

CSE207 Introduction to Networking Project 2 - Report

Chen Qianshan

1508670

qianshan.chen15@student.xjtlu.edu.cn

Program: CST

Lecturer: Charles Fleming

Due: Tuesday, 11:55 am, Dec.12 2017

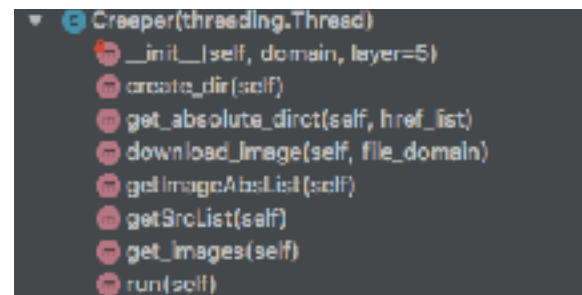
Image Crawler Application.

A one-page report specifying which parts of the project you implemented, known bugs, and citing any sources you take more than 10 lines of code from.

1. This is an object-oriented program. The main class Creeper is inheritance from threading.Thread, a multi-thread build in package of python.

2. There are 8 methods in Creeper class.

Shown in right. To initialize the class, **domain**, and **layer** are the two arguments domain is the aiming page waiting to be crawled, layer is the current layer, defaults as 5. To crawling each layer, the run() methods in the class will create and start a new Creeper class in one new thread. when the layer



3. Known Bugs:

- When crawling <https://www.baidu.com> the HTTP response body part cannot be decode with gzip or zlib.
- When the picture is big, it is possible to fail download it. Because of the socket connection timeout has been set to 2 seconds.
- User should enter the full [https/http://xxx.xxx](https://xxx.xxx) part for connection. When the web is a full ssl-site the program cannot automatically jump from http to https domain.

4. Citing source more than 10 lines

a. Header part: These codes are get from the chrome browser to simulate a browser:

```
se.send(b"GET " + path.encode() + b" HTTP/1.1\r\n")
se.send(b"Host: " + host.encode() + b"\r\n")
se.send(b"Connection:keep-alive\r\n")
se.send(b"Cache-Control:max-age=3600\r\n")
se.send(b"Accept:text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,*/*;q=0.8\r\n")
se.send(b"Upgrade-Insecure-Requests: 1\r\n")
se.send(
    b"User-Agent:Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/49.0.2623.110 Safari/537.35\r\n")
se.send(b"Accept-encoding:gzip, deflate, sdch\r\n")
se.send(b"Accept-Language:zh-CN,zh;q=0.8\r\n")
se.send(b"\r\n")
```

```
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/Cat_eating_a_rabbit.jpg.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/testImages/kittens-cat-cat-puppy-rush-45178.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub4/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/testImages/kittens-cat-cat-puppy-rush-45178.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/Cat_eating_a_rabbit.jpg.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/testImages/kittens-cat-cat-puppy-rush-45178.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub4/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/testImages/kittens-cat-cat-puppy-rush-45178.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub4/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/sub4/Cat_eating_a_rabbit.jpg.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/testImages/kittens-cat-cat-puppy-rush-45178.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/Cat_eating_a_rabbit.jpg.jpg
Download from http://cse.xjtlu.edu.cn/classes/CSE285/sub1/subsub1/, is completed!
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub4/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/sub4/Cat_eating_a_rabbit.jpg.jpg
At >>> cse.xjtlu.edu.cn/classes/CSE285/sub4/
Downloaded >>> http://cse.xjtlu.edu.cn/classes/CSE285/testImages/kittens-cat-cat-puppy-rush-45178.jpg
```