

Multi-threaded Web Crawler Built with C

Possible Files of Interest for Code Sample Review

- `findpng2.c` is the main driver program and is about 350 lines of code
- `stack.c` and `hash.c` (300 lines combined) define two memory-safe data structures
- Please feel free to look at any of the files you feel are appropriate. I included descriptions of each file below.
- Please also feel free to build using `make`

Project context

Purpose

Starting from a seed URL, recursively crawl to all other URLs linked on the web page of the seed URL, searching for PNGs. Stop when we find a specified number of PNGs or run out of URLs to visit. The program outputs all the pngs found in a `png_urls.txt` file at the end of the program. The user can also specify to output all the unique URLs the program crawled to a log file.

To speed up the search, the program launches multiple threads. The user may specify the number of runner threads that the program will launch. Each thread can process a different URL concurrent with other threads. The user may also specify the number of PNGs to find before stopping.

We use multiple threads primarily because network downloads can be slow compared to the CPU. Each thread synchronously waits for the network download of a page. The CPU can process other threads while this thread is blocked on the network download. If the program runs on one thread, the program cannot make progress while the one thread waits for its network download to complete.

Files

- `findpng2.c` : main driver program
 - takes input
 - runs the specified number of threads
 - the threads perform the crawl
 - handles synchronization between threads
 - handles clean-up of threads at the end of the program
- `stack.c` :
 - a memory-safe dynamic stack that holds strings
 - used for holding URLs to crawl (the `frontier`)
- `p_stack.c` :
 - a memory-safe dynamic stack that holds pointers
 - used in `hash.c` for holding pointers to memory for later deallocation
- `hash.c` :
 - a memory-safe hash set that holds strings as keys
 - used for holding visited URLs to prevent cycles in the crawling process
- `curl_xml.c` :
 - utility functions for downloading web pages using cURL
 - utility functions for processing XML pages using libxml
 - utility function for checking if a png file is a valid png
 - used for downloading web pages, searching for URLs listed on the pages, and checking the validity of found pngs

External libraries used

- cURL (<https://curl.se/libcurl/>)
- libxml (<https://github.com/GNOME/libxml2>)

Notes

- I wrote this program for the course, ECE 252 (Systems Programming and Concurrency).
 - I received 100% on this project (marked based on functionality, speed, lack of deadlocks, and proper memory management).
 - I received 100% on the course (60% final exam, 40% projects).
- The program is named `findpng2` because I have a few other versions that perform the same task with different techniques (e.g. using processes and interprocess communication for synchronization; or using asynchronous cURL).

Building

- run `make` in this directory

Usage

```
findpng2 [OPTION]... [ROOT_URL]
```

- options:
 - -t=NUM - the program will create NUM threads to crawl the web (default: 1)
 - -m=NUM - the program will find up to NUM unique PNG URLs (default: 50)
 - -v=LOGFILE - if specified, program will log the unique URLs visited in a file named LOGFILE
- output:
 - on terminal, `findpng2` execution time: S seconds

- the program will create a `png_urls.txt` file containing all the valid PNG URLs found
- if a log file is specified, the program will create a `<LOGFILE>` file containing all unique URLs visited
- for example, `./findpng2 -t 1 -m 1 -v test.txt https://www.cleanpng.com/static/img/logo.png` will launch 1 thread to find 1 png starting from the URL `https://www.cleanpng.com/static/img/logo.png` and output all visited URLs into `test.txt`. Note that since the seed URL is a png itself, the crawl will end immediately after the first visit.

Code Best Practice Notes

- `findpng2.c` uses global variables, which is not good practice; however, the course instructors recommended them for simplicity for the sake of the small project.
 - I could instead pass to the thread runner a pointer to a struct containing all the global variables.
- `findpng2.c`'s thread runner function exits with a `return` rather than `pthread_exit`, which is the recommended way to exit a thread. However, valgrind will report a small number of reachable blocks for threads that use `pthread_exit`. This is not a real issue in most applications. However, since a part of the marking scheme for this project was to achieve 0 reachable blocks at the end of program, my only option was to use a `return`.

Algorithm overview

findpng2.c

Global data structures

- `frontier`: a stack of discovered URLs that we have yet to process, shared by all threads
- `pngs`: a stack of all pngs found so far
- `visited`: a hash set of URLs we have visited (so we don't crawl repeat URLs)

Synchronization

- `frontier_mutex`: a lock for accessing `frontier` and other accounting variables related to the state of the crawl (whether if the crawl is done; number of threads waiting for the frontier to be non-empty; number of threads processing a URL)
 - used for the `frontier_empty` condition variable as well
- `pngs_mutex`: a lock for accessing `pngs`
- `visited_mutex`: a lock for accessing `visited`
- `frontier_empty`: a condition variable that threads will wait on when `frontier` is empty
 - when another thread adds to `frontier`, it will broadcast to wake up the sleeping threads
 - alternatively, a thread may broadcast when the program is finished (no more URLs we can recursively crawl or we have found `num_pngs_to_find` pngs) so that sleeping threads can wake up and exit

Thread runner

- The main function launches `t` number of these thread runners.
- The runner does the following in an infinite loop.
 - With lock `frontier_mutex`:
 - Check if our crawl is finished. If so, signal all threads to exit.
 - If `frontier` is empty and there are no threads currently processing a URL, we don't and won't have any more URLs to crawl: our crawl is finished.
 - If our crawl is finished, broadcast `frontier_empty` to wake up sleeping threads so they can exit.
 - Wait until there is a URL to crawl or our crawl is finished.
 - The thread does this by waiting on the condition variable `frontier_empty`.
 - The signal to wake up will come from a thread pushing a URL to `frontier` or a thread signalling our crawl is finished.
 - If our crawl is finished, exit loop.
 - Else, grab the first URL on `frontier` for processing.
 - With lock `visited_mutex`:
 - Check if we've already visited the URL.
 - If we have, don't progress further and return to the start of the loop.
 - If we haven't, add it to the `visited` hash set and continue on.
 - Download the URL's contents
 - If it is a HTML file, grab all URLs that it links to.
 - If it is a PNG file, determine if it is a valid png.
 - (If HTML file) With lock `frontier_mutex`:
 - Push all URLs found onto `frontier`.
 - If there are any sleeping threads waiting for a non-empty frontier, broadcast on `frontier_empty`.
 - (If PNG file) With lock `pngs_mutex`:
 - Push URL into `pngs`.
 - If we hit `num_pngs_to_find`:
 - With lock `frontier_mutex`:
 - Set a global variable indicating our crawl is finished, and broadcast on `frontier_empty`.
- When the runner exits the infinite loop, clean up any data structures used (memory deallocation, network libraries).

stack.c and p_stack.c

- Memory-safe stacks
- Resizing is done by allocating a larger chunk of memory, moving the old items over, and deallocating the old memory.

hash.c

- A hash set that holds strings as keys

Data structures

- `hmap` : `hsearch_data` struct used as the underlying hashmap
- `elements` : an array of pointers to keys currently in the hash set
 - used for linear access to keys, moving keys over on resize, and for deallocating the keys at destruction
- `ps` : a stack of pointers to keys that were used when searching `hmap`
 - used for deallocating the strings at destruction
 - note we don't deallocate the keys/strings immediately after the search since `hsearch_data` may hold on to the pointer to the key and use it in future operations after the search

Adding a key

- If full, resize.
- Create a copy of the passed-in key/string.
- Add the key (copy) to `hmap`.
- Add the key (copy) to `elements`.

Searching for a string

- Create a copy of the passed-in key/string.
- Search for the key (copy) in `hmap`.
- If we found the key:
 - Add the key (copy) to `ps`.
- Else:
 - Free memory for the key (copy).
- Why we add the key to `ps`:
 - If the key is found in `hmap`, `hmap` may use the memory of the passed-in string for future references instead of the string used when the key was initially added to `hmap`.
 - `hsearch` does not deallocate memory upon destruction. Thus, we keep track of `ps` : the strings to deallocate at destruction.

Resizing

- Destroy and reinitialize `hmap` with a larger size.
- Allocate a larger `elements`.
- Transfer keys from the old `elements` to the new `hmap` and `elements`.
- Deallocate old `elements`.

multithreaded/findpng2.c

```
1 #include "findpng2.h"
2
3 // I defined this for the code sample review
4 // so if you want to make and run the program,
5 // you can see the urls being crawled on command line
6 #define DEBUG_URL_PRINT
7
8 /* -- Global Variables -- */
9 // global collection of urls to be crawled by runner threads
10 STACK *frontier;
11 // pngs found (png urls)
12 STACK *pngs;
13 // urls visited
14 HSET *visited;
15 // whether we are done with the entire crawl
16 bool done;
17 // number of threads waiting for a non-empty frontier
18 size_t num_waiting_on_url;
19 // number of thread runners currently processing a url
20 size_t num_running;
21 // number of pngs to find before stopping
22 int num_pngs_to_find;
23 /* ----- */
24
25 /* -- Synchronization -- */
26 // condition variable for threads to wait on when the frontier is empty
27 pthread_cond_t frontier_empty;
28 // lock for frontier, done, num_waiting_on_url, and num_running;
29 // also used for frontier_empty
30 pthread_mutex_t frontier_mutex;
31 // lock for pngs stack
32 pthread_mutex_t pngs_mutex;
33 // lock for visited hash set
34 pthread_mutex_t visited_mutex;
35 /* ----- */
36
37 /**
38 * @brief initialize global variables and synchronization variables
39 */
40 void initialize_global()
41 {
42     frontier = malloc(sizeof(STACK));
43     memset(frontier, 0, sizeof(STACK));
44     init_stack(frontier, STACK_SIZE);
45
46     visited = malloc(sizeof(HSET));
47     memset(visited, 0, sizeof(HSET));
48     init_hset(visited, HMAP_SIZE);
```

```
49
50     pngs = malloc(sizeof(STACK));
51     memset(pngs, 0, sizeof(STACK));
52     init_stack(pngs, STACK_SIZE);
53
54     done = false;
55     num_waiting_on_url = 0;
56     num_running = 0;
57
58     pthread_cond_init(&frontier_empty, NULL);
59     pthread_mutex_init(&frontier_mutex, NULL);
60     pthread_mutex_init(&visited_mutex, NULL);
61     pthread_mutex_init(&pngs_mutex, NULL);
62 }
63
64 /**
65 * @brief cleanup global variables and synchronization variables
66 */
67 void cleanup_global()
68 {
69     cleanup_stack(frontier);
70     free(frontier);
71     frontier = NULL;
72
73     cleanup_hset(visited);
74     free(visited);
75     visited = NULL;
76
77     cleanup_stack(pngs);
78     free(pngs);
79     pngs = NULL;
80
81     pthread_cond_destroy(&frontier_empty);
82     pthread_mutex_destroy(&frontier_mutex);
83     pthread_mutex_destroy(&pngs_mutex);
84     pthread_mutex_destroy(&visited_mutex);
85 }
86
87 /**
88 * @brief runner function that crawls urls in the global frontier
89 * @param _ void*: not used; only defined to satisfy thread API
90 * @return NULL
91 * @details
92 * Any number of runner threads can be started.
93 * The runner function manages concurrency.
94 * The runner function assumes that all global structures and
95 * synchronization variables are initialized.
96 * The runner function will stop once there are no more urls to crawl
97 * or when we've found num_pngs_to_find pngs.
98 * The runner function does not clean up global variables.
```

```
99     */
100    void *runner(void *_)
101    {
102        /* -- Initialize cURL easy handle -- */
103        CURL *curl_handle = curl_easy_init();
104        if (curl_handle == NULL)
105        {
106            fprintf(stderr, "curl_easy_init: returned NULL\n");
107            exit(1);
108        }
109        /* ----- */
110
111        /* -- Defining variables used in the loop -- */
112        // response code from accessing url
113        long response_code;
114        // content type of data at url (e.g. HTML, PNG)
115        int content_type = DEFAULT_TYPE;
116        // the current url the thread is crawling
117        char *url_to_crawl = NULL;
118        // urls found on the web page visited; we will add these to the frontier
119        STACK *urls_found = NULL;
120        // if we have cleaned urls_found
121        bool cleaned_urls_found = false;
122        /* ----- */
123
124        while (true)
125        {
126            /* -- Cleanup structures from last iteration and re-initialize -- */
127            if (urls_found != NULL)
128            {
129                if (!cleaned_urls_found)
130                {
131                    cleanup_stack(urls_found);
132                    cleaned_urls_found = true;
133                }
134                free(urls_found);
135            }
136            urls_found = malloc(sizeof(STACK));
137            memset(urls_found, 0, sizeof(STACK));
138            init_stack(urls_found, 1);
139            cleaned_urls_found = false;
140
141            if (url_to_crawl != NULL)
142            {
143                free(url_to_crawl);
144                url_to_crawl = NULL;
145            }
146            /* ----- */
147
148            /* -- Check status of frontier and overall crawl -- */
```

```
149     pthread_mutex_lock(&frontier_mutex);
150 {
151     // If the crawl is finished, signal sleeping threads to
152     // wake up so they can exit
153     if (is_empty_stack(frontier) && num_running == 0)
154     {
155         done = true;
156         if (num_waiting_on_url > 0)
157         {
158             pthread_cond_broadcast(&frontier_empty);
159         }
160     }
161
162     // If there are no urls to crawl and the crawl is not done, wait
163     while (is_empty_stack(frontier) && !done)
164     {
165         ++num_waiting_on_url;
166         pthread_cond_wait(&frontier_empty, &frontier_mutex);
167         --num_waiting_on_url;
168     }
169
170     // If the crawl is finished, exit the loop
171     if (done)
172     {
173         pthread_mutex_unlock(&frontier_mutex);
174         break;
175     }
176
177     // Take the top url on the frontier
178     pop_stack(frontier, &url_to_crawl);
179
180     // Check if the url has been visited
181     pthread_mutex_lock(&visited_mutex);
182     {
183         // If the url has been visited, go back to the top of the loop
184         // (go to the next url in the frontier or if frontier is empty, wait)
185         if (search_hset(visited, url_to_crawl) == 1)
186         {
187             pthread_mutex_unlock(&visited_mutex);
188             pthread_mutex_unlock(&frontier_mutex);
189             continue;
190         }
191         // If the url has not been visited, mark it as visited.
192         // The thread will now process the url.
193         else
194         {
195             add_hset(visited, url_to_crawl);
196         }
197     }
198     pthread_mutex_unlock(&visited_mutex);
```

```
199         ++num_running;
200     }
201     pthread_mutex_unlock(&frontier_mutex);
202     /* ----- */
203
204 #ifdef DEBUG_URL_PRINT
205     printf("URL: %s\n", url_to_crawl);
206#endif
207
208     /* -- Crawl the url -- */
209     // download the contents at the url and process it
210     process_url(curl_handle, url_to_crawl, &content_type, urls_found, &response_code);
211     /* ----- */
212
213     /* -- Process url based on its contents -- */
214     if (is_processable_response(response_code))
215     {
216         // If the url was a HTML page, add all urls on that page to the frontier
217         if (content_type == HTML)
218         {
219             char *url_in_html = NULL;
220             while (pop_stack(urls_found, &url_in_html) == 0)
221             {
222                 // Add to the frontier and signal sleeping threads
223                 // (that a url is ready in frontier)
224                 pthread_mutex_lock(&frontier_mutex);
225                 {
226                     push_stack(frontier, url_in_html);
227                     if (num_waiting_on_url > 0)
228                     {
229                         pthread_cond_broadcast(&frontier_empty);
230                     }
231                 }
232                 pthread_mutex_unlock(&frontier_mutex);
233                 free(url_in_html);
234                 url_in_html = NULL;
235             }
236         }
237         // If the url was a valid PNG, add it to our collection of found pngs
238         else if (content_type == VALID_PNG)
239         {
240             pthread_mutex_lock(&pngs_mutex);
241             {
242                 push_stack(pngs, url_to_crawl);
243                 // If we've reached the maximum number of PNGs we want to find,
244                 // end the program
245                 if (num_elements_stack(pngs) >= num_pngs_to_find)
246                 {
247                     pthread_mutex_lock(&frontier_mutex);
248                     {
```

```
249         done = true;
250         pthread_cond_broadcast(&frontier_empty);
251     }
252     pthread_mutex_unlock(&frontier_mutex);
253 }
254 }
255 pthread_mutex_unlock(&pngs_mutex);
256 }
257 }
258 /* ----- */
259
260 /* -- The thread is no longer processing a url -- */
261 pthread_mutex_lock(&frontier_mutex);
262 {
263     --num_running;
264 }
265 pthread_mutex_unlock(&frontier_mutex);
266 /* ----- */
267 }
268
269 /* -- The thread is done all processing: clean up -- */
270 if (urls_found != NULL)
271 {
272     if (!cleaned_urls_found)
273     {
274         cleanup_stack(urls_found);
275     }
276     free(urls_found);
277 }
278
279 if (url_to_crawl != NULL)
280 {
281     free(url_to_crawl);
282     url_to_crawl = NULL;
283 }
284
285 curl_easy_cleanup(curl_handle);
286 /* ----- */
287
288 return NULL;
289 }
290
291 int main(int argc, char **argv)
292 {
293     /* -- command line inputs -- */
294     char *seed_url;
295     char *logfile = NULL;
296     size_t t = 1;
297     num_pngs_to_find = 50;
298 }
```

```
299     if (argc == 1)
300     {
301         printf("Usage: ./findpng2 OPTION[-t=<NUM> -m=<NUM> -v=<LOGFILE>] SEED_URL\n");
302         return -1;
303     }
304
305     seed_url = argv[argc - 1];
306
307     int c;
308     char *str = "option requires an argument";
309
310     while ((c = getopt(argc, argv, "t:m:v:")) != -1)
311     {
312         switch (c)
313         {
314             case 't':
315                 if (optarg == NULL)
316                 {
317                     t = 1;
318                     break;
319                 }
320                 t = strtoul(optarg, NULL, 10);
321                 if (t <= 0)
322                 {
323                     fprintf(stderr, "%s: %s > 0 -- 't'\n", argv[0], str);
324                     return -1;
325                 }
326                 break;
327             case 'm':
328                 if (optarg == NULL)
329                 {
330                     num_pngs_to_find = 50;
331                     break;
332                 }
333                 num_pngs_to_find = atoi(optarg);
334                 if (num_pngs_to_find < 0)
335                 {
336                     fprintf(stderr, "%s: %s >= 0 -- 'm'\n", argv[0], str);
337                     return -1;
338                 }
339                 break;
340             case 'v':
341                 if (optarg == NULL)
342                 {
343                     logfile = NULL;
344                     break;
345                 }
346                 logfile = malloc(sizeof(char) * FILE_PATH_SIZE);
347                 memset(logfile, 0, sizeof(char) * FILE_PATH_SIZE);
348                 strcpy(logfile, optarg);
```

```
349         break;
350     }
351 }
352 /* ----- */
353
354 /* -- initialize global variables and synchronization variables -- */
355 initialize_global();
356 /* ----- */
357
358 /* -- CURL global init -- */
359 curl_global_init(CURL_GLOBAL_DEFAULT);
360 /* ----- */
361
362 /* -- Initialize XML Parser -- */
363 xmlInitParser();
364 /* ----- */
365
366 /* -- Put the seed URL in the frontier -- */
367 push_stack(frontier, seed_url);
368 /* ----- */
369
370 /* -- Record time to be used for measuring speed -- */
371 double times[2];
372 struct timeval tv;
373 if (gettimeofday(&tv, NULL) != 0)
374 {
375     perror("gettimeofday");
376     exit(1);
377 }
378 times[0] = (tv.tv_sec) + tv.tv_usec / 1000000.;
379 /* ----- */
380
381 /* -- Create threads -- */
382 pthread_t *runners = malloc(t * sizeof(pthread_t));
383 memset(runners, 0, sizeof(pthread_t) * t);
384 if (runners == NULL)
385 {
386     perror("malloc\n");
387     exit(-1);
388 }
389 for (int i = 0; i < t; ++i)
390 {
391     pthread_create(&runners[i], NULL, runner, NULL);
392 }
393 /* ----- */
394
395 /* -- Wait for threads to finish -- */
396 for (int i = 0; i < t; ++i)
397 {
398     pthread_join(runners[i], NULL);
```

```
399     }
400     /* ----- */
401
402     /* -- Write to files -- */
403     // Write png urls
404     FILE *fpngs = fopen("./png_urls.txt", "w+");
405     if (fpngs == NULL)
406     {
407         fprintf(stderr, "Opening png file for write failed\n");
408         exit(1);
409     }
410     char *temp = NULL;
411     while (pop_stack(pngs, &temp) == 0)
412     {
413         fprintf(fpngs, "%s\n", temp);
414         free(temp);
415     }
416     fclose(fpngs);
417
418     // Write all urls visited into a log file if user desires
419     if (logfile != NULL)
420     {
421         char *logfile_name = malloc(sizeof(char) * FILE_PATH_SIZE);
422         memset(logfile_name, 0, sizeof(char) * FILE_PATH_SIZE);
423         sprintf(logfile_name, "./%s", logfile);
424         FILE *flogs = fopen(logfile_name, "w+");
425         free(logfile_name);
426         if (flogs == NULL)
427         {
428             fprintf(stderr, "Opening log file for write failed\n");
429             exit(1);
430         }
431         temp = NULL;
432         for (size_t i = 0; i < visited->cur_size; ++i)
433         {
434             fprintf(flogs, "%s\n", visited->elements[i]);
435             free(temp);
436         }
437         fclose(flogs);
438     }
439     free(logfile);
440     /* ----- */
441
442     /* -- Cleanup global variables and synchronization variables -- */
443     cleanup_global();
444     /* ----- */
445
446     /* -- Free threads -- */
447     free(runners);
448     /* ----- */
```

```
449
450     /* -- Clean up libraries used -- */
451     curl_global_cleanup();
452     xmlCleanupParser();
453     /* ----- */
454
455     /* -- Print time it took for crawl from the seed url -- */
456     if (gettimeofday(&tv, NULL) != 0)
457     {
458         perror("gettimeofday");
459         exit(1);
460     }
461     times[1] = (tv.tv_sec) + tv.tv_usec / 1000000.;
462     printf("findpng2 execution time: %.6f seconds\n", times[1] - times[0]);
463     /* ----- */
464
465     return 0;
466 }
```

multithreaded/findpng2.h

```
1 #include <stdbool.h>
2 #include "curl_xml.h"
3 #include "hash.h"
4 #include <pthread.h>
5
6 #define URL_SIZE 512
7 #define FILE_PATH_SIZE 512
8 #define STACK_SIZE 1024
9 #define HMAP_SIZE 1024
10
11 void *runner(void *args);
```

multithreaded/stack.c

```
1  /*
2  A dynamic stack holding strings
3  - note: this current implementation only resizes stack to increase size, never to decrease size
4  */
5
6 #include "stack.h"
7
8 /**
9  * @brief initialize stack with an initial size (capacity)
10 * @param p STACK*: a pointer to uninitialized memory
11 * @param stack_size size_t: initial capacity of the stack to be initialized
12 * @return 0 on success; 1 otherwise
13 */
14 int init_stack(STACK *p, size_t stack_size)
15 {
16     if (p == NULL || stack_size == 0)
17     {
18         return 1;
19     }
20
21     p->size = stack_size;
22     p->pos = -1;
23     p->items = (char **)malloc(stack_size * sizeof(char *));
24
25     memset(p->items, 0, stack_size * sizeof(char *));
26
27     return 0;
28 }
29
30 /**
31 * @brief push an item onto the stack; if the stack is full, resize the stack
32 * @param p STACK*: (pointer to) the stack the function will push item onto
33 * @param item char*: string to push onto stack
34 * @return 0 on success; 1 otherwise
35 */
36 int push_stack(STACK *p, char *item)
37 {
38     if (p == NULL)
39     {
40         return 1;
41     }
42
43     if (is_full_stack(p))
44     {
45         resize_stack(p);
46     }
47
48     ++(p->pos);
```

```
49     // strlen(item) + 1 for end of string character
50     p->items[p->pos] = malloc((strlen(item) + 1) * sizeof(char));
51     memset(p->items[p->pos], 0, (strlen(item) + 1) * sizeof(char));
52     strncpy(p->items[p->pos], item, strlen(item));
53
54     return 0;
55 }
56
57 /**
58  * @brief pop from the stack
59  * @param p STACK*: (pointer to) the stack the function will pop from
60  * @param p_item char**: pointer that will be populated with popped element
61  * @return 0 on success; 1 otherwise
62  * @note the caller is responsible for deallocating memory assigned to p_item
63 */
64 int pop_stack(STACK *p, char **p_item)
65 {
66     if ((p == NULL) || is_empty_stack(p))
67     {
68         return 1;
69     }
70
71     *p_item = malloc(sizeof(char) * (strlen(p->items[p->pos]) + 1));
72     memset(*p_item, 0, sizeof(char) * (strlen(p->items[p->pos]) + 1));
73     strncpy(*p_item, p->items[p->pos], strlen(p->items[p->pos]));
74     free(p->items[p->pos]);
75     p->items[p->pos] = NULL;
76     (p->pos)--;
77     return 0;
78 }
79
80 /**
81  * @brief check if the stack is full
82  * @param p STACK*: (pointer to) the stack to check
83  * @return true if full; false otherwise
84 */
85 bool is_full_stack(STACK *p)
86 {
87     if (p == NULL)
88     {
89         return 0;
90     }
91     return (p->pos == (p->size - 1));
92 }
93
94 /**
95  * @brief check if the stack is empty
96  * @param p STACK*: (pointer to) the stack to check
97  * @return true if empty; false otherwise
98 */
```

```
99  bool is_empty_stack(STACK *p)
100 {
101     if (p == NULL)
102     {
103         return 0;
104     }
105     return (p->pos == -1);
106 }
107
108 /**
109  * @brief resize stack to have greater capacity; maintain existing elements
110  * @param p STACK*: (pointer to) the stack to resize
111  * @return 0 on success; 1 otherwise
112 */
113 int resize_stack(STACK *p)
114 {
115     size_t old_size = p->size;
116     char **old_items = p->items;
117     p->size = (p->size) * STACK_RESIZE_FACTOR;
118     p->items = (char **)malloc((p->size) * sizeof(char *));
119     if (p->items == NULL)
120     {
121         return 1;
122     }
123
124     for (size_t i = 0; i < old_size; ++i)
125     {
126         p->items[i] = old_items[i];
127     }
128     for (size_t i = old_size; i < p->size; ++i)
129     {
130         p->items[i] = NULL;
131     }
132
133     free(old_items);
134     old_items = NULL;
135
136     return 0;
137 }
138
139 /**
140  * @brief returns number of elements currently in the stack
141  * @param p STACK*: (pointer to) the stack
142  * @return number of elements in the stack
143 */
144 size_t num_elements_stack(STACK *p)
145 {
146     return p->pos + 1;
147 }
148
```

```
149  /**
150  * @brief deconstruct stack: free all allocated memory
151  * @param p STACK*: (pointer to) the stack to deconstruct
152  * @return 0 on success; 1 otherwise
153  */
154  int cleanup_stack(STACK *p)
155  {
156      if (p == NULL || p->items == NULL)
157      {
158          return 0;
159      }
160      for (size_t i = 0; i < p->size; ++i)
161      {
162          if (p->items[i] != NULL)
163          {
164              free(p->items[i]);
165              p->items[i] = NULL;
166          }
167      }
168      free(p->items);
169      p->items = NULL;
170      return 0;
171 }
```

multithreaded/stack.h

```
1  /*
2  A dynamic stack holding strings
3  */
4
5 #include <stdio.h>
6 #include <string.h>
7 #include <stdlib.h>
8 #include <stdbool.h>
9
10 typedef struct stack
11 {
12     // max capacity of the stack
13     size_t size;
14     // position of the last item pushed onto the stack
15     size_t pos;
16     // memory for stack
17     char **items;
18 } STACK;
19
20 #define STACK_RESIZE_FACTOR 2
21
22 int init_stack(STACK *p, size_t stack_size);
23 bool is_full_stack(STACK *p);
24 bool is_empty_stack(STACK *p);
25 int push_stack(STACK *p, char *item);
26 int pop_stack(STACK *p, char **p_item);
27 int resize_stack(STACK *p);
28 size_t num_elements_stack(STACK *p);
29 int cleanup_stack(STACK *p);
```

multithreaded/p_stack.c

```
1  /*
2  A dynamic stack holding pointers
3  - note: this current implementation only resizes stack to increase size, never to decrease size
4  */
5
6 #include "p_stack.h"
7
8 /**
9  * @brief initialize stack with an initial size (capacity)
10 * @param p PSTACK*: a pointer to uninitialized memory
11 * @param stack_size size_t: initial capacity of the stack to be initialized
12 * @return 0 on success; 1 otherwise
13 */
14 int init_pstack(PSTACK *p, size_t stack_size)
15 {
16     if (p == NULL || stack_size == 0)
17     {
18         return 1;
19     }
20
21     p->size = stack_size;
22     p->pos = -1;
23     p->items = (void **)malloc(stack_size * sizeof(void *));
24
25     memset(p->items, 0, stack_size * sizeof(void *));
26
27     return 0;
28 }
29
30 /**
31 * @brief push an item onto the stack; if the stack is full, resize the stack
32 * @param p PSTACK*: (pointer to) the stack the function will push item onto
33 * @param item void*: pointer to push onto stack
34 * @return 0 on success; 1 otherwise
35 */
36 int push_pstack(PSTACK *p, void *item)
37 {
38     if (p == NULL)
39     {
40         return 1;
41     }
42
43     if (is_full_pstack(p))
44     {
45         resize_pstack(p);
46     }
47
48     ++(p->pos);
```

```
49     p->items[p->pos] = item;
50
51     return 0;
52 }
53
54 /**
55 * @brief pop from the stack
56 * @param p PSTACK*: (pointer to) the stack the function will pop from
57 * @param p_item void**: pointer that will be populated with popped element (which itself is a
58 pointer)
59 * @return 0 on success; 1 otherwise
60 */
61 int pop_pstack(PSTACK *p, void **p_item)
62 {
63     if (p == NULL || is_empty_pstack(p))
64     {
65         return 1;
66     }
67
68     *p_item = p->items[p->pos];
69     p->items[p->pos] = NULL;
70     (p->pos)--;
71     return 0;
72 }
73 /**
74 * @brief check if the stack is full
75 * @param p PSTACK*: (pointer to) the stack to check
76 * @return true if full; false otherwise
77 */
78 bool is_full_pstack(PSTACK *p)
79 {
80     if (p == NULL)
81     {
82         return 0;
83     }
84     return (p->pos == (p->size - 1));
85 }
86
87 /**
88 * @brief check if the stack is empty
89 * @param p PSTACK*: (pointer to) the stack to check
90 * @return true if empty; false otherwise
91 */
92 bool is_empty_pstack(PSTACK *p)
93 {
94     if (p == NULL)
95     {
96         return 0;
97     }
```

```
98     return (p->pos == -1);
99 }
100 /**
101 * @brief resize stack to have greater capacity; maintain existing elements
102 * @param p PSTACK*: (pointer to) the stack to resize
103 * @return 0 on success; 1 otherwise
104 */
105
106 int resize_pstack(PSTACK *p)
107 {
108     size_t old_size = p->size;
109     void **old_items = p->items;
110     p->size = (p->size) * PSTACK_RESIZE_FACTOR;
111     p->items = (void **)malloc((p->size) * sizeof(void *));
112     if (p->items == NULL)
113     {
114         return 1;
115     }
116
117     for (size_t i = 0; i < old_size; ++i)
118     {
119         p->items[i] = old_items[i];
120     }
121     for (size_t i = old_size; i < p->size; ++i)
122     {
123         p->items[i] = NULL;
124     }
125
126     free(old_items);
127     old_items = NULL;
128
129     return 0;
130 }
131
132 /**
133 * @brief returns number of elements currently in the stack
134 * @param p STACK*: (pointer to) the stack
135 * @return number of elements in the stack
136 */
137 size_t num_elements_pstack(PSTACK *p)
138 {
139     return p->pos + 1;
140 }
141
142 /**
143 * @brief deconstruct stack: free all allocated memory
144 * @param p PSTACK*: (pointer to) the stack to deconstruct
145 * @return 0 on success; 1 otherwise
146 */
147 int cleanup_pstack(PSTACK *p)
```

```
148 {  
149     if (p == NULL || p->items == NULL)  
150     {  
151         return 0;  
152     }  
153     for (size_t i = 0; i < p->size; ++i)  
154     {  
155         if (p->items[i] != NULL)  
156         {  
157             free(p->items[i]);  
158             p->items[i] = NULL;  
159         }  
160     }  
161     free(p->items);  
162     p->items = NULL;  
163     return 0;  
164 }
```

multithreaded/p_stack.h

```
1 #include <stdio.h>
2 #include <string.h>
3 #include <stdlib.h>
4 #include <stdbool.h>
5
6 typedef struct p_stack_struct
7 {
8     // max capacity of the stack
9     size_t size;
10    // position of the last item pushed onto the stack
11    size_t pos;
12    // memory for stack
13    void **items;
14 } PSTACK;
15
16 #define PSTACK_RESIZE_FACTOR 2
17
18 int init_pstack(PSTACK *p, size_t PSTACK_size);
19 bool is_full_pstack(PSTACK *p);
20 bool is_empty_pstack(PSTACK *p);
21 int push_pstack(PSTACK *p, void *item);
22 int pop_pstack(PSTACK *p, void **p_item);
23 int resize_pstack(PSTACK *p);
24 size_t num_elements_pstack(PSTACK *p);
25 int cleanup_pstack(PSTACK *p);
```

multithreaded/hash.c

```
1  /*
2  A hash set with strings as keys
3  - uses hsearch as the underlying hashmap
4  */
5
6 #include "hash.h"
7
8 /**
9  * @brief initialize hash set with an initial size (capacity)
10 * @param p HSET*: a pointer to uninitialized memory
11 * @param set_size size_t: initial capacity of the hash set to be initialized
12 * @return 0 on success; 1 otherwise
13 */
14 int init_hset(HSET *p, size_t set_size)
15 {
16     // create hsearch hash map
17     p->hmap = malloc(sizeof(struct hsearch_data));
18     memset(p->hmap, 0, sizeof(struct hsearch_data));
19     if (hcreate_r(set_size, p->hmap) == 0)
20     {
21         perror("hcreate\n");
22         return 1;
23     }
24
25     // array of keys in hash set
26     p->elements = (char **)malloc(sizeof(char *) * set_size);
27     for (int i = 0; i < set_size; ++i)
28     {
29         p->elements[i] = NULL;
30     }
31
32     p->cur_size = 0;
33     p->size = set_size;
34
35     // stack of pointers to strings used as arguments when searching for a string in hsearch;
36     // kept so we can deallocate them at cleanup
37     p->ps = malloc(sizeof(PSTACK));
38     memset(p->ps, 0, sizeof(PSTACK));
39     init_pstack(p->ps, 1);
40
41     return 0;
42 }
43
44 /**
45  * @brief check if hash set is at capacity
46  * @param p HSET*: (pointer to) the hash set to check
47  * @return true if full; false otherwise
48 */

```

```

49 bool is_full_hset(HSET *p)
50 {
51     return (p->size == p->cur_size);
52 }
53
54 /**
55 * @brief check if hash set is empty
56 * @param p HSET*: (pointer to) the hash set to check
57 * @return true if empty; false otherwise
58 */
59 bool is_empty_hset(HSET *p)
60 {
61     return (p->cur_size == 0);
62 }
63
64 /**
65 * @brief add key to the hash set; do nothing if key already exists
66 * @param p HSET*: (pointer to) the hash set to add key to
67 * @param key char*: key (string) to add
68 * @return 0 on success (no error); 1 otherwise (error)
69 */
70 int add_hset(HSET *p, char *key)
71 {
72     if (is_full_hset(p))
73     {
74         resize_hset(p);
75     }
76
77     // add key to hsearch
78     ENTRY item;
79     // note we create a copy of the parameter key since hsearch may
80     // continue to refer to the passed-in string's memory
81     item.key = malloc(strlen(key) * (sizeof(char) + 1));
82     memset(item.key, 0, strlen(key) * (sizeof(char) + 1));
83     strncpy(item.key, key, strlen(key));
84     item.data = NULL;
85     ACTION action = ENTER;
86     ENTRY *retval = NULL;
87     if (hsearch_r(item, action, &retval, p->hmap) == 0)
88     {
89         perror("hsearch_r\n");
90         return 1;
91     }
92
93     // add the new string to our array of keys
94     // (so we can transfer them on resize and free them at cleanup)
95     p->elements[p->cur_size] = item.key;
96     ++p->cur_size;
97
98     return 0;

```

```
99  }
100
101 /**
102 * @brief search for the key in the hash set
103 * @param p HSET*: (pointer to) the hash set to search
104 * @param key char*: key (string) to search
105 * @return 1 if the key is found; 0 if the key isn't found
106 */
107 int search_hset(HSET *p, char *key)
108 {
109     // search in hsearch
110     ENTRY item;
111     // note we create a copy of the parameter key since
112     // hsearch may continue to refer to the passed-in string's memory
113     item.key = malloc(strlen(key) * (sizeof(char) + 1));
114     memset(item.key, 0, strlen(key) * (sizeof(char) + 1));
115     strncpy(item.key, key, strlen(key));
116     item.data = NULL;
117     ACTION action = FIND;
118     ENTRY *retval;
119     hsearch_r(item, action, &retval, p->hmap);
120
121     // found key
122     if (retval != NULL)
123     {
124         // add to the pointer stack so we can deallocate at destruction
125         // we don't want to deallocate now, as hsearch may continue to refer
126         // to the string's memory
127         push_pstack(p->ps, item.key);
128         return 1;
129     }
130
131     // did not find key: we can deallocate the key string now
132     free(item.key);
133     item.key = NULL;
134     return 0;
135 }
136
137 /**
138 * @brief resize hash set to have greater capacity; maintain existing elements
139 * @param p HSET*: (pointer to) the hash set to resize
140 * @return 0 on success; 1 otherwise
141 */
142 int resize_hset(HSET *p)
143 {
144     // destroy hsearch
145     hdestroy_r(p->hmap);
146     free(p->hmap);
147     p->hmap = NULL;
148 }
```

```
149     // reinitialize hsearch
150     p->hmap = malloc(sizeof(struct hsearch_data));
151     memset(p->hmap, 0, sizeof(struct hsearch_data));
152     if (hcreate_r(HSET_RESIZE_FACTOR * (p->size), p->hmap) == 0)
153     {
154         perror("hcreate\n");
155         return 1;
156     }
157
158     // allocate resized array of keys
159     size_t old_size = p->size;
160     char **old_elements = p->elements;
161     p->size = (p->size) * HSET_RESIZE_FACTOR;
162     p->elements = (char **)malloc((p->size) * sizeof(char *));
163
164     // transfer old keys over into new array
165     p->cur_size = 0;
166     for (size_t i = 0; i < old_size; ++i)
167     {
168         add_hset(p, old_elements[i]);
169         free(old_elements[i]);
170         old_elements[i] = NULL;
171     }
172     for (size_t i = old_size; i < p->size; ++i)
173     {
174         p->elements[i] = NULL;
175     }
176
177     // deallocate old array of keys
178     free(old_elements);
179     old_elements = NULL;
180
181     return 0;
182 }
183
184 /**
185 * @brief deconstruct hash set: free all allocated memory
186 * @param p HSET*: (pointer to) the hash set to deconstruct
187 * @return 0 on success; 1 otherwise
188 */
189 int cleanup_hset(HSET *p)
190 {
191     for (size_t i = 0; i < p->size; ++i)
192     {
193         if (p->elements[i] != NULL)
194         {
195             free(p->elements[i]);
196             p->elements[i] = NULL;
197         }
198     }
199 }
```

```
199     free(p->elements);
200     p->elements = NULL;
201
202     hdestroy_r(p->hmap);
203     free(p->hmap);
204     p->hmap = NULL;
205
206     // clean up the pstack, which cleans up the extra string pointers created during search
207     cleanup_pstack(p->ps);
208     free(p->ps);
209
210     return 0;
211 }
```

multithreaded/hash.h

```
1  /*
2  A hash set with strings as keys
3  - uses Linux's hsearch as the underlying hashmap
4  */
5
6 #define _GNU_SOURCE
7 #include <stdio.h>
8 #include <stdlib.h>
9 #include <string.h>
10 #include <search.h>
11 #include "p_stack.h"
12
13 typedef struct hashmap
14 {
15     // current capacity (max number of keys)
16     size_t size;
17     // array of keys in the hash set
18     // - for linear access,
19     // - for moving keys over on resize,
20     // - and for deallocation at destruction
21     char **elements;
22     // hsearch structure
23     struct hsearch_data *hmap;
24     // current number of keys
25     size_t cur_size;
26     // stack of pointers to strings used for searching in hsearch;
27     // so we can deallocate them at cleanup
28     PSTACK *ps;
29 } HSET;
30
31 #define HSET_RESIZE_FACTOR 2
32
33 int init_hset(HSET *p, size_t set_size);
34 bool is_full_hset(HSET *p);
35 bool is_empty_hset(HSET *p);
36 int add_hset(HSET *p, char *key);
37 int search_hset(HSET *p, char *key);
38 int resize_hset(HSET *p);
39 int cleanup_hset(HSET *p);
```

multithreaded/curl_xml.c

```
1 #include "curl_xml.h"
2
3 /**
4  * @brief set options of curl easy handle
5  * @param curl_handle CURL*: (pointer to) already-initialized curl easy handle to configure
6  * @param ptr RECV_BUF*: (pointer to) user data needed by the curl write call back function
7  * @param url const char*: target url to crawl
8  * @return a valid CURL handle upon sucess; NULL otherwise
9  * note the caller is responsible for cleaning the returned curl handle
10 */
11 CURL *easy_handle_config(CURL *curl_handle, RECV_BUF *ptr, const char *url)
12 {
13     if (ptr == NULL || url == NULL)
14     {
15         return NULL;
16     }
17
18     // init user defined call back function buffer
19     if (recv_buf_init(ptr, BUF_SIZE) != 0)
20     {
21         return NULL;
22     }
23
24     // specify URL to get
25     curl_easy_setopt(curl_handle, CURLOPT_URL, url);
26
27     // register write call back function to process received data
28     curl_easy_setopt(curl_handle, CURLOPT_WRITEFUNCTION, write_cb_curl);
29     // user defined data structure passed to the call back function
30     curl_easy_setopt(curl_handle, CURLOPT_WRITEDATA, (void *)ptr);
31
32     // register header call back function to process received header data
33     curl_easy_setopt(curl_handle, CURLOPT_HEADERFUNCTION, header_cb_curl);
34     // user defined data structure passed to the call back function
35     curl_easy_setopt(curl_handle, CURLOPT_HEADERDATA, (void *)ptr);
36
37     // some servers require a user-agent field
38     curl_easy_setopt(curl_handle, CURLOPT_USERAGENT, CURL_USER_AGENT_FIELD);
39
40     // follow HTTP 3XX redirects
41     curl_easy_setopt(curl_handle, CURLOPT_FOLLOWLOCATION, 1L);
42     // continue to send authentication credentials when following locations
43     curl_easy_setopt(curl_handle, CURLOPT_UNRESTRICTED_AUTH, 1L);
44     // max number of redirects to follow sets to 5
45     curl_easy_setopt(curl_handle, CURLOPT_MAXREDIRS, 5L);
46     // supports all built-in encodings
47     curl_easy_setopt(curl_handle, CURLOPT_ACCEPT_ENCODING, "");
```

```
49     // Enable the cookie engine without reading any initial cookies
50     curl_easy_setopt(curl_handle, CURLOPT_COOKIEFILE, "");
51     // allow whatever auth the proxy speaks
52     curl_easy_setopt(curl_handle, CURLOPT_PROXYAUTH, CURLAUTH_ANY);
53     // allow whatever auth the server speaks
54     curl_easy_setopt(curl_handle, CURLOPT_HTTPAUTH, CURLAUTH_ANY);
55
56     return curl_handle;
57 }
58
59 /**
60 * @brief process a downloaded html page: get all urls from the page and push it onto stack
61 * @param curl_handle CURL*: (pointer to) curl handler that was used to access the url
62 * @param p_recv_buf RECV_BUF*: (pointer to) buffer that contains the received data
63 * @param content_type int*: (pointer to) int to be set with content type code
64 * @param stack STACK*: (pointer to) stack that will be populated with further urls to crawl
65 * @return 0 on success; non-zero otherwise
66 */
67 int process_html(CURL *curl_handle, RECV_BUF *p_recv_buf, int *content_type, STACK *stack)
68 {
69     *content_type = HTML;
70
71     int follow_relative_link = 1;
72     char *url = NULL;
73
74     curl_easy_getinfo(curl_handle, CURLINFO_EFFECTIVE_URL, &url);
75     find_http(p_recv_buf->buf, p_recv_buf->size, follow_relative_link, url, stack);
76     return 0;
77 }
78
79 /**
80 * @brief check if a png is a valid png
81 * @param buf uint8_t*: (pointer to) memory containing the supposed png image
82 * @param n size_t: size of the memory
83 * @return true if memory contains a valid png; false otherwise
84 * @details
85 * The check is derived from the png specification: https://www.w3.org/TR/png/
86 */
87 bool is_png(uint8_t *buf, size_t n)
88 {
89     if (n < 8)
90     {
91         return false;
92     }
93
94     if ((buf[0] == 0x89) &&
95         (buf[1] == 0x50) &&
96         (buf[2] == 0x4E) &&
97         (buf[3] == 0x47) &&
98         (buf[4] == 0x0D) &&
```

```
99         (buf[5] == 0x0A) &&
100        (buf[6] == 0x1A) &&
101        (buf[7] == 0x0A))
102    {
103        return true;
104    }
105
106    return false;
107}
108
109 /**
110 * @brief process a downloaded png: check if it's a valid png
111 * @param curl_handle CURL*: (pointer to) curl handler that was used to access the url
112 * @param p_recv_buf RECV_BUF*: (pointer to) buffer that contains the received data
113 * @param content_type int*: (pointer to) int to be set with content type code
114 * @return 0 on success; non-zero otherwise
115 */
116 int process_png(CURL *curl_handle, RECV_BUF *p_recv_buf, int *content_type)
117 {
118     // effective url
119     char *eurl = NULL;
120     curl_easy_getinfo(curl_handle, CURLINFO_EFFECTIVE_URL, &eurl);
121
122     if (is_png((uint8_t *)p_recv_buf->buf, p_recv_buf->size))
123     {
124         *content_type = VALID_PNG;
125     }
126     else
127     {
128         *content_type = INVALID_PNG;
129     }
130
131     return 0;
132 }
133
134 /**
135 * @brief process the downloaded content data
136 * @param curl_handle CURL*: (pointer to) curl handler that was used to access the url
137 * @param p_recv_buf RECV_BUF*: (pointer to) buffer that contains the received data
138 * @param content_type int*: (pointer to) int to be set with content type code
139 * @param stack STACK*: (pointer to) stack that will be populated with further urls to crawl
140 * @param response_code_p long*: (pointer to) int to be set with the response code
141 * @return 0 on success; non-zero otherwise
142 * @details
143 * if url points to a HTML page, populate stack with urls linked on the page
144 * if url points to a png, check if it's a valid png
145 * set the content type and response code via the appropriate pointers
146 */
147 int process_data(CURL *curl_handle, RECV_BUF *p_recv_buf, int *content_type, STACK *stack, long
*response_code_p)
```

```
148 {
149     CURLcode res;
150
151     // get response code, and return if it's a fail
152     long response_code;
153     res = curl_easy_getinfo(curl_handle, CURLINFO_RESPONSE_CODE, &response_code);
154     *response_code_p = response_code;
155     if (response_code >= BAD_REQUESTS)
156     {
157         return 1;
158     }
159
160     // get content, and handle differently depending on content type
161     char *ct = NULL;
162     res = curl_easy_getinfo(curl_handle, CURLINFO_CONTENT_TYPE, &ct);
163     if (res != CURLE_OK || ct == NULL)
164     {
165         return 2;
166     }
167     if (strstr(ct, CT_HTML))
168     {
169         return process_html(curl_handle, p_recv_buf, content_type, stack);
170     }
171     else if (strstr(ct, CT_PNG))
172     {
173         return process_png(curl_handle, p_recv_buf, content_type);
174     }
175
176     return 0;
177 }
178
179 /**
180 * @brief crawl specified url and process the downloaded data
181 * @param curl_handle CURL*: (pointer to) the curl handler that will be used to process the url
182 * @param seed_url char*: string containing the url to crawl
183 * @param content_type int*: (pointer to) int to be set with content type code
184 * @param stack STACK*: (pointer to) stack that will be populated with further urls to crawl
185 * @param response_code_p long*: (pointer to) int to be set with the response code
186 * @return 0 on success; non-zero otherwise
187 * @details
188 * if url points to a HTML page, populate stack with urls linked on the page
189 * if url points to a png, check if it's a valid png
190 * set the content type and response code via the appropriate pointers
191 */
192 int process_url(CURL *curl_handle, char *seed_url, int *content_type, STACK *stack, long
*response_code_p)
193 {
194     // set default response code to failure (if nothing fails, the code will be set later)
195     *response_code_p = INTERNAL_SERVER_ERRORS;
196 }
```

```
197     // configure the easy curl handle
198     char url[URL_LENGTH];
199     RECV_BUF recv_buf;
200     strcpy(url, seed_url);
201     curl_handle = easy_handle_config(curl_handle, &recv_buf, url);
202     if (curl_handle == NULL)
203     {
204         fprintf(stderr, "Curl configuration failed. Exiting...\n");
205         curl_global_cleanup();
206         abort();
207     }
208
209     // process the url
210     CURLcode res;
211     res = curl_easy_perform(curl_handle);
212     if (res != CURLE_OK)
213     {
214         recv_buf_cleanup(&recv_buf);
215         return 1;
216     }
217
218     // process the data from the url
219     process_data(curl_handle, &recv_buf, content_type, stack, response_code_p);
220
221     // clean up data buffer
222     recv_buf_cleanup(&recv_buf);
223     return 0;
224 }
225
226 /**
227 * @brief returns whether the response code is not an error (i.e. okay or redirect)
228 * @param response_code long: response_code in question
229 * @return true if response code is crawlable (not an error); false otherwise
230 */
231 bool is_processable_response(long response_code)
232 {
233     return (response_code >= OK_REQUESTS && response_code <= OK_REQUESTS + CODE_RANGE) ||
234     (response_code >= REDIRECT_REQUESTS && response_code <= REDIRECT_REQUESTS + CODE_RANGE);
235
236 /**
237 * @brief cURL header call back function to extract image sequence number from
238 *        http header data. An example header for image part n (assume n = 2) is:
239 *        X-Ece252-Fragment: 2
240 * @param p_recv char*: (pointer to) header data delivered by cURL
241 * @param size size_t: number of data elements
242 * @param nmemb size_t: size of one data element
243 * @param userdata void*: (pointer to) buffer containing user-defined data
244 *                    structure used for extracting sequence number
245 * @return total size of header data received
```

```
246 * @details
247 * cURL documentation: https://curl.se/libcurl/c/CURLOPT\_HEADERFUNCTION.html.
248 */
249 size_t header_cb_curl(char *p_recv, size_t size, size_t nmemb, void *userdata)
250 {
251     int realsize = size * nmemb;
252     RECV_BUF *p = userdata;
253
254     if (realsize > strlen(ECE252_HEADER) &&
255         strncmp(p_recv, ECE252_HEADER, strlen(ECE252_HEADER)) == 0)
256     {
257         // extract image sequence number
258         p->seq = atoi(p_recv + strlen(ECE252_HEADER));
259     }
260     return realsize;
261 }
262
263 /**
264 * @brief cURL write callback function that saves a copy of received data
265 * @param p_recv char*: (pointer to) memory that will be populated with received data
266 * @param size size_t: number of data elements
267 * @param nmemb size_t: size of one data element
268 * @param p_userdata void*: (pointer to) user data buffer that can be accessed outside CURL;
269 *                      this will be populated with the web page (or png file)
270 * @return total size of data received
271 * @details
272 * cURL documentation: https://curl.se/libcurl/c/CURLOPT\_WRITEFUNCTION.html
273 */
274 size_t write_cb_curl(char *p_recv, size_t size, size_t nmemb, void *p_userdata)
275 {
276     size_t realsize = size * nmemb;
277     RECV_BUF *p = (RECV_BUF *)p_userdata;
278
279     if (p->size + realsize + 1 > p->max_size)
280     {
281         // since received data is not 0 terminated, add one byte for terminating 0
282         size_t new_size = p->max_size + max(BUF_INC, realsize + 1);
283         char *q = realloc(p->buf, new_size);
284         if (q == NULL)
285         {
286             // out of memory
287             perror("realloc");
288             return -1;
289         }
290         p->buf = q;
291         p->max_size = new_size;
292     }
293
294     // copy data from libcurl into a buffer we can access later
295     memcpy(p->buf + p->size, p_recv, realsize);
```

```
296     p->size += realsize;
297     p->buf[p->size] = 0;
298
299     return realsize;
300 }
301
302 /**
303 * @brief initialize data structure used for downloading data via cURL
304 * @param ptr RECV_BUF*: a pointer to uninitialized memory
305 * @param max_size size_t: maximum expected size of data to download
306 * @return 0 on success; non-zero otherwise
307 * @details
308 * note that if the size of downloaded data is larger than the max_size,
309 * we will reallocate to accommodate
310 */
311 int recv_buf_init(RECV_BUF *ptr, size_t max_size)
312 {
313     void *p = NULL;
314
315     if (ptr == NULL)
316     {
317         return 1;
318     }
319
320     p = malloc(max_size);
321     if (p == NULL)
322     {
323         return 2;
324     }
325
326     ptr->buf = p;
327     ptr->size = 0;
328     ptr->max_size = max_size;
329     // a valid sequence number should be positive
330     ptr->seq = -1;
331     return 0;
332 }
333
334 /**
335 * @brief clean up data structure used for downloading data via cURL: deallocate memory
336 * @param ptr RECV_BUF*: (pointer to) RECV_BUF to clean
337 * @return 0 on success; non-zero otherwise
338 */
339 int recv_buf_cleanup(RECV_BUF *ptr)
340 {
341     if (ptr == NULL)
342     {
343         return 1;
344     }
345 }
```

```
346     if (ptr->buf != NULL)
347     {
348         free(ptr->buf);
349         ptr->buf = NULL;
350     }
351
352     ptr->size = 0;
353     ptr->max_size = 0;
354     return 0;
355 }
356
357 /**
358 * @brief clean up all cURL related data structures
359 * @param curl CURL*: (pointer to) curl easy handle to clean up
360 * @param ptr RECV_BUF*: a pointer to data structure used for downloading data via cURL
361 */
362 void cleanup(CURL *curl, RECV_BUF *ptr)
363 {
364     curl_easy_cleanup(curl);
365     curl_global_cleanup();
366     recv_buf_cleanup(ptr);
367 }
368
369 /**
370 * @brief get html document from data
371 * @param buf char*: (pointer to) memory containing document data
372 * @param size int: size of data
373 * @param url char*: url string of the html document
374 * @return document pointer if successful; NULL otherwise
375 */
376 xmlDocPtr mem_getdoc(char *buf, int size, const char *url)
377 {
378     int opts = HTML_PARSE_NOBLANKS | HTML_PARSE_NOERROR |
379                 HTML_PARSE_NOWARNING | HTML_PARSE_NONET;
380     xmlDocPtr doc = htmlReadMemory(buf, size, url, NULL, opts);
381
382     if (doc == NULL)
383     {
384         printf("%s\n", url);
385         fprintf(stderr, "Document not parsed successfully.\n");
386         return NULL;
387     }
388
389     return doc;
390 }
391
392 /**
393 * @brief get nodes on the xml page
394 * @param doc xmlDocPtr: xml document to get nodes from
395 * @param xpath xmlChar*: xpath
```

```
396 * @return nodes on success; NULL otherwise
397 */
398 xmlDocPtr getnodeset(xmlDocPtr doc, xmlChar *xpath)
399 {
400     xmlXPathContextPtr context;
401     xmlXPathObjectPtr result;
402
403     context = xmlXPathNewContext(doc);
404     if (context == NULL)
405     {
406         printf("Error in xmlXPathNewContext\n");
407         return NULL;
408     }
409     result = xmlXPathEvalExpression(xpath, context);
410     xmlXPathFreeContext(context);
411     if (result == NULL)
412     {
413         printf("Error in xmlXPathEvalExpression\n");
414         return NULL;
415     }
416     if (xmlXPathNodeSetIsEmpty(result->nodesetval))
417     {
418         xmlXPathFreeObject(result);
419         printf("No result\n");
420         return NULL;
421     }
422     return result;
423 }
424
425 /**
426 * @brief get all urls on the xml web page and push them onto stack
427 * @param buf char*: (pointer to) buffer that contains the HTML web page
428 * @param size int: size of the buffer
429 * @param follow_relative_links int: 1 if we're following relative links and 0 otherwise
430 * @param base_url const char*: base url of the page
431 * @param stack STACK*: (pointer to) stack that will be populated with further urls to crawl
432 * @return 0 on success; non-zero otherwise
433 */
434 int find_http(char *buf, int size, int follow_relative_links, const char *base_url, STACK
*stack)
435 {
436     int i;
437     xmlDocPtr doc;
438     xmlChar *xpath = (xmlChar *) "//a/@href";
439     xmlNodeSetPtr nodeset;
440     xmlXPathObjectPtr result;
441     xmlChar *href;
442
443     if (buf == NULL)
444     {
```

```
445     return 1;
446 }
447
448 doc = mem_getdoc(buf, size, base_url);
449
450 result = getnodeset(doc, xpath);
451 if (result)
452 {
453     nodeset = result->nodesetval;
454     for (i = 0; i < nodeset->nodeNr; i++)
455     {
456         href = xmlNodeListGetString(doc, nodeset->nodeTab[i]->xmlChildrenNode, 1);
457         if (follow_relative_links)
458         {
459             xmlChar *old = href;
460             href = xmlBuildURI(href, (xmlChar *)base_url);
461             xmlFree(old);
462         }
463         if (href != NULL && !strcmp((const char *)href, "http", 4))
464         {
465             char *temp = malloc((strlen((char *)href) + 1) * sizeof(char));
466             memset(temp, 0, (strlen((char *)href) + 1) * sizeof(char));
467             sprintf(temp, "%s", href);
468             push_stack(stack, (char *)temp);
469             free(temp);
470             temp = NULL;
471         }
472         xmlFree(href);
473     }
474     xmlXPathFreeObject(result);
475 }
476
477 xmlFreeDoc(doc);
478
479 return 0;
480 }
```

multithreaded/curl_xml.h

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <stdbool.h>
4 #include <string.h>
5 #include <sys/types.h>
6 #include <unistd.h>
7 #include <curl/curl.h>
8 #include <libxml/HTMLparser.h>
9 #include <libxml/parser.h>
10 #include <libxml>xpath.h>
11 #include <libxml/uri.h>
12 #include "stack.h"
13
14 #define SEED_URL "http://ece252-1.uwaterloo.ca/lab4/"
15 #define ECE252_HEADER "X-Ece252-Fragment: "
16 #define CURL_USER_AGENT_FIELD "ece252 lab4 crawler"
17 #define BUF_SIZE 1048576 /* 1024*1024 = 1M */
18 #define BUF_INC 524288 /* 1024*512 = 0.5M */
19
20 #define CT_PNG "image/png"
21 #define CT_HTML "text/html"
22 #define CT_PNG_LEN 9
23 #define CT_HTML_LEN 9
24 #define URL_LENGTH 256
25
26 #define OK_REQUESTS 200
27 #define REDIRECT_REQUESTS 300
28 #define BAD_REQUESTS 400
29 #define INTERNAL_SERVER_ERRORS 500
30 #define CODE_RANGE 99
31
32 #define DEFAULT_TYPE -1
33 #define HTML 0
34 #define VALID_PNG 1
35 #define INVALID_PNG 2
36
37 #define max(a, b) \
38 ({ __typeof__ (a) _a = (a); \
39   __typeof__ (b) _b = (b); \
40   _a > _b ? _a : _b; })
41
42 typedef struct recv_buf2
43 {
44     char *buf;          // memory to hold a copy of received data
45     size_t size;        // size of valid data in buf in bytes
46     size_t max_size;    // max capacity of buf in bytes
47     int seq;           // >=0 sequence number extracted from http header
48                           // <0 indicates an invalid seq number
```

```
49 } RECV_BUF;
50
51 xmlDocPtr mem_getdoc(char *buf, int size, const char *url);
52 xmlXPathObjectPtr getnodeset(xmlDocPtr doc, xmlChar *xpath);
53 int find_http(char *fname, int size, int follow_relative_links, const char *base_url, STACK *stack);
54 size_t header_cb_curl(char *p_recv, size_t size, size_t nmemb, void *userdata);
55 size_t write_cb_curl(char *p_recv, size_t size, size_t nmemb, void *p_userdata);
56 int recv_buf_init(RECV_BUF *ptr, size_t max_size);
57 int recv_buf_cleanup(RECV_BUF *ptr);
58 void cleanup(CURL *curl, RECV_BUF *ptr);
59 CURL *easy_handle_config(CURL *curl_handle, RECV_BUF *ptr, const char *url);
60 int process_data(CURL *curl_handle, RECV_BUF *p_recv_buf, int *content_type, STACK *stack, long *response_code_p);
61 int process_png(CURL *curl_handle, RECV_BUF *p_recv_buf, int *content_type);
62 bool is_png(uint8_t *buf, size_t n);
63 int process_url(CURL *curl_handle, char *seed_url, int *content_type, STACK *stack, long *response_code_p);
64 bool isprocessable_response(long response_code);
```