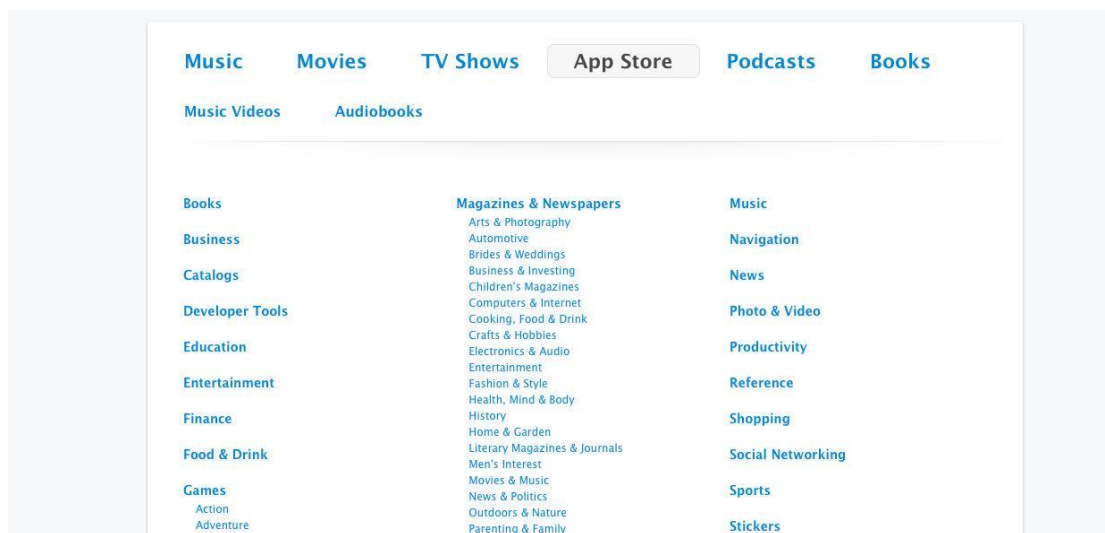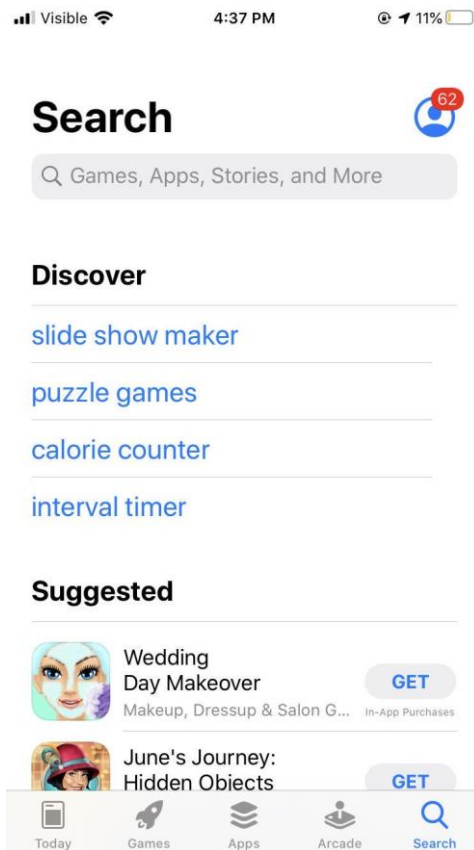# CS6200-IR-Final Project Report

Author: Yuchen Xie

## Background and Problems

As we can see from graph1 which is the official site of apps.apple.com, there is no search user interface for the applications in itunes store. The official site is where we can see all the applications' complete data, each one of them is displayed on apple's official site, but we can not search on these apps. In other words, the search interface on itunes and mobile app store does not have transparent progress to users.

Complete data is available on the official site, and the project's goals is to build a search engine on these data.



1 *Official Site of apple*
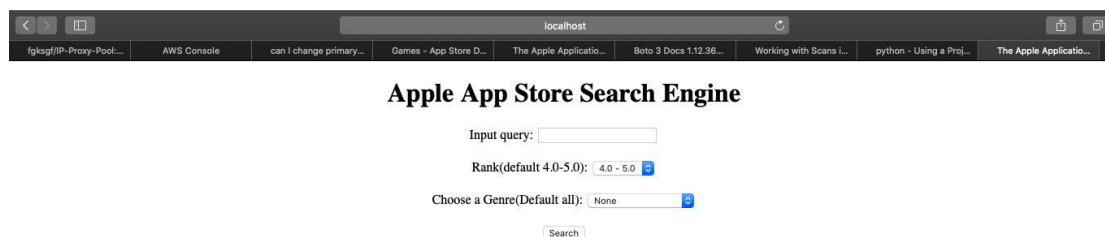
**2** *Itunes search page on mobile phone*

On the mobile page, the users only have a input search interface, no other settings can be made, like genre setting or ranking settings. When we are trying to search some game and other tools in photos & videos popping out would make people frustrating. In the meantime, we do not know how apple ranks the results that it's returning to us, some applications got returned top of the list with no reason, maybe developers can trade higher recommendation ranks with apple, there are always some promotion on the page, we don't know, but the ranking is clearly vague.

When we are actually using the search function apple has offered us, sometimes we forget about the name of the application, we would try to type in as much characters as we remember, but the search can only be

based on names of the apps and other obvious aspects like zombie games or dungeon games. We want to search all the applications by all information, we want to search on developers, description and genres. Thus, the goal of the project is to come up with a more comprehensive ad-hoc search engine, plus some of the vertical search function.

For this search engine, I want it to has as much data as possible, with no other influence, the engine would rank the return results only on its similarity and relevance to the query and previous settings.

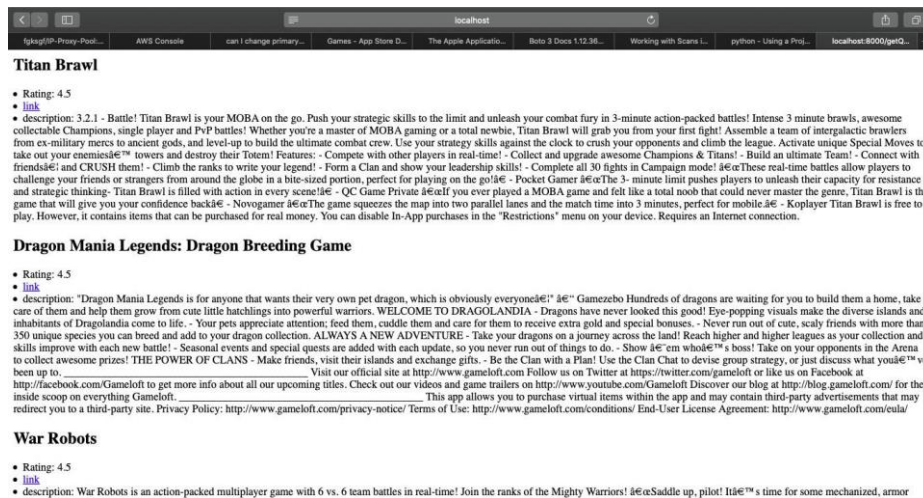## Overview of the search engine



**3** *Search Engine page*

The page of my search engine is showed above, it has one query input, which would allow users to type in whatever they want.

In addtion, there are two bar tools that would allow users to select which genre and which rank interval the users want to set before search. If the genres and rankings has been put, the search tends to cost less time,

since the data set is smaller, and there is less things to compute.



**4** *Search result*

The search results of the engine have simple structure, it would return all the information that users want, actually what I chose to show on the result page is limited, in storage I have all kinds of data like how many people rated this app, what age should be appropriate for the application, but the users only need some information for them to see. In such case, I chose to show Title, ranking and description. The link here can lead users to the application's actual apple store page.
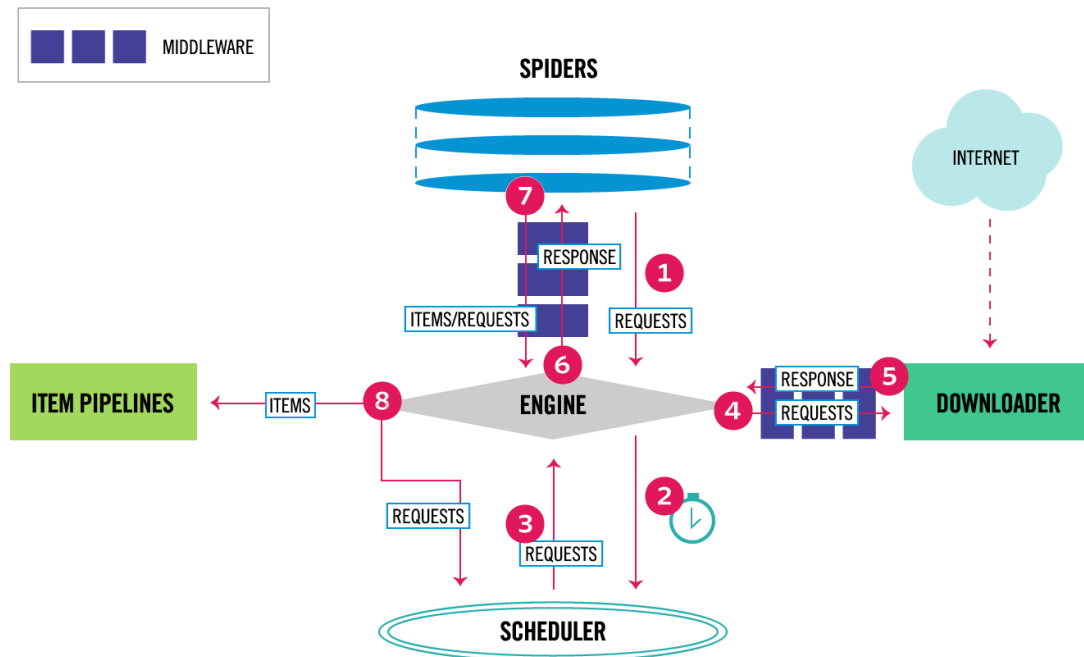
## How it works

The whole project can be split into 4 parts to demonstrate: **Crawler**, **Data storage and synchronize**, **Ranking** and **web server**. Let me explain the project by detail.

## 1. Crawler

For crawler I use scrapy as crawler framework, which has some

advantages like default concurrency, download delays and rotating IP proxy pool. With customized middlewares and pipelines, users can optimize their crawlers' performance and alter them according to actual usage.



5 *Scrapy Architecture*

Due to the IP banning policy toward robots of Apple, I need to change the download delay to be 10 seconds per page, if I download the pages faster, like 300 pages per minutes, my IP would be banned for 3 to 4 hours, which is unacceptable. Apparently there's a limit number of pages on Apple that each IP can open in a fixed period.

For the crawler's each request, I would disguise the crawler with customized user-agent:

```python
def parseItemTop(self, response):

    item = response.meta['item']
    currentUrl = response.url

    data = response.xpath('//div[@id="selectedcontent"]//li//a')
    headers = {'User-Agent': 'Mozilla/5.0 (X11; Linux x86_64; rv:48.0) Gecko/20100101 Firefox/48.0'}

    for letterNum in range(26):
        currentLetter = chr(letterNum + 65)
        for pageNum in range(1, 21):
            params = {'letter':currentLetter,'page':pageNum}
            yield Request(
                url = currentUrl + '?letter='+currentLetter+'&page=' + str(pageNum) + '#page',
                callback = self.parseClassPage,
                meta={'item':item}
                )
    return
```

6 *Crawler-User agent*

After downloading the page, I used scrapy built-in function xpath

to parse the response text by selecting elements with specific id or class.

```python
def parseSpecificPage(self, response):

    item = response.meta['item']

    # The name of the app
    name = response.xpath('//h1[@class="product-header__title app-header__title"]//text()').get()
    name = name.strip(' \n')
    if not name:
        name = "No name?"

    # The developer of the app
    developerSelec = response.xpath('//h2[@class="product-header__identity app-header__identity"]//a//text()')
    developer = ''
    if developerSelec:
        developer = developerSelec[0].get()
        developer = developer.strip(' \n\t')
    if not developer:
        developer = 'Unknown'

    # The rating of the app
    ratingSelec = response.xpath('//figcaption[@class="we-rating-count star-rating__count"]//text()')
    rating, rateNum = "0", "0"
    if ratingSelec:
        rating, rateNum = ratingSelec.get().split(',')
```

7 *Xpath example*

When the parsing is done, the crawler would pass an *item* to

pipeline, then it would be stored into database.

```
pageItem = SpecificPageItem()
pageItem['description'] = description
pageItem['name'] = name
pageItem['url'] = response.url
pageItem['rating'] = rating
pageItem['rateNum'] = rateNum
pageItem['developer'] = developer
pageItem['classification'] = response.meta['item']['name']
pageItem['uid'] = response.url.split('/')[-1][2:]
yield pageItem
```

**8** *Scrapy Item*

We can see from graph 7, there are 8 elements in the item, which would be the features/keys of the data that we store into the database.

There are several things need to be set previously in setting of scrapy, like DOWNLOAD_DELAY, ROBOTSTXT_OBEY and concurrency etc., if pipelines or items are being implemented, we need to set them active in setting as well.

## 2. Data Storage and Synchronization

When the crawler has finished downloading the web page and parsing it, it will be passed to pipeline of scrapy, in the pipeline, I would do the tokenization, and then store all I got into database: AWS Dynamodb.



**9** *AWS Dynamodb*

The biggest reason I am using Dynamodb is, cheap. It has some free tier and credit for students. Other consideration like it is a relational database which can prevent duplicate by setting the primary key which would overwrite the records with the same key attribute; It has user interreact page which can offer us better way to monitor the tables; Dynamodb provides two different ways to retrieve the data, scan and query, that suits different scenarios; Convenient API in Python, etc.

Before putting the whole data into the database, I will preprocess the data, which is stemming and tokenizing, since I am using cosine similarity to rank the return result, it will improve the efficiency if I parse the tokens at the time when crawler is waiting to download another page, this action would save some time. In stemming and tokenizing I am using nltk to perform such actions.

```python
import boto3

from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
from collections import Counter

class AppleapplicationPipeline(object):

    def open_spider(self, spider):
        dynamodb = boto3.resource('dynamodb', region_name='us-east-2')
        self.table = dynamodb.Table("apple")
        self.table_token = dynamodb.Table("apple-tokens")
        stoplist = ' '.join(open('stoplist.txt', 'r').readlines())
        self.stop_words = set(stoplist.split())
        self.ps = PorterStemmer()

    def process_item(self, item, spider):
        res = self.table.put_item (Item = dict(item))

        sentence = (item['description'] + item['name']).strip()
        tokens = word_tokenize(sentence)
        tokens = [self.ps.stem(token) for token in tokens]
        token_cnt = Counter(tokens)
        final_token = ''
        for token in token_cnt:
            if token in self.stop_words:
                continue
            else:
                final_token += (token)
                final_token += ":"
                final_token += str(token_cnt[token])
                final_token += ","
        token_item = {'id':item['uid'], 'tokens':final_token}

        res1 = self.table_token.put_item(Item = dict(token_item))
        return item
```

**10** *Dynamodb put item*

Due to the big dataload on the official site of Apple, the crawler is consistent running, the data is continually being put into the database. When the search engine is running, if the engine retrieves the data online from the database, it would not only be incredibly slow, but also it can result in dirty read and conflicts.

Thus, I would daily synchronize the data into local storage files in order to accelerate the search and compute. (You can see searchEngineServer / refresh_app.py). All data are store in searchEngineServer/data_genre/*.

The data that I would query actually looks like this:



**11** *Actuall Data*

The Data stored in local are in simpler form, only contains import attributes which would help us to identify the data, like name and id, and other attributes which would help us to rank the data like the ranking on app store. This option would make data stored locally with smaller size and computed more conveniently.

## 3. Ranking

For the results ranking, I am using the simplest measures to rank them, TF-IDF and cosine similarity.

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│ Parse Query  │ ───▶ │ Get Related  │ ───▶ │ Compute TF-IDF│
│              │      │ Data         │      │ Weights      │
└──────────────┘      └──────────────┘      └──────────────┘
                                                    │
                                                    ▼
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│ Data Return  │ ◀─── │ Get the scores│ ◀─── │ Compute Cosine│
│              │      │ for each data │      │ Similarity   │
│              │      │ entry,        │      │              │
│              │      │ rank them     │      │              │
└──────────────┘      └──────────────┘      └──────────────┘
```

Process are as above. Actual code about data computation can be viewed in searchEngineServer/parser.py. Code about data retrieve from local storage can be viewed in searchEngineServer/dydb.py. Stemming and Tokenizing are based on nltk, a python package, tf-idf and cosine care computed by raw code without other packages.

The query would be consisting of a text query, a rank interval (If set), a targeted genre (If set). According to the specific requirement of the query, the ranking function would first request the data from local storage, and then start to compute the scores of each data entry. After the computation is done, the results would be ranked by their relevance firstly, and then sorted by their actual store ranking secondly, which would give users better experience.

The process is relatively easy but functional. The tokens of each data

entry are parsed by one application's name, description and developers, which can ensure better search upon attributes other than names.

The returned data are in simple form, only has name, id and rank, so that the computation and local storage pressure would be lighter. After we ranked the results, we have to retrieve the data entry in complete form which has all attributes from database, then we can return them to the server.

## 4. Web Server

The server is based on Python's baseHTTPServer, which is a simple http server packet supported by python2, no longer a separate package in Python3.

There's no much to talk about the server. The query is passed to backend with POST method, after all the process being done, the results would be returned and showed directly on the web page.

# Conclusion

This project let me build the search engine from data crawling, although I thought about using other search engine like Elastic Search and Solr, I decided not to apply other techniques since we have already done the tokenizing and ranking assignments, I think those can really help the final project.

I have learned great knowledge about IR and search engine from this project and this course, I appreciate the help I got from TA Michelle and Gustavo, and Professor Alonso and Professor Ricardo has offered great deal of help, I deeply thank you for educating me such wonderful course.

Github link of the project:

*https://github.com/Yuchen112211/CS6200-IR-AppleSearchEngine*