

COMP6331 Computer Networks
Assignment 2 - Web Crawler
u6743886 Jinpei Chen

Report:

Item	What to report	Results							
1	The total number of distinct URLs found on the site (including any errors and redirects)	46 (including off-site URLs) 43 (excluding off-site URLs)							
2	The number of html pages	32							
3	The number of non-html objects on the site (e.g. images)	9							
4	The smallest html pages, and size	comp3310.ddns.net:7880/B/23.html, size: 1443							
5	The largest html pages, and size	comp3310.ddns.net:7880/C/307.html, size: 7261							
6	The oldest modified page, and its date/timestamps	comp3310.ddns.net:7880A/10.html, modified time: Tue, 01 Jan 2019 05:05:00 GMT							
7	The most-recently modified page, and its date/timestamps	comp3310.ddns.net:7880C/30.html, modified time: Sun, 05 May 2019 01:01:00 GMT							
8	A list of invalid URLs (not) found (404) (including off site invalid URLs)	<table><tr><td>comp3310.ddns.net:7880/F/10.html</td></tr><tr><td>comp3310.ddns.net:7880/F/20.html</td></tr><tr><td>comp3310.ddns.net:7880/F/30.html</td></tr><tr><td>comp3311.ddns.net:7880/B/207.html</td></tr><tr><td>comp3310.ddns.net/A/19.html</td></tr><tr><td>comp3310.ddns.net/B/29.html</td></tr><tr><td>comp3310.ddns.net/C/39.html</td></tr></table>	comp3310.ddns.net:7880/F/10.html	comp3310.ddns.net:7880/F/20.html	comp3310.ddns.net:7880/F/30.html	comp3311.ddns.net:7880/B/207.html	comp3310.ddns.net/A/19.html	comp3310.ddns.net/B/29.html	comp3310.ddns.net/C/39.html
comp3310.ddns.net:7880/F/10.html									
comp3310.ddns.net:7880/F/20.html									
comp3310.ddns.net:7880/F/30.html									
comp3311.ddns.net:7880/B/207.html									
comp3310.ddns.net/A/19.html									
comp3310.ddns.net/B/29.html									
comp3310.ddns.net/C/39.html									

9	A list of on-site redirected URLs found (30x) and where they redirect to	<table><tr><th>Source</th><th>Redirect to</th></tr><tr><td>comp3310.ddns.net:7880/C/3 A.html</td><td>comp3310.ddns.net/A/19.html</td></tr><tr><td>comp3310.ddns.net:7880/A/1 A.html</td><td>comp3310.ddns.net/B/29.html</td></tr><tr><td>comp3310.ddns.net:7880/B/2 A.html</td><td>comp3310.ddns.net/C/39.html</td></tr></table>	Source	Redirect to	comp3310.ddns.net:7880/C/3 A.html	comp3310.ddns.net/A/19.html	comp3310.ddns.net:7880/A/1 A.html	comp3310.ddns.net/B/29.html	comp3310.ddns.net:7880/B/2 A.html	comp3310.ddns.net/C/39.html				
Source	Redirect to													
comp3310.ddns.net:7880/C/3 A.html	comp3310.ddns.net/A/19.html													
comp3310.ddns.net:7880/A/1 A.html	comp3310.ddns.net/B/29.html													
comp3310.ddns.net:7880/B/2 A.html	comp3310.ddns.net/C/39.html													
10	A list of off-site URLs found (either 30x redirects or html references), and whether those sites are valid	<table><tr><th>Off-site URLs</th><th>Flag</th></tr><tr><td>comp3310.ddns.net/A/19.html</td><td>Invalid</td></tr><tr><td>comp3310.ddns.net/B/29.html</td><td>Invalid</td></tr><tr><td>comp3310.ddns.net/C/39.html</td><td>Invalid</td></tr><tr><td>comp3311.ddns.net:7880/B/207.ht ml</td><td>Invalid</td></tr><tr><td>www.canberratimes.com.au</td><td>Valid</td></tr></table>	Off-site URLs	Flag	comp3310.ddns.net/A/19.html	Invalid	comp3310.ddns.net/B/29.html	Invalid	comp3310.ddns.net/C/39.html	Invalid	comp3311.ddns.net:7880/B/207.ht ml	Invalid	www.canberratimes.com.au	Valid
Off-site URLs	Flag													
comp3310.ddns.net/A/19.html	Invalid													
comp3310.ddns.net/B/29.html	Invalid													
comp3310.ddns.net/C/39.html	Invalid													
comp3311.ddns.net:7880/B/207.ht ml	Invalid													
www.canberratimes.com.au	Valid													