

Kernel Data Structures

Evolution of Kernel Scheduling

Department of Computer Science
and Engineering

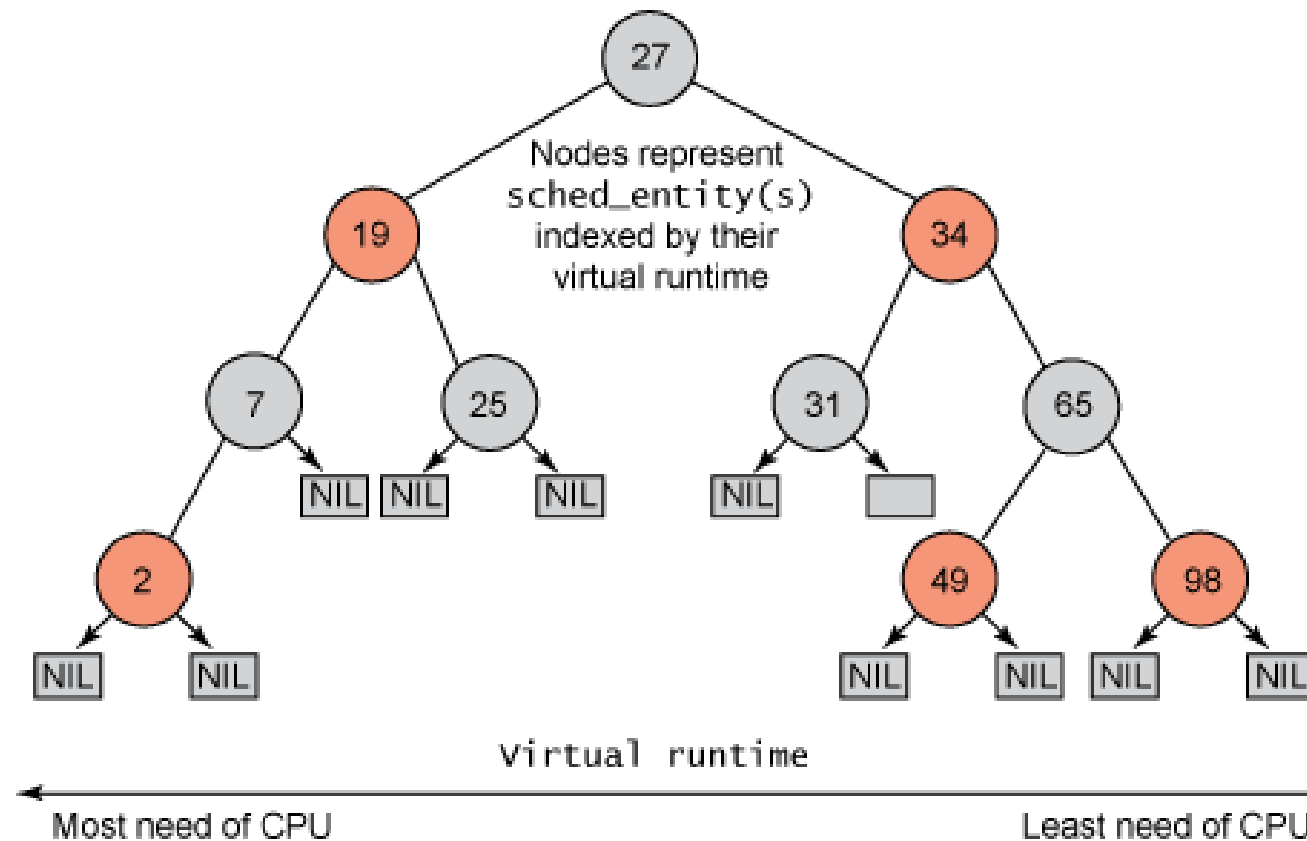


INDIAN INSTITUTE OF TECHNOLOGY
KHARAGPUR

Sandip Chakraborty
sandipc@cse.iitkgp.ac.in



The Evolution of Linux Process Schedulers



The Genesis Scheduler

- Introduced in 1991 within Kernel 0.01
 - Fairly simple, minimum design
 - Does not aim to support massive multiprocessing systems

The Genesis Scheduler

- Introduced in 1991 within Kernel 0.01
 - Fairly simple, minimum design
 - Does not aim to support massive multiprocessing systems
- A single process queue – the scheduler iterates over the entire queue to select a task to run
 - Queue size is not very long; default NR_TASKS of kernel set to 32

The Genesis Scheduler – The Code

```
1 void schedule(void) {
2   int i,next,c;
3   struct task_struct ** p;
4
5   /* check alarm, wake up any interruptible tasks
6      that have got a signal */
7   for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
8     if (*p) {
9       if ((*p)->alarm && (*p)->alarm < jiffies) {
10         (*p)->signal |= (1<<(SIGALRM-1));
11         (*p)->alarm = 0;
12       }
13       if ((*p)->signal && (*p)->state==TASK_INTERRUPTIBLE)
14         (*p)->state=TASK_RUNNING;
15     }
16 }
```

The Genesis Scheduler – The Code

```
1 void schedule(void) {
2   int i,next,c;
3   struct task_struct ** p;
4
5   /* check alarm, wake up any interruptible tasks
6      that have got a signal */
7   for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
8     if (*p) {
9       if ((*p)->alarm && (*p)->alarm < jiffies) {
10         (*p)->signal |= (1<<(SIGALRM-1));
11         (*p)->alarm = 0;
12       }
13       if ((*p)->signal && (*p)->state==TASK_INTERRUPTIBLE)
14         (*p)->state=TASK_RUNNING;
15     }
16 }
```


Iterate over all the tasks in the task array



The Genesis Scheduler – The Code

```
1 void schedule(void) {
2   int i,next,c;
3   struct task_struct ** p;
4
5   /* check alarm, wake up any interruptible tasks
6      that have got a signal */
7   for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
8     if (*p) {
9       if ((*p)->alarm && (*p)->alarm < jiffies) {
10         (*p)->signal |= (1<<(SIGALRM-1));
11         (*p)->alarm = 0;
12       }
13       if ((*p)->signal && (*p)->state==TASK_INTERRUPTIBLE)
14         (*p)->state=TASK_RUNNING;
15     }
16 }
```

Check whether an alarm has been raised before the current jiffies (the interrupt timer counter)



The Genesis Scheduler – The Code


```
1 void schedule(void) {
2   int i,next,c;
3   struct task_struct ** p;
4
5   /* check alarm, wake up any interruptible tasks
6      that have got a signal */
7   for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
8     if (*p) {
9       if ((*p)->alarm && (*p)->alarm < jiffies) {
10         (*p)->signal |= (1<<(SIGALRM-1));
11         (*p)->alarm = 0;
12       }
13       if ((*p)->signal && (*p)->state==TASK_INTERRUPTIBLE)
14         (*p)->state=TASK_RUNNING;
15     }
16 }
```



**Raise the corresponding
signal and reset the alarm**

The Genesis Scheduler – The Code

```
1 void schedule(void) {
2   int i,next,c;
3   struct task_struct ** p;
4
5   /* check alarm, wake up any interruptible tasks
6      that have got a signal */
7   for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
8     if (*p) {
9       if ((*p)->alarm && (*p)->alarm < jiffies) {
10         (*p)->signal |= (1<<(SIGALRM-1));
11         (*p)->alarm = 0;
12       }
13       if ((*p)->signal && (*p)->state==TASK_INTERRUPTIBLE)
14         (*p)->state=TASK_RUNNING;
15     }
16 }
```



**If the task is waiting, then
move it to the running state
(ready to run)**

The Genesis Scheduler – The Code

```
17  /* this is the scheduler proper: */
18  while (1) {
19      c = -1;
20      next = 0;
21      i = NR_TASKS;
22      p = &task[NR_TASKS];
23      while (--i) {
24          if (!*--p)
25              continue;
26          if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27              c = (*p)->counter, next = i;
28      }
29      if (c) break;
```

**Iterate until a runnable
process with largest unused
chunk of timeslice is found**



The Genesis Scheduler – The Code

```
17  /* this is the scheduler proper: */
18  while (1) {
19      c = -1;
20      next = 0;
21      i = NR_TASKS;
22      p = &task[NR_TASKS];
23      while (--i) {
24          if (!*--p)
25              continue;
26          if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27              c = (*p)->counter, next = i;
28      }
29      if (c) break;
```



Iterate backwards

The Genesis Scheduler – The Code

```
17  /* this is the scheduler proper: */
18  while (1) {
19      c = -1;
20      next = 0;
21      i = NR_TASKS;
22      p = &task[NR_TASKS];
23      while (--i) {
24          if (!*--p)
25              continue;
26          if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27              c = (*p)->counter, next = i;
28      }
29      if (c) break;
```



Not a valid task, continue the iteration

The Genesis Scheduler – The Code

```
17  /* this is the scheduler proper: */
18  while (1) {
19      c = -1;
20      next = 0;
21      i = NR_TASKS;
22      p = &task[NR_TASKS];
23      while (--i) {
24          if (!*--p)
25              continue;
26          if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27              c = (*p)->counter, next = i;
28      }
29      if (c) break;
```

A runnable task and have the maximum available timeslice

The Genesis Scheduler – The Code

```
17  /* this is the scheduler proper: */
18  while (1) {
19      c = -1;
20      next = 0;
21      i = NR_TASKS;
22      p = &task[NR_TASKS];
23      while (--i) {
24          if (!*--p)
25              continue;
26          if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27              c = (*p)->counter, next = i;
28      }
29      if (c) break;
```

Come out of the loop if you have seen a process with maximum available timeslice

The Genesis Scheduler – The Code

```
30     for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
31         if (*p)
32             (*p)->counter = ((*p)->counter >> 1) + (*p)->priority;
33     }
34     switch_to(next);
35 }
```

Iterate over the task array if
all the processes have taken
the allocated timeslice



The Genesis Scheduler – The Code


```
30     for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
31         if (*p)
32             (*p)->counter = ((*p)->counter >> 1) + (*p)->priority;
33     }
34     switch_to(next);
35 }
```



**Reallocate timeslice equals to its
previous timeslice divided by 2 plus its
priority**

The Genesis Scheduler – The Code

```
30     for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
31         if (*p)
32             (*p)->counter = ((*p)->counter >> 1) + (*p)->priority;
33     }
34     switch_to(next);
35 }
```




**Trigger a context switch to
execute the next task (task
having the maximum
available timeslice)**

```
26         if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27             c = (*p)->counter, next = i;
```

The Genesis Scheduler – The Code

```
30     for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
31         if (*p)
32             (*p)->counter = ((*p)->counter >> 1) + (*p)->priority;
33     }
34     switch_to(next);
35 }
```

Running complexity?




**Trigger a context switch to
execute the next task (task
having the maximum
available timeslice)**

```
26         if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27             c = (*p)->counter, next = i;
```

The Genesis Scheduler – The Code

```
30     for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
31         if (*p)
32             (*p)->counter = ((*p)->counter >> 1) + (*p)->priority;
33     }
34     switch_to(next);
35 }
```



Trigger a context switch to
execute the next task (task
having the maximum
available timeslice)

Running complexity?

$O(n)$

```
26     if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
27         c = (*p)->counter, next = i;
```

O(n) Scheduler – Kernel 2.4

- Conceptually same as the Genesis scheduler
 - Update is in the metric used for selecting the next process – **Goodness** of a process
- Goodness of a process is calculated as the number of clock-ticks a task had left plus some weight based on the task's priority; returns integer values
 - -1000: Never select this task to run
 - Positive number: The goodness value
 - +1000: A real time process

O(n) Scheduler – Kernel 2.4

- Conceptually same as the Genesis scheduler
 - Update is in the metric used for selecting the next process – *Goodness* of a process
- Goodness of a process is calculated as the number of clock-ticks a task had left plus some weight based on the task's priority; returns integer values
 - -1000: Never select this task to run
 - Positive number: The goodness value
 - +1000: A real time process

O(n) Scheduler – Kernel 2.4

- Conceptually same as the Genesis scheduler
 - Update is in the metric used for selecting the next process – **Goodness** of a process
- Goodness of a process is calculated as the number of clock-ticks a task had left plus some weight based on the task's priority; returns integer values
 - -1000: Never select this task to run
 - Positive number: The goodness value
 - +1000: A real time process
- **Cons:** A lot of time goes for goodness calculation

O(1) Scheduler – Kernel 2.6

Introduced a number of new features --

- **Global priority scale** (0-139), separation between real-time (0-99) and normal (100-139) processes based on priority values

O(1) Scheduler – Kernel 2.6

Introduced a number of new features --

- **Global priority scale** (0-139), separation between real-time (0-99) and normal (100-139) processes based on priority values
- **Early preemption:** When a new task enters in the TASK_RUNNING state, it can preempt the currently running process if having lower priority

O(1) Scheduler – Kernel 2.6

Introduced a number of new features --

- **Global priority scale** (0-139), separation between real-time (0-99) and normal (100-139) processes based on priority values
- **Early preemption:** When a new task enters in the TASK_RUNNING state, it can preempt the currently running process if having lower priority
- Static priority for real-time tasks, dynamic priority for normal tasks

O(1) Scheduler – Kernel 2.6

Introduced a number of new features --

- **Global priority scale** (0-139), separation between real-time (0-99) and normal (100-139) processes based on priority values
- **Early preemption:** When a new task enters in the TASK_RUNNING state, it can preempt the currently running process if having lower priority
- Static priority for real-time tasks, dynamic priority for normal tasks
- Dynamic priority for the normal tasks is decided depending on its interactivity in the past – how frequently the process runs CPU-bound instructions

Moving from $O(n)$ to $O(1)$

- **The Good:** Timeslices are allocated in proportion to the priority of the tasks
 - higher priority tasks get more timeslice to complete the job

Moving from $O(n)$ to $O(1)$

- **The Good:** Timeslices are allocated in proportion to the priority of the tasks
 - higher priority tasks get more timeslice to complete the job
- **The Bad:** Priorities are not fixed all the time
 - New tasks with higher priority can arrive – preempt the currently running process?
 - Process priority can get updated dynamically (Kernel 2.6)

Moving from $O(n)$ to $O(1)$

- **The Good:** Timeslices are allocated in proportion to the priority of the tasks
 - higher priority tasks get more timeslice to complete the job
- **The Bad:** Priorities are not fixed all the time
 - New tasks with higher priority can arrive – preempt the currently running process?
 - Process priority can get updated dynamically (Kernel 2.6)
- **The Ugly:** Recompute the amount of timeslices to be allocated to each task after the current epoch is complete – **$O(n)$**
 - Iterate over the entire runqueue
 - Recompute each task's priority and the timeslice to be allocated

$O(1)$ Scheduler – Extract MAX Priority Task in $O(1)$

- Reorganize the runqueue data structure

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



Timeslice available

[0]
[1]
[2]
...
[99]
[100]
...
[139]

Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

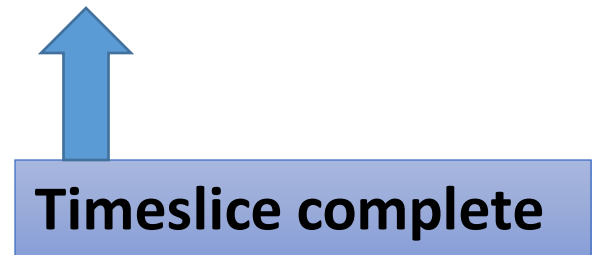
- Reorganize the runqueue data structure

Active Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

Expired Array


[0]
[1]
[2]
...
[99]
[100]
...
[139]



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1] 
[2]
...
[99]
[100]
...
[139]

Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1] 1
[2]
...
[99] 2
[100]
...
[139]

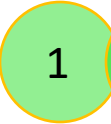
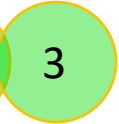

Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1]  
[2]
...
[99] 
[100]
...
[139]

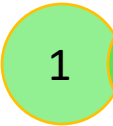
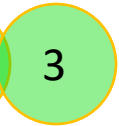


Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1]  
[2]
...
[99] 
[100]
...
[139] 

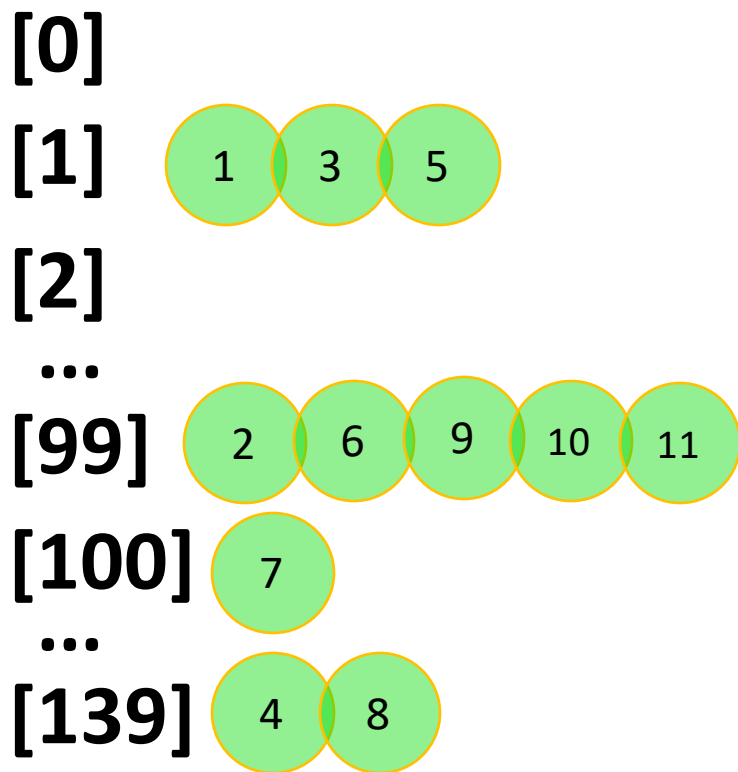
Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



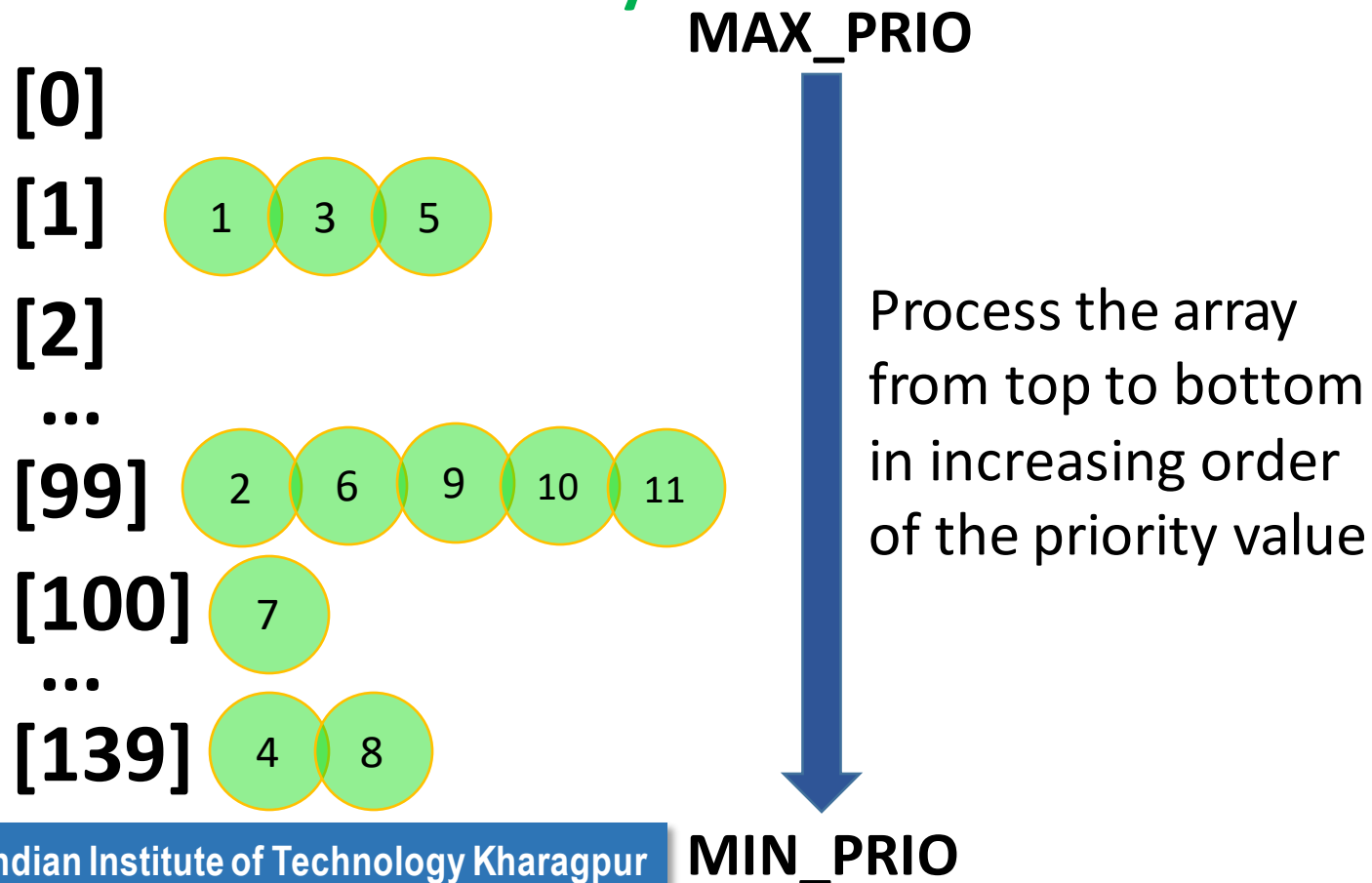
Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



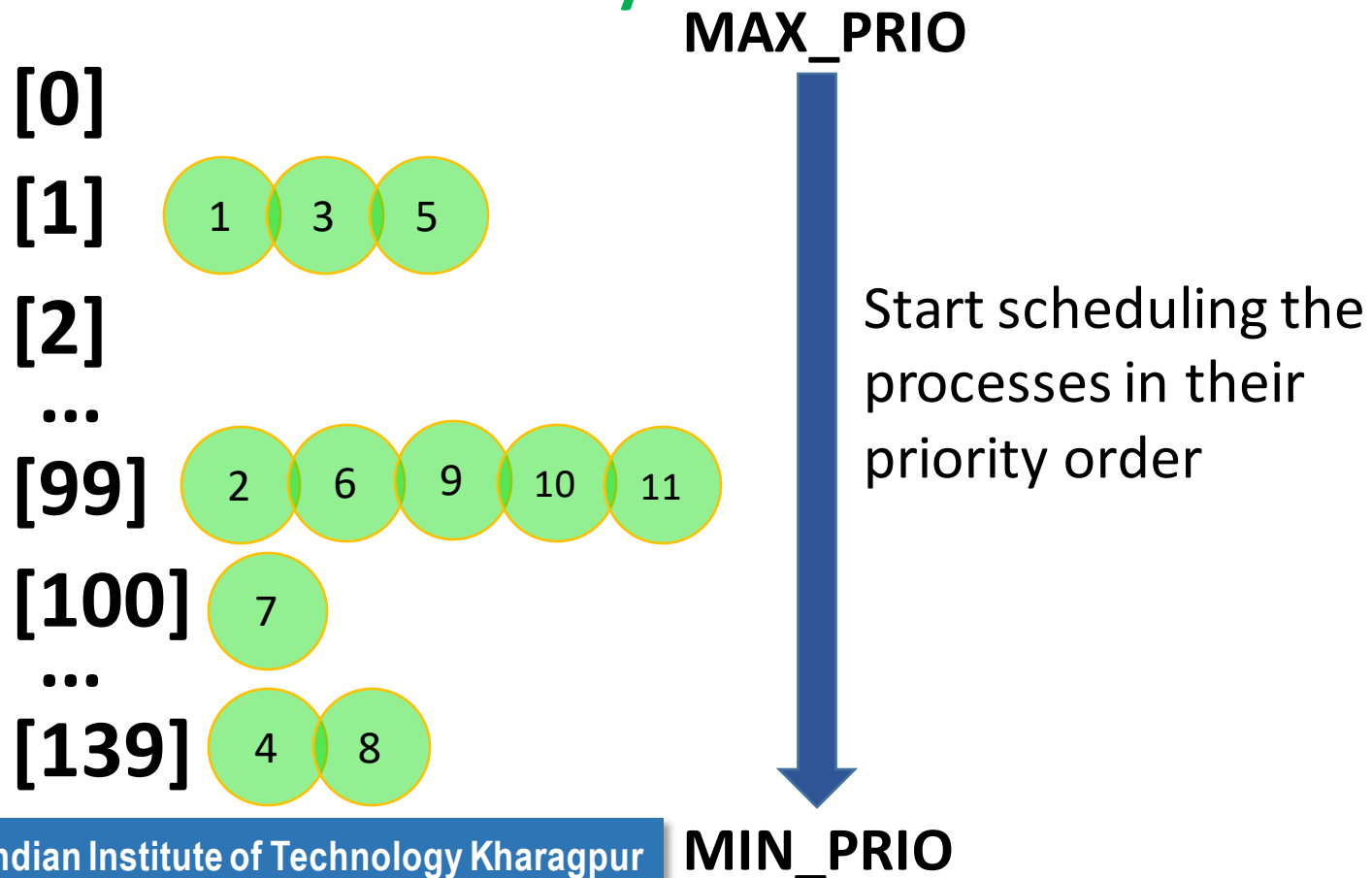
Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



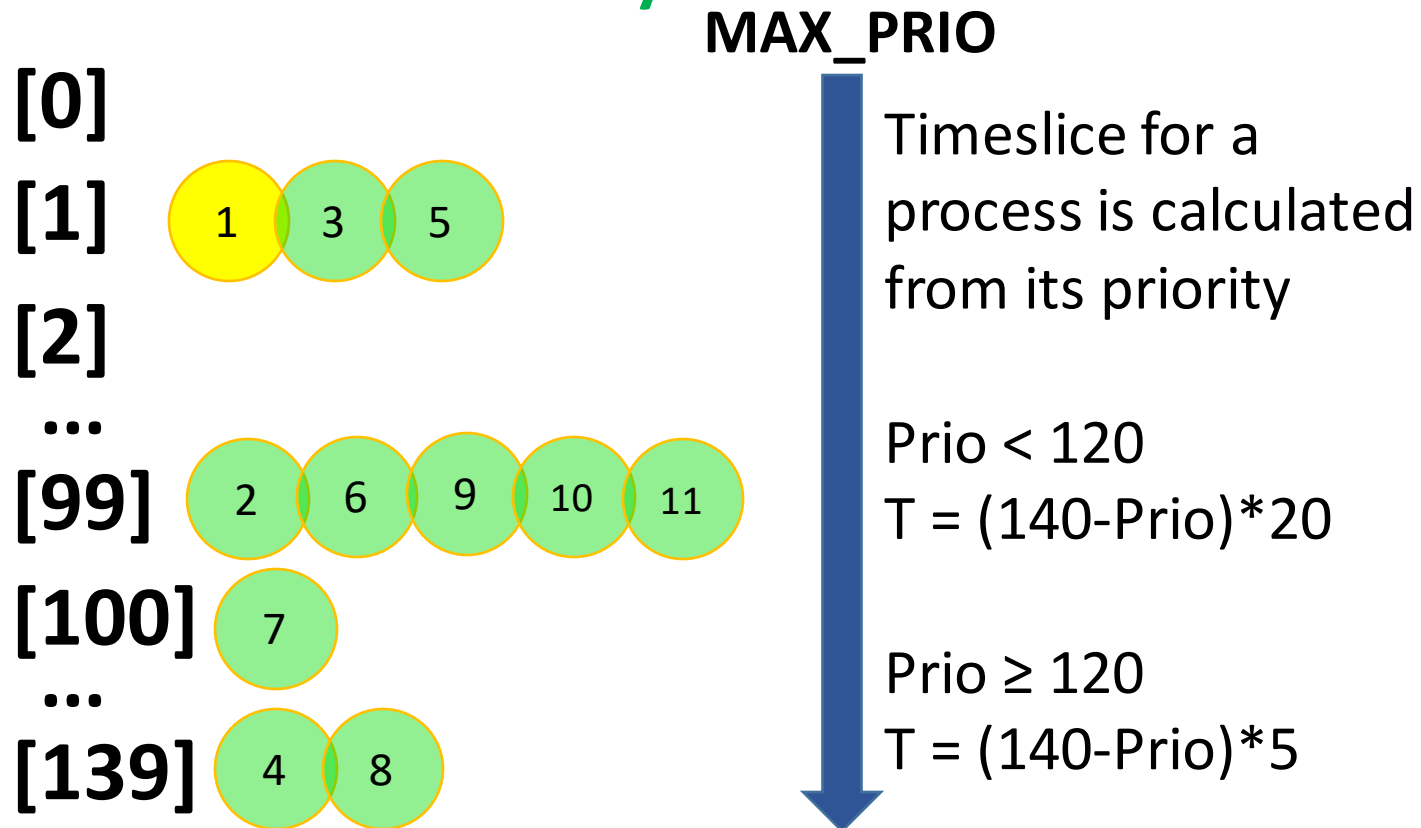
Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



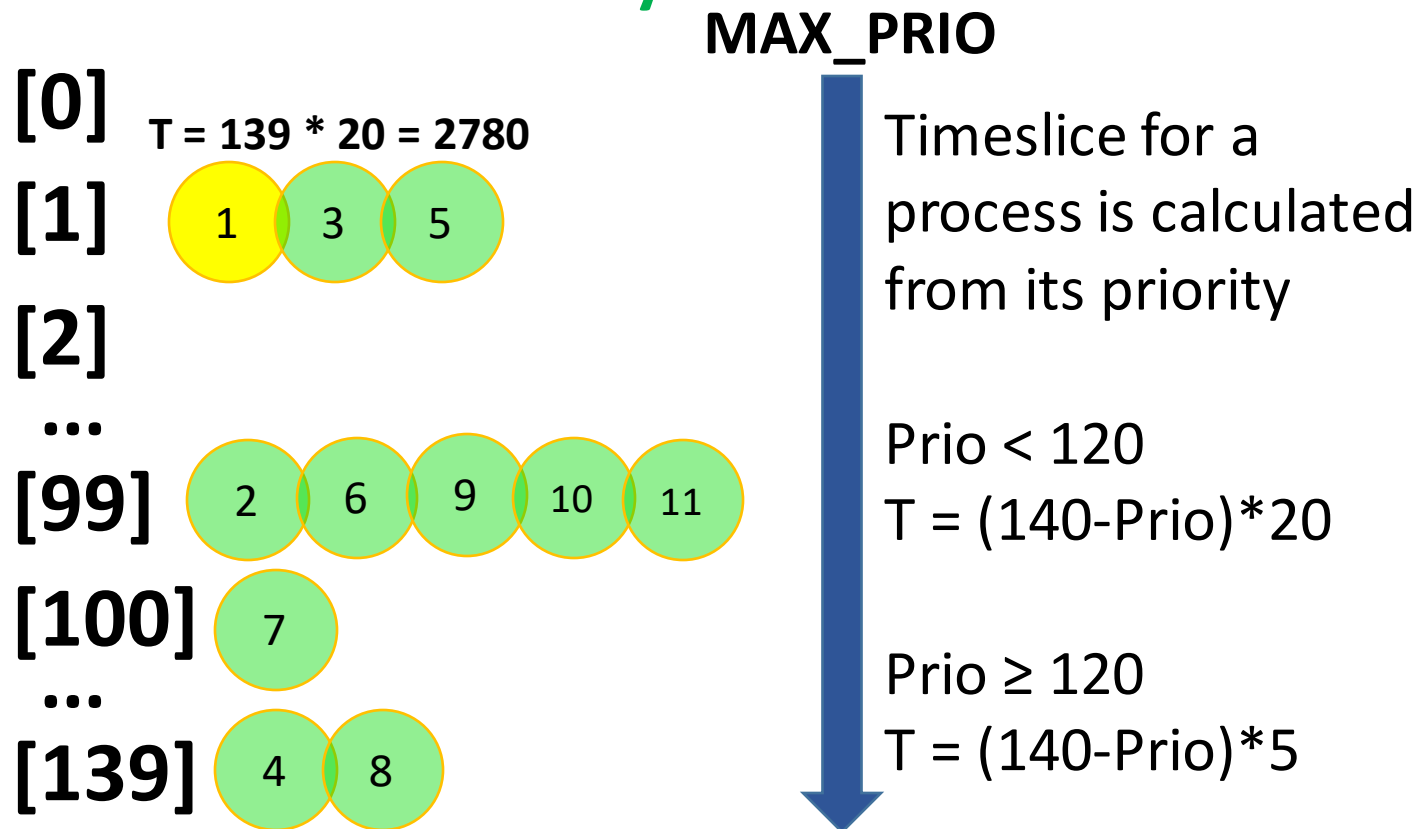
Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



Expired Array

[0]

[1]

[2]

...

[99]

[100]

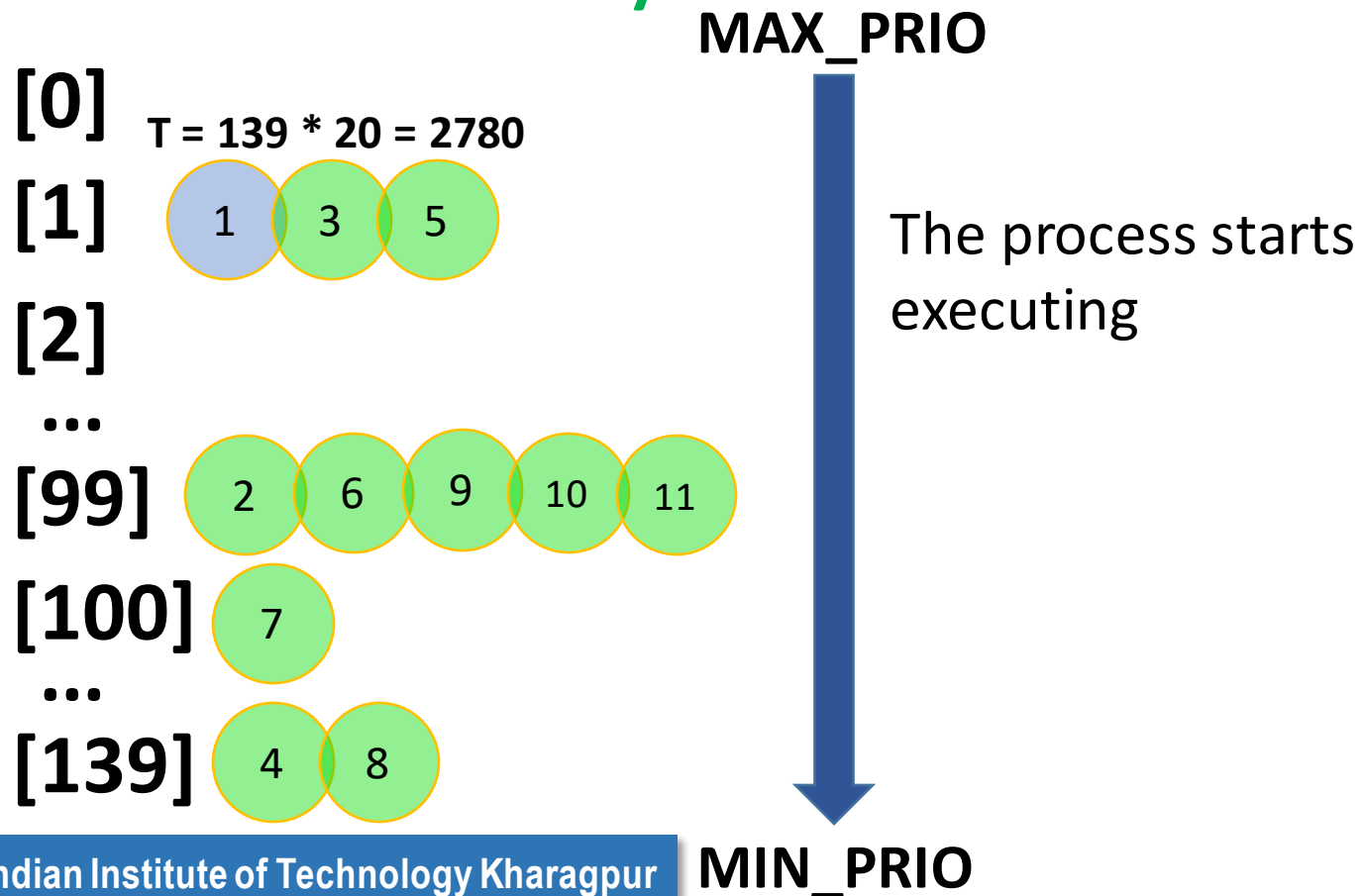
...

[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



Expired Array

[0]

[1]

[2]

...

[99]

[100]

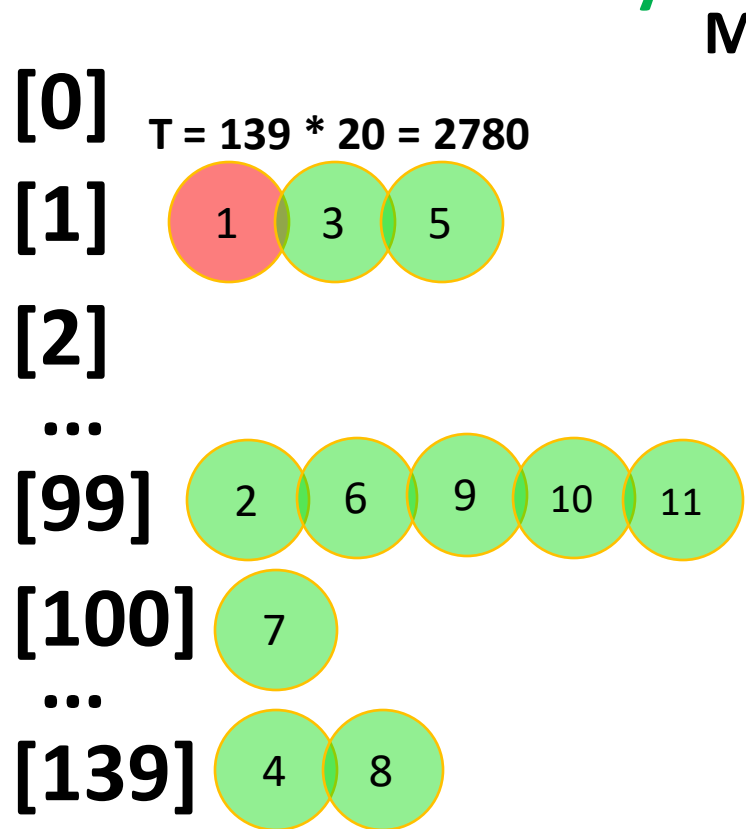
...

[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

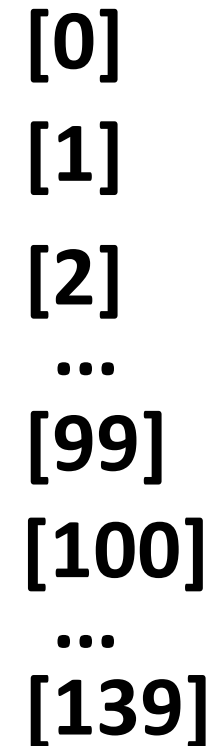


MAX_PRIO

The process starts
executing

Timer Interrupts at
2780 clock tick

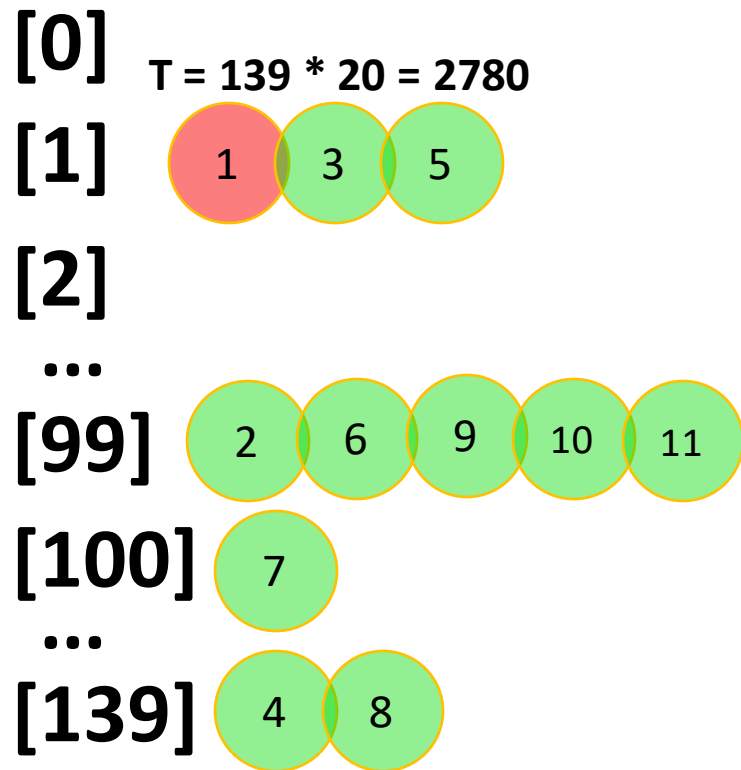
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

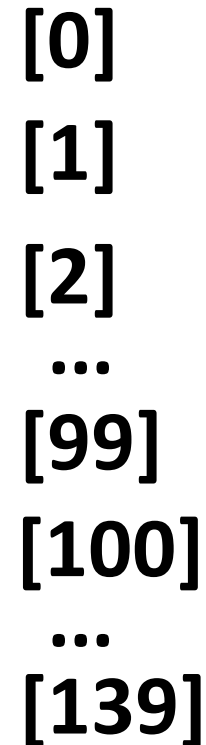


MAX_PRIO

The priority of the process is recalculated

- Niceness
- Interactivity
- Ageing

Expired Array

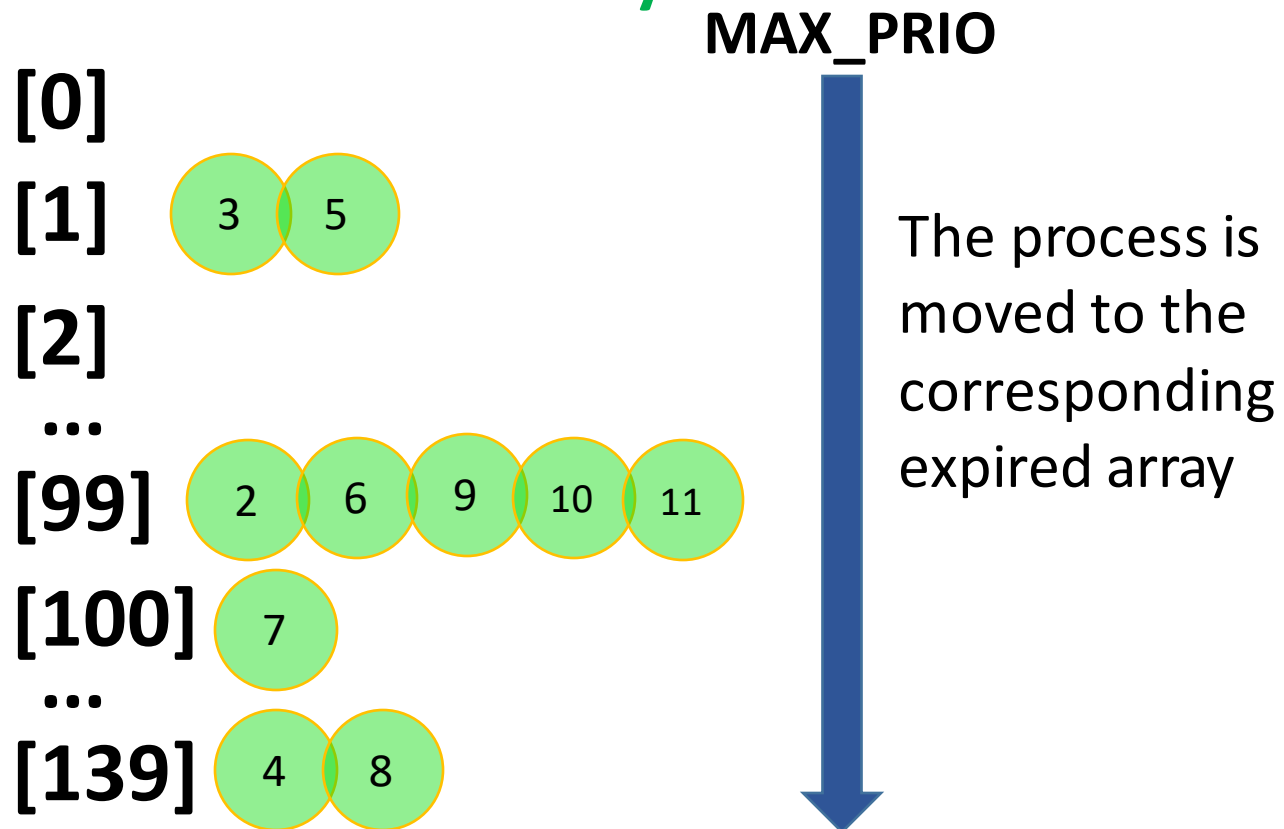


MIN_PRIO

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



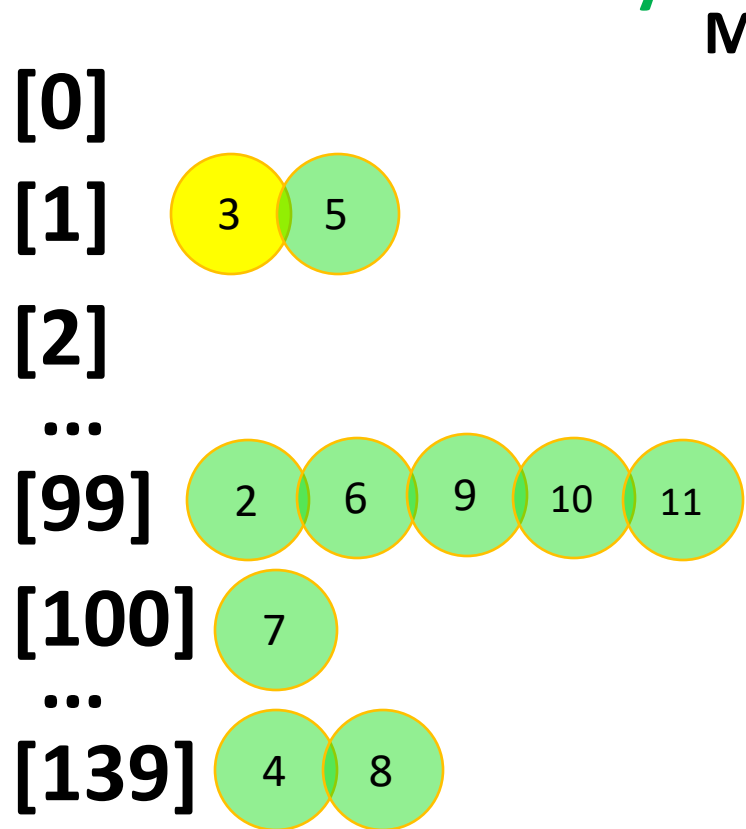
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



MAX_Prio

Context switch to
the next process for
the same priority
runqueue

MIN_Prio

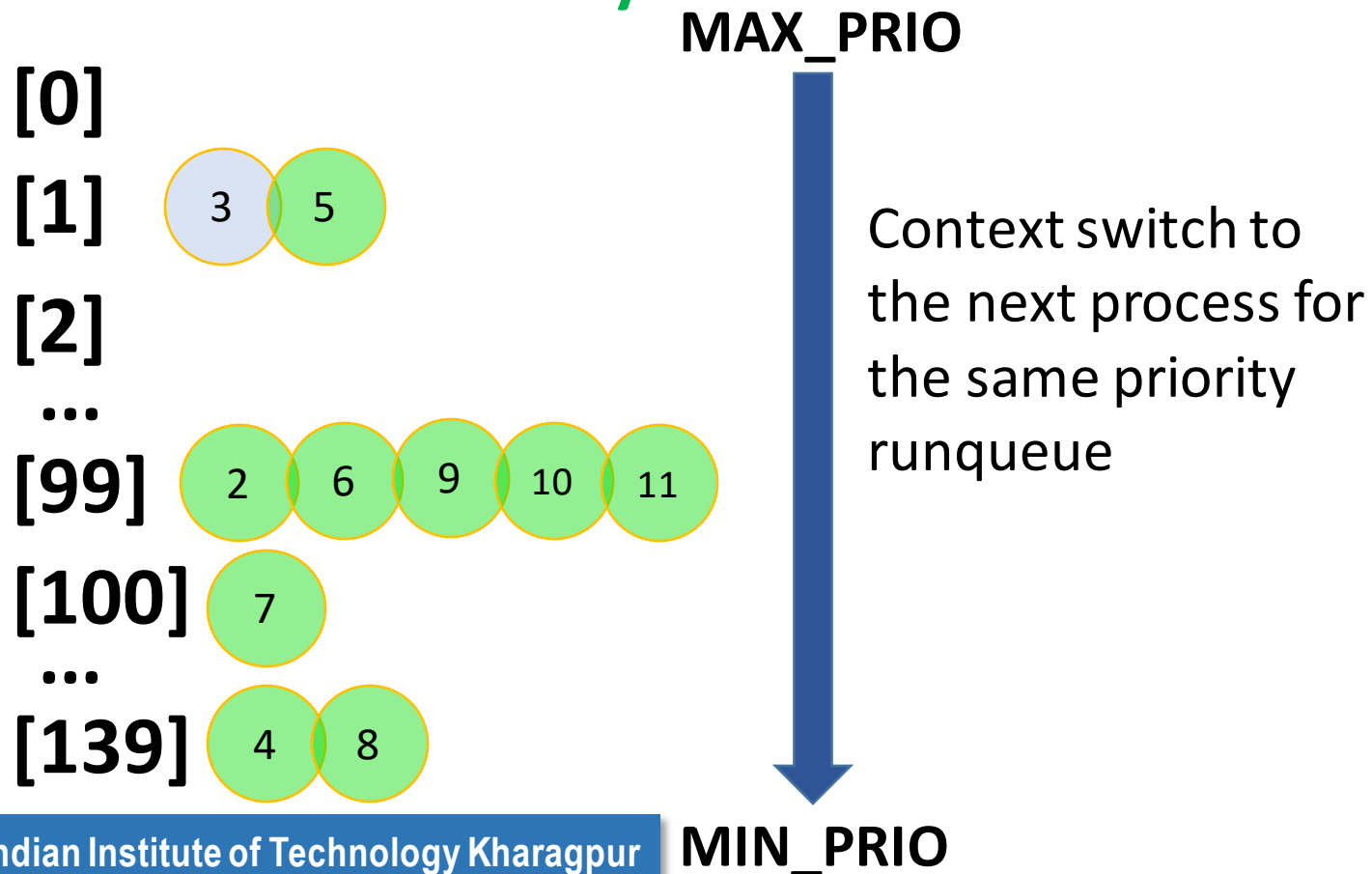
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



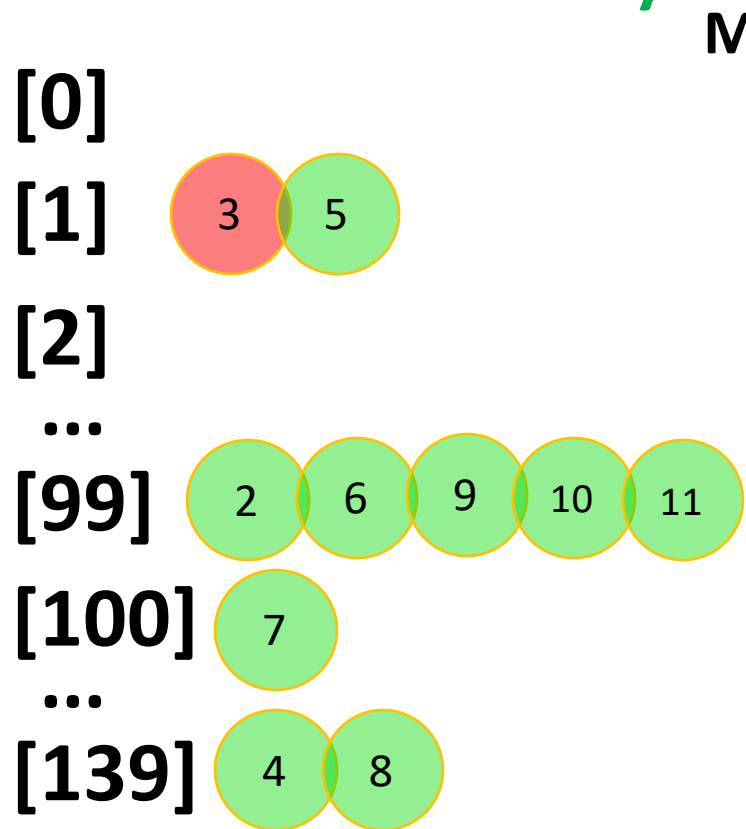
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



MAX_Prio

Context switch to
the next process for
the same priority
runqueue

MIN_Prio

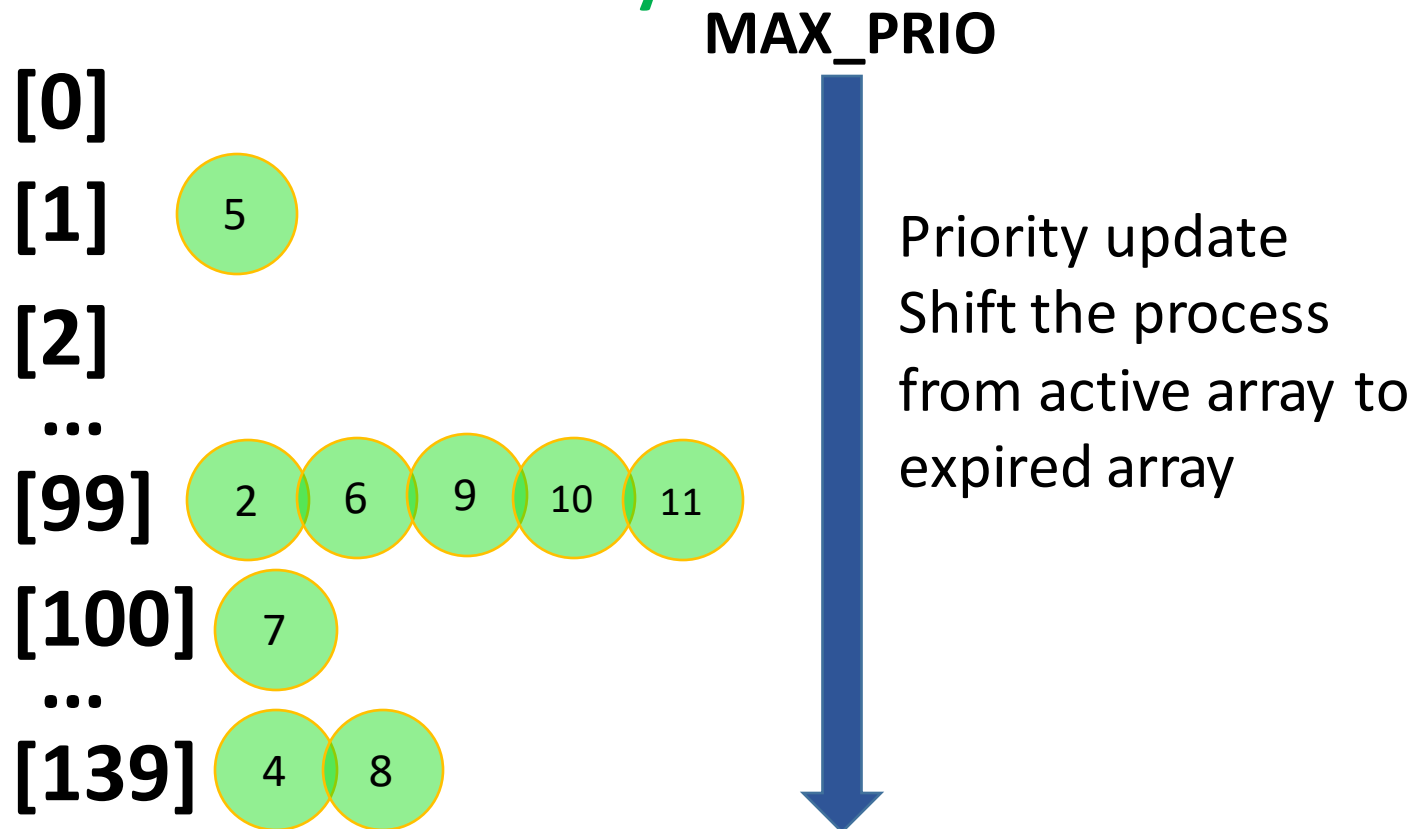
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



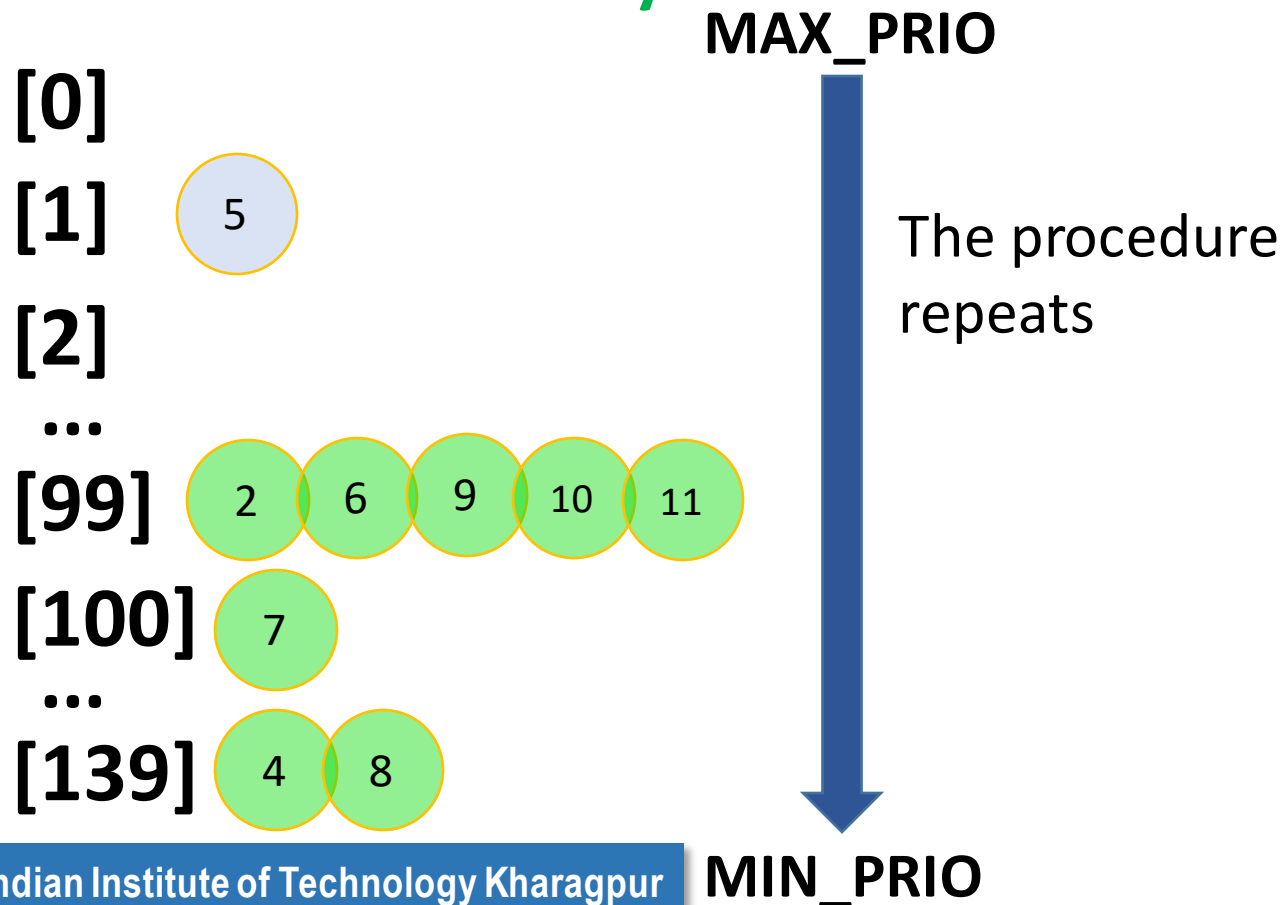
Expired Array



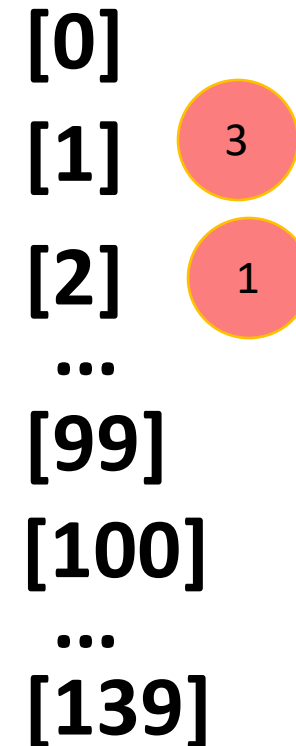
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



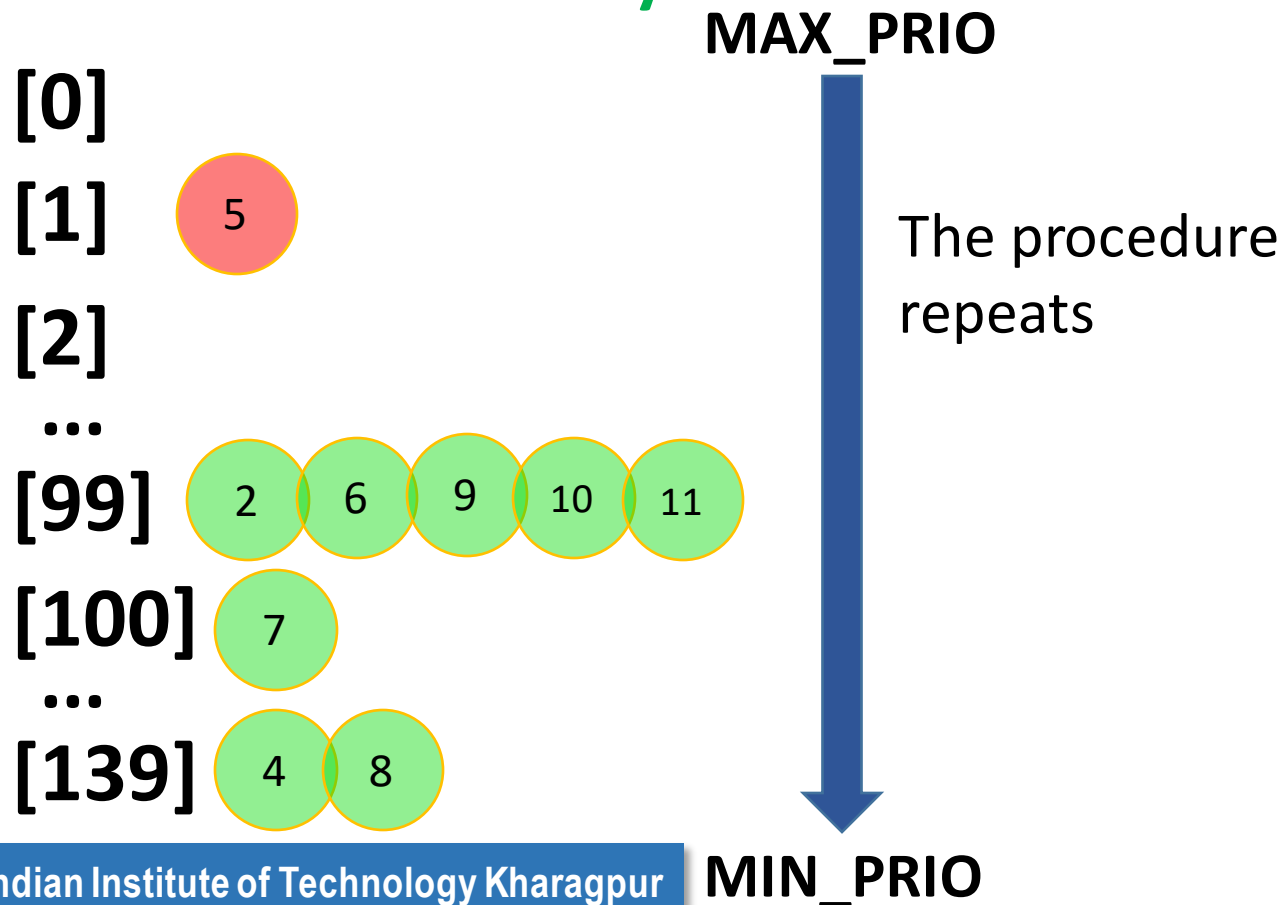
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



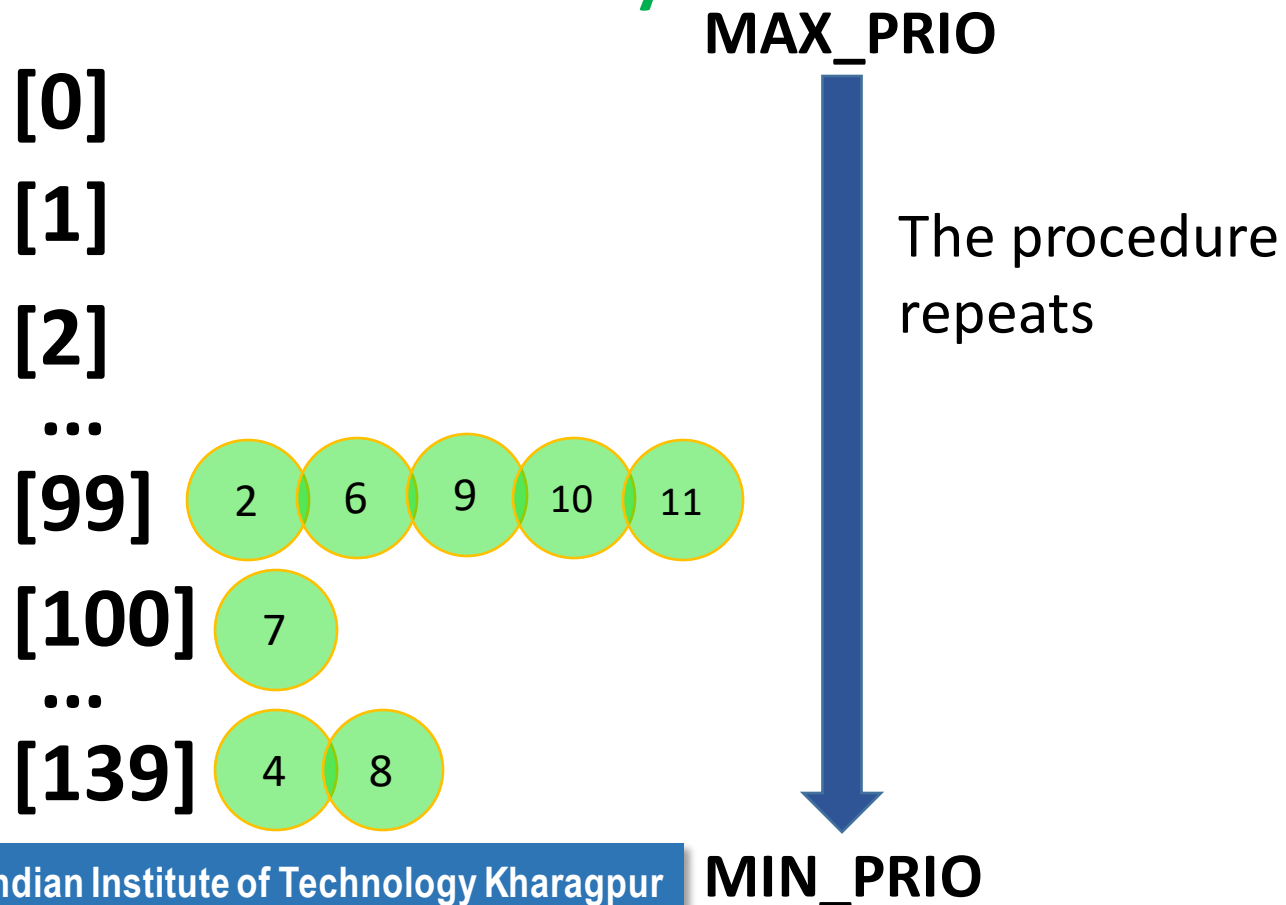
Expired Array



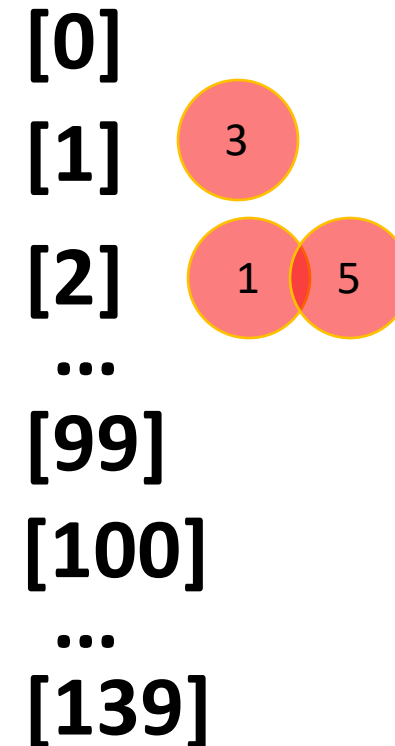
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



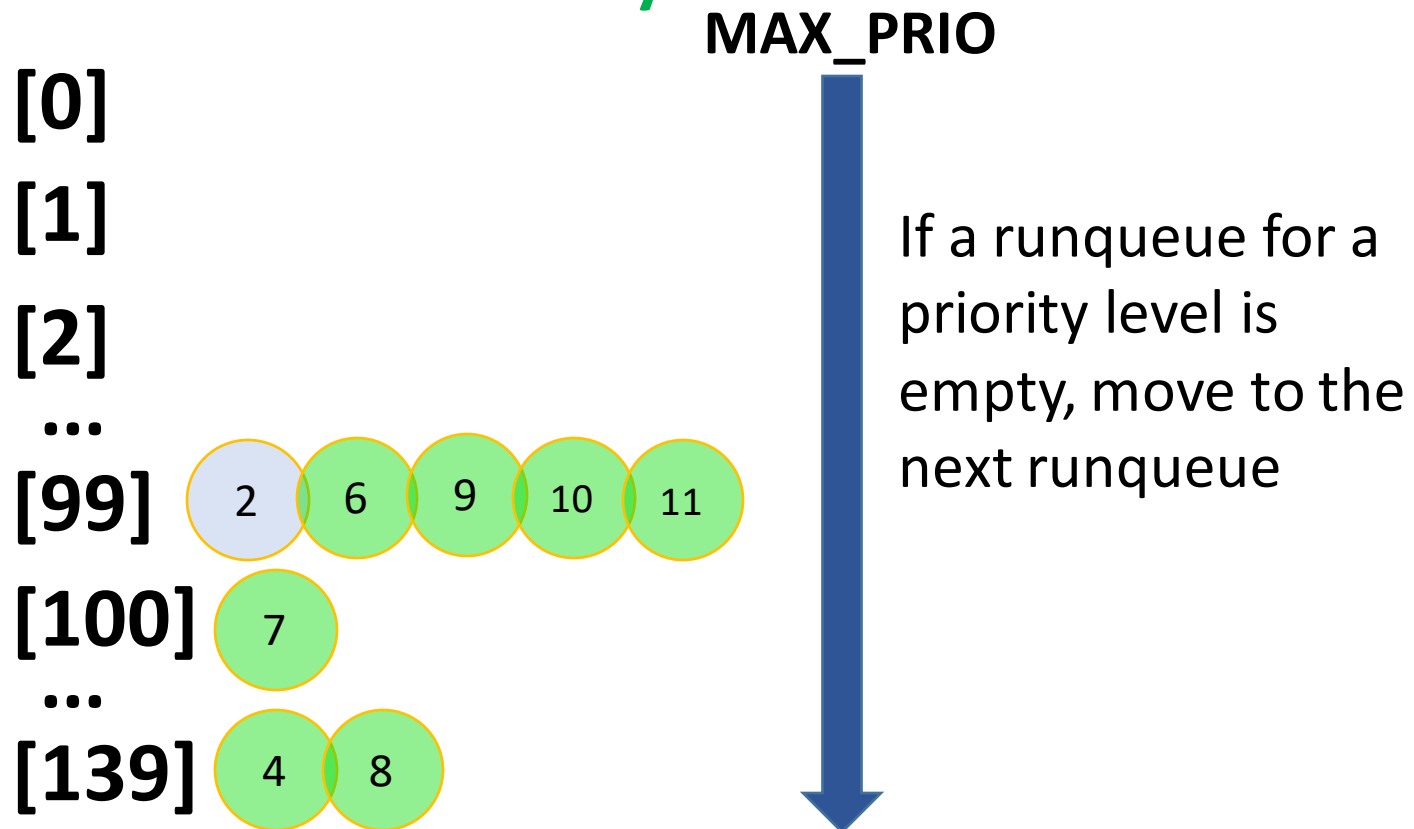
Expired Array



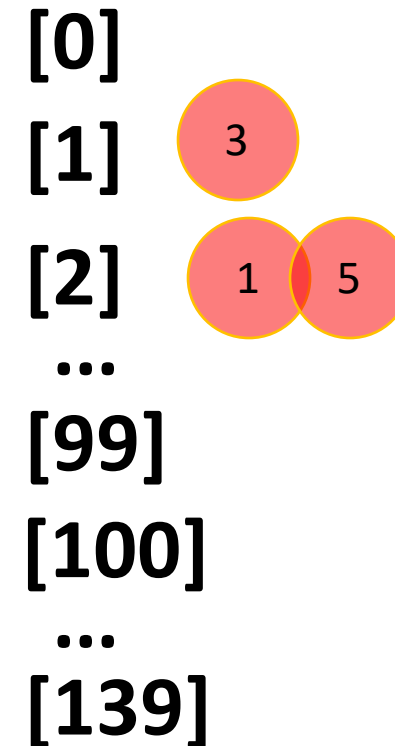
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



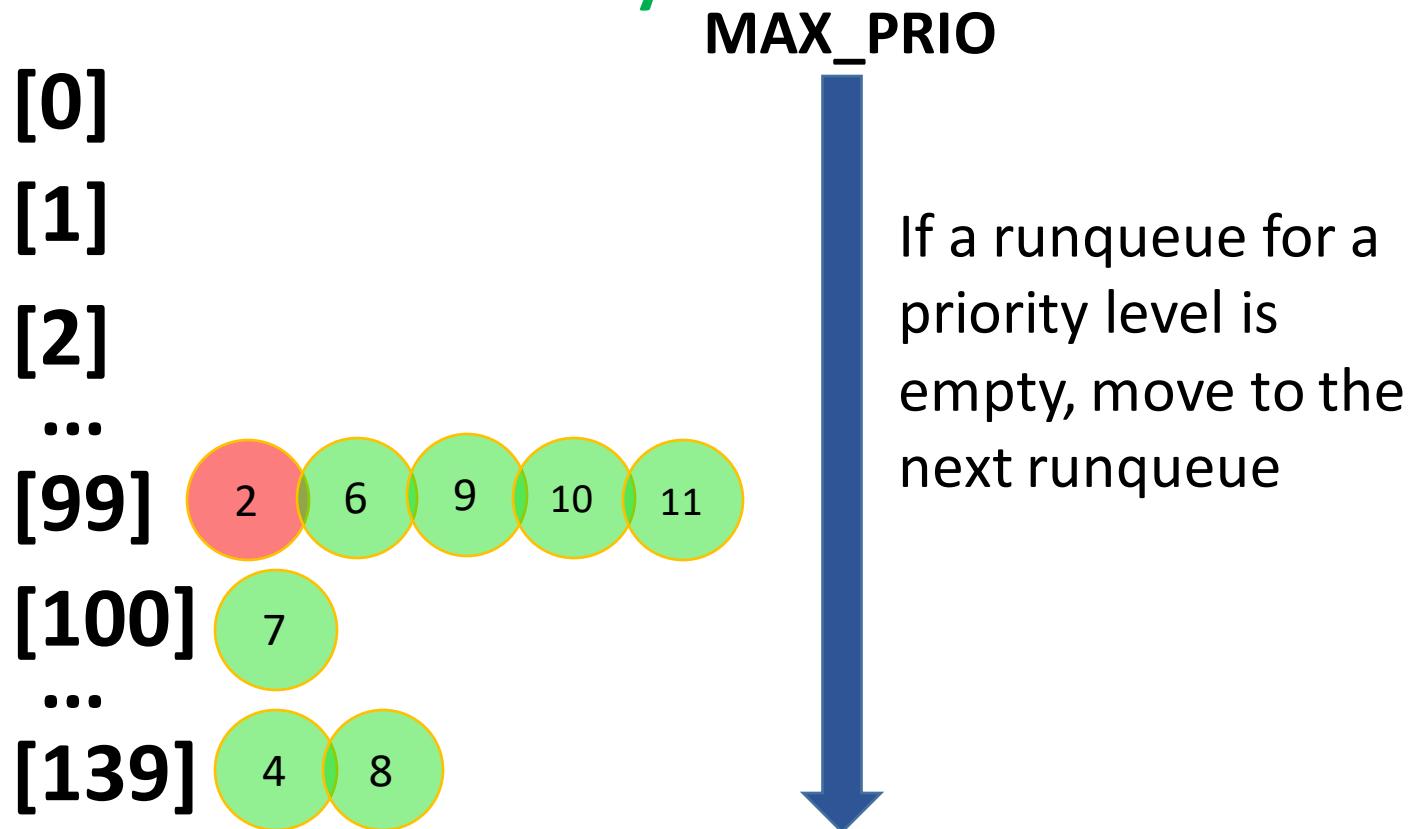
Expired Array



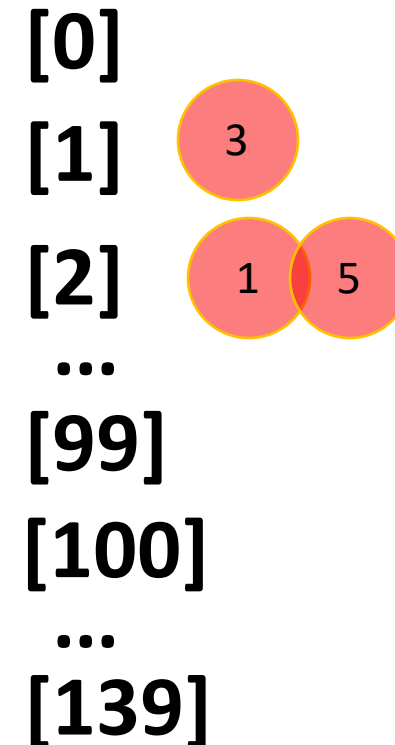
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



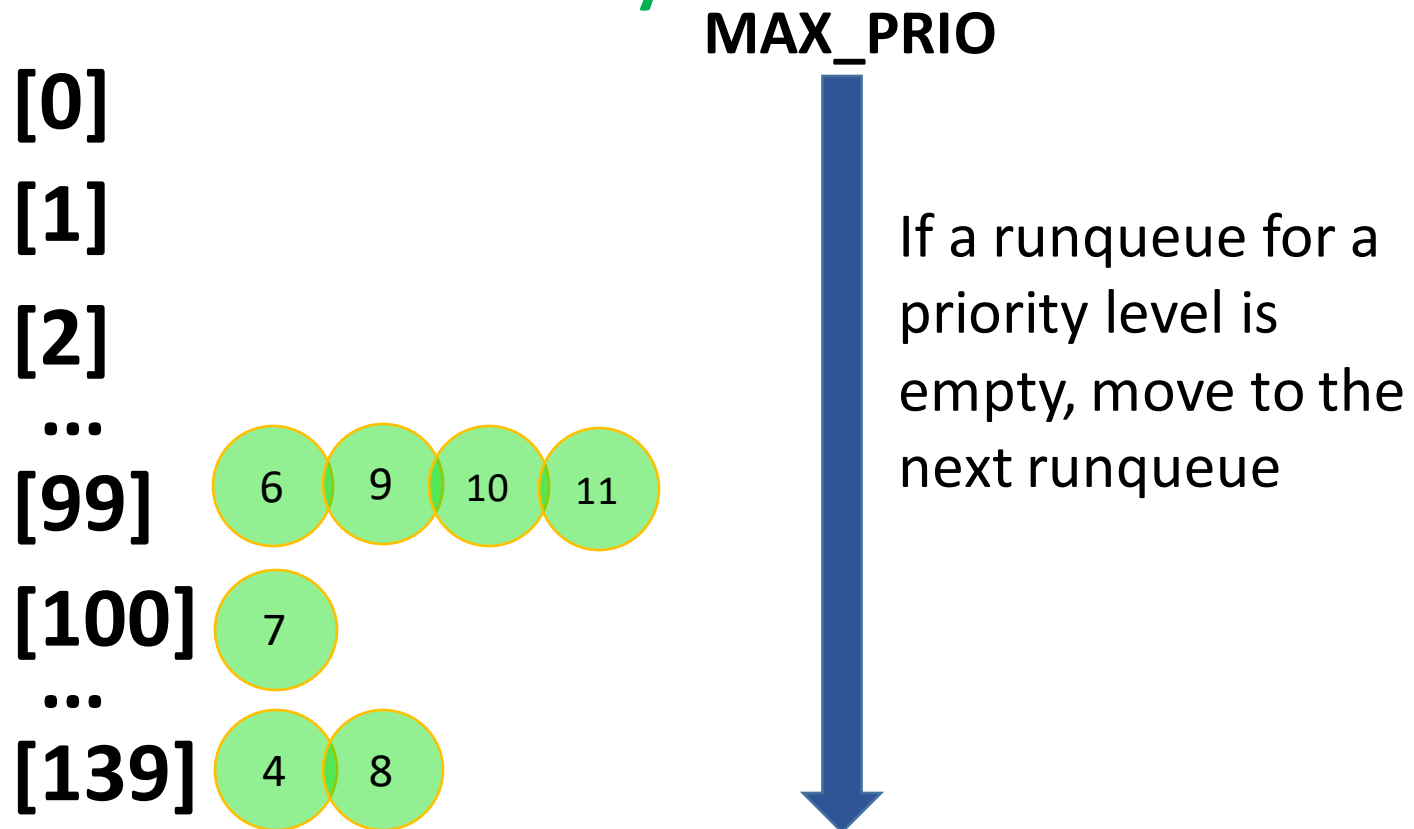
Expired Array



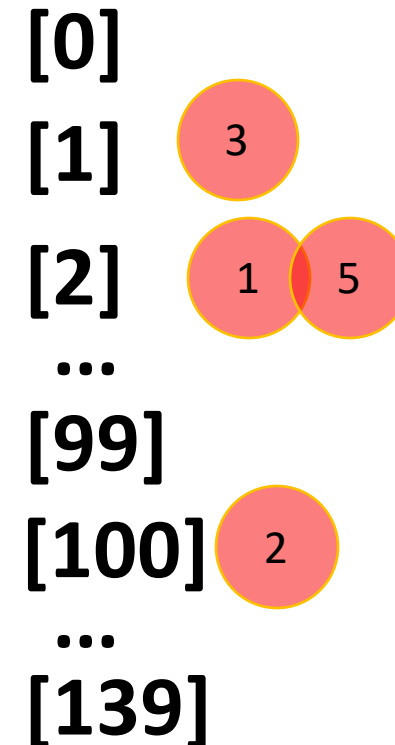
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



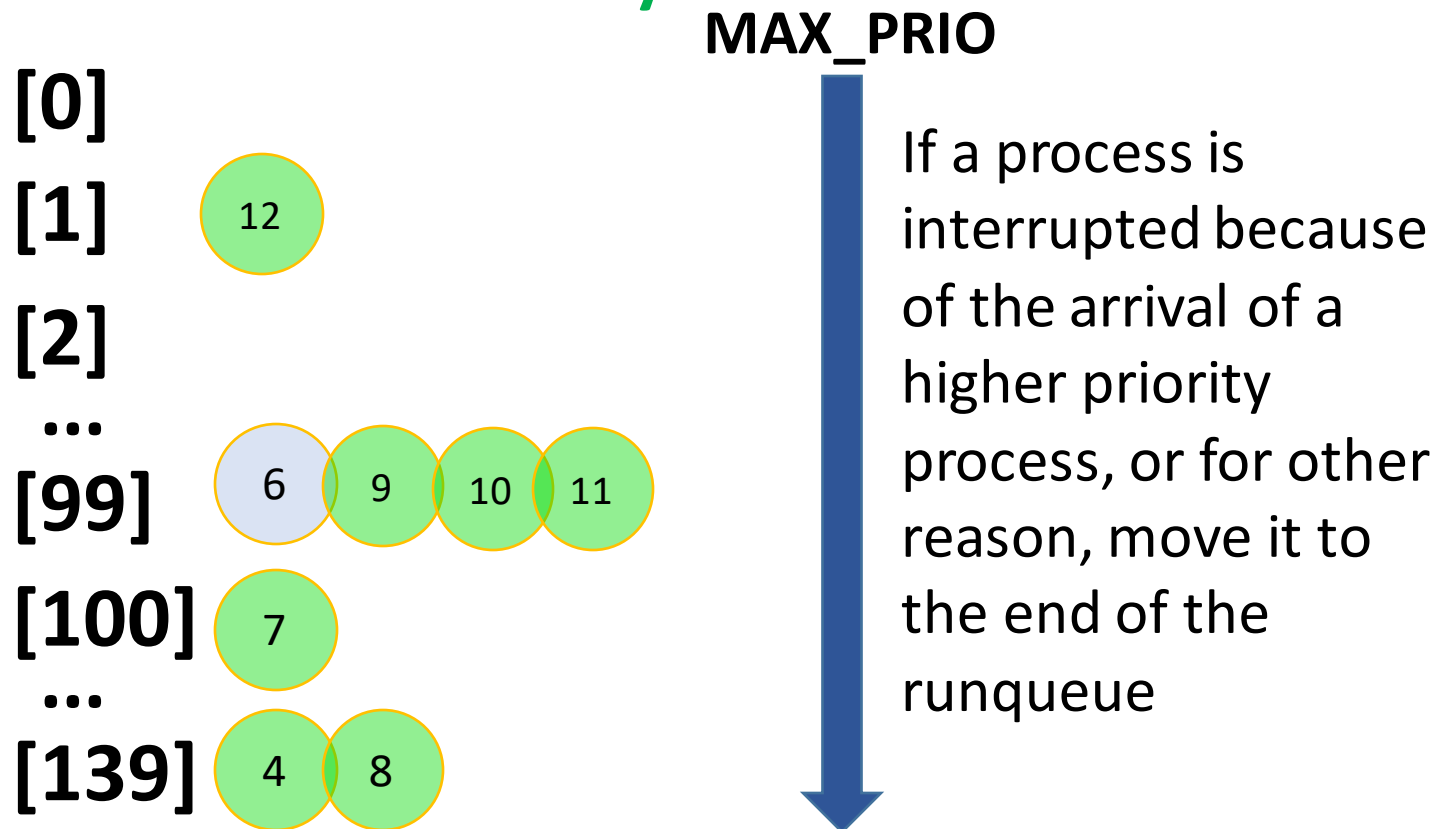
Expired Array



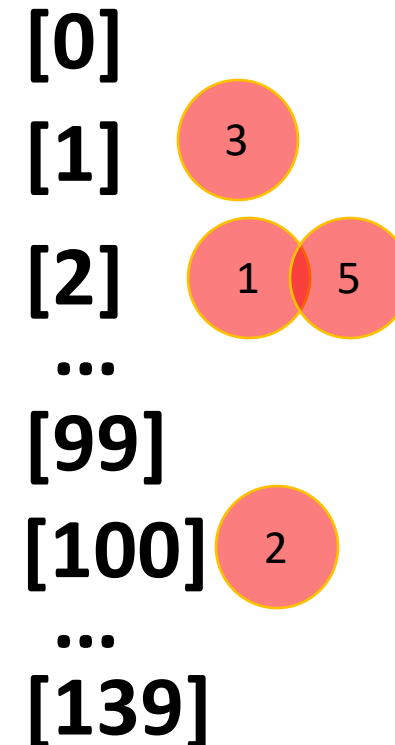
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



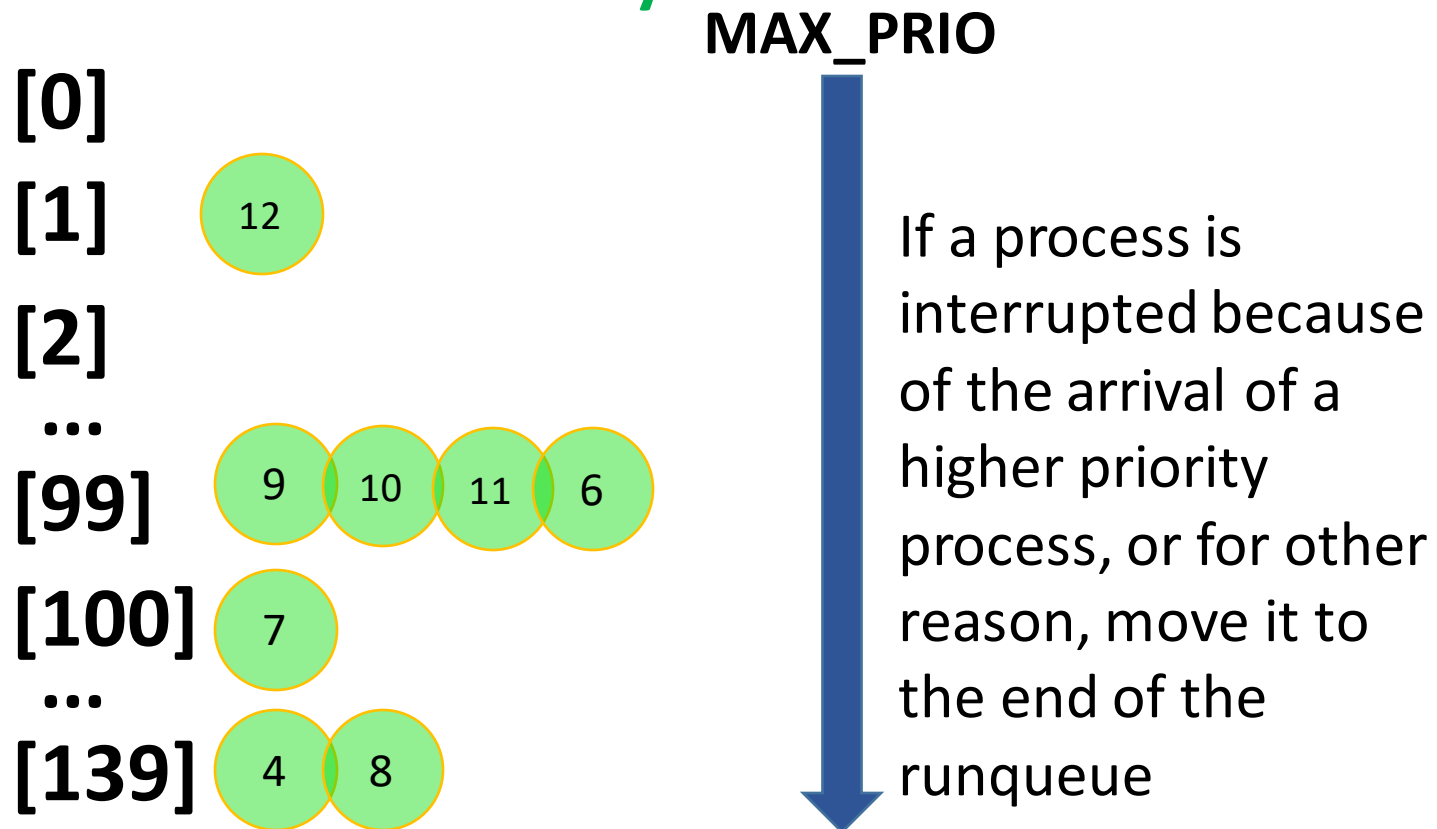
Expired Array



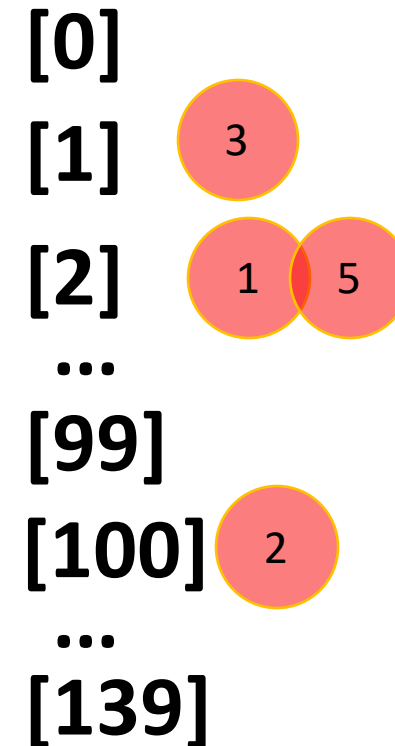
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



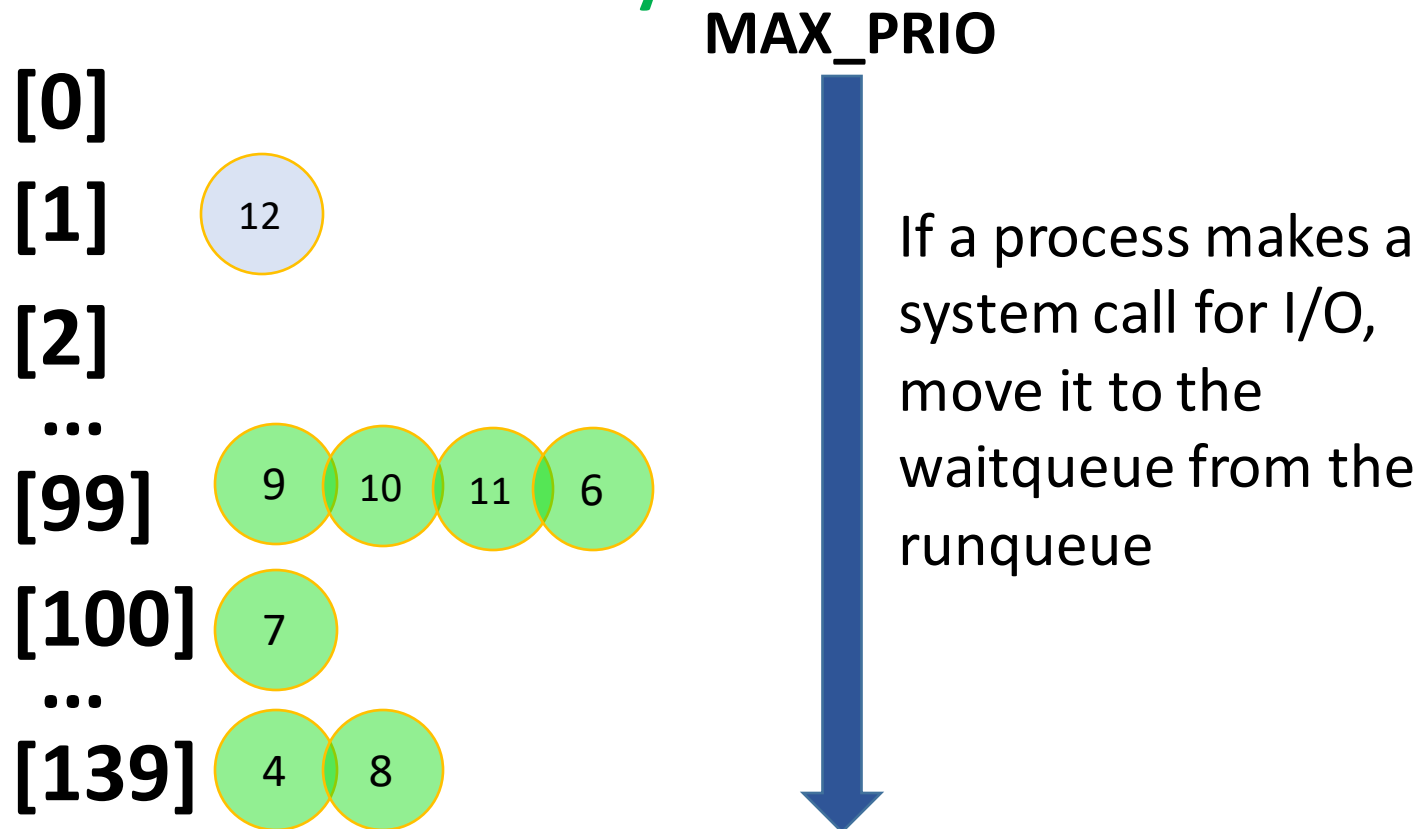
Expired Array



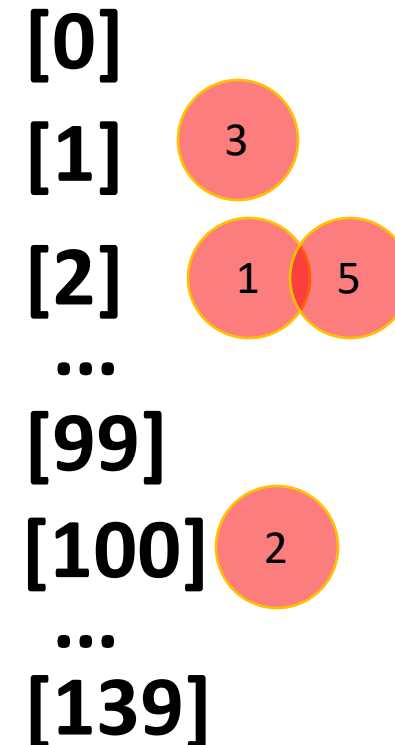
O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

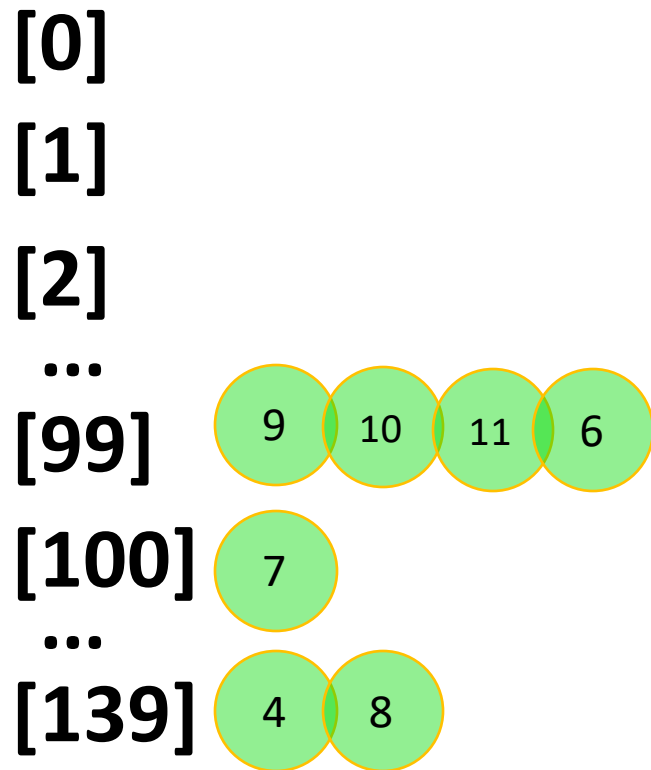
- Reorganize the runqueue data structure

Wait Queue

12

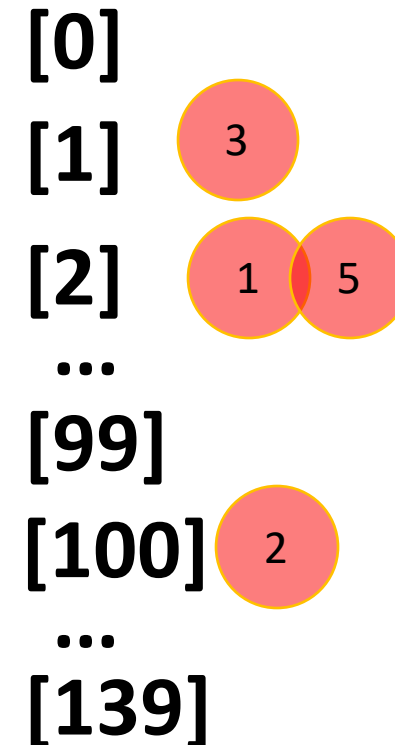
Active Array

MAX_Prio



If a process makes a system call for I/O, move it to the waitqueue from the runqueue

Expired Array

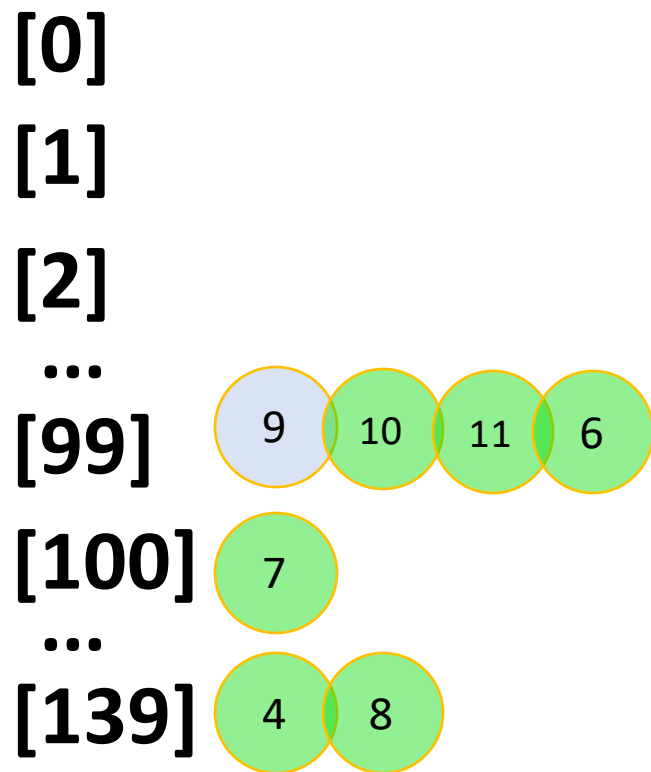


MIN_Prio

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



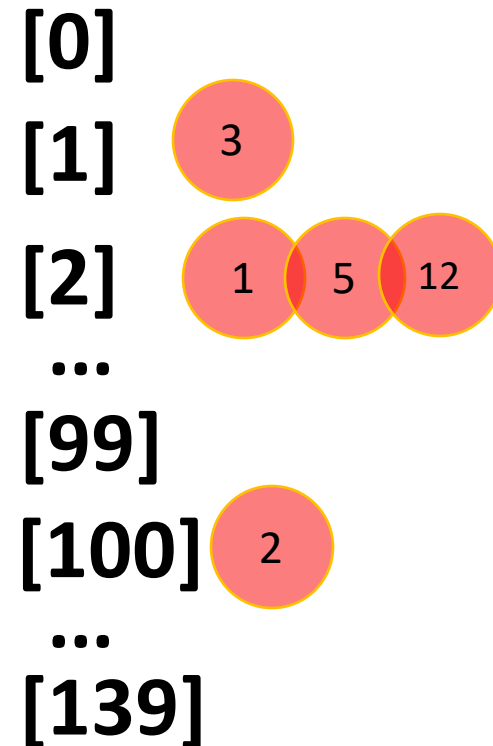
MAX_PRIO

Once the I/O is complete, move the task to the expired array (or active array when the task needs immediate scheduling)



MIN_PRIO

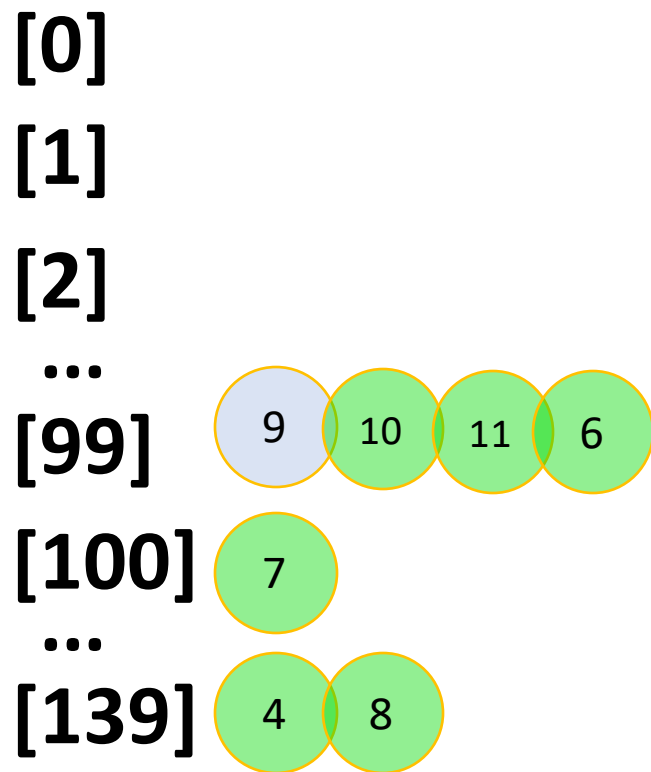
Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



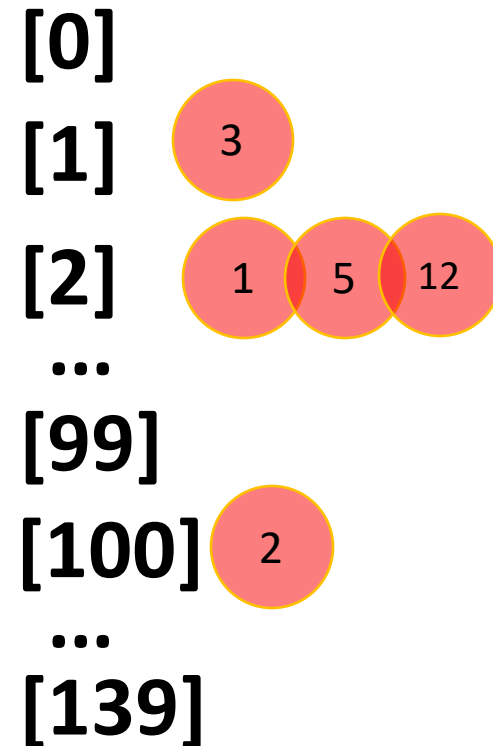
MAX_PRIO

Continue the
execution of all the
processes from the
active array



MIN_PRIO

Expired Array



O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

MAX_PRIO

Continue the
execution of all the
processes from the
active array



MIN_PRIO

Expired Array

[0]
[1] 3
[2] 1 5 12
...
[99] 7 10 6
[100] 2 9 4 11
...
[139] 8

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

MAX_PRIO

Make the Expired
Array as the
Active array



MIN_PRIO

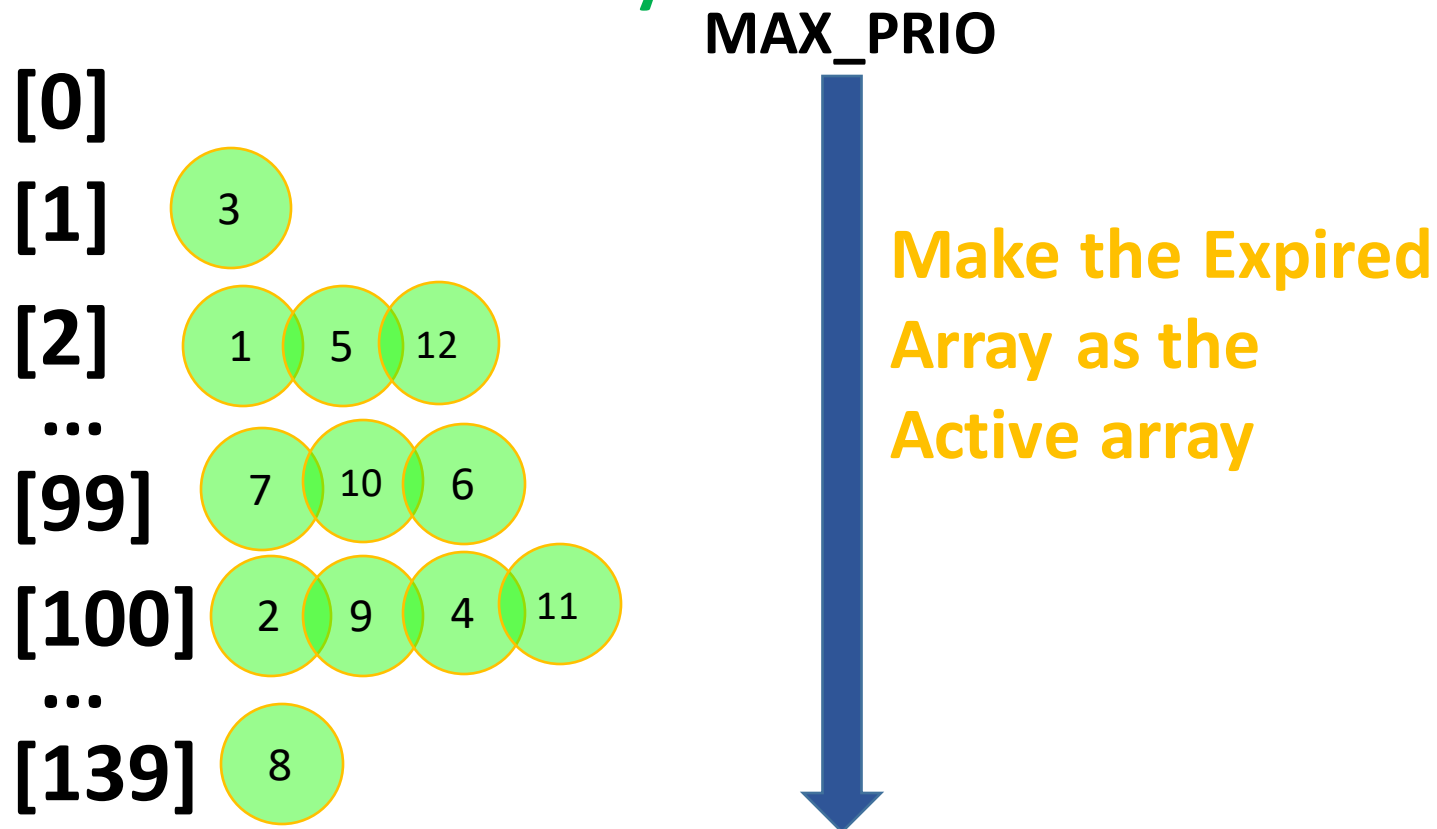
Expired Array

[0]
[1] 3
[2] 1 5 12
...
[99] 7 10 6
[100] 2 9 4 11
...
[139] 8

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array



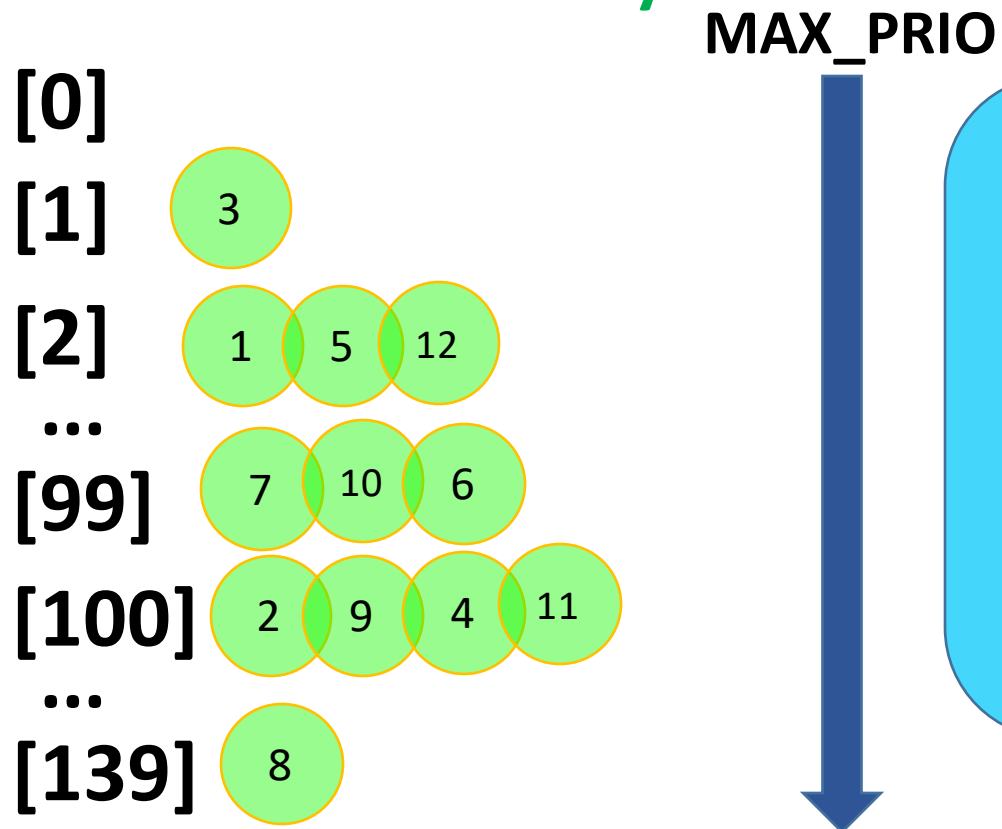
Expired Array

[0]
[1]
[2]
...
[99]
[100]
...
[139]

O(1) Scheduler – Extract MAX Priority Task in O(1)

- Reorganize the runqueue data structure

Active Array

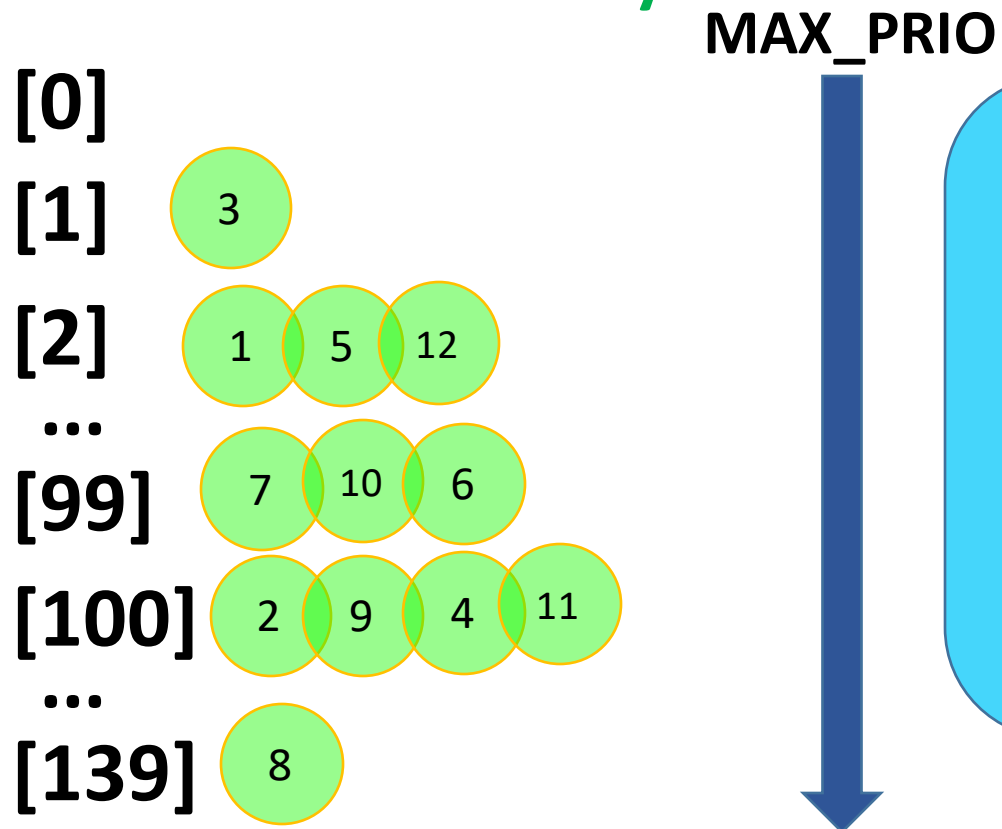


You do not need to iterate over all the processes to update their priority, at the end of the current epoch

O(1) Scheduler – Extract MAX Priority Task in O(1)

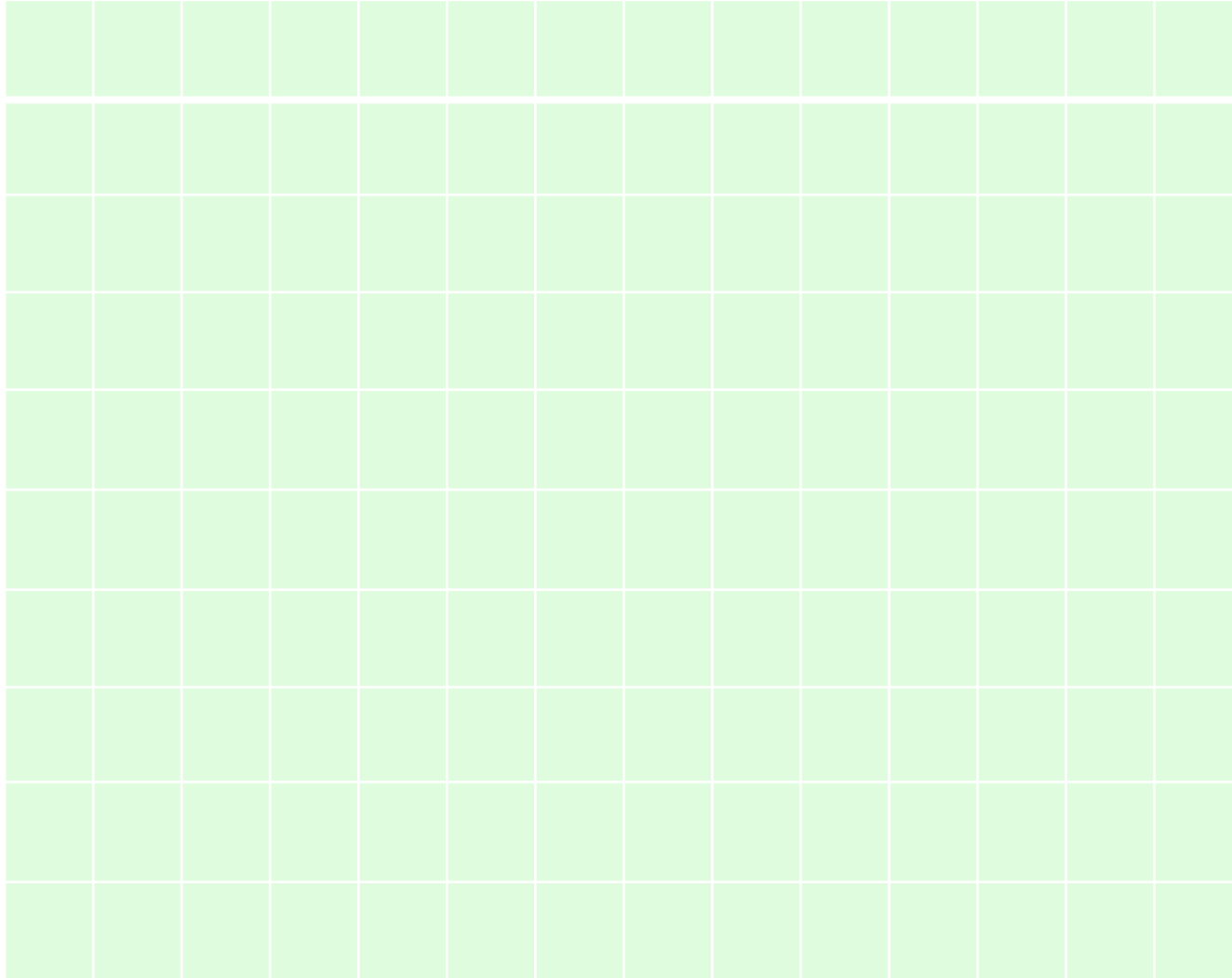
- Reorganize the runqueue data structure

Active Array



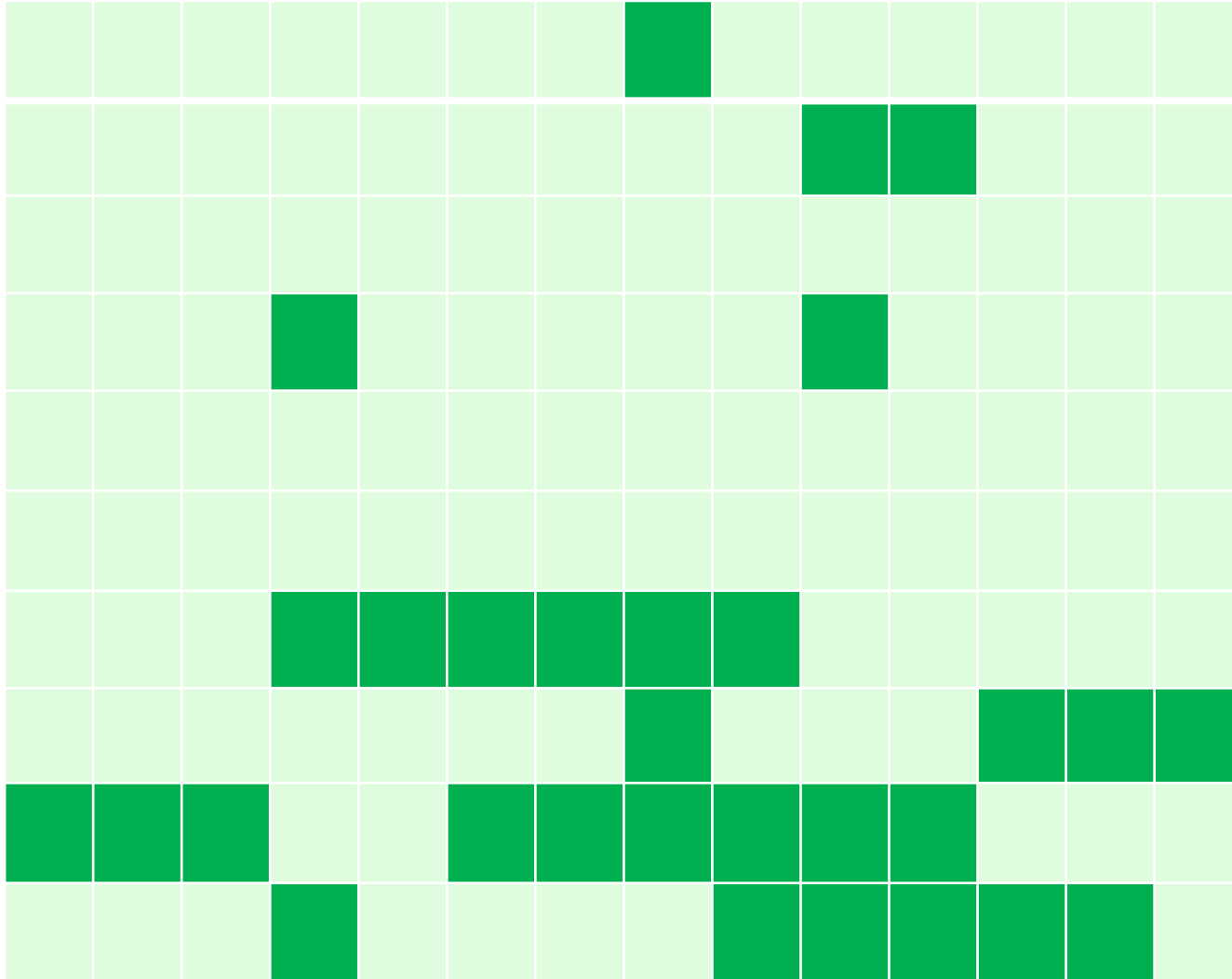
Only thing that remains:
How do we check that a higher priority process has arrived in the Active Array?

Priority Bitmap



```
struct prio_array {  
  
    /* number of tasks */  
    int nr_active;  
  
    /* priority bitmap */  
    unsigned long bitmap[BITMAP_SIZE];  
  
    /* priority queues */  
    struct list_head queue[MAX_PRIO];  
  
};
```

Priority Bitmap



```
struct prio_array {  
  
    /* number of tasks */  
    int nr_active;  
  
    /* priority bitmap */  
    unsigned long bitmap[BITMAP_SIZE];  
  
    /* priority queues */  
    struct list_head queue[MAX_PRIO];  
  
};
```

Update Priority based on Interactivity

- Dynamically increase the priority level of the interactive processes (processes which does a lot of I/Os) for better user experience

Update Priority based on Interactivity

- Dynamically increase the priority level of the interactive processes (processes which does a lot of I/Os) for better user experience
- **Sleep Ratio:**
 - Number of clock ticks spent while getting executed in CPU / Number of clock ticks while blocked in the wait queue

Update Priority based on Interactivity

- Dynamically increase the priority level of the interactive processes for better user experience
- **Sleep Ratio:**
 - Number of clock ticks spent while getting executed in CPU / Number of clock ticks while blocked in the wait queue
 - **Mostly Sleeping:** I/O Bound (Interactive)
 - **Mostly Running:** CPU Bound (Batch)

O(1) Scheduler

- **Pros**

- Better scalability
- Faster performance
- Reduced context switching time
- Incorporated interactivity metric

O(1) Scheduler

- **Pros**

- Better scalability
- Faster performance
- Reduced context switching time
- Incorporated interactivity metric

- **Cons**

- Complex heuristic to mark a process as interactive or batch – didn't work well in practice
- The interactivity may change over time – initially CPU bound, then I/O bound
- 140*2 runqueues for each processor – complex logic is required for runqueue locking
- Complicated codebase – difficult for debugging purpose

