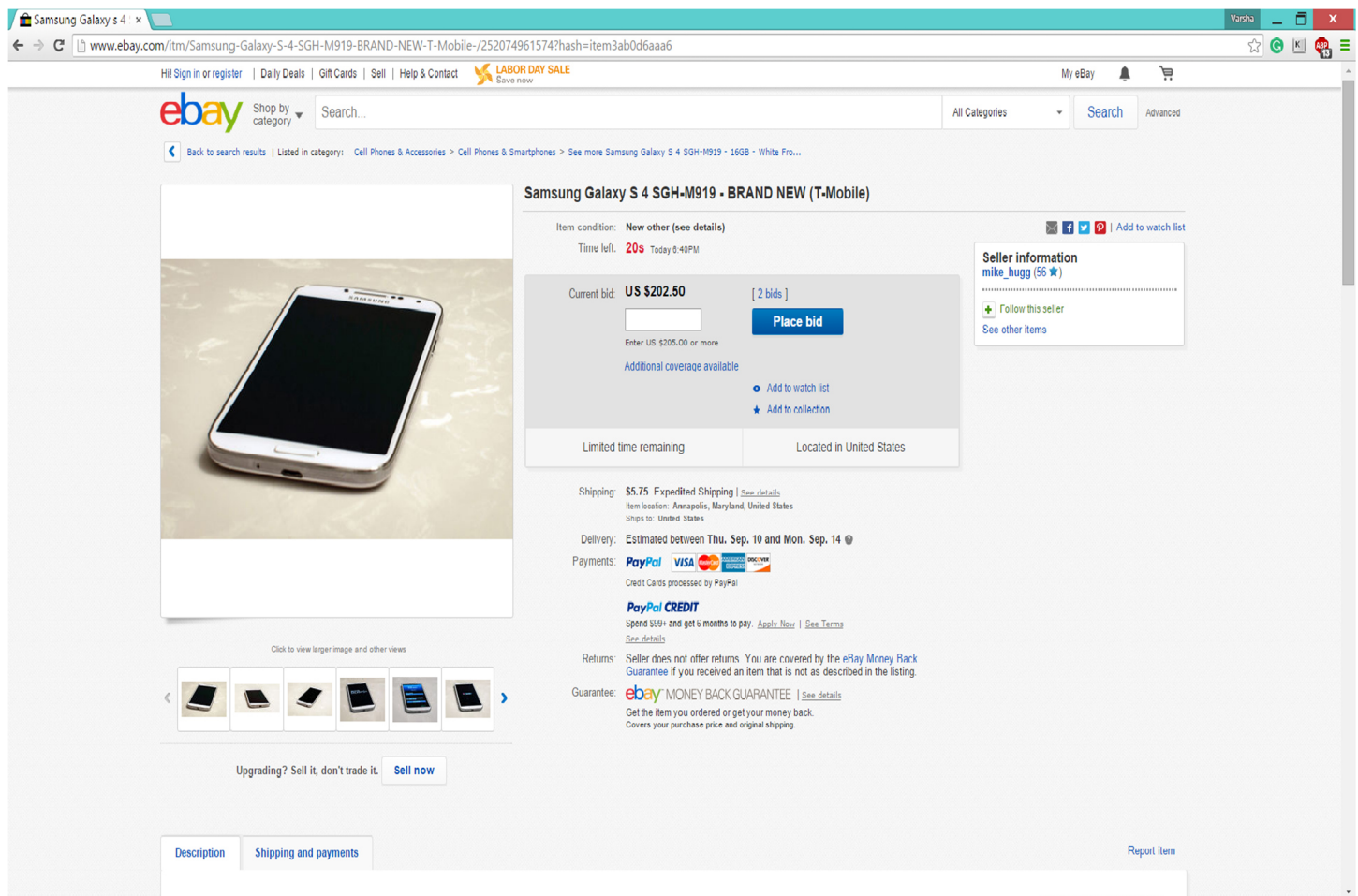
**Technologies used: Java, jsoup.**



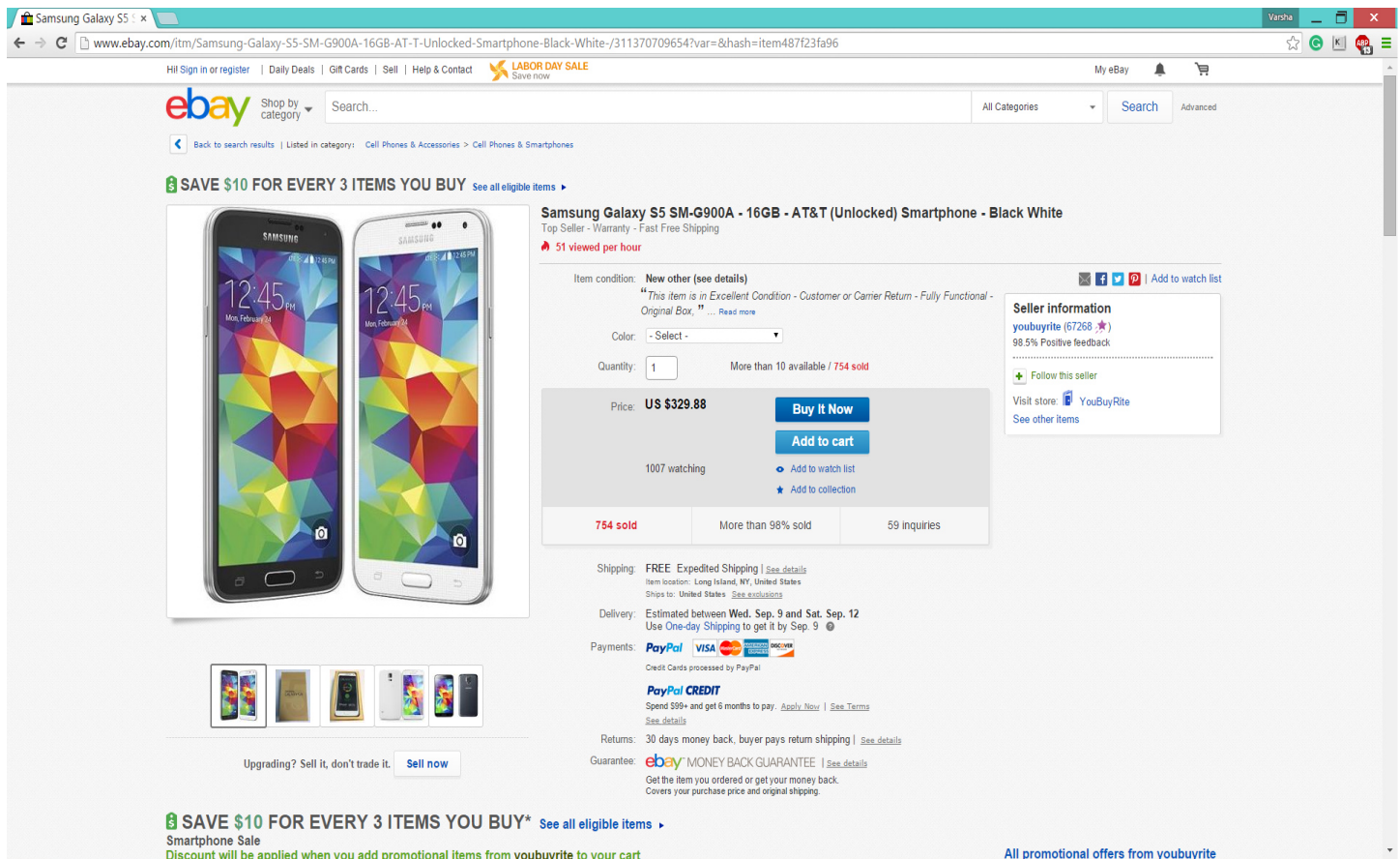**Figure1.** A webpage corresponding to the query I selected

**Figure2.** A second webpage corresponding to the query

## QUESTION 2: How does your crawler avoid a crawler trap (infinite loop)?

**Answer:**

- I started with some seed URLs and kept adding outgoing links from those pages to a HashSet, say 'pagesToVisit'.
- I also add pages already crawled to another HashSet, say 'pagesVisited'.
- I wrote a function to retrieve the next URL to be crawled from the HashSet('pagesToVisit') but made sure it is not crawled already by checking if 'pagesVisited' already contains the URL.
- So I never visit/crawl a webpage that has already been crawled. This is how I avoided *crawler trap*(Infinite Loop)

**QUESTION 3: What is the seed URL(s)?**

    **Answer:** The seed URL is: http://www.ebay.com/sch/samsung-galaxy

**QUESTION 4: How did you manage to only collect the webpages respecting the template(s) in Task 2? How did you discard irrelevant pages?**

    **Answer:**

- I observed that ebay has a way/pattern in naming the file and also storing it. All webpages that are about a particular product instead of a list of items are from a particular folder called "itm" and all URLs of webpages that match the query usually have the words "Samsung" or "Galaxy" in them.
- I used Jsoup to connect to the webpages and extract HTML.
- The price information of the products is usually embedded in the content attribute of a meta tag with a property attribute having the value "og:Description".
- I made sure that the webpages are from the "itm" folder and also the products are Samsung cell phones by parsing the URL.
- I then parsed the content attribute value of the meta tag to make sure that the price of the product complies with the query.
- I also made sure that the crawler does not go out of the domain and visit webpages on sites other than ebay by checking the domain name of each URL using the getHost() method in the URI class(Java)

**QUESTION 5: If you were not able to collect 1000 pages, explain what the problem was.**

    **Answer:** I was able to collect 1000 webpages.

## *ADDITIONAL:

- I have saved the webpages as 1.html, 2.html and so on.
- This does not provide much information about the URLs of the webpages. So I saved the URLs of the pages to a file named "urlList.txt".
- I am including this file along with the other folders I am supposed to submit. (This is saved in the source folder)
- The name of the source code file is Specific.java
  There is another file called trial.java in the source folder. Please ignore it.