

# Web Crawler Design Grading Sheet

Name: \_\_\_\_\_

TA: \_\_\_\_\_

Score    Max Possible

## Data Structures

Detailed description of the following data structures, including justification for each choice

- \_\_\_\_\_                      **Data structure for tracking unprocessed pages**
- \_\_\_\_\_                      **Data structure for tracking already processed pages**
- \_\_\_\_\_                      **Data structure for managing stop words**
- \_\_\_\_\_                      **Data structure for storing word index**
- \_\_\_\_\_                      **Data structure for implementing robot filter**

## Class Responsibilities

- \_\_\_\_\_                      Drive the overall crawling process
- \_\_\_\_\_                      Store the URL and description for a page
- \_\_\_\_\_                      Store index that maps words to pages
- \_\_\_\_\_                      Keep track of yet-to-be indexed pages
- \_\_\_\_\_                      Keep track of already indexed pages
- \_\_\_\_\_                      Load and store stop words
- \_\_\_\_\_                      Robot Filter
- \_\_\_\_\_                      Distinguish between HTML and non-HTML links
- \_\_\_\_\_                      Distinguish between in-scope and out-of-scope links
- \_\_\_\_\_                      Resolve relative URLs
- \_\_\_\_\_                      Parses words, links, and description from HTML pages
- \_\_\_\_\_                      Populate word index
- \_\_\_\_\_                      Generate HTML index pages

## Algorithms

Top-level code for the following algorithms:

_____	Main driver for the crawling process
_____	HTML parser
_____	HTML index page generation

## Design Quality

_____	Cohesive classes and methods
_____	Effective information hiding
_____	Effective class, method, and variable names
_____	Clear, easy-to-read document

\_\_\_\_\_ **Total**