

Assignment 4 – Processing CSV Data with Threads

Description:

The purpose of this assignment is to deal with parsing CSV files, error handling, handling large datasets, summarizing data with specific criteria, utilizing vectors, implementing threads with thorough timing understanding, and preventing race conditions to ensure data integrity.

Approach:

Mostly follow the “steps” provided in the README, beginning with figuring out how to parse the CSV file.

Parsing

First, I need to figure out how to parse the CSV file. The simple approach here is to simply split on a comma. However, there are other considerations to take care of, such as fields enclosed in quotes, and newlines within fields. This means I can’t simply just read the file line-by-line and split on the comma. Also, I think it would be best to avoid `strtok` completely because it mutates the original string, in addition, it would not handle blank fields e.g., “`a`, , `b`” because it would replace all instances of the delimiter (,) with `\0`.

Here’s a very rough pseudocode of how I would approach parsing the CSV file.

- for each character:
 - if it’s a quotation mark
 - if there is, check for another quote
 - then, this is an escaped quote
 - skip the next quote
 - if it’s a comma
 - if there was an opening quote
 - then, this character is part of the field
 - otherwise
 - this is the end of the column
 - if it’s a newline
 - if there was an opening quote
 - then, the character is part of the field
 - otherwise
 - this is the end of the row
 - if it’s none of the above

- then, it's just a "regular" character

Processing

Looking at the screenshot in the provided README, we see that each row is separated by the call type, and for each district, we separate the response ("dispatch" and "on-scene") times into different buckets, with a special column for the aggregate.

Since the table is listed by the call type, I would store the parse data and list them by the call types. Beginning with a struct that represents it:

```
typedef struct ResponseTime
{
    ResponseType type;
    int under_2_mins;
    int mins_3_5;
    int mins_6_10;
    int over_10_mins;
} ResponseTime;

typedef struct Subfield
{
    char *name;
    ResponseTime *responseTimes[2];
} Subfield;

typedef struct CallType
{
    char *name;
    int total_count;
    Subfield *call_total;
    Subfield **subfields;
} CallType;
```

- for each row:
 - get the call type
 - get the four timestamps
 - parse and calculate deltas for dispatch and on-scene times
 - if any one of these fields do not exist, ignore the row
 - find the call type from the struct:
 - if it exists, retrieve it
 - increment the total for the call type
 - find and update the response time of the neighborhood
 - otherwise, create the call type and initialize everything
 - set the call type name
 - set the total to one

- initialize the array for neighborhood, and create them using the values passed from argv
- update the response times

Threading

- make a function that the threads call that call csvnext
- have a mutex lock on when the file is changed
 - only on the getline, not the processing
- have another mutex lock for when the data structure is changed
 - only when values are updated, not read

Analysis:

nearly half time at 2 threads. Not fully half, because there is a mutex lock which leads to code being ran serialized. Adding more threads did not lead to a decrease in time. This is caused by the VM only having 2 cores. The threads need a core to run, so at most we can only see close to half time in the configuration of the VM. If there were 4 cores then we can expect to see close to a quarter time.

Issues and Resolutions:

Issue one: Not initializing the `tm` struct when using `strptime` to parse the time fields from the CSV. Since the provided CSV files contains two different types of timestamps, I set up the helper function to handle two formats. Not initializing the struct would cause the value return from `mkttime` to be seemingly randomized. In addition, one of the format (in the 5,000-row file) do not contain seconds in the timestamp. When using `strptime` to parse, the seconds field would be set to the current computer's time. In order to make it consistent, I needed to set the seconds field to zero, `tm.tm_sec = 0`.

```
student@student: ~/asmt/csc415-assignment-4-csvthreads-...
1704777780 1704777780
1/8/24 21:23 1/8/24 21:23
1704777780 1704777780
DELTAS: Dispatch 0.000000 | Onscene -3600.000000
PASSING CALL
1/8/24 9:33 1/8/24 9:33
1704735180 1704735180
1/8/24 9:33 1/8/24 9:33
1704735180 1704735180
DELTAS: Dispatch 0.000000 | Onscene -3600.000000
PASSING CALL
1/8/24 13:24 1/8/24 13:24
1704749040 1704749040
1/8/24 13:24 1/8/24 13:24
1704749040 1704749040
DELTAS: Dispatch 0.000000 | Onscene -3600.000000
Name: TEST_TYPE
Subfield 0: BAYVIEW
Subfield 1: MISSION
Subfield 2: TENDERLOIN
ptr: 0xaaaae7b39930
Name: TEST_TYPE
Total Time was 0.000000042 seconds
student@student:~/asmt/csc415-assignment-4-csvthreads-mosguinz$
```

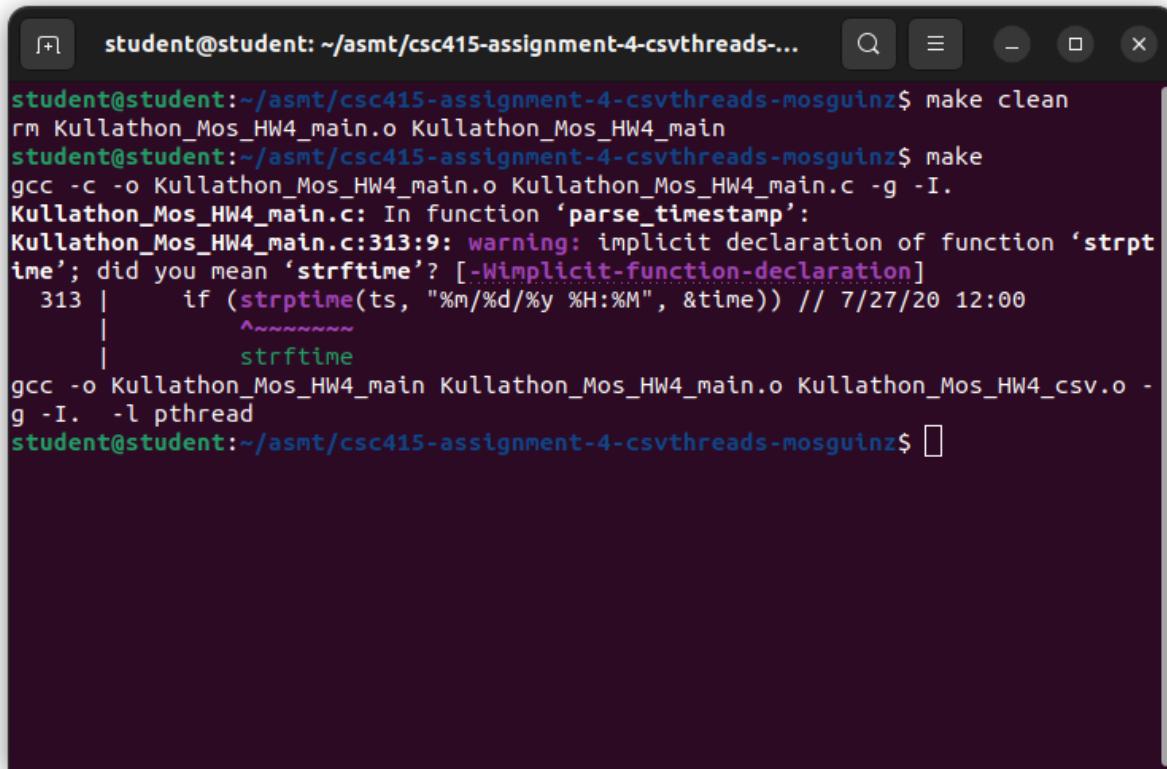
Issue two: total call count for each type being “doubled.” I had a helper function `update_calltype` to update the number of calls for each call type, for a particular response time (i.e. either on-scene or dispatch). This function is being invoked twice, once to update the call type for on-scene, and another for dispatch. I was also updating the *total* number of calls inside of this helper, which meant that they were being double counted. I fixed this by moving the incrementation of the total number of calls outside of this helper.

Issue three: Memory leak occurring when running for a certain number of rows. For the most part, I “centralize” all of my `malloc` and `free` logic into their respective functions e.g., `init_calltype` and `free_calltype` in order to properly keep track and ensure that everything was being freed at the end. However, I did not take into account where the lines being read from the CSV file were being freed at the end of the loop. When a row is determined to be “garbage,” the loop terminates early but doesn’t free the row that is read.

Issue four: Structs containing string fields are producing garbage data. I initially did not heap-allocate any of the string fields. I simply assign all string values after creating the struct for each row. This did not seem to be an issue at first when simply printing in the loop for debugging purposes. However, once I fixed the memory leak above by freeing each line, the

string fields are no longer valid since they are `char*`. Instead, I had to `strdup` the all such fields, and of course update my `free_*` and `init_*` helpers.

Screen shot of compilation:



A screenshot of a terminal window titled "student@student: ~/asmt/csc415-assignment-4-csvthreads-mosguinz...". The window shows the following command-line session:

```
student@student:~/asmt/csc415-assignment-4-csvthreads-mosguinz$ make clean
rm Kullathon_Mos_HW4_main.o Kullathon_Mos_HW4_main
student@student:~/asmt/csc415-assignment-4-csvthreads-mosguinz$ make
gcc -c -o Kullathon_Mos_HW4_main.o Kullathon_Mos_HW4_main.c -g -I.
Kullathon_Mos_HW4_main.c: In function 'parse_timestamp':
Kullathon_Mos_HW4_main.c:313:9: warning: implicit declaration of function 'strptime'; did you mean 'strftime'? [-Wimplicit-function-declaration]
  313 |     if (strptime(ts, "%m/%d/%y %H:%M", &time)) // 7/27/20 12:00
      |     ^
      |     strftime
gcc -o Kullathon_Mos_HW4_main Kullathon_Mos_HW4_main.o Kullathon_Mos_HW4_csv.o -g -I. -l pthread
student@student:~/asmt/csc415-assignment-4-csvthreads-mosguinz$
```

Screen shot(s) of the execution of the program:

See below

1 thread - 1.54 seconds

2 threads - 1.04 seconds

The screenshot shows a terminal window titled "OperatingSystemsSP24ARM" with the command "student@student: ~/amkt/csc415-assignment-4-csvthreads-mosguinz" running. The output is a large table of data with many columns and rows, representing performance metrics over time. The table includes columns for CPU usage (e.g., %CPU, %idle), memory (e.g., %MEM, %SWAP), disk activity (e.g., %DISK, %IO), and other system statistics. The data is presented in a grid format with horizontal and vertical lines separating the columns and rows.

Time	%CPU	%idle	%MEM	%SWAP	%DISK	%IO	Other Metrics
00:00:00	1.00	98.99	10.00	0.00	0.00	0.00	Initial values
00:00:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:00:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:00:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:00:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:00:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:01:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:01:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:01:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:01:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:01:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:01:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:02:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:02:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:02:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:02:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:02:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:02:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:03:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:03:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:03:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:03:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:03:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:03:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:04:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:04:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:04:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:04:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:04:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:04:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:05:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:05:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:05:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:05:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:05:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:05:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:06:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:06:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:06:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:06:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:06:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:06:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:07:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:07:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:07:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:07:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:07:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:07:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:08:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:08:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:08:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:08:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:08:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:08:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:09:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:09:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:09:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:09:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:09:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:09:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:10:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:10:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:10:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:10:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:10:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:10:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:11:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:11:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:11:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:11:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:11:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:11:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:12:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:12:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:12:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:12:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:12:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:12:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:13:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:13:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:13:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:13:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:13:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:13:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:14:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:14:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:14:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:14:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:14:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:14:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:15:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:15:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:15:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:15:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:15:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:15:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:16:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:16:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:16:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:16:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:16:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:16:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:17:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:17:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:17:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:17:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:17:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:17:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:18:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:18:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:18:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:18:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:18:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:18:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:19:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:19:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:19:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:19:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:19:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:19:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:20:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:20:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:20:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:20:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:20:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:20:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:21:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:21:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:21:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:21:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:21:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:21:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:22:00	1.00	98.99	10.00	0.00	0.00	0.00	
00:22:10	1.00	98.99	10.00	0.00	0.00	0.00	
00:22:20	1.00	98.99	10.00	0.00	0.00	0.00	
00:22:30	1.00	98.99	10.00	0.00	0.00	0.00	
00:22:40	1.00	98.99	10.00	0.00	0.00	0.00	
00:22:50	1.00	98.99	10.00	0.00	0.00	0.00	
00:23:00	1.00	98.99	10.00	0.00	0.00	0.00	

4 threads - 1.34 seconds

Category	ID	Count	Time
CITIZEN ARREST	270	68	34
CITIZEN STANDBY	1120	190	203
CIVIL PROCESS	4	4	0
DAM PROTEST	77	17	37
EMERGENCY BACKUP	77	3	0
EXPLOSIVE FOUND	17	3	0
FIGHT / VIOLENCE	270	180	20
FORCED ENTRY	12	1	0
FRAUD	150	22	2
GANG VANDALISM	199	28	29
GRAND THEFT	199	28	29
HIGHWAY ACCIDENT	198	24	10
HOME INVADERS	1	1	0
HOME CHECK	204	155	20
HOMEOWNER COMPLAINT	303	3	0
HOMICIDE	303	3	0
HOSTAGE SITUATION	466	236	147
INJURY / VIOLENCE	58	33	194
JOKE	30	7	10
KIDNAPING	981	663	10
MEET WITNESS	199	130	265
MISSING PERSON / DISAPPEARANCE	97	32	28
MISSING PERSON / CRIME	872	261	355
MENTALLY DISTURBED	226	24	24
MESSING AROUND	526	51	187
MISSING ADULT	9	5	5
MISSING CHILD	10	6	2
MISSING PERSON	1618	74	65
MISSING PERSON / CRIME	2108	210	96
PANIC ALARM	218	94	117
PASSIVE AGGRESSION	2108	210	96
PERSON BREAKING IN	12	21	45
PERSON DUMPING	12	4	0
PERSON KIDNAPING	251	189	137
PERSON KIDNAPING	301	138	145
PERSON KIDNAPING / CRIME	1	1	0
PERSON KIDNAPING / VICTIM	112	219	187
PERSON TRANSPORT	129	90	671
PROSTITUTION / POLICING	168	52	68
PURSE SNATCH	124	5	46
PUBLIC HEALTH VIOLATION	27	1	0
RECORDS REQUEST	10	1	0
RESISTING ARREST	24	14	20
ROBBERY	264	69	1
ROLLING DRUNK PERSON	98	7	12
SILENT PERSON	10	10	10
SHOT SPOTTING	145	141	7
SOCIAL DISTANCING	2	2	0
SILENT HOLUP ALARM	97	38	54
STOLEN PROPERTY	63	38	32
STABBING	3	1	0
STOLEN PROPERTY	38	7	1
STOLEN PROPERTY	38	3	26
STRONGARM ROBBERY	164	32	68
SUSPICIOUS PERSON	3305	1170	146
SUSPICIOUS PERSON	3305	1170	146
THREATS / HARASSMENT	965	37	146
TRESPASS	759	367	383
TRAF VIOLATION CITY	579	488	456
TRAF VIOLATION STATE	582	177	187
TRAFFIC HAZARD	2759	264	204
TRESPASSER	246	246	1998
UNLAWFUL IDENTIFICATION	758	163	230
VANDALISM	83	1	4
VEHICLE ALARM	203	123	123
WELL BEING CHECK	3045	522	495

8 threads - 1.33 seconds