

به نام خدا

نام : **عماد**

نام خانوادگی : **آقاجانی**

شماره دانشجویی : ۸۸۵۲۱۳۴۴

استاد درس: دکتر شریفی

# Linux Memory Management

فهرست مندرجات :

- (۱) توضیحات
- (۲) مکانیزم **Paging**
- (۳) فضای آدرس فیزیکی و منطقی حافظه
- (۴) مدیریت **Page Frame** در لینوکس
- (۵) انواع مدیریت حافظه های استفاده شده در لینوکس
  - a. تخصیص پیوسته
  - i. الگوریتم **Buddy**
  - ii. **Slab Allocvator**
  - b. تخصیص ناپیوسته
- (۶) فضای آدرس دهی پردازنده ها (Processes)
- (۷) رسیدگی کننده **Page Fault (Handler)**
- (۸) نحوه اختصاص حافظه به پردازنده جدید
- (۹) مدیریت حافظه **Heap**
- (۱۰) منابع تحقیقاتی

## توضیحات

با عرض سلام و خسته نباشید

تحقیق پیش رو چکیده ای از مطالعه و یادگیری بنده از منابع مختلف است .

علت کار اضافی انجام شده روی این تکلیف صرفا یادگیری و علاقه شخصی و ایجاد یک منبع مطالعاتی چکیده در رابطه با نحوه مدیریت حافظه در لینوکس برای (حداقل) خودم بوده که هنوز هم قصد کاملتر کردنشو دارم .

همچنین در صورت تمایل از ارایه حضوری این تکلیف در کلاس درس، استقبال میکنم .

در ضمن چون هدف نحوه مدیریت حافظه بوده ، نه تعریف **Heap Memory** و **Page File** و **External Fragmentation** ، فرض بر اینه که این اطلاعات رو داریم .

در مورد متن هم باید بگم که الکی نخواستم با کلمات جمله سازی کنم و برای همین به سبک **note** نوشته شده .

با تشکر از توجه شما

آقاجانی (=)

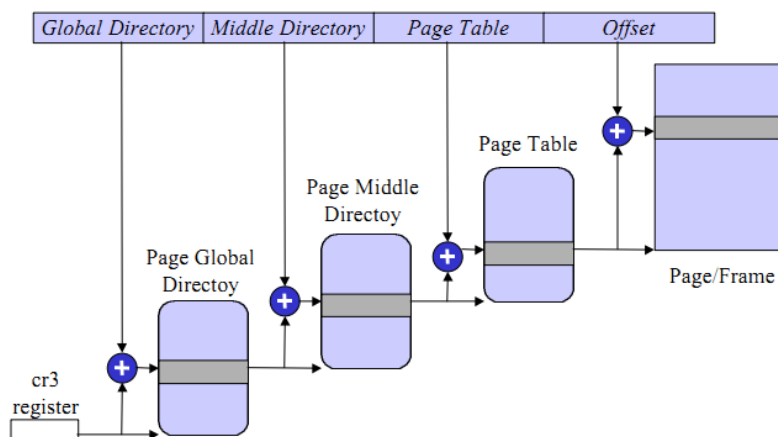
## مکانیزم Paging

- پردازنده هایی از قبیل Intel x86 توانایی پشتیبانی از شیوه Segment را دارا می باشند.
- ولی در لینوکس از شیوه Segment استفاده نمیشود، چراکه :
  - (۱) استفاده از مقدار ثبات Segment یکسان برای تمام پردازنده ها (Process)، مدیریت حافظه را آسان میکند.
  - (۲) حفظ قابلیت استفاده بر روی سایر پردازنده ها (Portability).
- در شیوه Paging استفاده شده در لینوکس :
  - (۱) اندازه هر Page File برابر 4K.
  - (۲) یک Page Table سه لایه برای مدیریت ۶۴ بیت خط آدرس.
  - (۳) در پردازنده هایی از قبیل x86 :
    - تنها به دو لایه Page Table نیاز میباشد.
    - Page Table توسط سخت افزار پشتیبانی میشود.
    - امکان بهره گیری از TLB(Translation lookaside buffer) فراهم آمده است.

## فضای آدرس (Address Space) منطقی / فیزیکی

- شمای آدرس های منطقی در لینوکس :

(For x86 processors, Middle Directory is 0 bits)



- حدود ۲ مگابایت اول **حافظه فیزیکی** بمنظور معماری کامپیوتر جاری و اطلاعات و ... سیستم عامل رزرو شده است و مابقی برای انجام **Paging** در دسترس میباشد.

- فضای منطقی** آدرس های یک پرده به دو قسمت تقسیم شده است :

(۱) که شامل 0X00000000 تا PAGE\_OFFSET-1 میشود و میتوان چه در **User mode** و چه در **Kernel mode**

آن را آدرس دهی کرد

(۲) شامل PAGE\_OFFSET تا 0xffffffff میشود و تنها میتوان در **Kernel mode** آن را آدرس دهی کرد

▪ PAGE\_OFFSET معمولاً 0xc0000000 میباشد

## مدیریت Page Frame در لینوکس

- کرنل وضعیت کنونی هر Page Frame را در آرایه ای از struct page ها ، بنام mem\_map ، نگه میدارد.
  - همچنین کرنل تعداد استفاده از یک Page Frame را نگه میدارد . مقدار \* بمعنای آزاد بودن Page Frame و عدد بزرگتر از \* بمعنای تعداد استفاده از آن Page Frame میباشد.
  - همچنین پرچم (Flag) هایی از قبیل Dirty, Locked, Referenced و غیره نیز برای Page Frame ها وجود دارد.
  - کرنل Page Frame ها را بوسیله دستورات زیر گرفته و یا آزاد میکند :  

```
__get_free_pages ( gfp_mask, order )  
free_pages ( addr, order )
```
  - در تئوری، شیوه Paging نیاز به تخصیص حافظه پیوسته را از بین برده است - بعضی از عملیات مانند DMA در روند Paging خلل ایجاد میکنند .
  - از طرف دیگر باید توجه داشت که شیوه تخصیص حافظه پیوسته Page Table کرنل را تغییر نمیدهد که در نتیجه باعث حفظ TLB میگردد و زمان دسترسی را کاهش میدهد !
  - در نتیجه لینوکس از شیوه ای بمنظور تخصیص پیوسته Page Frame (allocating contiguous page frames) استفاده میکند .
- در این روش مشکل External Fragmentation نیز وجود نخواهد داشت.

## انواع مدیریت حافظه های استفاده شده در لینوکس :

(۱) تخصیص حافظه پیوسته (contiguous memory allocation)

(۲) تخصیص حافظه ناپیوسته (noncontiguous memory allocation)

### تخصیص حافظه پیوسته (contiguous memory allocation) [۴و۵و۶]

- در این شیوه مدیریت حافظه که یک شیوه سنتی محسوب میشود، تخصیص حافظه به پردازنده ها از بلاک های متوالی حافظه صورت میگیرد.
- اشکال این شیوه : مشکل External Fragmentation .
- ویژگی این روش : در حافظه هایی با حجم بالا مشکل های مرتبط با کش (Cache) و تاخیر دسترسی (Access Latency) وجود ندارد.

### الگوریتم Buddy [۷و۸]

- تمامی page frame های موجود به ۱۰ لیست از بلاک ها تقسیم بندی میشوند که شامل لیست هایی با ۱، ۲، ۴، ۸، ۱۶، ۳۲، ۶۴، ۱۲۸، ۲۵۶، ۵۱۲ page frame پیوسته میشود.
- آدرس اولین page frame در یک بلاک ضربی از سایز آن دسته میباشد. بعنوان مثال در یک بلاک ۱۶ تایی، آدرس اولین page frame برابر  $16 \times 2^{12}$  خواهد بود.
- الگوریتم برای تخصیص مثلاً یک بلاک ۱۲۸ تایی از page frame های پیوسته :
  - ابتدا به دنبال یک بلاک خالی در لیست ۱۲۸ تایی بگرد
  - اگر بلاک خالی پیدا نشد، در لیست ۲۵۶ تایی دنبال یک بلاک خالی بگرد
  - اگر بلاکی پیدا کرد، کرنل ۱۲۸ تا از ۲۵۶ page frame موجود را تخصیص میدهد و ۱۲۸ تای باقیمانده را در لیست ۱۲۸ تایی ها قرار میدهد
  - اگر در لیست های بزرگتر پیدا شد، انقدر بلاک را خرد کن ( بصورت باینری) تا حافظه مورد نیاز را اختصاص بدهی و بلاک های بوجود آمده را در لیست مناسب قرار بده
  - اگر هیچ بلاکی قابل تخصیص نیست یک پیغام خطا بازگردان

- یک مثال از نحوه تخصیص حافظه (Buddy System) :

Step	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K	64K		
1	2 <sup>4</sup>																	
2.1	2 <sup>3</sup>								2 <sup>3</sup>									
2.2	2 <sup>2</sup>				2 <sup>2</sup>				2 <sup>3</sup>									
2.3	2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>									
2.4	2 <sup>0</sup>	2 <sup>0</sup>	2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>									
2.5	A: 2 <sup>0</sup>	2 <sup>0</sup>	2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>									
3	A: 2 <sup>0</sup>	2 <sup>0</sup>	B: 2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>									
4	A: 2 <sup>0</sup>	C: 2 <sup>0</sup>	B: 2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>									
5.1	A: 2 <sup>0</sup>	C: 2 <sup>0</sup>	B: 2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>3</sup>									
5.2	A: 2 <sup>0</sup>	C: 2 <sup>0</sup>	B: 2 <sup>1</sup>		D: 2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>3</sup>									
6	A: 2 <sup>0</sup>	C: 2 <sup>0</sup>	2 <sup>1</sup>		D: 2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>3</sup>									
7.1	A: 2 <sup>0</sup>	C: 2 <sup>0</sup>	2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>1</sup>		2 <sup>3</sup>									
7.2	A: 2 <sup>0</sup>	C: 2 <sup>0</sup>	2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>									
8	2 <sup>0</sup>		C: 2 <sup>0</sup>		2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>							
9.1	2 <sup>0</sup>		2 <sup>0</sup>		2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>							
9.2	2 <sup>1</sup>			2 <sup>1</sup>		2 <sup>2</sup>				2 <sup>3</sup>								
9.3	2 <sup>2</sup>				2 <sup>2</sup>				2 <sup>3</sup>									
9.4	2 <sup>3</sup>								2 <sup>3</sup>									
9.5	2 <sup>4</sup>																	

This allocation could have occurred in the following manner

- The initial situation.
- Program A** requests memory 34K, order 0.
  - No order 0 blocks are available, so an order 4 block is split, creating two order 3 blocks.
  - Still no order 0 blocks available, so the first order 3 block is split, creating two order 2 blocks.
  - Still no order 0 blocks available, so the first order 2 block is split, creating two order 1 blocks.
  - Still no order 0 blocks available, so the first order 1 block is split, creating two order 0 blocks.
  - Now an order 0 block is available, so it is allocated to A.
- Program B** requests memory 66K, order 1. An order 1 block is available, so it is allocated to B.
- Program C** requests memory 35K, order 0. An order 0 block is available, so it is allocated to C.
- Program D** requests memory 67K, order 1.
  - No order 1 blocks are available, so an order 2 block is split, creating two order 1 blocks.
  - Now an order 1 block is available, so it is allocated to D.
- Program B** releases its memory, freeing one order 1 block.
- Program D** releases its memory.
  - One order 1 block is freed.
  - Since the buddy block of the newly freed block is also free, the two are merged into one order 2 block.
- Program A** releases its memory, freeing one order 0 block.
- Program C** releases its memory.
  - One order 0 block is freed.
  - Since the buddy block of the newly freed block is also free, the two are merged into one order 1 block.
  - Since the buddy block of the newly formed order 1 block is also free, the two are merged into one order 2 block.
  - Since the buddy block of the newly formed order 2 block is also free, the two are merged into one order 3 block.
  - Since the buddy block of the newly formed order 3 block is also free, the two are merged into one order 4 block.

As you can see, what happens when a memory request is made is as follows:

- If memory is to be **allocated**

1. Look for a memory slot of a suitable size (the minimal  $2^k$  block that is larger or equal to that of the requested memory)
  1. If it is found, it is allocated to the program
  2. If not, it tries to make a suitable memory slot. The system does so by trying the following:
    1. Split a free memory slot larger than the requested memory size into half
    2. If the lower limit is reached, then allocate that amount of memory
    3. Go back to step 1 (look for a memory slot of a suitable size)
    4. Repeat this process until a suitable memory slot is found
- If memory is to be **freed**
  1. Free the block of memory
  2. Look at the neighboring block - is it free too?
  3. If it is, combine the two, and go back to step 2 and repeat this process until either the upper limit is reached (all memory is freed), or until a non-free neighbour block is encountered

- وقتی یک بلاک آزاد می‌گردد، کرنل سعی بر ادغام جفت های مشابه (free Buddy Block) با سایز  $b$  به یک  $2b$  دارد.
- دو بلاک buddy تلقی میشوند اگر :
  - (۱) اگر هم سایز باشند
  - (۲) هر دو در آدرس فیزیکی متوالی واقع شده باشند
  - (۳) آدرس فیزیکی page اول از بلاک اول ضریبی از  $2b \times 2^{12}$  باشد
- عملیات ادغام متوالیا تکرار میشود.
- لینوکس از دو نوع Buddy System متفاوت برای مدیریت حافظه استفاده میکند.
  - یکی برای page frame های مناسب برای DMA (بعنوان مثال page frame های با آدرس هایی کمتر از 16MB).
  - دیگری برای سایر.
- هر Buddy System مبتنی است بر ..
  - آرایه توصیفگر frame ها ، mem\_map
  - آرایه ۱۰ تایی از free\_area\_struct، هر درایه برای یک سایز (۱، ۲، ...، ۵۱۲) (هر درایه یک لینک لیست دوطرفه از بلاک های page frame با سایز مربوطه خود میباشد).
  - ده بیت‌مپ (Bitmap) برای نگهداری وضعیت تخصیص بلاک های موجود.
- الگوریتم Buddy برای مدیریت حجم های بالای حافظه بخوبی عمل میکند. اما در رابطه با بخش های کوچک حافظه با مشکل Internal Fragmentation مواجه است !
- برای حل این مشکل از نسخه ۲.۲ لینوکس slab allocator معرفی گردید

### مکانیزم Slab Allocator [۹ و ۱۰]

- ایده اصلی : نگهداری اشیا (objects) بتازگی استفاده شده در کش
- در این شیوه :
  - قسمت های مختلف حافظه بعنوان یک اشیا همراه با اطلاعات و متدها (مانند سازنده و مخرب و ...) دیده میشوند.
  - اشیا بصورت دسته بندی شده به کش فرستاده میشوند
  - یک مجموعه ویژه در کش برای عملیات های سیستم عامل در نظر گرفته میشود



- Slab cache شامل هیچ یا چند slab میشود. slab یک یا چند page frame پیوسته از Buddy System میباشد.
- اشیا بوسیله دستور `kmem_cache_alloc(cachep)` تخصیص می یابند ( `cachep` یک اشاره گر به کش میباشد).
- اشیا بکمک دستور `kmem_cache_free(cachep, objp)` آزاد میشوند.
- یک گروه از کش های عمومی وجود دارد که با سایز های گسترده ای بین ۳۲ تا ۱۳۱۰۷۲ بایت توزیع شده است.
- برای گرفتن شی ای از این کش ها از دستور `kmalloc(size, flags)` استفاده میشود.
- برای آزاد کردن شی ای از این کش ها عمومی از دستور `kfree(objp)` استفاده میشود.

### تخصیص حافظه ناپیوسته (noncontiguous memory allocation) [۶]

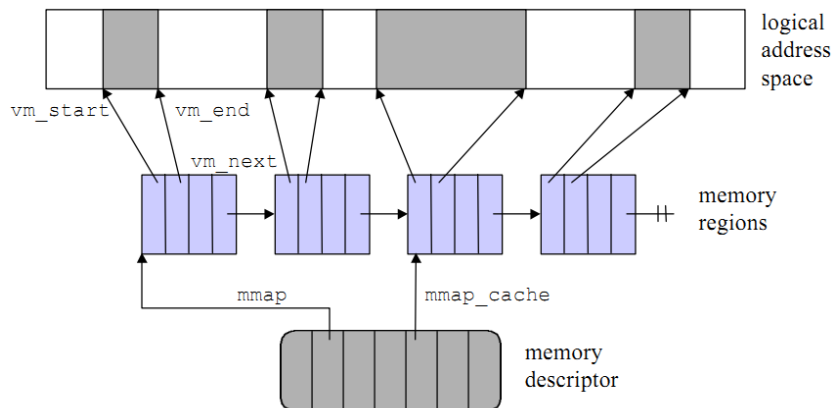
- در این شیوه از یک حافظه مجازی (Virtual Memory) برای پیاده سازی شیوه سنتی " تخصیص حافظه ناپیوسته " استفاده میگردد.
- یک فضا در فضای آدرس مجازی بین `VMALLOC_START` و `VMALLOC_END` رزرو شده است. مکان `VMALLOC_START` به میزان حافظه فیزیکی بستگی دارد ولی محدوده همیشه حداقل `VMALLOC_RESERVE` میباشد که بر روی یک **X86** برابر ۱۲۸ مگابایت میباشد.
- لینوکس از تخصیص حافظه ناپیوسته عموماً پرهیز میکند. اما گاهی برای اختصاص حافظه به گونه ای از درخواست های محدود استفاده از این روش کارآمدتر میشود.
- لینوکس بیشتر آدرس های بالاتر از `PAGE_OFFSET` را بمنظور تخصیص ناپیوسته استفاده میکند.
- برای تخصیص و آزاد کردن حافظه به روش ناپیوسته از دستورات `vmalloc(size)` و `vfree(addr)` استفاده میگردد.

مرور : انواع تخصیص حافظه در کرنل لینوکس

- کرنل حافظه پویا خود را به یکی از ۳ روش زیر میگیرد :
  - (۱) استفاده از `__get_free_pages()` برای گرفتن **pages** از Buddy System
  - (۲) `kmem_cache_alloc()` یا `kmalloc()` برای استفاده از تخصیص گر **slab** برای گرفتن اشیا خاص یا عمومی
  - (۳) `vmalloc()` برای گرفتن حافظه ناپیوسته

## فضای آدرس دهی پردازش ها

- فضای آدرس (Address Space) یک پردازش، شامل تمام آدرس های منطقی قابل دسترس توسط پردازش میشود.
- فضای آدرس هر پردازش مستقل از دیگری است (مگر اینکه مشترک باشد).
- کرنل عملیات تخصیص آدرس منطقی برای پردازش ها را در یک بازه بنام "مناطق حافظه" (Memory Regions) انجام میدهد. این ناحیه یک مقدار اولیه و یک طول دارد که مضربی از ۴۰۹۶ می باشد.
- برای کرنل، دو دسته درخواست حافظه اتفاق می افتد :
  - (۱) درخواست های غیر ضروری :
    - a. نیازی به تمام Page ها همیشه نیست
    - b. تخصیص حافظه ممکن است برای لحظه ای امکان پذیر نباشد
  - (۲) درخواست های نامطمئن
    - a. کرنل باید آمادگی روبرو شدن با خطا را داشته باشد
- در نتیجه کرنل لینوکس با توجه به پردازش، عملیات تخصیص حافظه را گوناگون انجام میدهد.
- شرایطی که پروسه ها حافظه جدید میگیرند :
  - (۱) ایجاد یک پروسه جدید (fork())
  - (۲) بارگزاری یک برنامه جدید (execve())
  - (۳) نگاشت حافظه به فایل (mmap())
  - (۴) رشد پشته (Stack)
  - (۵) ایجاد حافظه اشتراکی (shmat())
  - (۶) گسترش heap (malloc())
- تمامی اطلاعات مرتبط با فضای آدرس یک پردازش در بخش توصیفگر حافظه (memory descriptor) (mm\_struct) قرار گرفته است که از طریق فیلد mm از توصیفگر پردازش (process descriptor) قابل دسترس میباشد. بعضی از اطلاعاتی که قابل دسترس میباشد :
- (۱) یک اشاره گر به بالاترین مرحله Page Table یا همان Page Global Directory از طریق pgd
  - (۲) تعداد page frame های تخصیص یافته از طریق rss
  - (۳) ...
- توصیفگر حافظه از کش اختصاص گر slab (slab allocator cache) و با دستور mm\_alloc() تخصیص می یابد.
- شمایی انتزاعی از توصیفگر حافظه، مناطق حافظه و فضای آدرس منطقی



- برای تخصیص و آزاد سازی یک آدرس مجازی، کرنل بترتیب از دستورات `Do_mmap()` و `do_munmap()` استفاده میکند.

## رسیدگی کننده Page Fault

- زمانی که یک پردازش از کرنل تقاضای حافظه بیشتر میکند، پردازش تنها فضای آدرس منطقی بیشتر میگیرد و نه حافظه فیزیکی بیشتر.
- زمانی که پردازش قصد استفاده و دسترسی به فضای آدرس منطقی جدید خود را میکند، یک **Page Fault** رخ میدهد تا به کرنل بگوید که واقعا به این حافظه نیاز دارد.
- رسیدگی کننده **Page Fault**، آدرس منطقی تقاضا شده را با ناحیه ای از حافظه که در اختیار پردازش هست مقایسه میکند تا یا متوجه خطا در تقاضا شود (عدم امکان دسترسی پردازش به آن آدرس مجازی) و یا متوجه نیاز پردازش به اختصاص از حافظه فیزیکی شود. همچنین آدرس تقاضا شده، ممکن است بدلایلی از روی حافظه فیزیکی به روی دیسک سخت منتقل شده باشد.

## نحوه اختصاص حافظه به پردازش جدید

- زمانی که کرنل یک پردازش جدید ایجاد میکند، به آن یک فضای آدرس کاملاً جدید نمیدهد و در ابتدا پردازش جدید از فضای آدرس پردازش پدر خود بصورت اشتراکی استفاده میکند. کرنل در این زمان از **page frame** های مشترک محافظت میکند.
  - سپس اولین مرتبه ای که پردازش پدر یا پردازش جدید قصد نوشتن بر روی **page frame** های مشترک را داشته باشند یک استثنا (**exception**) رخ میدهد.
  - کرنل این استثنا را دریافت کرده و یک نسخه کپی از آن **frame**، برای پردازش ای که قصد نوشتن داشت، ایجاد میشود.

## مدیریت حافظه Heap

- هر پردازش میتوانند بر روی **heap** خود، حافظه پویا جدید بدست بیاورد.
- دستورات زبان C بمنظور انجام تغییرات در **Heap**: **malloc()**, **calloc()** و **free()**.
- تمامی این دستورات از دستور اصلی **brk()** مشتق شده اند.
  - این تنها دستوری میباشد که یک فراخوانی سیستمی (**System Call**) را انجام میدهد.
  - این دستور مستقیماً میزان حافظه **Heap** را تغییر میدهد.
  - این دستور در حقیقت عمل تخصیص و رها سازی فضای آدرس منطقی را برعهده دارد.
- زمانی که پردازش یک **page frame** میگیرد، عمل تخصیص قطعه حافظه های کوچک (مانند: **malloc ( sizeof(char)\*50**) در **user space** صورت میگیرد.

## منابع تحقیقاتی

منبع اصلی:

[www.inf.fu-berlin.de/lehre/SS01/OS/Lectures/Lecture14.pdf](http://www.inf.fu-berlin.de/lehre/SS01/OS/Lectures/Lecture14.pdf)

سایر منابع:

- [1] <http://tldp.org/LDP/tlk/mm/memory.html>
- [2] [http://en.wikipedia.org/wiki/Translation\\_lookaside\\_buffer](http://en.wikipedia.org/wiki/Translation_lookaside_buffer)
- [3] <http://www.cs.umd.edu/~meesh/cmsc411/website/saltz/cs412/lect3.html>
- [4] <http://www.techopedia.com/definition/3769/contiguous-memory-allocation>
- [5] <http://siber.cankaya.edu.tr/ozdogan/OperatingSystems/ceng328/node180.html>
- [6] <http://www.kernel.org/doc/gorman/html/understand/understand010.html>
- [7] [http://en.wikipedia.org/wiki/Buddy\\_memory\\_allocation](http://en.wikipedia.org/wiki/Buddy_memory_allocation)
- [8] [www.cs.purdue.edu/homes/hosking/690M/p421-peterson.pdf](http://www.cs.purdue.edu/homes/hosking/690M/p421-peterson.pdf)
- [9] [en.wikipedia.org/wiki/Slab\\_allocation](http://en.wikipedia.org/wiki/Slab_allocation)
- [10] <http://www.kernel.org/doc/gorman/html/understand/understand011.html>