

# COL334 - Assignment 2

## HTTP

Akshay Kumar Gupta  
2013CS50275

Barun Patra  
2013CS10773

Haroun Habeeb  
2013CS10225

We made a python script ( Q4.py ) on the two Object Tree files, namely, www.nytimes.com.objt and www.vox.com.objt. All data collection was performed on the same network, i.e, Airtel Broadband

### Q4a. Overview

The script takes in as input an **Object Tree** or a **HAR file**, the **maximum allowed TCP connections** per domain, and the **maximum number of objects** per TCP connection. If the input is a **HAR file**, it is first converted to an **Object Tree** and then processed.

Each layer of the **Object Tree** is processed sequentially, i.e objects of parents are downloaded before the objects of the children. At every layer, objects belonging to a domain are passed to the thread.

In the thread, at any given time,  $\text{TCP connections} \leq (\text{TCP connections})_{max}$  are opened in parallel to the domain in consideration, with each TCP connection requesting and downloading  $\text{objects} \leq (\text{objects})_{max}$ . Downloading across domains is done in parallel on separate threads.

### Q4c.

The parameters obtained from Q3c. were :

- **nytimes:**

- TCP Connections:
- Number of Objects:
- Browser Time:
- Downloader Time:

- **vox:**

- TCP Connections:

- Number of Objects:
- Browser Time:
- Downloader Time:

In conclusion,

#### Q4d. Running time for different parameters

#### Q4e. Improving Running Time:

Some points of improvement are as follows:

- **Managing Data Distribution:** In the implemented downloader, objects belonging to the same domain are distributed to different threads, taking into consideration the maximum objects a connection can handle at any given time. A better way to carry out this distribution would be to ensure every TCP connection gets (more or less) an even split in terms of the size of data being downloaded on the same, to ensure maximum benefit of parallelization. The file sizes can be taken from the HAR file.
- **Scheduling of Threads:** In the implemented version of the downloader, for a given domain, at most  $K$  ( $K = (\text{TCP connections})_{max}$ ) threads are spawned, each establishing its own TCP connection. Then, the downloader waits for all the threads to join, before continuing, in order to ensure that at most  $K$  connections are opened at any given time to the server. This can be improved by not waiting for the thread join. Instead, at the end of each thread execution, another thread can be spawned then to handle the next batch of requests. This can be used, in conjunction with the aforementioned improvement to ensure the newly spawned thread takes up an optimal number of objects instead of the maximum allowed number of objects, so as to, again, maximize the benefits of parallelization.
- **Number of TCP connections opened to a domain:** The HAR file can be used to infer the maximum number of TCP connections that a server allows in parallel(i.e if the server places a cap on the maximum allowed connections). Using this information, we can dynamically allocate the value of  $K$  for different domains for the case of large files; so that the overhead of creating new threads and making a new TCP connection does not outweigh the benefit of parallelization.
- **Fractional Requests:** In the case that the server sets a cap on the bandwidth allocated to a single TCP connection, large objects can be requested in chunks on separate TCP connections, with each TCP connection requesting for a fraction of the file.