

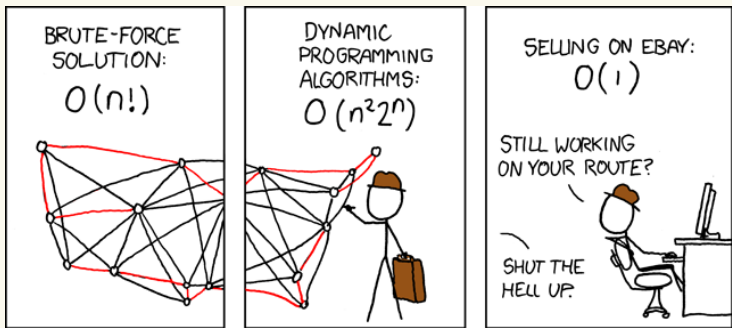
Linux Kernel Programming

Kernel Data Structures

Pierre Olivier

Systems Software Research Group @ Virginia Tech

February 7, 2017



► Source: <https://xkcd.com/399/>

Outline

- 1 Linked lists
- 2 Queues
- 3 Maps
- 4 Binary trees
- 5 The right data structure for the right problem

Introduction

- ▶ The kernel has efficient implementations of:
 - ① **Lists** (singly/doubly linked): `include/linux/list.h`
 - ② **Queues**: `include/linux/kfifo.h`
 - ③ **Maps**: `include/linux/idr.h`
 - ④ **Binary trees** (*red-black trees*): `include/linux/rbtree.h`
- ▶ Do not reinvent the wheel!

Outline

- 1 **Linked lists**
- 2 Queues
- 3 Maps
- 4 Binary trees
- 5 The right data structure for the right problem

Linked lists

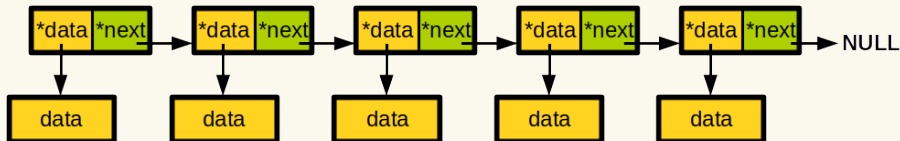
Singly linked list

```

1 struct my_list_element
2 {
3     void                *data;
4     struct my_list_element *next;
5 };

```

- ▶ void pointer to point on generic data
 - ▶ Can also contain data directly for a non-generic version
- ▶ Pointer to the next element

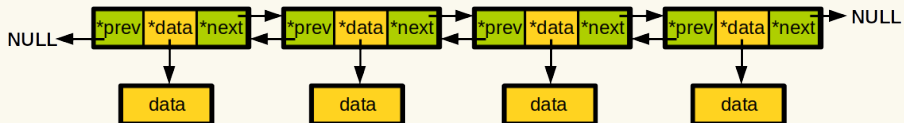


Linked lists

Doubly linked list

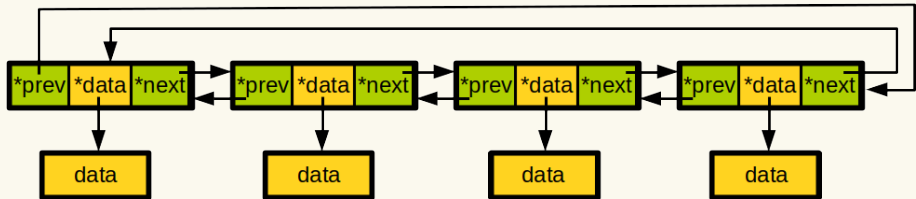
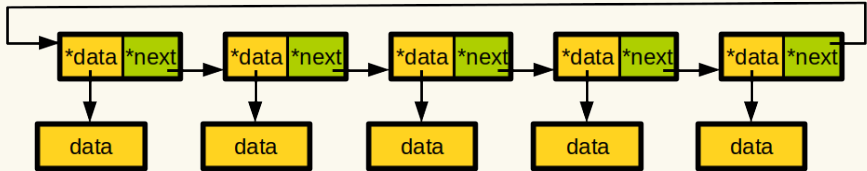
```
1 struct my_list_element
2 {
3     void *data;
4     struct my_list_element *prev;
5     struct my_list_element *next;
6 };
```

- ▶ Additional `prev` pointer
- ▶ Allows backward traversal



Linked lists

Circular lists



- ▶ Linked lists are iterated sequentially
 - ▶ Inappropriate when random (direct) access to a specific element is needed

Linked lists

Linux implementation: standard approach vs Linux

- ▶ Linux implements linked list a bit differently than the standard approach
- ▶ Standard approach → add `next/prev` pointers to a data structure:

```
1 struct car
2 {
3     unsigned int    max_speed;
4     unsigned int    drive_wheel_num;
5     double          price_in_dollars;
6     struct car      *prev;           /* we add this ... */
7     struct car      *next;         /* ... and this */
8 };
```

Linked lists

Linux implementation: `struct list_head`

```
1 struct list_head
2 {
3     struct list_head *next;
4     struct list_head *prev;
5 };
```

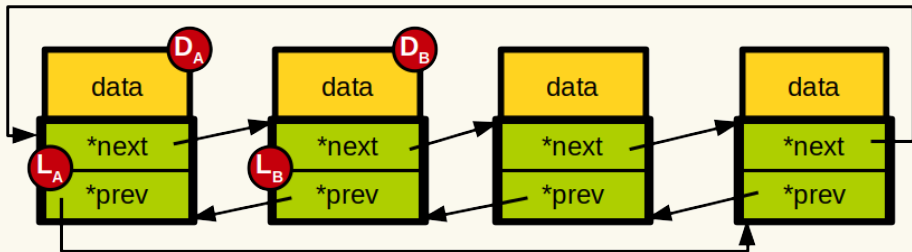
- ▶ Current implementation was introduced in Linux 2.1
- ▶ **struct list_head** as the central data structure

▶ `list_head` is **embedded** in the structure we want to link:

```
1 struct car
2 {
3     unsigned int    max_speed;
4     unsigned int    drive_wheel_num;
5     double          price_in_dollars;
6     struct list_head list;          /* we add this */
7 };
```

Linked lists

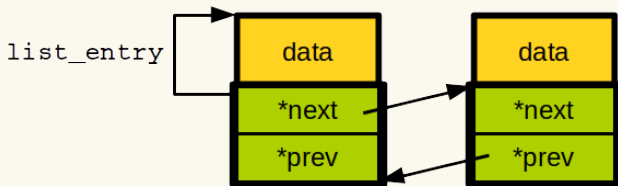
Linux implementation: `struct list_head` (2)



Linked lists

Linux implementation: `list_entry`

- ▶ The kernel provides a *generic* API to manipulate such lists
 - ▶ ex: `list_add(struct list_head *new, struct list_head *head)`
 - ▶ Manipulates `struct list_head` objects
- ▶ **How to get access to the containing data structure given a `struct list_head`?**
 - ▶ Use `list_entry`



Linked lists

Linux implementation: `list_entry` (2)

```
1 /* let's assume we have a pointer car_list_ptr to a struct list_head embedded into a struct
2  * car data object */
3 struct car *amazing_car = list_entry(car_list_ptr, struct car, list)
```

```
1 #define list_entry(ptr, type, member) \
2     container_of(ptr, type, member)
```

```
1 #define container_of(ptr, type, member) ({ \
2     const typeof( ((type *)0)->member ) *__mptr = (ptr); \
3     (type *) ( (char *)__mptr - offsetof(type,member) );})
```

Linked lists

The container_of macro

```

1 #define container_of(ptr, type, member) ({ \
2     const typeof( ((type *)0)->member ) *__mptr = (ptr);
3     (type *) ( (char *)__mptr - offsetof(type,member) );})

```

```

1 /* call to list_entry expands into: */
2 struct car *amazing_car = container_of(car_list_ptr, struct car, list);

```

```

1 /* next expansion: */
2 amazing_car = ({
3     const typeof( ((struct car *)0)->list ) *__mptr = car_list_ptr;
4     (struct car *) ( (char *)__mptr - offsetof(struct car, list));
5 });

```

```

1 /* last expansion: */
2 amazing_car = ({
3     const struct list_head *__mptr = car_list_ptr;
4     (struct car *) ( (char *)__mptr - 0x10);
5 });

```

Linked lists

Defining a linked list

▶ Previous example:

```

1 struct car
2 {
3     unsigned int    max_speed;
4     unsigned int    drive_wheel_num;
5     double          price_in_dollars;
6     struct list_head list;
7 };

```

▶ Static (compile-time) definition:

```

1 struct car my_car =
2 {
3     .max_speed = 150,
4     .drive_wheel_num = 2,
5     price_in_dollars = 10000.0,
6     .list = LIST_HEAD_INIT()
7 }

```

▶ Dynamic (runtime) definition, most commonly used:

```

1 struct car *my_car =
2     kmalloc(sizeof(*my_car), GFP_KERNEL);
3 my_car->max_speed = 150;
4 my_car->drive_wheel_num = 2;
5 my_car->price_in_dollars = 10000.0;
6 INIT_LIST_HEAD(&my_car->list);

```

▶ Canonical pointer representing the list as a whole:

```

1 LIST_HEAD(my_car_list); /* my_car_list is a struct list_head */

```

Linked lists

Adding/deleting a node to/from a list

- ▶ **list_add**
(**struct list_head *new, struct list_head *head**)
 - ▶ Add the node right after the `head` node
- ▶ **list_add_tail**
(**struct list_head *new, struct list_head *head**)
 - ▶ Add the node at the end of the list, i.e. before the `head` node
- ▶ **list_del(struct list_head *entry)**
 - ▶ Remove the element from the list
 - ▶ You still have to take care of the memory deallocation if needed

```
1 list_add(&my_car->list, &my_car_list);  
2 list_add_tail(&my_car->list, &my_car_list);  
3 list_del(&my_car->list);
```


Linked lists

Moving/Splicing nodes

- ▶ **list_move**
(`struct list_head *list, struct list_head *head`)
- ▶ **list_move_tail**
(`struct list_head *list, struct list_head *head`)
 - ▶ Move a node from one list to another one
- ▶ **list_empty(struct list_head *head)**
 - ▶ Returns nonzero if the list is empty
- ▶ **list_splice**
(`struct list_head *list, struct list_head *head`)
 - ▶ Insert the list pointed by `list` after the element `head`

Linked lists

Iterating over a list

▶ `list_for_each()`, `list_for_each_entry()`

```
1 /* Temporary variable needed to iterate: */
2 struct list_head p;
3 /* This will point on the actual data structures (struct car) during the iteration: */
4 struct car *current_car;
5
6 list_for_each(p, &my_car_list)
7 {
8     current_car = list_entry(p, struct car, list);
9     printk(KERN_INFO "Price: %lf\n", current_car->price_in_dollars);
10 }
11
12 /* Simpler: use list_for_each_entry */
13 list_for_each_entry(current_car, &my_car_list, list)
14 {
15     printk(KERN_INFO "Price: %lf\n", current_car->price_in_dollars);
16 }
```

▶ `list_for_each_entry_reverse()`

▶ Iterate backwards

Linked lists

Removing while iterating

```
1 /* This will point on the actual data structures (struct car) during the iteration: */
2 struct car *current_car, *next;
3
4 list_for_each_entry_safe(current_car, next, my_car_list, list)
5 {
6     printk(KERN_INFO "Price: %lf\n", current_car->price_in_dollars);
7     list_del(current_car->list);
8     kfree(current_car); /* if this was dynamically allocated through kmalloc */
9 }
```

- ▶ For each iteration, `next` points to the next node
 - ▶ Can safely remove the current node
 - ▶ Otherwise: → *use-after-free* bug

Linked lists

Linked lists: where are they used in the kernel?

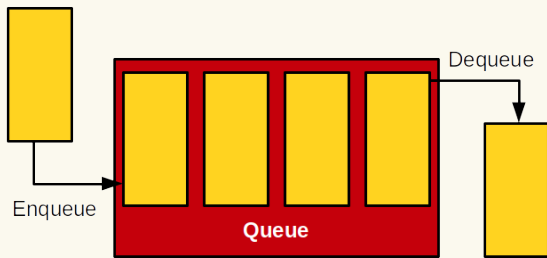
- ▶ **Kernel code makes extensive use of linked lists:**
 - ▶ Linking threads that share a common PID;
 - ▶ Linking the superblocks of all partitions sharing a common file system type;
 - ▶ Linking processes in a CPU run queue
 - ▶ etc.

Outline

- 1 Linked lists
- 2 Queues**
- 3 Maps
- 4 Binary trees
- 5 The right data structure for the right problem

Queues

Definition



- ▶ Producer/consumer programming model
- ▶ **FIFO**: First-In-First-Out
- ▶ Implemented in Linux through `struct kfifo`

Queues

Queues: creation

```
1 struct kfifo my_queue;
2 char *buffer;
3 int ret;
4
5 /* Version 1: dynamic allocation + initialization of a 1024 bytes sized queue */
6 ret = kfifo_alloc(&my_queue, 1024, GFP_KERNEL);
7 if(ret)
8     return ret; /* can fail! return error */
9
10 /* Version 2: initialization of a (dynamically) pre-allocated buffer to be used as a queue */
11 buffer = kmalloc(1024, GFP_KERNEL);
12 ret = kfifo_init(&my_queue, buffer, 1024);
13 if(ret)
14     return ret;
15
16 /* Version 3: static declaration */
17 DECLARE_KFIFO(another_queue, 1024); /* type of another queue is struct kfifo */
18 INIT_KFIFO(another_queue);
```

- ▶ Size should be a power of 2

Queues

Enqueueing/dequeueing

▶ Prototypes:

```

1 /* (these are actually macros in recent kernel versions) */
2 unsigned int kfifo_in(struct kfifo *fifo, const void *from, unsigned int len);
3 unsigned int kfifo_out(struct kfifo *fifo, void *to, unsigned int len);
4 unsigned int kfifo_out_peek(struct kfifo *fifo, void *to, unsigned int len);

```

▶ Enqueueing:

```

1 struct car
2 {
3     unsigned int max_speed;
4     unsigned int drive_wheel_num;
5     double      price_in_dollars;
6 };

```

```

1 unsigned int ret;
2 struct car car_to_add = {100, 2, 10000.0};
3
4 ret = kfifo_in(&fifo, &car_to_add,
5     sizeof(struct car));
6 if(ret != sizeof(struct car))
7     /* Not enough space left in the queue */

```

▶ Dequeueing:

```

1 struct car amazing_car;
2 unsigned int ret = kfifo_out(&fifo, &amazing_car, sizeof(struct car));

```

- ▶ Use `kfifo_out_peek` to access the head of the queue without removal

Queues

Queue size/reset/destroy

▶ Information on queue size - prototypes:

```
1 /* Let's assume we have struct kfifo my_kfifo */
2 unsigned int buffer_total_size_in_bytes = kfifo_size(&my_kfifo);
3 unsigned int bytes_used = kfifo_len(&my_kfifo);
4 unsigned int bytes_free = kfifo_avail(&my_kfifo);
5 int empty = kfifo_is_empty(&my_kfifo);
6 int full = kfifo_is_full(&my_kfifo);
```

▶ Reset a queue (removes all content):

```
1 kfifo_reset(&my_kfifo); /* returns void */
```

▶ Free a queue previously allocated through `kfifo_alloc()`

```
1 kfifo_free(&my_kfifo); /* returns void */
```

▶ Sample of kernel queues usage:

- ▶ In linux sources, `samples/kfifo`

Queues

Queues usage in the kernel

- ▶ List of free blocks for the SmartMedia flash driver
- ▶ Used in the message queue driver for TI OMAP processors to buffer messages
- ▶ etc.

Outline

- 1 Linked lists
- 2 Queues
- 3 Maps**
- 4 Binary trees
- 5 The right data structure for the right problem

Maps

Definition & initialization

- ▶ A **map** maps keys to values, supporting 3 main operations:
 - ▶ add(key, value)
 - ▶ remove(key)
 - ▶ value = lookup(key)
- ▶ Linux implementation indexes content using a **binary search tree**
 - ▶ Keys must support the operation \leq
- ▶ Linux does not implement a general purpose map
 - ▶ The implementation named **idr**,
maps integers (keys) to pointers (values)
 - ▶ These integers are named Unique Identification Numbers (UIDs)
- ▶ Initialization:

```
1 /* Statically */
2 struct idr my_map;
3 idr_init(&my_map);
```

```
1 /* Dynamically */
2 struct idr *my_map_ptr = kmalloc(sizeof(
3     struct idr), GFP_KERNEL);
4 idr_init(my_map_ptr);
```

Maps

New UID allocation

- ▶ 3-steps process
- ▶ Prototypes:

```
1 /* 1. Pre-allocate the memory for the UID allocation request */
2 void idr_preload(gfp_t gfp_mask);
3 /* 2. Actual allocation request */
4 int idr_alloc(struct idr *idr, void *ptr, int start, int end, gfp_t gfp_mask);
5 /* 3. idr_preload disables preemption, needs to re-enable it: */
6 void idr_preload_end(void);
```

- ▶ Note that the interface to add a new UID has changed since the textbook publication
 - ▶ Simplified, removed the need for looping
 - ▶ <https://lwn.net/Articles/536293/>

Maps

Insertion: full example

```
1 int id;
2
3 idr_preload(GFP_KERNEL);
4 id = idr_alloc(&the_idr, &some_data, 10, 550, GFP_KERNEL);
5 idr_preload_end();
6 if(id == -ENOSPC)
7     /* error, no id available in the requested range */
8 else if(id == -ENOMEM)
9     /* error, could not allocate memory */
```

- ▶ UID range constraints provide more control on allocated UIDs
 - ▶ Ex: loop device driver indexes a loop device partitions based on their *minor number*

Maps

UID lookup/removal, map destruction

▶ Prototypes:

```
1 void *idr_find(struct idr *idp, int id);
2 void idr_remove(struct idr *idp, int id);
3 void idr_destroy(struct idr *idp);
```

▶ UID lookup:

```
1 struct car *my_car = idr_find(&my_map, id); /* returns NULL on error */
2 if(!my_car)
3     return -EINVAL; /* not found */
```

▶ UID removal:

```
1 idr_remove(&my_map, id);
```

▶ Map destruction:

```
1 idr_destroy(&my_map);
```

Maps

Maps usage in the kernel

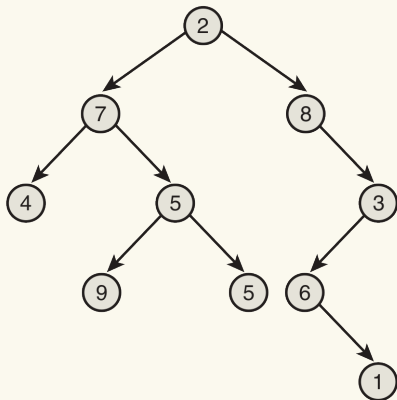
- ▶ Indexing loop devices
 - ▶ (File acting like a virtual block device (ex: disk images, ISOs, etc.))
- ▶ Index permission data structures for IPCs in a namespace
 - ▶ (OS level virtualization)
- ▶ Index Performance Monitoring Unit events
- ▶ etc.

Outline

- 1 Linked lists
- 2 Queues
- 3 Maps
- 4 Binary trees**
- 5 The right data structure for the right problem

Binary trees

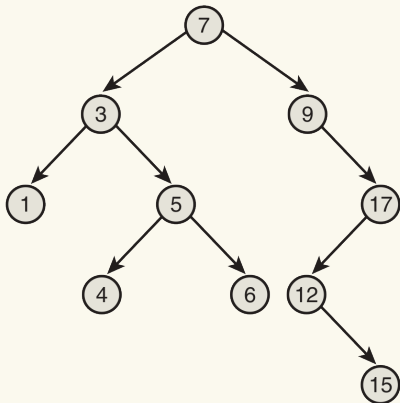
Definition



- ▶ **Binary Tree**
- ▶ Nodes have zero, one or two children
- ▶ Root has 0 parent, other nodes have one

Binary trees

Definition (2)



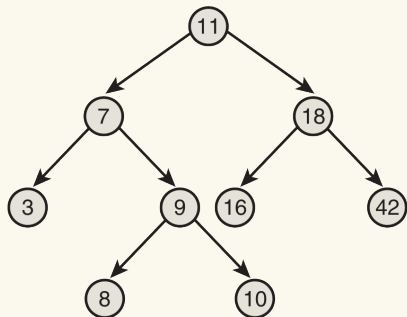
▶ Binary Search Tree

▶ Ordered:

- ▶ Left children $<$ parent
- ▶ Right children $>$ parent
- ▶ Search and in-order traversal are efficient

Binary trees

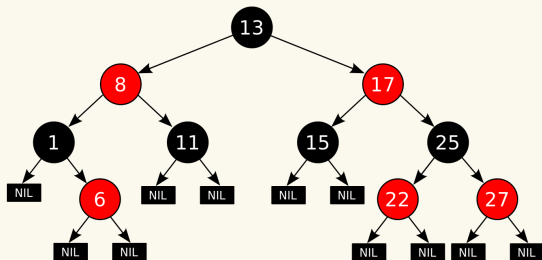
Definition (3)



- ▶ **Balanced Binary Search Tree**
- ▶ Depth of all leaves differs by at most 1
 - ▶ Puts a boundary on the worst case search operation

Binary trees

Definition (4)



- ▶ Linux implements **Red-Black Trees**
- ▶ Nodes: red or black
- ▶ Leaves: black, no data
- ▶ Non-leaves: two children
- ▶ Red nodes have two black children

- ▶ Path from one node to one of its leaves has same amount of black nodes as the shortest path to any of its other leaves
- ▶ Properties are maintained during tree modifications:
 - ▶ Red-black trees are **self-balanced**
 - ▶ Try to stay (semi-)balanced with modifications
 - ▶ Efficient insert operations

Binary trees

Tree creation and search

▶ Creation:

```
1 struct rb_root my_tree_root = RB_ROOT;
```

▶ Search routing must be implemented by the programmer:

```
1 static struct zswap_entry *zswap_rb_search(struct
      rb_root *root, pgoff_t offset)
2 {
3     struct rb_node *node = root->rb_node;
4     struct zswap_entry *entry;
5
6     while (node) {
7         entry = rb_entry(node, struct zswap_entry,
8             rbnode);
9         if (entry->offset > offset)
10            node = node->rb_left;
11        else if (entry->offset < offset)
12            node = node->rb_right;
13        else
14            return entry;
15    }
16    return NULL;
}
```

```
1 struct zswap_entry {
2     struct rb_node rbnode;
3     pgoff_t offset;
4     int refcount;
5     unsigned int length;
6     struct zswap_pool *pool;
7     unsigned long handle;
8 };
```

▶ Use `rb_entry` to get the data structure from the corresponding indexing node

Binary trees

Insertion/deletion

```
1 static int zswap_rb_insert(struct rb_root *root, struct zswap_entry *entry,
2     struct zswap_entry **dupentry)
3 {
4     struct rb_node **link = &root->rb_node, *parent = NULL;
5     struct zswap_entry *myentry;
6
7     while (*link) {
8         parent = *link;
9         myentry = rb_entry(parent, struct zswap_entry, rbnode);
10        if (myentry->offset > entry->offset)
11            link = &(*link)->rb_left;
12        else if (myentry->offset < entry->offset)
13            link = &(*link)->rb_right;
14        else {
15            *dupentry = myentry;
16            return -EEXIST;
17        }
18    }
19    rb_link_node(&entry->rbnode, parent, link);
20    rb_insert_color(&entry->rbnode, root);
21    return 0;
22 }
```

▶ Deletion:

- ▶ `rb_erase(struct rb_node *node, struct rb_root *root)`

Binary trees

Rbtrees: where are they used in the kernel?

- ▶ Rbtrees usage in the kernel:
 - ▶ Processes runqueues for the CFS (default) Linux scheduler
 - ▶ Indexing file (inode) fragments for the CEPH filesystem
 - ▶ Indexing memory areas in a process address space
 - ▶ etc.

Outline

- 1 Linked lists
- 2 Queues
- 3 Maps
- 4 Binary trees
- 5 The right data structure for the right problem

The right data structure for the right problem

- Linked lists:**
- ▶ Sequential iteration over all data is needed
 - ▶ There is an unknown number of elements
-

- Queues:**
- ▶ Useful with producer/consumer pattern
 - ▶ When it's OK to work with a fixed size buffer
-

- Maps:**
- ▶ Need to map a unique integer to a pointer
-

- Red-black trees:**
- ▶ Large amount of data, efficient search

- ▶ Other data structures in the kernel:

- ▶ Radix trees [2]
- ▶ Bitmaps [1]
- ▶ etc.

Bibliography I

- [1] [Bit arrays and bit operations in the linux kernel.](https://0xax.gitbooks.io/linux-insides/content/DataStructures/bitmap.html)
[https://0xax.gitbooks.io/linux-insides/content/DataStructures/bitmap.html.](https://0xax.gitbooks.io/linux-insides/content/DataStructures/bitmap.html)
Accessed: 2017-02-07.
- [2] [Lwn - trees i: Radix trees.](https://lwn.net/Articles/175432/)
[https://lwn.net/Articles/175432/.](https://lwn.net/Articles/175432/)
Accessed: 2017-02-07.