# EPFL

## École Polytechnique Fédérale de Lausanne

# An Operating System Kernel Frame Allocator in Rust

## Semester Project

Author: David Desboeufs
david.desboeufs@epfl.ch

Supervisors:
Edouard Bugnion
edouard.bugnion@epfl.ch

Charly Castes
charly.castes@epfl.ch

3rd January 2023

# Contents

# 1   INTRODUCTION

We aim to develop a simple frame allocator for x86-64 systems, also known as AMD64, which is the x86 architecture's 64-bit variant. We want to develop a simple algorithm that can be proved using symbolic execution, such that it avoids path explosion. Moreover, unbounded loops are not permitted in the algorithm.

The algorithm is written in Rust, which offers several advantages. Rust is designed to be memory safe (e.g. buffer overflows and data races) as it uses strong type system and a system of ownership and borrowing [1]. Furthermore, when compared to C code, it offers good performance [2].

Intel® x86-64 page tables come in three sizes: 4Kb, 2Mb (big page) and 1Gb (huge page) [3]. All of these sizes are used by our frame allocator.

Internal and external fragmentation should be minimal in a good allocator. Internal fragmentation is defined as allocating a larger block of memory than is required [4]. External fragmentation is defined as the inability to allocate memory despite sufficient memory. This is due to several small blocks of free memory being distributed across the memory space [4].

# 2   ALGORITHM DESCRIPTION

## 2.1   TREE REPRESENTATION

The buddy allocation technique inspired our frame allocator [5]. This consists of recursively splitting blocks in two to obtain the smallest possible block to satisfy a request.

Because blocks are divided into two parts, each one has a buddy. When a block is deallocated, it is merged with his buddy if this one is free. This technique allows us to efficiently allocate and deallocate power of two size blocks of memory.

Buddy allocation technique can be illustrated using a full binary tree [6], with node values indicating whether or not a node is free (1 for free, 0 for used). Figure 1 shows the occupation tree, where two blocks are allocated: 8Kb and 4Kb, respectively (in red). We can note that a parent block is free if his two children are also free (logical AND tree).
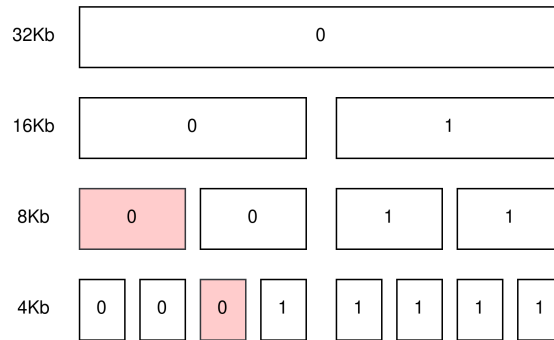


**FIGURE 1**
32Kb Buddy Allocator Representation

Except when the top level is free, it does not provide any information about available 8Kb and 4Kb blocks. For example, if we want to allocate a 4Kb block, we cannot begin at the top of the tree, instead, we must to go through each 4Kb blocks to find the first one that is free.

In order to know which child is free, each node has two bits instead of one, one bit per child. On the right, Figure 2 depicts the new binary tree. Note that a node now corresponds to two blocks rather than a single block. When we see "01" in a node, we know that the right child is available. However, this new architecture does not address the issue of locating a free 4Kb starting from the top level.
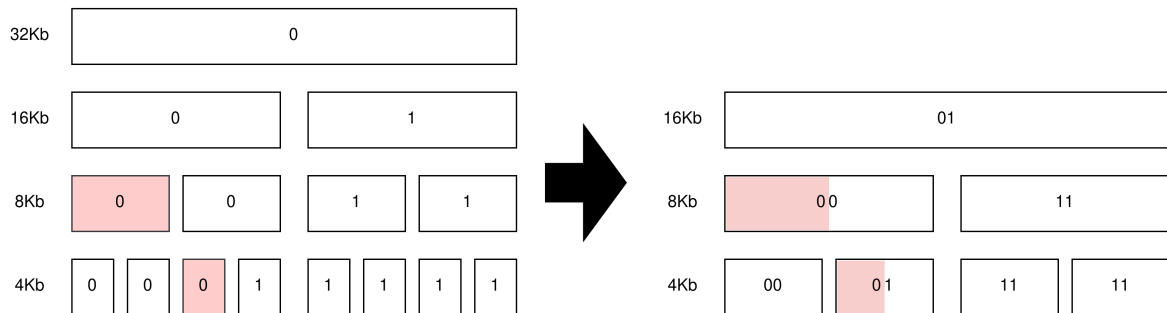


**FIGURE 2**
2 Bits per Node

One solution is to create an OR tree rather than an AND tree, which means that we have a 1 (free) if at least one child is free. Figure 3 shows the new tree that can be used efficiently to find a free 4Kb blocks beginning at the top. To find a free 4Kb block we must traverse three nodes, which correspond to tree's height.
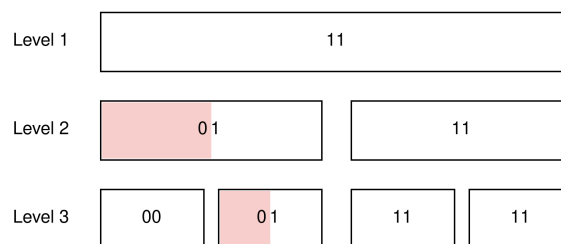


**FIGURE 3**
4Kb Tree

Because this tree can only be used to allocate 4Kb blocks, we must create a separate tree for each block size. Figure 4 depicts the two additional OR trees which allow to allocate 8Kb blocks (on the left) and 16Kb block (on the right). These three data structures allow us to allocate the three block sizes in a simple manner, for each tree we begin at the level 1, and if it contains a 1, we have at least one free block of the desired size in memory. Because there is only 4Kb memory left on the left side, we see that both trees have level 1 equal to "01".
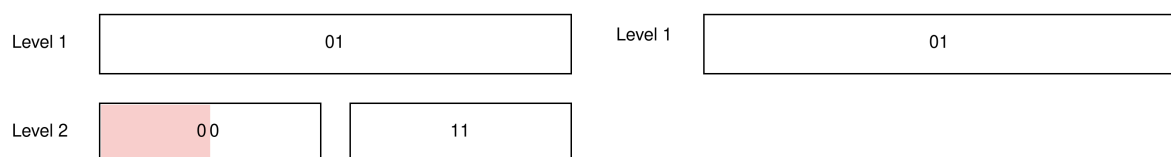


**FIGURE 4**
8Kb and 16Kb Trees

We can easily allocate blocks of 4Kb, 8Kb, and 16Kb using the previously created trees. When an allocation is completed, the three trees must be updated to reflect the changes. Observe that some changes are unnecessary, for example, when a 8Kb block is allocated, all levels of the 4Kb tree are updated. The level 3 update is irrelevant and can be ignored, when we set the parent bit (level 2) to 0 we will never visit the children. This reduces the number of updates we must perform at each allocation. Figure 5 illustrates the new three data structures, which have the same configuration as before (8Kb and 4Kb allocated).
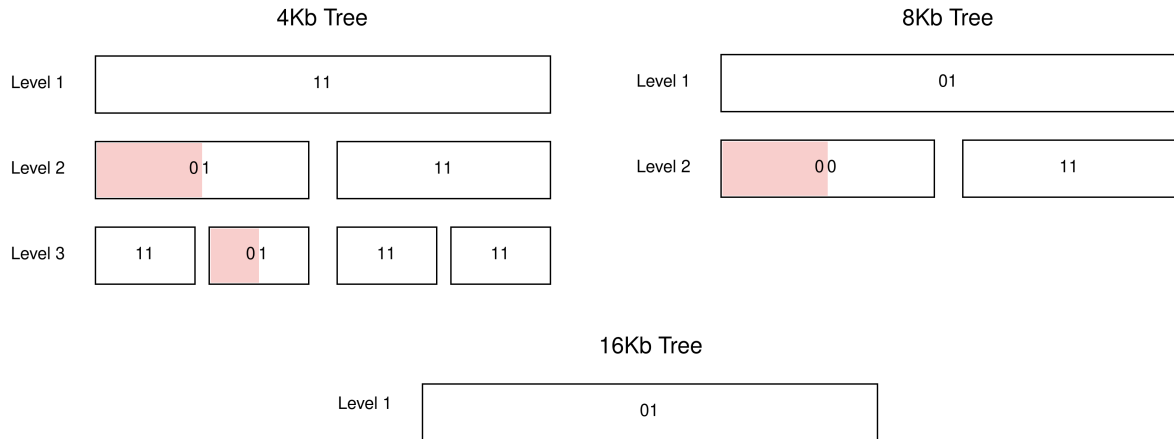


**FIGURE 5**
4Kb, 8Kb and 16Kb Trees

When we perform an allocation, we first search through the dedicated tree corresponding to the block size. If a block is available, we set the bit of the chosen block to 0 of the bottom level first. After that, we update the upper levels to maintain a consistent OR tree. To accomplish this, we check recursively if we need to change the upper bit, if the buddy is already used, then both are used, we set the parent bit to 0. We do this for all upper level.

In the two other trees, we set to 0 the following bits: level 1 for 16Kb block allocation, level 2 for 8Kb and 4Kb block allocation (we ignore unnecessary updates). Updates are done recursively in the same manner as for the corresponding tree to the block size.

The same logic applies to block deallocation, except that we must set upper block to 1 (free) in order to have a consistent OR tree.

## 2.2 AMD64 PAGE TABLE

We must address the allocation of 4Kb, 2Mb and 1Gb pages. With the previously described structure, we need a 4Kb tree of height 18 ($\log_2(512^2)$). However, there is a factor of 512 between 4Kb and 2Mb, and the same between 2Mb and 1Gb. Here comes the suggestion of having 512 children per node rather than the previously stated two. Each node has 512 bits (one bit per child). The height of the 4Kb tree is three now. Figure 6 portrays the representation of a single node and its 512 children.
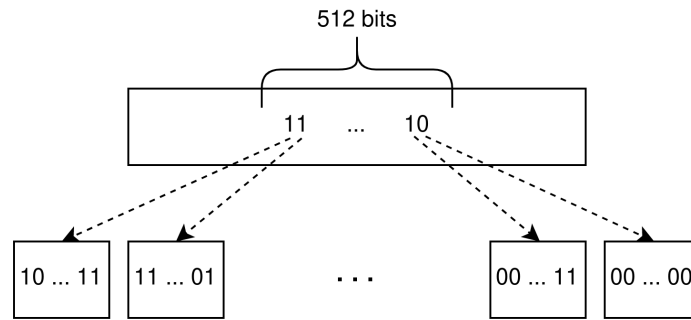
**FIGURE 6**
512-ary Node Representation

We eventually end up with three trees, one for each page size. We have a maximum capacity of 512Gb because each node contains 512 bits ($512^3 * 4Kb = 512Gb$).

## 2.3 FIND A FREE PAGE

As we begin at the top level, we must decide which child to choose from the ones are free. A simple technique is to always choose the rightmost one (search from right to left). Section 4.2 illustrates the effect of changing the direction of search on external fragmentation.

The option that minimizes external fragmentation is the following: search from right to left for 4Kb and 2Mb trees and search from left to right for the 1Gb tree. This comes from the fact that when a page of 1Gb is freed, a page of 4Kb or 2Mb may take its place, causing external fragmentation. This can be avoided if 1Gb page is allocated beginning from the other side.

# 3 RUST IMPLEMENTATION

Figure 7 illustrates the three functions dedicated to allocation and figure 8 shows the three functions dedicated to deallocation.

```rust
/**
 * Allocate 4kb page
 * return None if allocation fails
 */
pub fn allocate_frame(&mut self) -> Option<usize> {…


/**
 * Allocate 2Mb page
 * return None if allocation fails
 */
pub fn allocate_big_page(&mut self) -> Option<usize> {…


/**
 * Allocate 1Gb page
 * return None if allocation fails
 */
pub fn allocate_huge_page(&mut self) -> Option<usize> {…
```

**FIGURE 7**
Allocation Functions

```
/**
 * Deallocate frame
 * nothing is done if frame was not previously allocated
 */
pub fn deallocate_frame(&mut self, frame_id: usize) {…

/**
 * Deallocate big page
 * nothing is done if page was not previously allocated
 */
pub fn deallocate_big_page(&mut self, frame_id: usize) {…

/**
 * Deallocate huge page
 * nothing is done if frame was not previously allocated
 */
pub fn deallocate_huge_page(&mut self, frame_id: usize) {…
```

**FIGURE 8**
Deallocation Functions

## 3.1   TREE STORAGE

Trees are stored in simple arrays, figure 9 shows how they are flattened (binary tree as an example),
where node content corresponds to its array index [7]. Given a node $index$, the left child corresponds to
$2 * index + 1$ and the right child to $2 * index + 2$. This is similar for the 512-ary tree, except we multiply
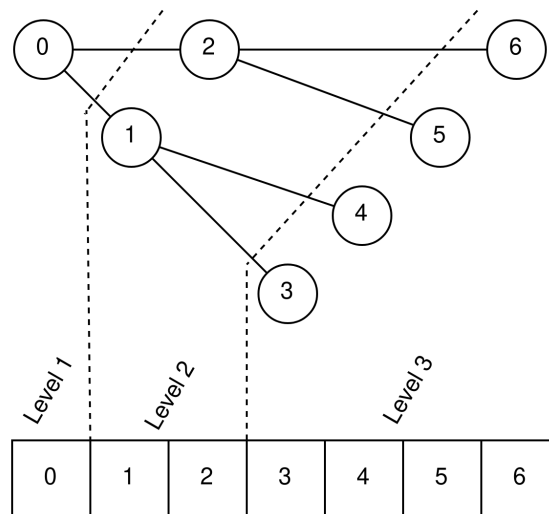by 512 and add from 1 to 512 to iterate over all the children.



**FIGURE 9**
Flatten Binary Tree

6

## 3.2 FIND FIRST FREE PAGE EFFICIENTLY

In this section, we will look at how to efficiently find the first 1 in a sequence of bits. We have a sequence of 512 bits, which must be divided into 8 parts of 64 bits each, as we use 64 bits machine. Figure 10 shows the simplest method in Rust as well as its optimized assembly code on the right.

```rust
pub extern fn find_first_one(input: u64) -> u64 {
    let mut temp = input;
    for i in 0..64 {
        if temp & 1 == 1{
            return i;
        }
        temp >>= 1;
    }
    return 0;
}
```

```asm
example::find_first_one:
        xor     eax, eax
.LBB0_1:
        test    dil, 1
        jne     .LBB0_10
        test    dil, 2
        jne     .LBB0_9
        test    dil, 4
        jne     .LBB0_8
        test    dil, 8
        jne     .LBB0_7
        shr     rdi, 4
        add     rax, 4
        cmp     rax, 64
        jne     .LBB0_1
        xor     eax, eax
        ret
.LBB0_9:
        add     rax, 1
.LBB0_10:
        ret
.LBB0_8:
        add     rax, 2
        ret
.LBB0_7:
        add     rax, 3
        ret
```

**FIGURE 10**
Simple For Loop

We observe that the compiler uses loop and rolling optimization technique, but we still have a linear search. In fact, the Intel instruction set includes instructions for determining the index of the first bit set to 1, *Bit Scan Forward* (BSF) and *Bit Scan Reverse* (BSR) [8]. Figure 11 depicts BSF's Rust code as well as its assembly code on the right.

```rust
pub extern fn bsf(input: u64) -> u64 {
    let mut pos: u64;
    unsafe {
        asm! {
            "bsf {pos}, {input}",
            input = in(reg) input,
            pos = out(reg) pos,
            options(nomem, nostack),
        };
    };
    return pos;
}
```

```asm
example::bsf:

        bsf     rax, rdi

        ret
```

**FIGURE 11**
Bit Scan Forward

Instead of directly using assembly code in Rust, we can use the *trailing_zeros* Rust method. This method returns the number of trailing zeros, which correspond to index of the first one. Figure 12 shows the *trailing_zeros* method code as well as its assembly code one the right. We observe that it is similar to the previous assembly code, with an additional check performed.

```
pub extern fn trailing_zeros(input: u64) -> u32 {
    input.trailing_zeros()
}
```

```
example::trailing_zeros:
        test    rdi, rdi
        je      .LBB1_1
        bsf     rax, rdi
        ret
.LBB1_1:
        mov     eax, 64
        ret
```

**FIGURE 12**
Trailing Zeros Method

Figure 13 illustrates average time of 1000 BSF and simple loop method iterations. We notice both BSF and *trailing_zeros* run in constant time.
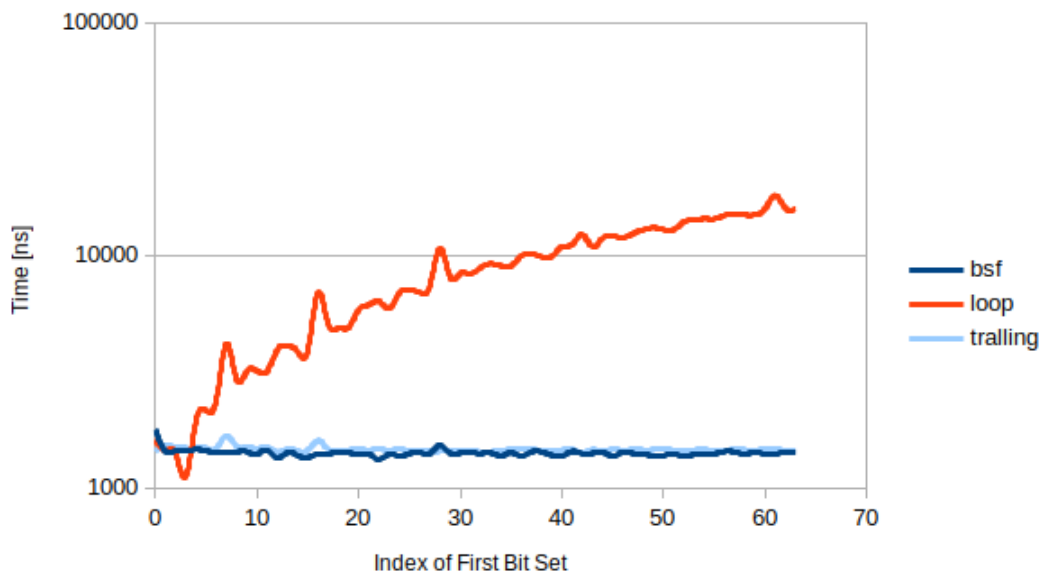


**FIGURE 13**
Benchmark BSF

## 3.3 SAFE DEALLOCATION

Deallocation is safe in the sense that it has no undefined behaviour when given an invalid address, it does nothing instead. Figure 14 illustrates the code diagram that enables us to determine the type of allocated page given an address. For example, when an address is not aligned, we can discard it: 2Mb page must have an address that is a multiple of $512$ and 1Gb page must have an address that is a multiple of $512^2$. Other cases are solved by inspecting the three data structures.

Furthermore, to distinguish which block type is allocated, we use the principle described in figure 5: children are not automatically set to 0.



**FIGURE 14**
Code Diagram To Find Page Type

## 3.4 PSEUDO-CODE

The six pseudo-codes to allocate and deallocate pages are shown below. The strategy to allocate a frame can be summarized as searching for an available frame in the corresponding tree. Then we set the necessaries bits to used in the three trees.

The strategy for deallocation can be summarized as determining whether or not the given address is valid. Then, in the three trees, set the necessary bits to free.

### 3.4.1 ALLOCATE 4KB PAGE

---

**Algorithm 1** $allocate\_frame()$

---

▷ Search free 4Kb page
$l1\_idx \leftarrow search\_first\_bit\_set(Tree4Kb, 0, LEVEL\_1)$
**if** $l1\_idx$ not valid **then**
    **return** None                                               ▷ No available memory
**end if**

$l2\_child\_node \leftarrow compute\_node\_child(l1\_idx, LEVEL\_2)$
$l2\_idx \leftarrow search\_first\_bit\_set(Tree4Kb, l2\_child\_node, LEVEL\_2)$

$l3\_child\_node \leftarrow compute\_node\_child(l2\_idx, LEVEL\_3)$
$l3\_idx \leftarrow search\_first\_bit\_set(Tree4Kb, l3\_child\_node, LEVEL\_3)$

▷ 4Kb tree: set bits to 0
$set\_bit\_to\_zero(Tree4Kb, l3\_idx, LEVEL\_3)$
**if** 512 bits of $l3\_child\_node$'s 4Kb tree are used **then**              ▷ XOR operation
    $set\_bit\_to\_zero(Tree4Kb, l2\_idx, LEVEL\_2)$
    **if** 512 bits of $l2\_child\_node$'s 4Kb tree are used **then**          ▷ XOR operation
        $set\_bit\_to\_zero(Tree4Kb, l1\_idx, LEVEL\_1)$
    **end if**
**end if**

▷ 2Mb tree: set bits to 0
$set\_bit\_to\_zero(Tree2Mb, l2\_idx, LEVEL\_2)$
**if** 512 bits of $l2\_child\_node$'s 2Mb tree are used **then**              ▷ XOR operation
    $set\_bit\_to\_zero(Tree2Mb, l1\_idx, LEVEL\_1)$
**end if**

▷ 1Gb tree: set bit to 0
$set\_bit\_to\_zero(Tree1Gb, l1\_idx, LEVEL\_1)$

**return** $Some((l1\_idx << 18) + (l2\_idx << 9) + l3\_idx)$

---

### 3.4.2 ALLOCATE 2MB PAGE

**Algorithm 2** $allocate\_big\_page()$

  ▷ Search free 2Mb page
$l1\_idx \leftarrow search\_first\_bit\_set(Tree2Mb, 0, LEVEL\_1)$
**if** $l1\_idx$ not valid **then**
    **return** None                                   ▷ No available memory
**end if**

$l2\_child\_node \leftarrow compute\_node\_child(l1\_idx, LEVEL\_2)$
$l2\_idx \leftarrow search\_first\_bit\_set(Tree2Mb, l2\_child\_node, LEVEL\_2)$

  ▷ 2Mb tree: set bits to 0
$set\_bit\_to\_zero(Tree2Mb, l2\_idx, LEVEL\_2)$
**if** 512 bits of $l2\_child\_node$'s 2Mb tree are used **then**             ▷ XOR operation
    $set\_bit\_to\_zero(Tree2Mb, l1\_idx, LEVEL\_1)$
**end if**

  ▷ 4Kb tree: set bits to 0
$set\_bit\_to\_zero(Tree4Kb, l2\_idx, LEVEL\_2)$         ▷ Note that $LEVEL\_3$ is not modified
**if** 512 bits of $l2\_child\_node$'s 4Kb tree are used **then**             ▷ XOR operation
    $set\_bit\_to\_zero(Tree4Kb, l1\_idx, LEVEL\_1)$
**end if**

  ▷ 1Gb tree: set bit to 0
$set\_bit\_to\_zero(Tree1Gb, l1\_idx, LEVEL\_1)$

  **return** $Some((l1\_idx << 18) + (l2\_idx << 9))$

### 3.4.3 ALLOCATE 1GB PAGE

**Algorithm 3** $allocate\_huge\_page()$

  ▷ Search free 1Gb page
$l1\_idx \leftarrow search\_first\_bit\_set(Tree1Gb, 0, LEVEL\_1)$
**if** $l1\_idx$ not valid **then**
    **return** None                                   ▷ No available memory
**end if**

  ▷ 1Gb tree: set bit to 0
$set\_bit\_to\_zero(Tree1Gb, l1\_idx, LEVEL\_1)$

  ▷ 4Kb tree: set bit to 0
$set\_bit\_to\_zero(Tree4Kb, l1\_idx, LEVEL\_1)$

  ▷ 2Mb tree: set bit to 0
$set\_bit\_to\_zero(Tree2Mb, l1\_idx, LEVEL\_1)$

  **return** $Some((l1\_idx << 18))$

### 3.4.4 DEALLOCATE 4KB PAGE

---

**Algorithm 4** $deallocate\_frame(frame\_id)$

---

**if** frame was not previously allocated **then**
    **return**
**end if**

▷ Extract level indices
$l3\_idx \leftarrow frame\_id \ \& \ 0x1FF$
$l2\_idx \leftarrow (frame\_id \ >> \ 9) \ \& \ 0x1FF$
$l1\_idx \leftarrow (frame\_id \ >> \ 18) \ \& \ 0x1FF$

▷ Free the 3 Level of 4Kb Tree
$set\_bit\_to\_one(Tree4Kb, \ l1\_idx, \ LEVEL\_1)$
$set\_bit\_to\_one(Tree4Kb, \ l2\_idx, \ LEVEL\_2)$
$set\_bit\_to\_one(Tree4Kb, \ l3\_idx, \ LEVEL\_3)$

▷ if all 4Kb are free, free upper level for 2Mb
**if** $are\_all\_free(Tree4Kb, \ l3\_idx, \ LEVEL\_3)$ **then**
    $set\_bit\_to\_one(Tree2Mb, \ l1\_idx, \ LEVEL\_1)$
    $set\_bit\_to\_one(Tree2Mb, \ l2\_idx, \ LEVEL\_2)$
**end if**

▷ if all 2Mb are free, free 1Gb
**if** $are\_all\_free(Tree2Mb, \ l2\_idx, \ LEVEL\_2)$ **then**
    $set\_bit\_to\_one(Tree1Gb, \ l1\_idx, \ LEVEL\_1)$
**end if**

---

### 3.4.5 DEALLOCATE 2MB PAGE

---

**Algorithm 5** $deallocate\_big\_page(frame\_id)$

---

**if** frame was not previously allocated **then**
    **return**
**end if**

▷ Extract level indices
$l2\_idx \leftarrow (frame\_id >> 9) \& 0x1FF$
$l1\_idx \leftarrow (frame\_id >> 18) \& 0x1FF$

▷ Free the 2 Level of 2Mb Tree
$set\_bit\_to\_one(Tree2Mb, l1\_idx, LEVEL\_1)$
$set\_bit\_to\_one(Tree2Mb, l2\_idx, LEVEL\_2)$

▷ Free the 2 upper Level of 4Kb Tree
$set\_bit\_to\_one(Tree4Kb, l1\_idx, LEVEL\_1)$
$set\_bit\_to\_one(Tree4Kb, l2\_idx, LEVEL\_2)$

▷ if all 2Mb are free, free 1Gb
**if** $are\_all\_free(Tree2Mb, l2\_idx, LEVEL\_2)$ **then**
    $set\_bit\_to\_one(Tree1Gb, l1\_idx, LEVEL\_1)$
**end if**

---

### 3.4.6 DEALLOCATE 1GB PAGE

---

**Algorithm 6** $deallocate\_huge\_page(frame\_id)$

---

**if** frame was not previously allocated **then**
    **return**
**end if**

▷ Extract level index
$l1\_idx \leftarrow (frame\_id >> 18) \& 0x1FF$

▷ Free all the Level 1
$set\_bit\_to\_one(Tree1Gb, l1\_idx, LEVEL\_1)$
$set\_bit\_to\_one(Tree4Kb, l1\_idx, LEVEL\_1)$
$set\_bit\_to\_one(Tree2Mb, l1\_idx, LEVEL\_1)$

---

# 4 EVALUATION

In this section we evaluate two metrics: the memory overhead of the additional structure and the measure of external fragmentation given a predefined scenario.

## 4.1 MEMORY OVERHEAD

Figure 15 depicts memory usage in relation to total memory available, we see that overhead accounts for $\sim 0.004\%$ of total available memory. The 4Kb tree with its tree levels is the largest contributor to overhead.
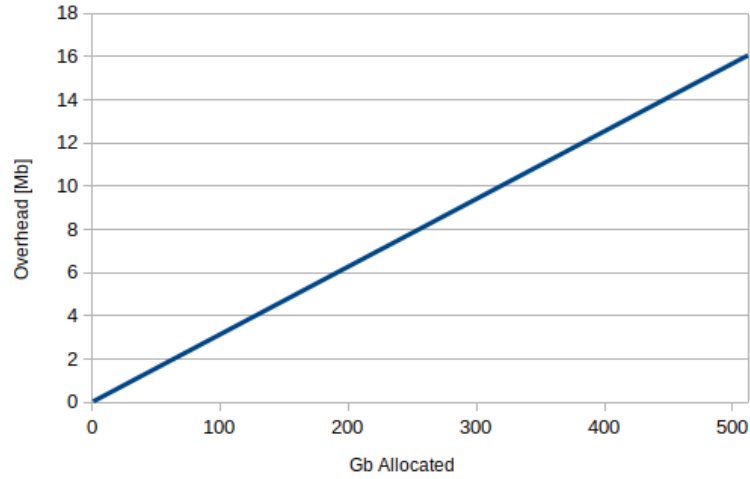


**FIGURE 15**
Memory Overhead

## 4.2 EXTERNAL FRAGMENTATION MEASUREMENT

Here is the benchmark for measuring the allocator's internal fragmentation. The goal is to allocate and deallocate pages of varying sizes at random while measuring fragmentation. We want to reach a memory usage of 70%. To achieve this goal we use an inverse cumulative distributive Poisson function to determine the probability of doing an allocation, as illustrated in figure 16. Note that we only perform allocation until memory utilisation exceeds about 50%.
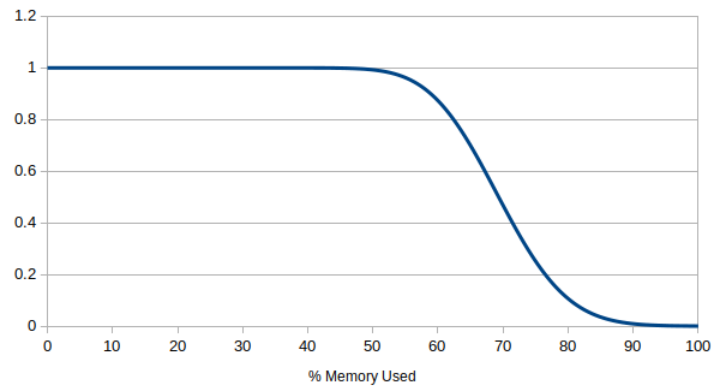


**FIGURE 16**
Inverse Poisson Cumulative Distribution

14

The goal is to have one third occupations for each page size, probabilities are shown below:

$$P(allocate/deallocate\ 4Kb) = \frac{512^2}{512^2 + 512 + 1}$$
$$P(allocate/deallocate\ 2Mb) = \frac{512}{512^2 + 512 + 1}$$
$$P(allocate/deallocate\ 1Gb) = \frac{1}{512^2 + 512 + 1}$$

Figures 17 and 18 illustrate the outcome of two billion successive random allocation and deallocation operations. Figure 17 shows the number of allocated blocks for each size, as well as the blocks that can only be allocated to 2Mb and 4Kb. One can see the external fragmentation on the top of the figure with the $4Kb\_free$ and $2Mb\_free$, indicating that this portion of the memory cannot be allocated to 1Gb page. The goal is to have a large majority of $1Gb\_free$ for unused blocks. One $1Gb\_free$ can be used to allocate each of the three page sizes. The $4Kb\_free$ memory zone can only be used to allocate 4Kb page. There is at most approximately 15% fragmentation, which seems reasonable.

Figure 18 illustrates the spatial representation of memory blocks. We observe 1Gb blocks are placed on the opposite side, as mentioned in section 2.3, search direction for 1Gb pages is reversed.
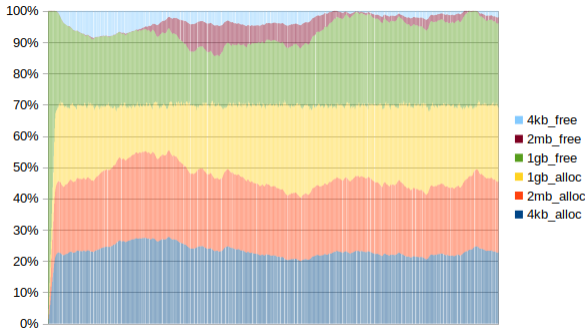


**FIGURE 17**
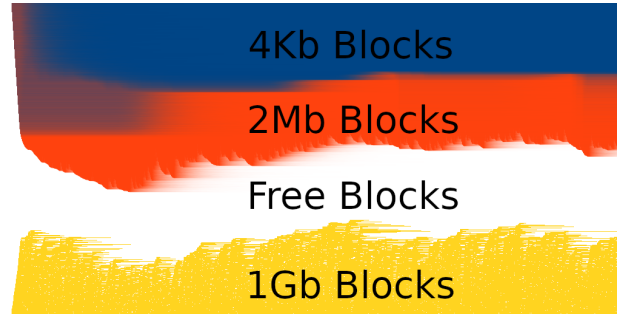Quantitative Memory Representation



**FIGURE 18**
Spatial Memory Representation

# 5   CONCLUSION

We have an efficient algorithm that allows us to allocated three different blocks sizes in few operations, in a simple manner. Furthermore, the code contains no unbounded loops. The code contains $\sim 500$ lines of code (LoC).

The *bit scan forward* x86 assembly instruction improves allocator performance. Memory overhead is kept to 0.004%. Because allocations are done in a dummy manner, external fragmentation may result. Furthermore, without the use of an additional memory structure, deallocations are safe.

## 5.1   PERSPECTIVES

There is no comparison with other allocator such as the x86 page table allocator from Linux, there are too many dependencies to simply extract the standalone allocator's code. The benchmark only checks behaviour of algorithm against external fragmentation, it is not a real-world benchmark.

One can also measure performance of our allocator, in terms of how many cycles are required to complete each task.

# REFERENCES

[1] *Rust*. 2022. URL: https://www.rust-lang.org/ (visited on 26th Dec. 2022).

[2] Yi Lin et al. 'Rust as a Language for High Performance GC Implementation'. In: *Proceedings of the 2016 ACM SIGPLAN International Symposium on Memory Management*. ISMM 2016. Santa Barbara, CA, USA: Association for Computing Machinery, 2016, pp. 89–98. ISBN: 9781450343176. DOI: 10.1145/2926697.2926707. URL: https://doi.org/10.1145/2926697.2926707.

[3] Intel®. 'Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3 Part 1'. In: 2022. Chap. 4.1.1.

[4] Greg Gagne Abraham Silberschatz Peter Baer Galvin. *Operating System Concepts*. 2013. ISBN: 978-1-118-06333-0.

[5] J.M. Chang and E.F. Gehringer. 'A high performance memory allocator for object-oriented systems'. In: *IEEE Transactions on Computers* 45.3 (1996), pp. 357–366. DOI: 10.1109/12.485574.

[6] William D. McQuain. *Full and Complete Binary Trees*. 2009. URL: https://courses.cs.vt.edu/~cs3114/Fall10/Notes/T03a.BinaryTreeTheorems.pdf (visited on 26th Dec. 2022).

[7] Romolo Marotta et al. 'A Non-blocking Buddy System for Scalable Memory Allocation on Multi-core Machines'. In: *2018 IEEE International Conference on Cluster Computing (CLUSTER)*. 2018, pp. 164–165. DOI: 10.1109/CLUSTER.2018.00034.

[8] Intel®. 'Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 2 Part 1'. In: 2022. Chap. 3-125.