

Page up and page down

Johan Montelius

HT2021

1 Introduction

In this tutorial you will implement an array allocation scheme that on the surface behaves very much like the array framework that you implemented in the segmentation exercise. This time you will use paging and to understand why you first must recap the problems with segmentation.

The problem with segmentation is that you will end up with a fragmented memory. Many small free memory areas that together could make up a large part of the available memory but they are all too small to be very useful. We could of course implement some compaction scheme but this is not always easy and you don't want to sit around waiting for defragmentation (remember Windows?).

To mitigate this problem we instead use a paging scheme. You will see that it is a bit trickier and if we did not have hardware support in the CPU we would not be able to use it efficiently.

2 Memory and arrays

Let's start as before and define a memory that we will use and what an array structure will look like. The general idea is that the memory will be divided into a sequence of *frames* and that each array will hold a *page table* that holds one entry per *page*. To make the system more flexible we define two macros that defines how many frames we will have in memory and the *size* of pages/frames (they are of course of the same size since a page should fit into a frame).

```
#define FRAMES 64
```

```
#define SIZE 16
```

We have here chosen 64 and 16 but could of course have chosen anything, the size is however preferably a power of two. The *virtual addresses* will be broken down into an *offset* (where in a page) and *index* (which page). Since the page size is 16 we need four bits as an offset so we create a macro that does a bit-wise *and operation* to select the four least significant bits. Another macro will shift an integer four bits to select the higher bits as page number.

```
#define Offset(addr) (addr & 0b1111)
```

```
#define PageNr(addr) (addr >> 4)
```

Note that we we do these operations we take for granted how an integer (that will be our address) is represented binary.

We also define two values, *FREE* and *TAKEN*, that we will use in a structure that keeps track of which frames that are used. The memory itself is represented as an array of integers.

```
typedef enum available {FREE, TAKEN} available;

available framemap[FRAMES];

int memory[FRAMES*SIZE];
```

So now to the actual array; the array is represented as an array of frame numbers i.e. the page table. We also keep track of how large this structure is since we both want to check that we're not addressing outside of the array bit also since we want to return all frames when we delete the array. The page table will use -1 to indicate that no frame has yet been allocated to the page.

```
typedef struct array {
    int pages;
    int *pagetable;    // an array of frame numbers (or -1 if not allocated)
} array;
```

Before you continue you should have a clear picture of where we are going. We will have a huge memory (well not that huge) that is divided into frames. We will know, by looking in the frame map if a frame is taken or not. An address in an array is broken down into a page number and an offset. A page number is translated using the page table into frame number. Draw this on a paper, once we start to manipulate these data structures you need to be able to visualize what we're doing.

3 Creating an array

When an array is created we need to find free frames to allocate to the array. We therefore start by implementing a procedure that will search through all frames and select the first one that is free. The frame is marked as taken and the frame number is returned. If no free frame is found we return -1 to signal that we have no more frames to offer.

```
int find_free() {
    for(int i = 0; i < FRAMES; i++) {
        if( framemap[i] == FREE ) {
            :
            :
        }
    }
}
```

```

    }
    return -1;
}

```

We're now ready to allocate an array of a given size. We first allocate the array itself and the page table and then allocate the frames that we need. If we fail to find a free frame we need to return everything and why not then use a delete procedure that we will define later.

Fill in the dotted lines and you should soon have this up and running.

```

array *allocate(int size) {
    int pages = size / SIZE;
    int rem = size % SIZE;
    if( rem > 0)
        pages += 1;

    array *new = (array*) ...
    int *pagetable = (int*) ...

    new->pages = pages;
    new->pagetable = pagetable;

    for(int i = 0; i < pages; i++) {
        new->pagetable[i] = -1;    // no frame yet allocated
    }

    printf("allocate array, frames: ");
    for(int i = 0; i < pages; i++) {
        int f = ...
        if( f == -1 ) {
            delete(new);
            return NULL;
        }
        printf("%d ", f);
        :
    }
    printf("\n");
    return new;
}

```

To delete an array we simply reverse what `find_free()` does and deallocate the data structures.

```

void delete(array *arr) {
    int pages = arr->pages;
    printf("delete array, freeing frames: ");

```

```

    for( int i = 0; i < pages; i++) {
        if( arr->pagetable[i] != -1) {
            printf(" %d", arr->pagetable[i]);
            ... = FREE;
        }
    }
    printf("\n");
    free(...);
    free(...);
    return;
}

```

One more wrapper function that will create a new array if possible. If we fail we have nothing else to do but to fail the whole computation. No garbage collector in the world would save us since no matter how we restructure the arrays we will not be able to create more free frames.

```

array *create(int size) {
    array *new = allocate(size);
    if(new == NULL) {
        printf("out of memory\n");
        exit(-1);
    }
    return new;
}

```

We're now ready to implement the set and get operations and this will be slightly more complicated compared to the segmentation exercise. We first need to divide the "address" into the page number and the offset. Once this is done we check that the number is a valid page. If the page is ok, we retrieve the frame number and access the memory location.

```

void set(array *arr, int pos, int val) {
    printf("set: arr %p pos %d val %d\n", arr, pos, val);
    int offset = Offset(pos);
    int page = PageNr(pos);

    if( page >= arr->pages ){
        printf("segmentation fault\n");
        exit(1);
    }

    int frame = arr->pagetable[page];
    printf("set: page %d offset %d frame %d\n", page, offset, frame);

    memory[frame*SIZE + offset] = val;
}

```

```

    return;
}

int get(array *arr, int pos) {
    printf("get: %p pos %d\n", arr, pos);
    int offset = Offset(pos);
    int page = PageNr(pos);

    if( page >= arr->pages ){
        printf("segmentation fault\n");
        exit(1);
    }

    int frame = arr->pagetable[page];
    printf("get: page %d offset %d frame %d\n", page, offset, frame);

    return memory[frame*SIZE + offset];
}

```

Will we catch all faulty addressing of the array? What if we allocate an array of size 40 and have a page size of 16 - what will happen when we address position 50?

4 A first run

I think you are ready to do a small test run of your system. Let's write a small benchmark to see if things work.

```

void bench() {

    array *a = create(20);
    array *b = create(40);

    set(a, 10, 110);
    set(a, 18, 118);

    set(b, 8, 208);
    set(b, 36, 212);

    printf(" a[10] + a[18] = %d\n", get(a,10) + get(a, 18));
    printf(" b[8] + b[36] = %d\n", get(b,8) + get(b, 36));

    delete(a);
    delete(b);
}

```

Also try to create larger arrays than what we know we have memory for and accessing outside of an array. You could try to access an array with a negative index and this will likely result in a crash. Update your implementation to catch this error and print a nice error message.

5 Be lazy

When you know that you have things up and running it's time to do a trick that most operating systems do. We will be as lazy as possible and only allocate frames if they are actually needed.

Take a look at the `allocate()` procedure; do we really need to allocate all frames directly. Could we wait until we see that they are actually needed?

If we want to keep track of if a frame has been allocated or not we can extend the page table to hold more information. We could extend the page table entry to be a small structure that holds a status information but why not be lazy. What if we leave all entries in the page table with the value `-1`?

How about this:

```
array *allocate(int size) {
    int pages = size / SIZE;
    int rem = size % SIZE;
    if( rem > 0)
        pages += 1;

    array *new = (array*) ...
    int *pagetable = (int*) ...

    new->pages = pages;
    new->pagetable = pagetable;

    for(int i = 0; i < pages; i++) {
        new->pagetable[i] = -1;    // no frame yet allocated
    }

    return new;
}
```

Now we must be very careful when we write to the array. If the page table entry returns `-1` we quickly need to find a free frame, insert the frame number into the page table and then continue as if nothing has happened. You will be able to do this with just a few lines of code.

```
:
if( frame == -1 ) {
    printf("page fault ... ");
```

```

    frame = ...
    if( frame == -1 ) {
        printf("out of memory\n");
        exit(-1);
    }
    printf("ok\n");
    arr->pagetable[page] = frame;
}
:

```

When we fix the `get()` procedure we realize that reading from a not yet allocated page could return zero. There is no need to allocate a frame and then do a read operation, the page has never been written to.

If you think that this is only a fun trick you're wrong. This is what an operating system does every time we request more memory. It will set up the page table correctly but will not allocate any frames unless we actually write to the pages. As you will see we can do even more tricks by delaying operations until they are actually needed.

6 Lazy copy

The lazy strategy can also be used when we copy an array. Why not try to delay the copying procedure until it is actually needed. The idea is to create a *lazy copy* of an array; the two array structures should share the frames. We should still be able to read from either array; only when we write to any of the two array structures will we create a copy but then of course only of the page that is written to.

This will require some more book-keeping so let's extend our page table to hold some thing that can hold more information. Let's define a page table entry and then let the array hold a proper page table.

```

typedef enum pte_status {ALLOCATED, LAZY, SHARED} pte_status;

typedef struct pt_entry {
    int frame;
    pte_status status;
    struct pt_entry *copy; // who else shares the frame
} pt_entry;

typedef struct array {
    int pages;
    pt_entry *pagetable;
} array;

```

This is much more interesting, an entry could now be either properly

allocated, a *lazy* allocation that we should fix (the -1 that we used before) or a *shared* frame. If it is a shared frame we also have a pointer to the page table entry that is the lazy copy. This will be a bit tricky so buckle up.

First you should go through your code (or why not create a copy and work on the copy) and update the code so that it works with the new representation of page tables entries. When we before treated it as frame number or -1 we must now look inside the data structure to figure out what to do.

If you look at this updated version of `allocate()` you will be able to update also `delete()`, `set()` and `get()`.

```

:
array *new = (array*) malloc(sizeof(array));
pt_entry *pagetable = (pt_entry*) malloc(sizeof(pt_entry)*pages);

new->pages = pages;
new->pagetable = pagetable;

for(int i = 0; i < pages; i++) {
    pt_entry *entry = &new->pagetable[i];
    entry->copy = entry;    // a trick
    entry->status = LAZY;
}
:

```

If you can run your previous benchmarks you should be fine. Now for the tricky part, how do we implement `copy()` and what changed do we have to do to the implementation of `set()`, `get()` and `delete()`?

The idea is like this, a page table entry could be in either of three states:

- **ALLOCATED** : a frame has been allocated for the page and the array in the only user of the frame.
- **LAZY**: a frame has not yet been allocated but will be as soon as a set or get procedure is called.
- **SHARED**: a frame has been allocated but the page is shared by two or more arrays, as long as we read from the page no harm done but as soon as we do a write operation we need to create a copy.

We write “two or more” since we of course we need to handle the case when we have taken a copy of a copy. This will of course complicate things since we need keep track of which other arrays share the page and when there is only one left it should treat it as its own allocated frame.

To keep track of who else shares the frame we link all page table entries that share a frame in a circular list. If an entry is the only entry that holds

a reference to the frame, i.e. ALLOCATED, the copy reference is a circular pointer to the entry it self (this is a trick that will do our coding easier).

This sounds complicated and it is, make some drawing of what things could look like. Also write down in your own words what should be done when we create a copy of a page table entry.

- ALLOCATED : ... now shared ... SHARED
- LAZY : this is easy LAZY
- SHARED : ... could also share ... linked in a circular ...

When you have done some drawings you're ready to copy an array; this skeleton code should give you a head start.

```
array *copy(array *orig) {

    array *copy = (array*)malloc(sizeof(array));
    pt_entry *pagetable = (pt_entry*)malloc(sizeof(pt_entry) * orig->pages)

    int pages = orig->pages;

    copy->pages = pages;
    copy->pagetable = pagetable;

    for(int i = 0; i < pages; i++) {
        pt_entry *orig_entry = &orig->pagetable[i];
        pt_entry *copy_entry = &copy->pagetable[i];

        switch(orig_entry->status) {
            case LAZY:
                copy_entry->status = ...
                copy_entry->copy = ... // circular
                break;
            case ALLOCATED:
            case SHARED:
                copy_entry->frame = ...
                orig_entry->status = ...
                copy_entry->status = ...
                // linking in the circular structure
                copy_entry->copy = ...
                orig_entry->copy = ...
                break;
        }
    }
}
```

```

    return copy;
}

```

Notice that we do the same thing if the original entry is allocated or shared. If the original entry has an allocated frame we simply mark this as a copy and create the initial circular structure knowing that an allocated array has a self reference in the *copy field*,

That was not that complicated and the reason is that the complicated copying is delayed until we do a `set()` operation. When we do a set operation we need to take care of the situation when two or more arrays share the same page; if an entry is referring to an allocated page or a lazy page the situations is as before. This is what we need to insert:

```

:
if( entry->status == SHARED ) {
    int f = find_free();
    if( f == -1 ) {
        delete(arr);
        exit(-1);
    }

    for(int i = 0; i < SIZE; i++) {
        memory[f*SIZE + i] = memory[entry->copy->frame*SIZE + i];
    }

    entry->frame = ...
    entry->status = ...

    // find the entry that is previous to entry
    pt_entry *prev = entry;
    while(prev->copy != entry) {
        prev = prev->copy;
    }
    // it should now point to ...
    prev->copy = ...

    // and if it is pointing to itself ...
    if(prev->copy == prev) {
        prev->status = ...
    }
}
:

```

Ahh, not that complicated (it took me an hour) after all? The `get()` procedure does not need any changes since reading from an allocated or

shared frame will be the same thing. You might want to write this as a switch statement to make it obvious but nothing special needs to be done.

The `delete()` procedure needs to be updated since we could now delete an array that holds a shared frame. The procedure is then similar to the last part in the set operation i.e. we unlink the entry from the list of copies and if there is only one copy left that entry is changed to allocated.

The lazy copying that you have now implemented is what the operating system does every time you do a `fork()`. The code area is of course read only and will be shared until either process calls `exec()`. The data areas will be shared until written to but only the pages that are written to will be copied.