

# Operating systems – Assignment 2

## Scheduling

Lennart Jern  
CS: ens16ljn

**Teacher**  
Ahmed Aley

December 17, 2016

# 1 Introduction

The Linux kernel provides a number of different scheduling policies that can be used to fine tune the performance of certain applications. In this report, five different schedulers are evaluated using an artificial, CPU intensive, task. Three of the tested schedulers are “normal”, while the last two are “real-time” schedulers, meaning that they provide higher priority for their processes than the normal ones do.

The work load consists of a simple program, called `work`, that sums over a part of Grandi’s series<sup>1</sup> ( $1 - 1 + 1 - 1 + \dots$ ), using a specified number of threads. The source code for `work` can be found in listing 1. Since the task is easy to parallelize, only require minimal memory access and no disk access, it should be comparable to CPU intense tasks like compression and matrix calculations.

## 2 Method

A Bash script (`timer.sh`, listing 2) was used to time the complete task 10 times for each scheduler, for thread counts ranging from 1 to 10. See code listing 2 for the code. Additionally, each thread keeps track of the time taken from the start of its execution until it is finished, and prints this information to `stdout`, which is forwarded to data files by the bash script.

All this data was then processed by a simple Python program, `stats`, in order to calculate the median, minimum and maximum run time for each scheduler and thread count; for both the total run times and the individual thread times. This program can be found in listing 3. In addition to the statistical calculations, `stats` also produces some figures to describe the data. The figures and calculations are mostly done using the python libraries Pandas<sup>2</sup> and Matplotlib<sup>3</sup>.

It should be noted here that the processes running with real-time schedulers were run with maximum priority. Since the other schedulers does not accept a priority setting, other than the nice value, these were left untouched.

All tests were run on my personal computer with the specifications seen in table 1.

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Grandi's\\_series](https://en.wikipedia.org/wiki/Grandi's_series)

<sup>2</sup><http://pandas.pydata.org/>

<sup>3</sup><http://matplotlib.org/>

| Component   | Specification                    |
|-------------|----------------------------------|
| OS:         | Fedora 25                        |
| Kernel:     | Linux 4.8.13-300.fc25.x86_64     |
| CPU:        | Intel Core i5-2500K CPU @ 3.7GHz |
| RAM:        | 7965MiB                          |
| GCC:        | 6.2.1                            |
| Bash:       | 4.3.43                           |
| Python:     | 3.5.2                            |
| Pandas:     | 0.18.1                           |
| Matplotlib: | 1.5.3                            |

Table 1: Test system specification.

### 3 Results

For the total run time of the process, there does not seem to be much of a difference between the five schedulers (see fig. 1). This is perhaps not too surprising considering that the system was practically idle, except for the work load process, so this process naturally got all resources the system could offer.

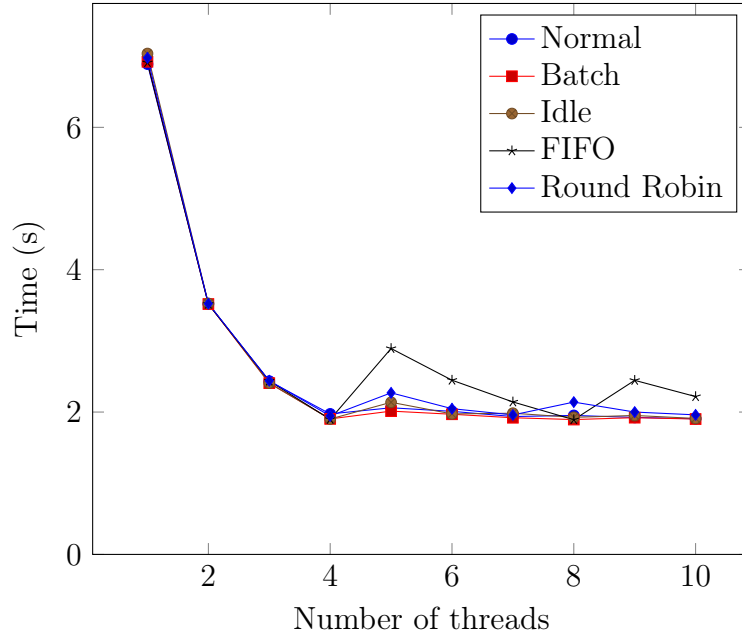


Figure 1: The median time required to finish the complete task.

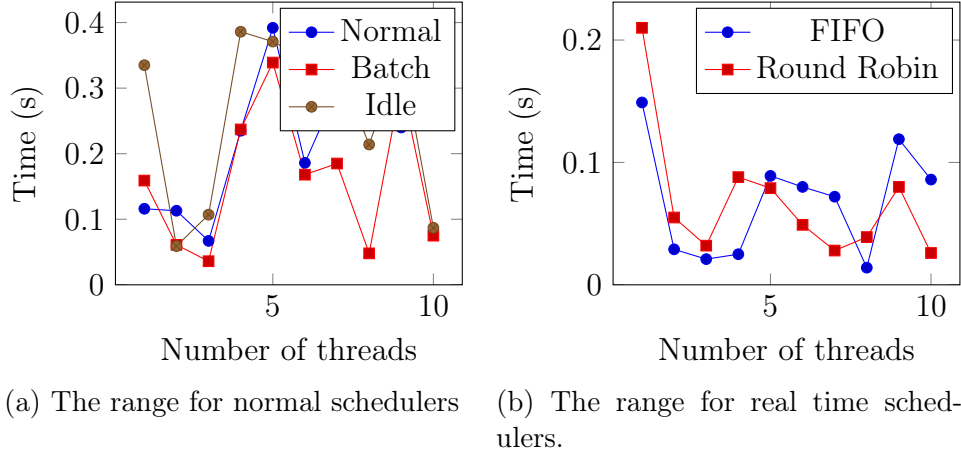


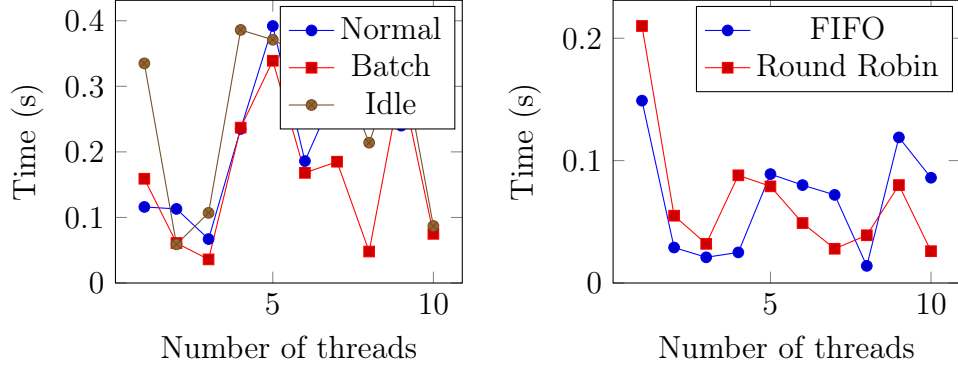
Figure 2: A comparison of response times for the complete task, between the different schedulers.

More interesting is the range (difference between maximum and minimum) of response times for the total run time of the process (fig. 2). This reveals a clear difference between the real-time and normal schedulers, where the real-time ones are clearly more predictable for two or more threads.

The same trend can also be seen for the thread response times in fig. 3. Furthermore, there is a small difference in median run time for the schedulers, when looking at the individual thread times, for a thread count between 6 and 8. This can be seen in fig. 4, where the Batch scheduler have a lower median time than the others. However, the difference seem to be somewhat periodic and vanishes almost completely with 9 threads. The Round Robin scheduler also follows this pattern but with a shorter period. With the exception of Batch, Round Robin beats the other schedulers for thread counts of 5, 6, 9 and 10.

The differences are even clearer when looking at the box charts for the individual thread times in figures 5 and 6 where the range between the maximum and minimum run times can also be seen.

To further analyze the behavior of the schedulers, we can take a look at the density plots for the thread times in figures 7 and 8. There are only small differences when running with two threads. But with four threads instead, we can see that the Batch scheduler have two distinct “batches” while RR is much closer to have all threads finish in one lump. I am quite surprised to see the FIFO curve so flat and spread out at four threads, since they should all be able to run together and finish simultaneously.



(a) The range in response times for normal schedulers (b) The range in response times for real time schedulers.

Figure 3: A comparison of response times per thread between the different schedulers.

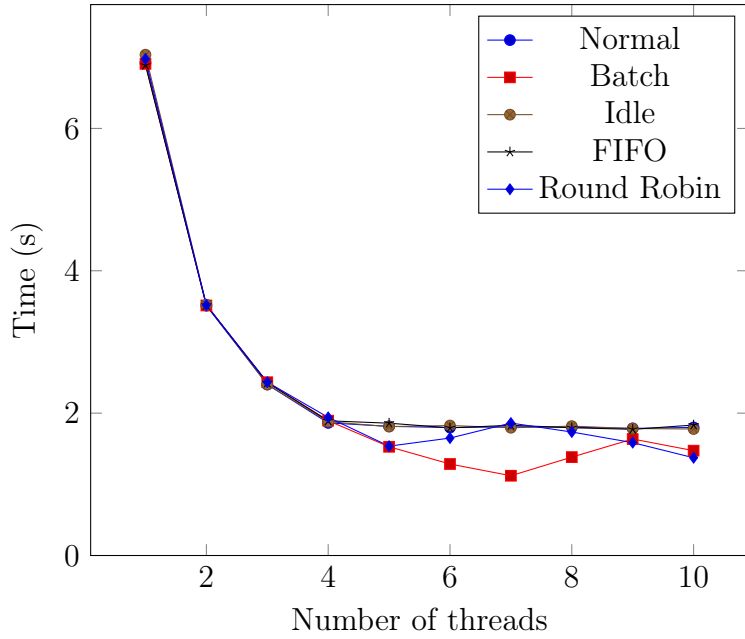
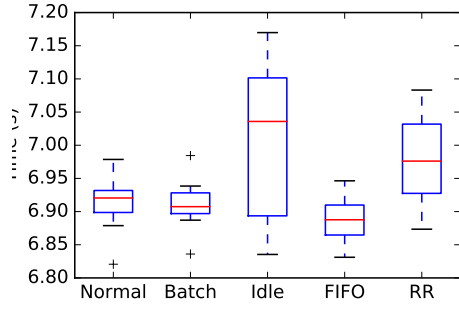
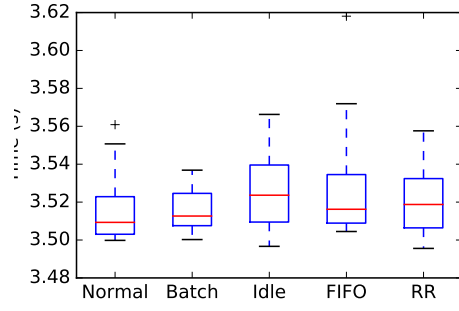


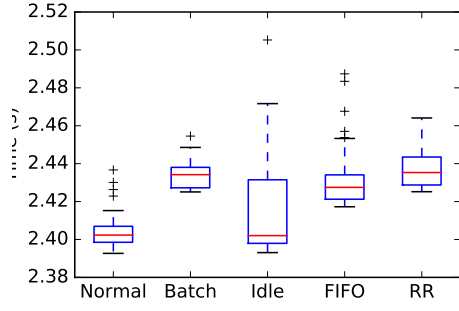
Figure 4: The median time per thread.



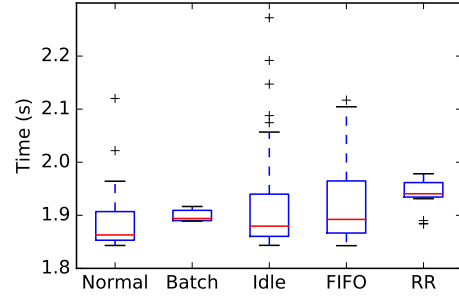
(a) Thread count: 1



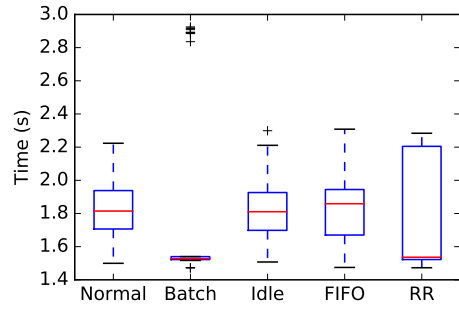
(b) Thread count: 2



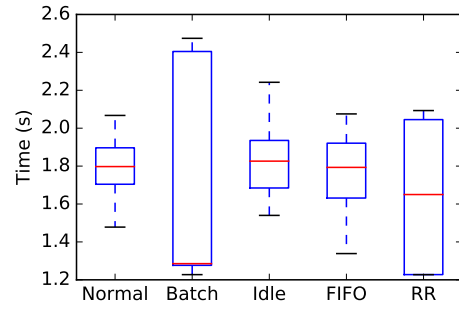
(c) Thread count: 3



(d) Thread count: 4

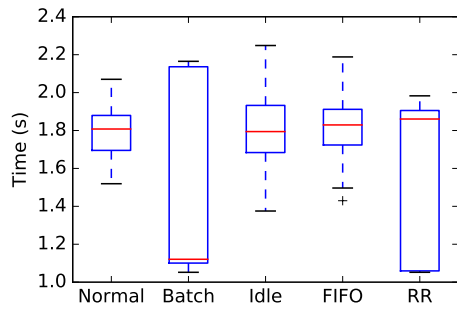


(e) Thread count: 5

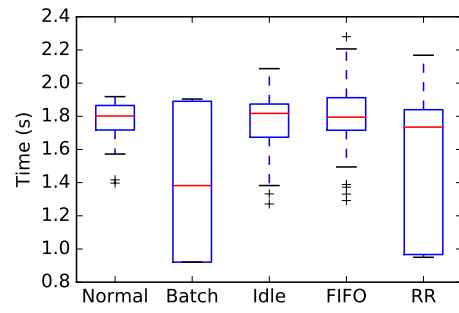


(f) Thread count: 6

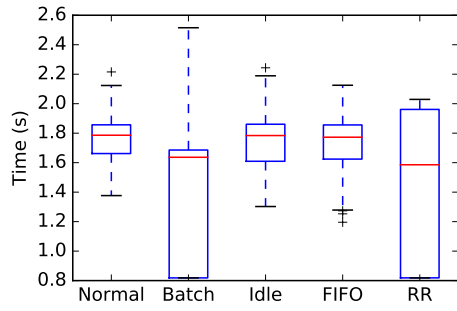
Figure 5: Box charts for thread counts 1 to 6.



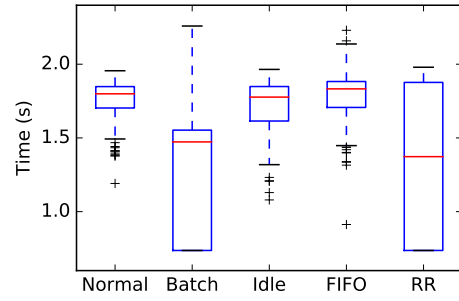
(a) Thread count: 7



(b) Thread count: 8

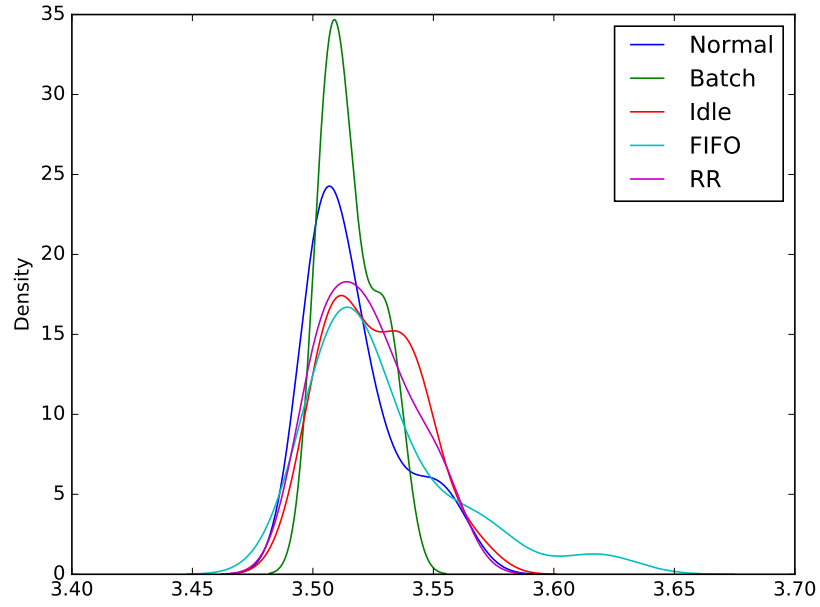


(c) Thread count: 9

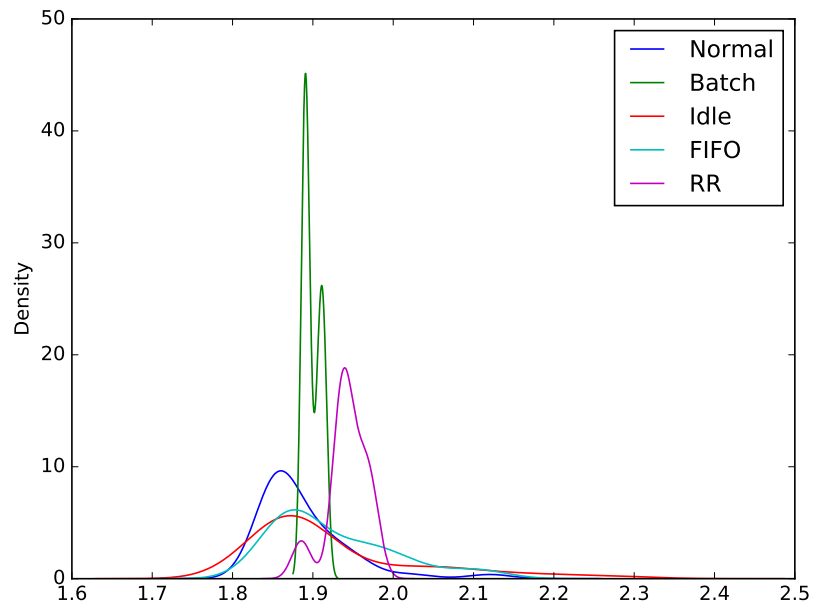


(d) Thread count: 10

Figure 6: Box charts for thread counts 7 to 10.



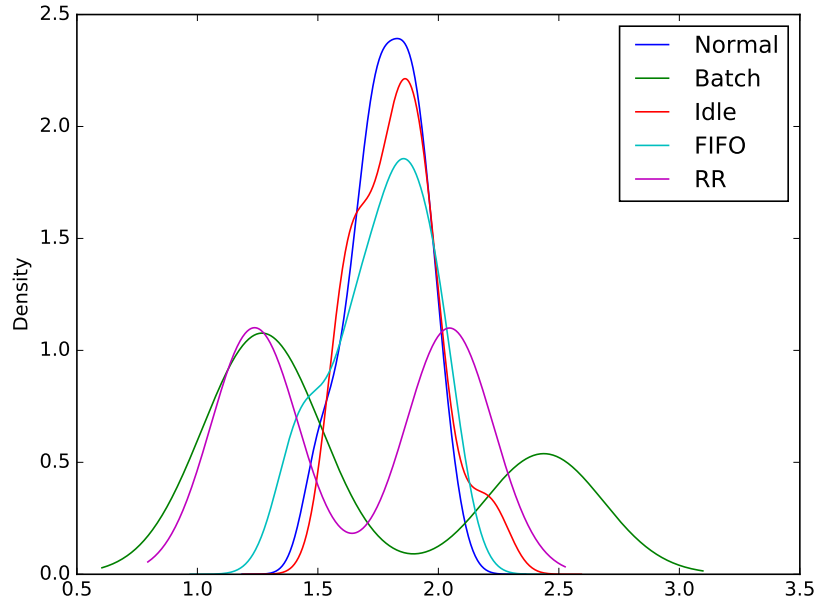
(a) Thread count: 2



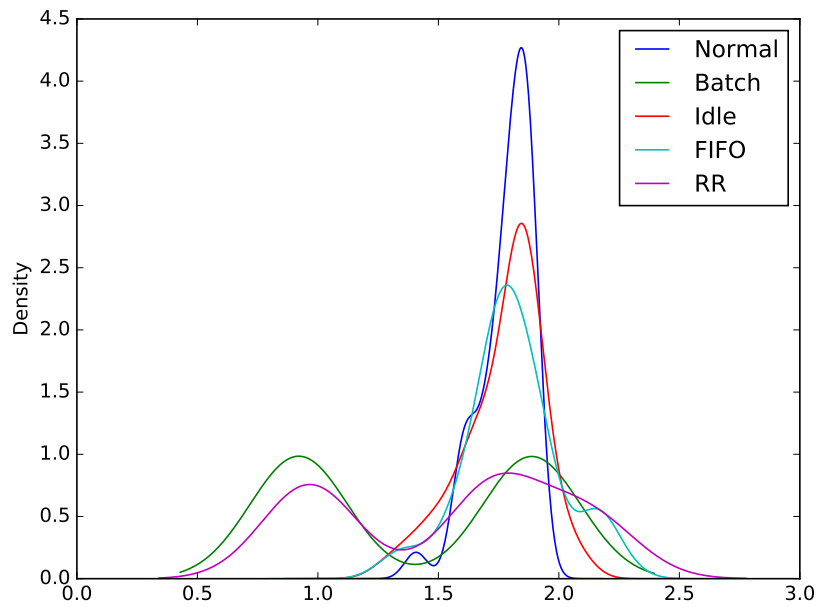
(b) Thread count: 4

Figure 7: Density plots for the different schedulers when running with 2 and 4 thread





(a) Thread count: 6



(b) Thread count: 8

Figure 8: Density plots for the different schedulers when running with 6 and 8 thread

At six and eight threads, a pattern emerges where the Normal, Idle and FIFO schedulers finishes all their threads close together, while the Batch and RR schedulers both have two groups of threads finishing together. With a higher resolution we would probably see two peaks for FIFO also, a small bump is already visible in the plots for six and eight threads.

The results for Round Robin follow the theory nicely, as it becomes more and more flat as the thread count increases (compare the figures for six and eight threads), since it tries to distribute the CPU time as evenly as possible. The Batch scheduler avoids switching between threads, and should thus get one batch of four and one of two threads finishing together when running with six threads total (since the processor has four cores). With eight threads, there should be two equal batches instead, and this is exactly what we can see in figures 8a and 8b. FIFO should in theory get two distinct peaks, but the resolution in the data does not seem to be high enough for this to appear.

## 4 Final thoughts and lessons learned

It is quite clear that it is possible to spend a considerable amount of time just analyzing schedulers. The results are intriguing as they show clear differences in the behavior between the schedulers, and with several different tasks to compare, even more patterns would surely emerge.

A lesson learned is that one should think carefully about when and where to start and stop the timers. If measuring just the whole task, this is quite easy, but for the individual threads it gets more tricky. The thread that starts the other threads must also start the timers since the working thread may not get to start immediately. Similarly, the working threads must stop their own timers, since there might be a pause between the threads finishing and actually getting joined.

## A Code listings

Listing 1: work.c

```
1  /**
2   * work.c
3   *
4   * Just a silly "do something that takes time" program.
5   * It tries to calculate the sum of Grandi's series (1-1+1-1+1-1...).
6   * As long as the length it is summing over is even the sum should always be 0.
7   *
8   * Author: Lennart Jern (ens16ljn)
9   *
10  * Call sequence: work [-p <policy>] [-j <number of jobs>]
11  * The policy is given by a single char according to this:
12  * n - Normal
13  * b - Batch
14  * i - Idle
15  * f - FIFO
16  * r - RR
17  */
18
19 #include <stdio.h>
20 #include <stdlib.h>
21 #include <time.h>           // timing
22 #include <errno.h>
23 #include <pthread.h>        // threading
24 #include <sys/types.h>      // pid
25 #include <unistd.h>         // pid, getopt
26 #include <linux/sched.h>    // scheduling policies
27
28 // Length of sequence to sum
29 #define LENGTH 2147483400
30 #define ONE_OVER_BILLION 1E-9
31
32 typedef struct work_load {
33     int nworkers;
34     long data_length;
35     char scheduler;
36 } WorkLoad;
37
38 typedef struct work_packet {
39     long index;
40     long length;
41     long result;
42     struct timespec start;
43 } Packet;
44
45 WorkLoad *get_work_load();
46 void *work(void *data);
47 int get_grandi(int index);
48 long calculate_sum(long index, long length);
49 void run_workers(WorkLoad *wl);
50 void print_scheduler();
51 void set_scheduler(WorkLoad *wl);
52 void set_settings(WorkLoad *wl, int argc, char *argv[]);
53
54 int num_policies = 5;
55 char c_policies[] = {'n', 'b', 'i', 'f', 'r'};
56 char *str_policies[] = {"Normal", "Batch", "Idle", "FIFO", "RR"};
57 int policies[] = {SCHED_NORMAL, SCHED_BATCH, SCHED_IDLE,
```

```

58         SCHED_FIFO, SCHED_RR};
59
60 int main(int argc, char *argv[]) {
61
62     WorkLoad *wl = get_work_load();
63
64     set_settings(wl, argc, argv);
65
66     // Set scheduler
67     set_scheduler(wl);
68
69     // Print scheduler to make sure it is set coorectly
70     print_scheduler();
71
72     run_workers(wl);
73
74     free(wl);
75     printf("Done\n");
76 }
77
78 /**
79  * get_work_load - initialize the work load and return a pointer to it
80  * @return pointer to allocated memory
81  */
82 WorkLoad *get_work_load() {
83     WorkLoad *wl;
84     // Allocate memory for work load
85     wl = malloc(sizeof(WorkLoad));
86     // Initialize work load
87     wl->nworkers = 1;
88     // wl->data_length = 1073741824; // 2^30
89     wl->data_length = LENGTH;
90     return wl;
91 }
92
93 /**
94  * work - a silly attempt to calculate the limit of Grandi's series
95  * @param packet the part of the work load to work on
96  * @return nothing
97  */
98 void *work(void *packet) {
99     Packet *pkt = (Packet *)packet;
100     long sum = 0;
101     // Starting time
102     struct timespec start = pkt->start;
103     // Time when finished
104     struct timespec end;
105
106     // Do the actual work
107     sum = sum + calculate_sum(pkt->index, pkt->length);
108
109     // Get the time when finished
110     clock_gettime(CLOCK_REALTIME, &end);
111     // Calculate time it took
112     double time_taken = (end.tv_sec - start.tv_sec)
113         + (end.tv_nsec - start.tv_nsec)
114         * ONE_OVER_BILLION;
115     printf("%lf\n", time_taken);
116 }
117
118 /**
119  * get_grandi - calculate the i:th number of Grandi's series

```

```

120  * @param index    index of the number you want to know
121  * @return         1 if index is even, -1 otherwise
122  */
123  int get_grandi(int index) {
124      if (index % 2 == 0) {
125          return 1;
126      } else {
127          return -1;
128      }
129  }
130
131  /**
132   * calculate_sum - sum Grandi's series from index over a given length
133   * @param index    index to start from
134   * @param length    how many numbers to sum over
135   * @return         the sum
136   */
137  long calculate_sum(long index, long length) {
138      long sum = 0;
139
140      for (int i = index; i < index+length; i++) {
141          sum = sum + get_grandi(i);
142      }
143      return sum;
144  }
145
146  /**
147   * run_workers - start work threads and wait for them to finish
148   */
149  void run_workers(WorkLoad *wl) {
150      int num = wl->nworkers;
151      long total_sum = 0;
152
153      // Allocate memory for all packets
154      Packet *pkt[num];
155      for (int i = 0; i < num; i++) {
156          pkt[i] = malloc(sizeof(Packet));
157          if (!pkt[i]) {
158              perror("malloc");
159          }
160          pkt[i]->result = 0;
161      }
162
163      // create threads
164      pthread_t threads[num];
165      long len = wl->data_length;
166      long p_len = len / num;
167      int i;
168      for (i = 0; i < num-1; i++) {
169          pkt[i]->index = i * len / num;
170          pkt[i]->length = p_len;
171          // Set the starting time
172          clock_gettime(CLOCK_REALTIME, &pkt[i]->start);
173          if (pthread_create(&threads[i], NULL, work, (void *)pkt[i]) != 0) {
174              perror("Could not create thread");
175          }
176      }
177      // Run this thread instead of just waiting for the others
178      pkt[i]->index = i * len / num;
179      pkt[i]->length = p_len;
180      clock_gettime(CLOCK_REALTIME, &pkt[i]->start);
181      work((void *)pkt[i]);

```

```

182
183     total_sum += pkt[i]->result;
184     free(pkt[i]);
185
186     // Join the threads
187     for (int i = 0; i < num-1; i++) {
188         pthread_join(threads[i], NULL);
189         total_sum += pkt[i]->result;
190         free(pkt[i]);
191     }
192
193     printf("Sum is %d\n", total_sum);
194 }
195
196 /**
197  * print_scheduler - print the current scheduler
198  * @param pid the pid of the process
199  */
200 void print_scheduler() {
201     pid_t pid = getpid();
202     int schedlr = sched_getscheduler(pid);
203
204     char *schedlr_name;
205     switch (schedlr) {
206         case SCHED_NORMAL:
207             schedlr_name = "Normal/Other";
208             break;
209         case SCHED_BATCH:
210             schedlr_name = "Batch";
211             break;
212         case SCHED_IDLE:
213             schedlr_name = "Idle";
214             break;
215         case SCHED_FIFO:
216             schedlr_name = "FIFO";
217             break;
218         case SCHED_RR:
219             schedlr_name = "RR";
220             break;
221         default:
222             schedlr_name = "Unknown";
223     }
224     printf("Scheduler: %s\n", schedlr_name);
225 }
226
227 /**
228  * set_scheduler - update scheduler to reflect the given WorkLoad
229  * @param wl the work load
230  */
231 void set_scheduler(WorkLoad *wl) {
232     struct sched_param param;
233     pid_t pid = getpid();
234     int policy = SCHED_NORMAL;
235
236     for (int i = 0; i < num_policies; i++) {
237         if (wl->scheduler == c_policies[i]) {
238             policy = policies[i];
239             break;
240         }
241     }
242
243     // Set the priority

```

```

244     param.sched_priority = sched_get_priority_max(policy);
245
246     if (sched_setscheduler(pid, policy, &param) != 0) {
247         perror("Set scheduler");
248     }
249 }
250
251 /**
252  * set_settings - parse arguments and set the settings for the work load
253  * @param wl     the work load to update
254  * @param argc   argument count
255  * @param argv   array of arguments
256  */
257 void set_settings(WorkLoad *wl, int argc, char *argv[]) {
258     // Two possible options: j(obs) and p(policy)
259     char *optstr = "j:p:";
260     int opt;
261     char policy = 'n';
262     int num_threads = 1;
263     int policy_ok = 0;
264     int threads_ok = 0;
265
266     // Parse flags
267     while ((opt = getopt(argc, argv, optstr)) != -1) {
268         char *end;
269         switch (opt) {
270             case 'p':
271                 policy = *optarg;
272                 break;
273             case 'j':
274                 errno = 0;
275                 num_threads = strtol(optarg, &end, 10);
276                 if (errno != 0) {
277                     perror("strtol");
278                 }
279                 break;
280             default:
281                 printf("Option %c not supported\n", opt);
282         }
283     }
284
285     // Check the parsed options
286     for (int i = 0; i < num_policies; i++) {
287         if (policy == c_policies[i]) {
288             policy_ok = 1;
289             break;
290         }
291     }
292
293     // Boundaries for the number of threads
294     if (num_threads <= 100 && num_threads > 0) {
295         threads_ok = 1;
296     }
297
298     // Set values if they are safe, or set defaults
299     if (policy_ok) {
300         wl->scheduler = policy;
301     } else {
302         wl->scheduler = 'n';
303     }
304
305     if (threads_ok) {

```

```

306     wl->nworkers = num_threads;
307 } else {
308     wl->nworkers = 1;
309 }
310 }

```

Listing 2: timer.sh

```

1  #!/bin/bash
2
3  # timer.sh
4  #
5  # A timer script to measure the differences between schedulers/policies
6  #
7  # Author: Lennart Jern (ens16ljn@cs.umu.se)
8
9  for THREADS in $(seq 1 10)
10 do
11     DATA="Normal,Batch,Idle,FIFO,RR"
12     echo "Running with_${THREADS}_threads"
13     # Time the commands 10 times
14     for i in $(seq 1 10)
15     do
16         LINE=""
17         # For the policies n(ormal) b(atch) and i(dle)
18         for POLICY in n b i f r
19         do
20             # Set policy and number of threads
21             FLAGS="-p$POLICY-_$THREADS"
22             COMMAND="./work_${FLAGS}>>../data/threads$THREADS$POLICY.log"
23             # Run the command and store the time
24             t="$(sh-c_"TIMEFORMAT='%5R'; time $COMMAND">>&1)"
25             # Build the line
26             if [ "$POLICY" = "n" ]; then
27                 LINE="$t"
28             else
29                 LINE="$LINE,$t"
30             fi
31         done
32         DATA=$DATA$'\n'$LINE
33         # A little progress report
34         echo "Run_${i}_done."
35     done
36
37     # Write data to a file
38     echo "$DATA" > "../data/data$THREADS.csv"
39     chown lennart ../data/threads*
40     chown lennart ../data/data*
41
42 done

```

Listing 3: stats.py

```

1  """
2  stats.py
3
4  Process the data produced by timer.sh by calculating the
5  medians, max values and min values for each scheduler and thread count.
6  Also collects the data about individual thread times, and produces similar
7  statistics for them.

```



```

8  Lastly, produces a set of plots describing the data.
9
10 Author: Lennart Jern (ens16ljn@cs.umu.se)
11 """
12
13 import pandas as pd
14 import re
15 import matplotlib.pyplot as plt
16
17 def total_stats():
18     """
19     Read the data from files and calculate statistical values and make plots.
20     """
21     # The data file names are of the form data<thread count>.csv
22     base = "../data/data"
23     thread_base = "../data/threads"
24     ext = ".csv"
25     header = ("Normal", "Batch", "Idle", "FIFO", "RR")
26     # Data frames to store the results in
27     med = pd.DataFrame(columns=header)          # Medians (total runtime)
28     mx = pd.DataFrame(columns=header)          # Max (total runtime)
29     mn = pd.DataFrame(columns=header)          # Min (total runtime)
30     thread_med = pd.DataFrame(columns=header)  # Medians (threads)
31     thread_mx = pd.DataFrame(columns=header)   # Max (threads)
32     thread_mn = pd.DataFrame(columns=header)   # Min (threads)
33
34     # For each number of threads
35     for i in range(1,11):
36         # Build the file name
37         f = base + str(i) + ext                # Total run times
38         thr_f = thread_base + str(i) + ext     # Thread times
39         # Read the time data
40         df = pd.read_csv(f)
41         thr_df = pd.read_csv(thr_f)
42
43         # Calculate some statistical properties
44         med.loc[i] = df.median()
45         mx.loc[i] = df.max()
46         mn.loc[i] = df.min()
47         thread_med.loc[i] = thr_df.median()
48         thread_mx.loc[i] = thr_df.max()
49         thread_mn.loc[i] = thr_df.min()
50
51         # Plot and save some nice figures
52         # Density curves for thread count 2, 4, 6 and 8
53         if (i == 2 or i == 4 or i == 6 or i == 8):
54             ax = thr_df.plot.kde()
55             ax.set_xlabel("Time(s)")
56             fig = ax.get_figure()
57             fig.savefig("density"+str(i)+".pdf")
58
59         # Box plots for all thread counts
60         ax = thr_df.plot.box(figsize=(4.5,3))
61         ax.set_ylabel("Time(s)")
62         fig = ax.get_figure()
63         fig.savefig("box"+str(i)+".pdf")
64
65     # Calculate ranges
66     rng = mx-mn
67     thr_rng = thread_mx-thread_mn
68
69

```

```

70     # Write everything to files
71     data_frames = [med, mx, mn, thread_med, thread_mx, thread_mn, rng, thr_rng]
72     base = "../data/"
73     names = ["medians", "max", "min", "thread_medians", "thread_max",
74             "thread_min", "range", "thread_range"]
75     for frm, name in zip(data_frames, names):
76         frm.to_csv(base+name + ext, index_label="Threads", float_format="%.5f")
77
78
79 def collect_thread_times(file_name):
80     """Read thread times from a file."""
81     f = open(file_name)
82     times = []
83     # Regular expression to find floats
84     time = re.compile("(\\d+\\.\\d+)")
85
86     for line in f:
87         match = time.match(line)
88
89         if (match):
90             t = float(match.group(1))
91             times.append(t)
92
93     return times
94
95
96 def thread_stats():
97     """Collect timing information about all threads and store in csv files"""
98     threads = [i for i in range(1, 11)]
99     schedulers = ["n", "b", "i", "f", "r"]
100    base = "../data/threads"
101    ext = ".log"
102    header=("Normal", "Batch", "Idle", "FIFO", "RR")
103
104    # Collect all times for one thread count in one file
105    for t in threads:
106        times = {key: [] for key in schedulers}
107        for s in schedulers:
108            f = get_file_name(t, s)
109            times[s] = collect_thread_times(f)
110        # Write to file
111        df = pd.DataFrame(times)
112        df.to_csv(get_csv_name(t), index=False, header=header)
113
114
115 def get_file_name(threads, scheduler):
116     """Get the file name for the data regarding <scheduler> and <threads>"""
117     base = "../data/threads"
118     ext = ".log"
119     return base + str(threads) + scheduler + ext
120
121 def get_csv_name(threads):
122     """For a specific number of threads: Get name of file to write data to"""
123     base = "../data/threads"
124     ext = ".csv"
125     return base + str(threads) + ext
126
127 # Collect thread timings
128 thread_stats()
129 # Get statistics and plots
130 total_stats()

```